

# **(Hyper)Safety Certification of Neural Network Surrogates for Aircraft Braking Distance Estimation**

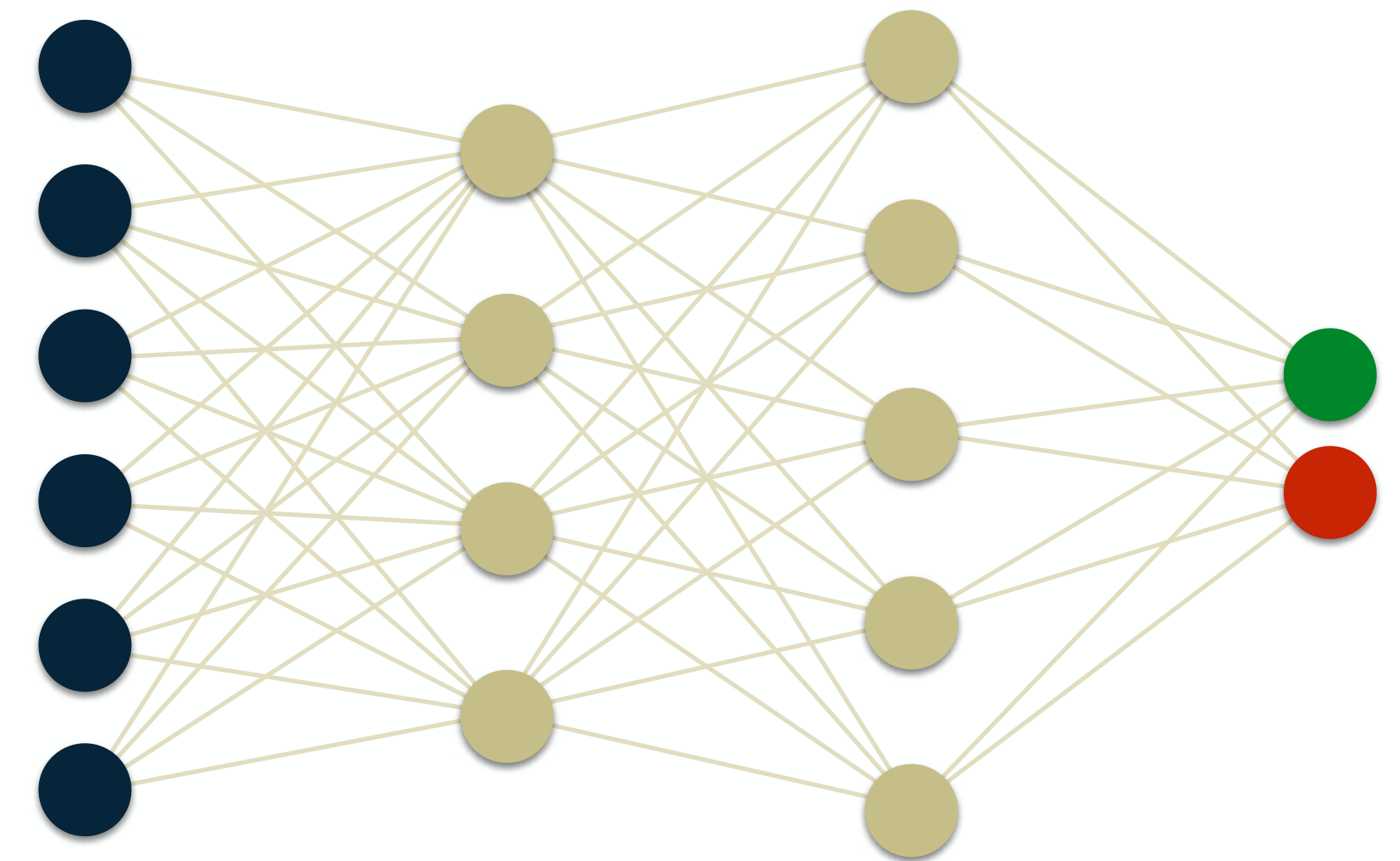
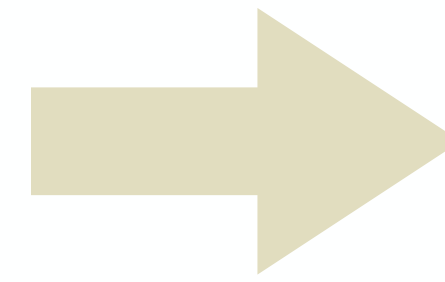
**Tech Talk: Formal Verification of AI**

**Caterina Urban**

Inria & École Normale Supérieure | Université PSL

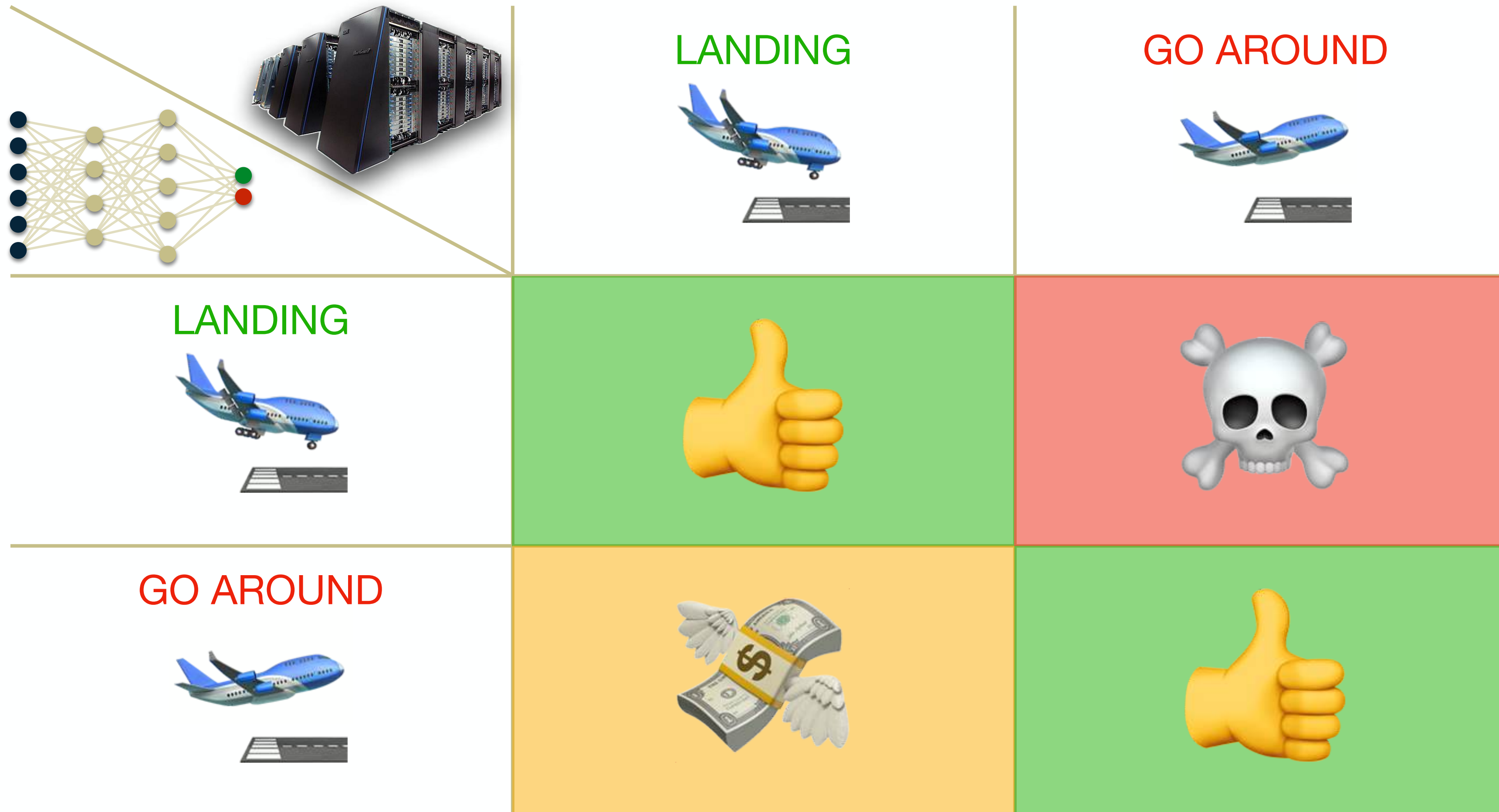
# Neural Network Surrogates

Less Computing Power and Less Computing Time



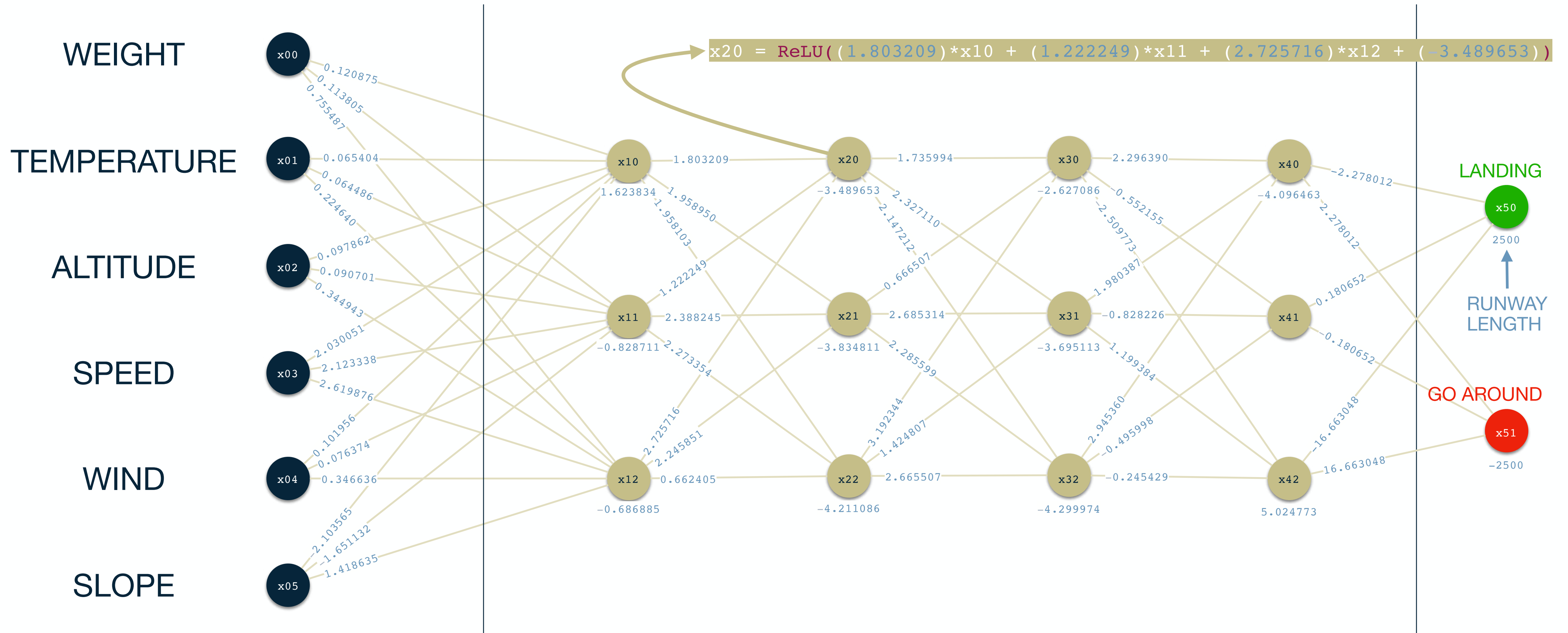
# Runway Overrun Warning

## Safety of Neural Network Surrogate



# Runway Overrun Warning

## Toy Example



# Runway Overrun Warning

## Toy Example

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Neural Network Verification

# Neural Network Explainability

# Neural Network Verification

# Neural Network Explainability

# **(Hyper)Safety Certification of Neural Network Surrogates for Aircraft Braking Distance Estimation**

**= by means of Abstract Interpretation-Based Static Analysis**



# Abstract Interpretation

SOFTWARE



€ 2.25    € 3



€ 2.95

€ 3



€ 3.65    € 4



€ 5.35

€ 6

ABSTRACTION



PROPERTY OF INTEREST

SOUNDNESS



€ 3 +  
 € 3 +  
 € 4 +  
 € 6  
 -----  
 € 16



€ 2.25 +  
 € 2.95 +  
 € 3.65 +  
 € 5.35  
 -----  
 € 14.20



FALSE ALARM

COMPLETENESS

# Abstract Interpretation

## 3-Step Recipe

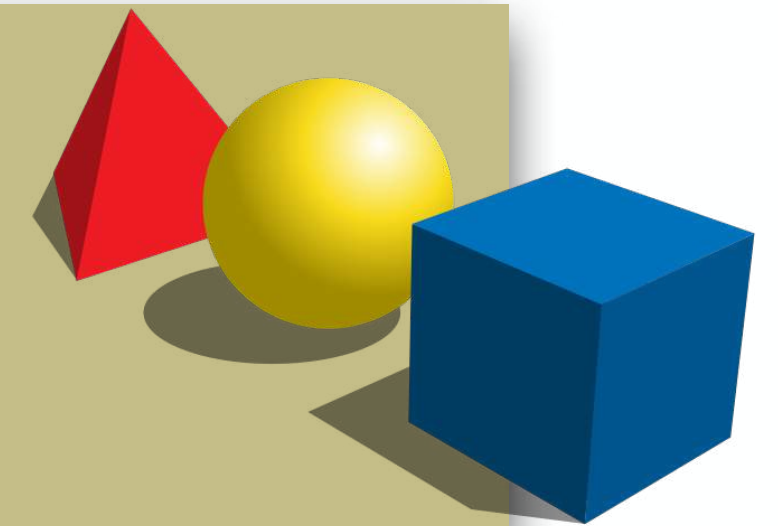
**practical tools**

targeting specific programs



**abstract semantics, abstract domains**

algorithmic approaches to decide program properties



**concrete semantics**

**mathematical models** of the program behavior



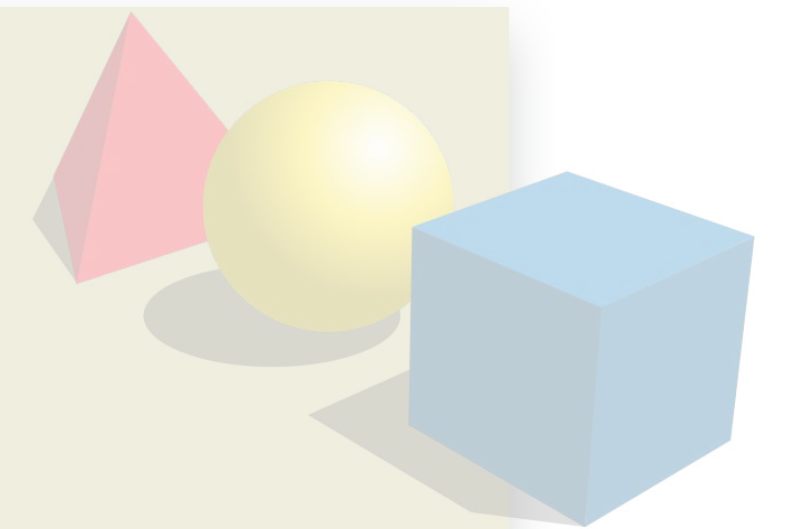
# Abstract Interpretation

## 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



# Forward Pass Trace Semantics

```

x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())

```

```

x10 = ReLU((0.120875)*x00 + (0.065404)*x01 +
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 +
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 +

```

```

x20 = ReLU((1.803209)*x10 + (1.222249)*x11 +
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 +
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 +

```

```

x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.
x31 = ReLU((2.927110)*x20 + (2.685314)*x21 + (1.4
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.66

```

```

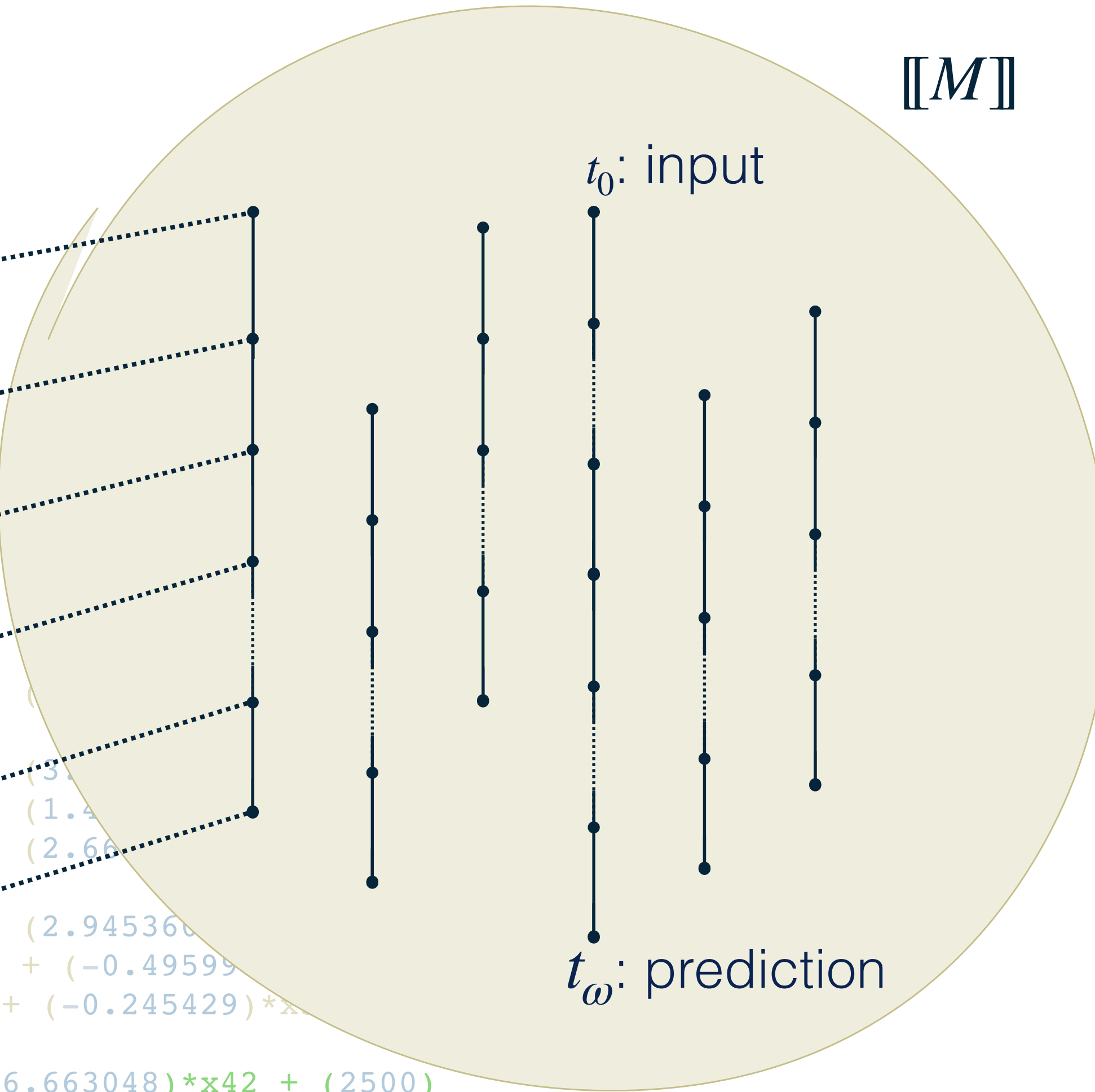
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.49599
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x

```

```

x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)

```

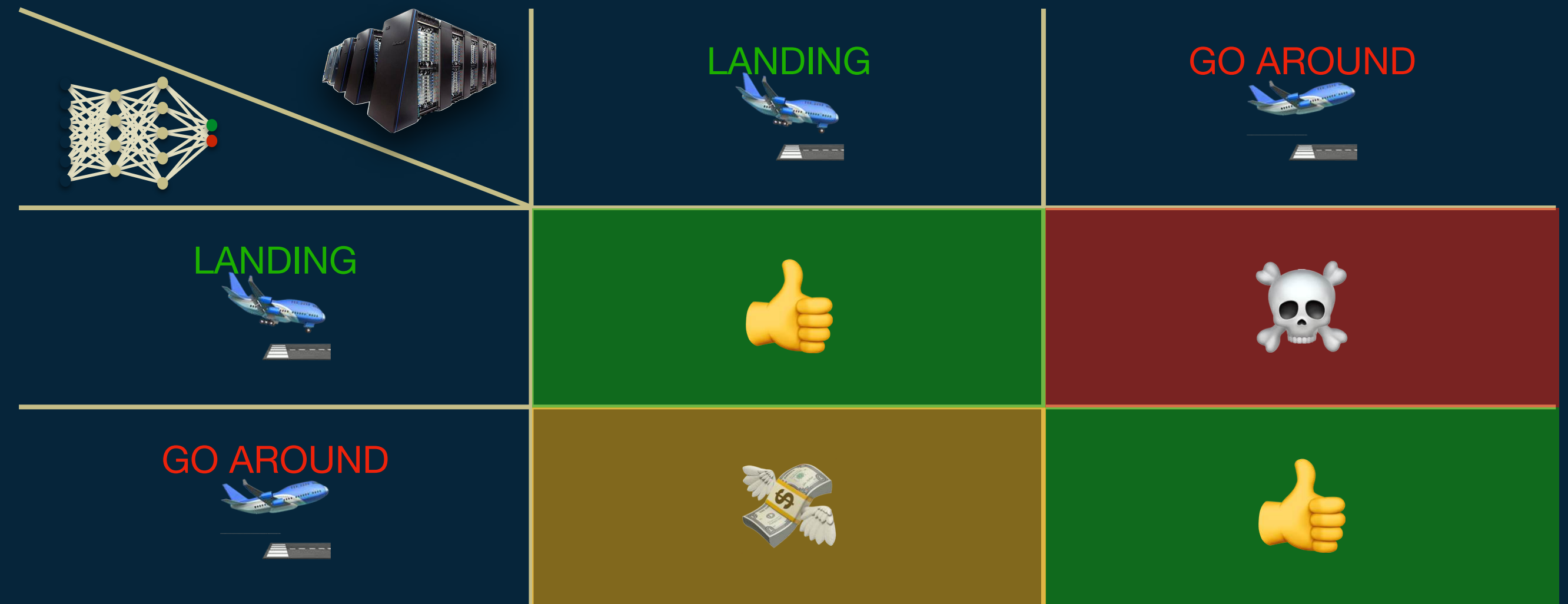


```

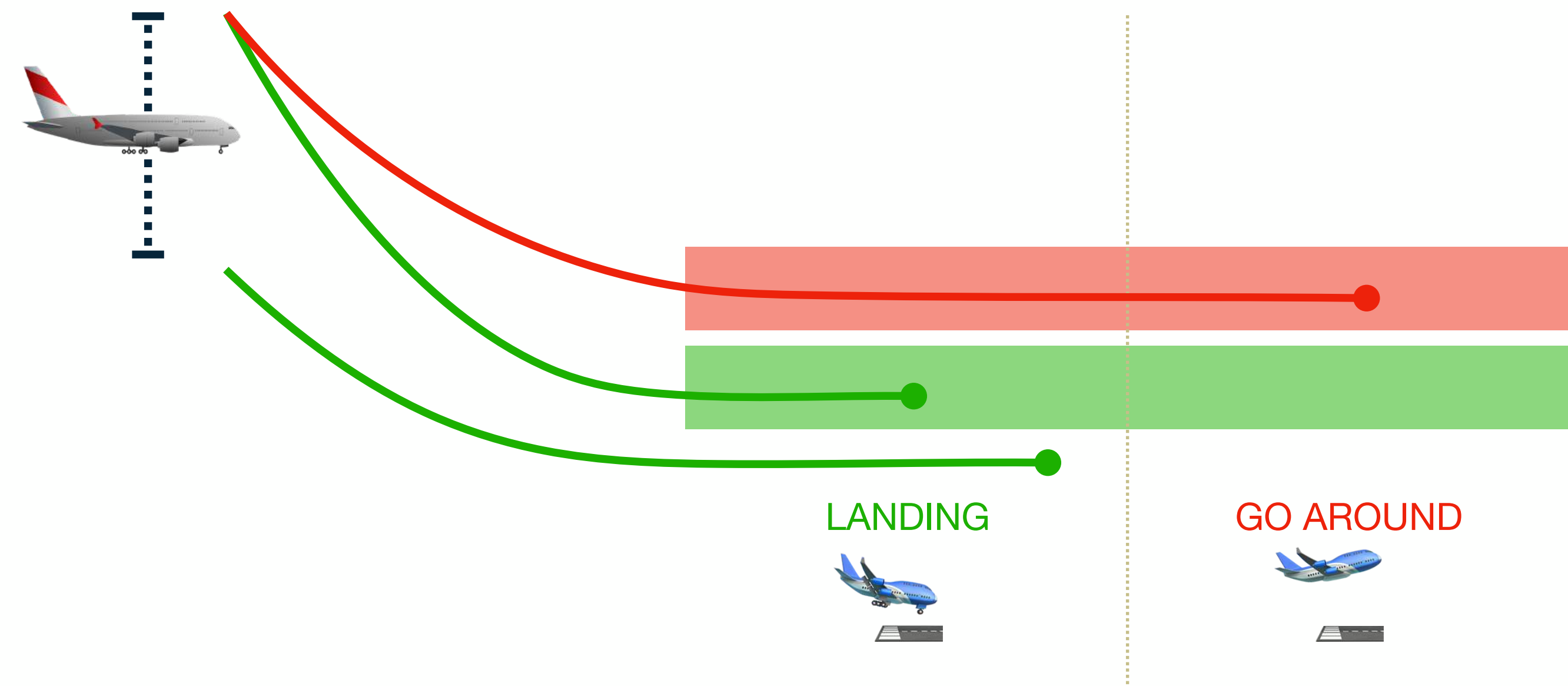
.103565)*x05 + (1.623834)
651132)*x05 + (-0.828711)
18635)*x05 + (-0.686885)

```

# Safety

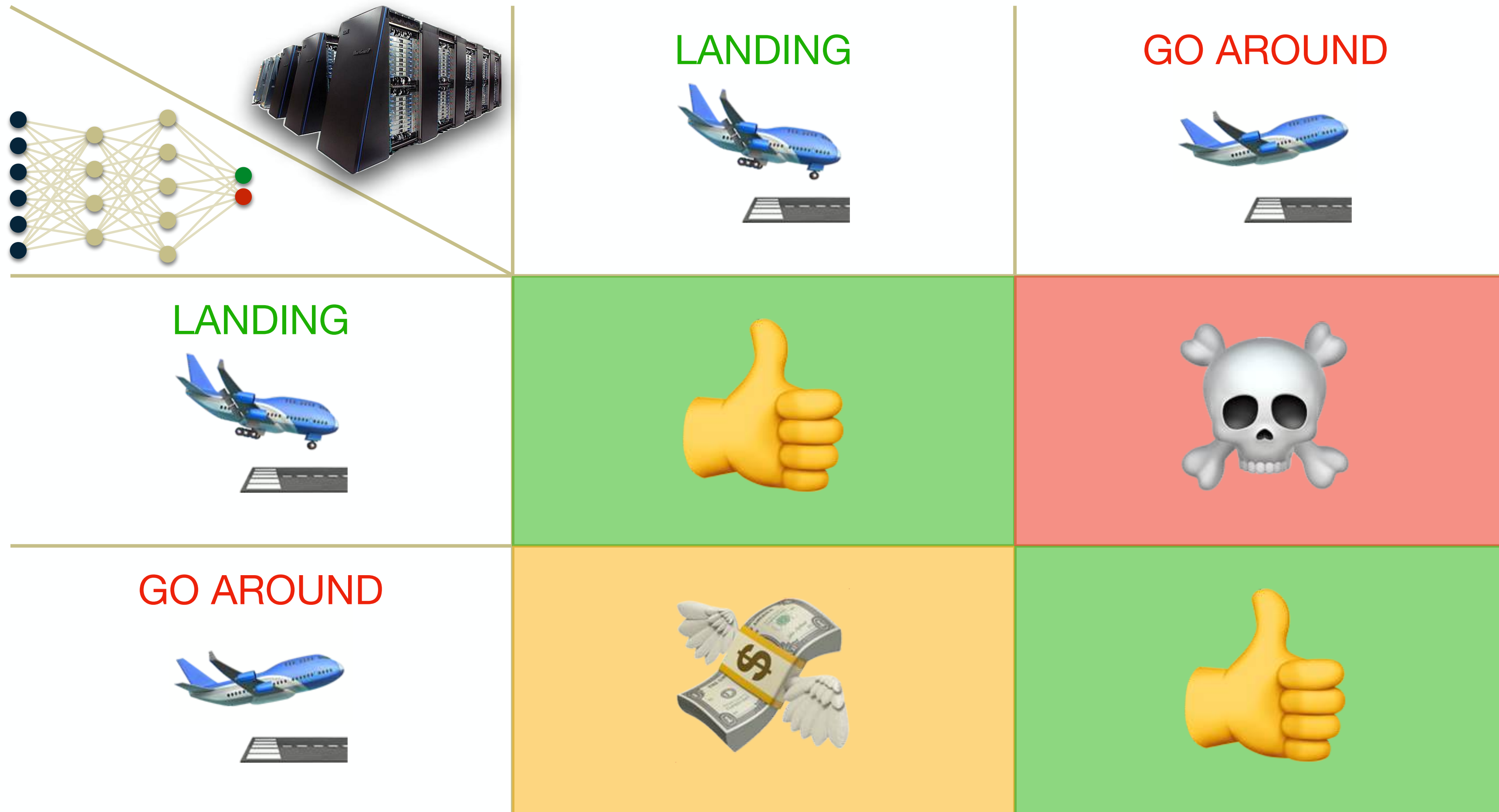


# Hypersafety



# Runway Overrun Warning

## Safety of Neural Network Surrogate



# Safety Verification

## Extensional Properties

**I**: input specification

**O**: output specification

$$\mathcal{S} \stackrel{\text{def}}{=} \left\{ t \mid t_0 \models \mathbf{I} \Rightarrow t_\omega \models \mathbf{O} \right\}$$

$\mathcal{S}$  is the set of all forward pass traces that **satisfy** the specification

Theorem

$$M \models \mathcal{S} \Leftrightarrow \llbracket M \rrbracket \subseteq \mathcal{S}$$

Corollary

$$M \models \mathcal{S} \Leftarrow \llbracket M \rrbracket \subseteq \llbracket M \rrbracket^\dagger \subseteq \mathcal{S}$$

# Safety Verification

## Example

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

**I:**

```
-1 ≤ x00 ≤ 1
-1 ≤ x01 ≤ 1
-1 ≤ x02 ≤ 1
-1 ≤ x03 ≤ 1
-1 ≤ x04 ≤ 1
-1 ≤ x05 ≤ 1
```

**O:**

```
x50 > x51
```



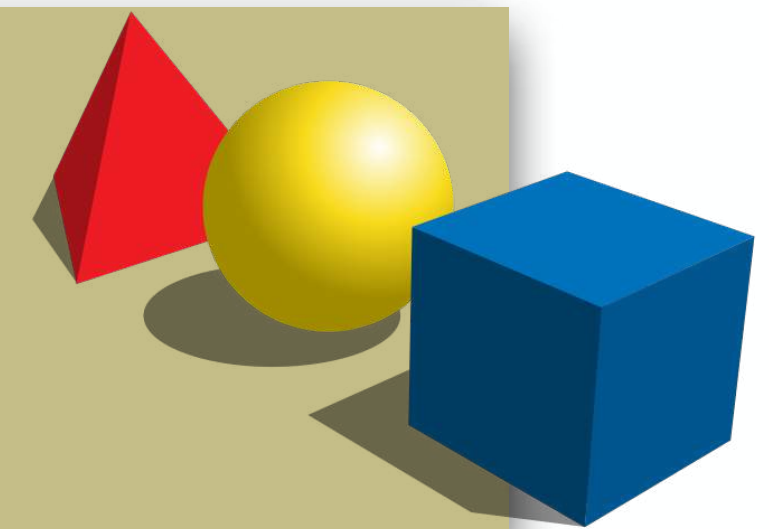
# Abstract Interpretation

## 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



# Safety Verification

## Static Forward Analysis

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

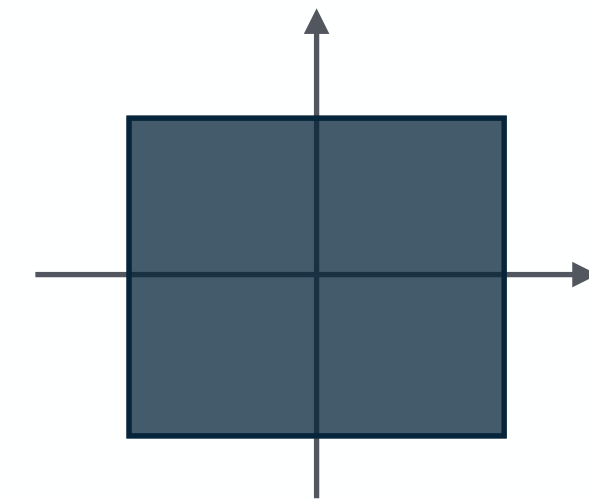
```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

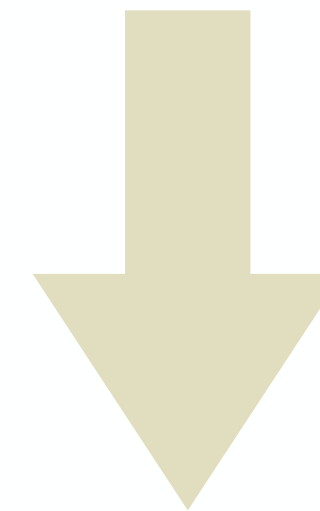
```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

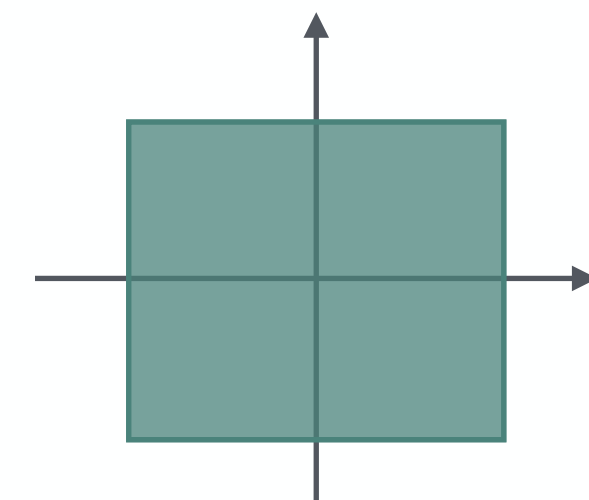
```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```



① start from an **abstraction** of all possible inputs



② proceed **forwards abstracting** the neural network computations



③ check output for **inclusion** in **expected output**:  
included → **safe**  
otherwise → **alarm**

# Safety Verification

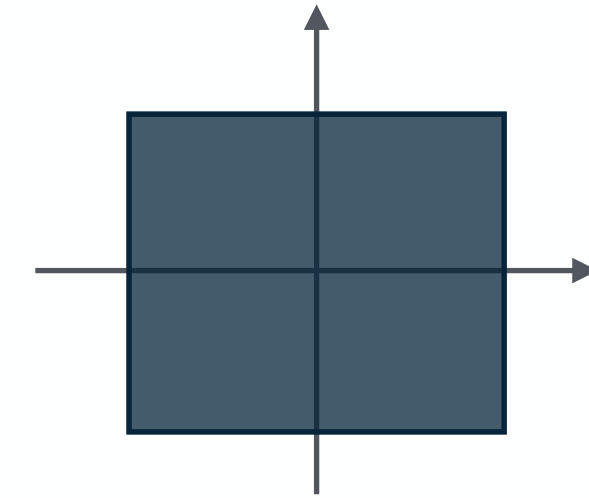
$$x \mapsto [a, b]$$
$$a, b \in \mathbb{R}$$

## Abstraction #1: Boxes Abstract Domain

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

I:

x00: [-1, 1]
x01: [-1, 1]
x02: [-1, 1]
x03: [-1, 1]
x04: [-1, 1]
x05: [-1, 1]



① start from an **abstraction** of all possible inputs

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

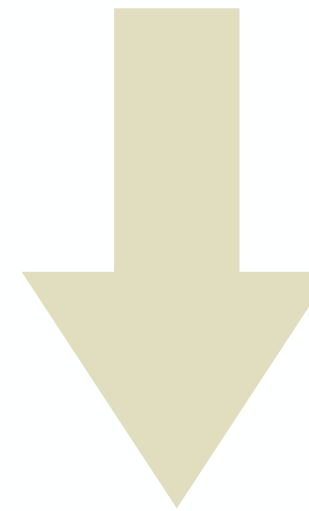
$$x \mapsto [a, b]$$
$$a, b \in \mathbb{R}$$

## Abstraction #1: Boxes Abstract Domain

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

I:

x00: [-1, 1]
x01: [-1, 1]
x02: [-1, 1]
x03: [-1, 1]
x04: [-1, 1]
x05: [-1, 1]



② proceed **forwards**  
**abstracting** the neural  
network computations

```
x10' = (0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834)
```

```
x10 -> [-2.90, 6.14]
```

```
x10 = ReLU(x10')
```

```
x10 -> [0, 6.14]
```

```
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
```

```
x11 -> [0, 3.29]
```

```
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x12 -> [0, 5.02]
```

⋮

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
```

```
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

$$x \mapsto [a, b]$$

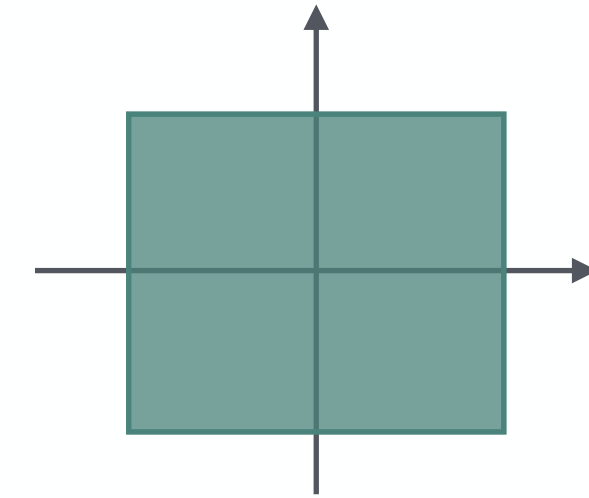
$$a, b \in \mathbb{R}$$

## Abstraction #1: Boxes Abstract Domain

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

I: 

x00: [-1, 1]
x01: [-1, 1]
x02: [-1, 1]
x03: [-1, 1]
x04: [-1, 1]
x05: [-1, 1]



③ check output for **inclusion**  
in **expected output**:  
**included** → **safe**  
otherwise → **alarm**

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

x10 -> [0, 6.14]    x11 -> [0, 3.29]    x12 -> [0, 5.02]

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

⋮ 

x20 -> [0, 25.30]	x21 -> [0, 27.34]	x22 -> [0, 18.63]
x30 -> [0, 118.99]	x31 -> [0, 155.15]	x32 -> [0, 162.18]

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

x40 -> [0, 1054.08]    x41 -> [0, 0]    x42 -> [0, 191.11]

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

O:  $x50 - x51 \sqsubset [0, \infty]$

$x50 - x51 \rightarrow [-6171.35, 5000]$



# Abstract Interpretation



# Safety Verification

## Abstraction #2: Symbolic Abstract Domain [Li19]

$$x \mapsto \begin{cases} E \\ [a, b] \quad a, b \in \mathbb{R} \end{cases}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$$\mathbf{I}: x_{00}: \begin{cases} x_{00} \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} x_{01} \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} x_{02} \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} x_{03} \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} x_{04} \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} x_{05} \\ [-1, 1] \end{cases}$$

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

## Abstraction #2: Symbolic Abstract Domain [Li19]

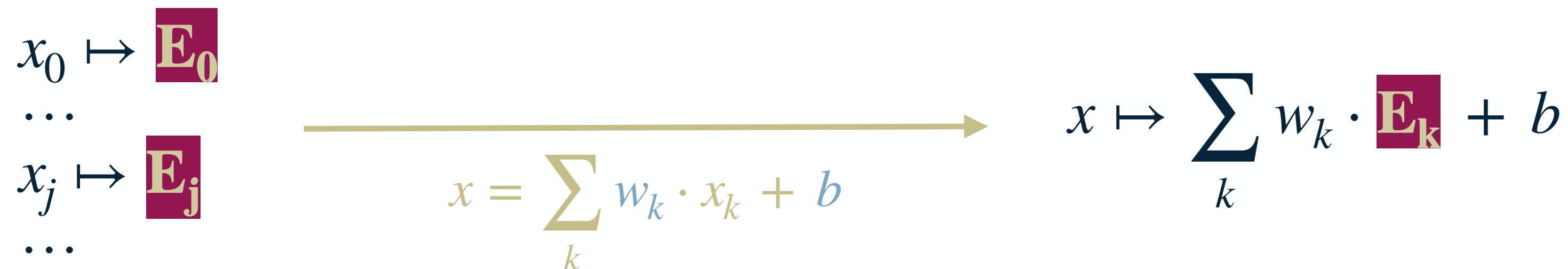
$$x \mapsto \begin{cases} E \\ [a, b] \quad a, b \in \mathbb{R} \end{cases}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$$\mathbf{I}: x_{00}: \begin{cases} x_{00} \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} x_{01} \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} x_{02} \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} x_{03} \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} x_{04} \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} x_{05} \\ [-1, 1] \end{cases}$$

```
x10' = (0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834)
```

$$x_{10}': \begin{cases} 0.12 * x_{00} + 0.07 * x_{01} + 0.10 * x_{02} + 2.03 * x_{03} + 0.10 * x_{04} - 2.10 * x_{05} + 1.62 \\ [-2.90, 6.14] \end{cases}$$



```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```



# Safety Verification

## Abstraction #2: Symbolic Abstract Domain [Li19]

$$x \mapsto \begin{cases} E \\ [a, b] \quad a, b \in \mathbb{R} \end{cases}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

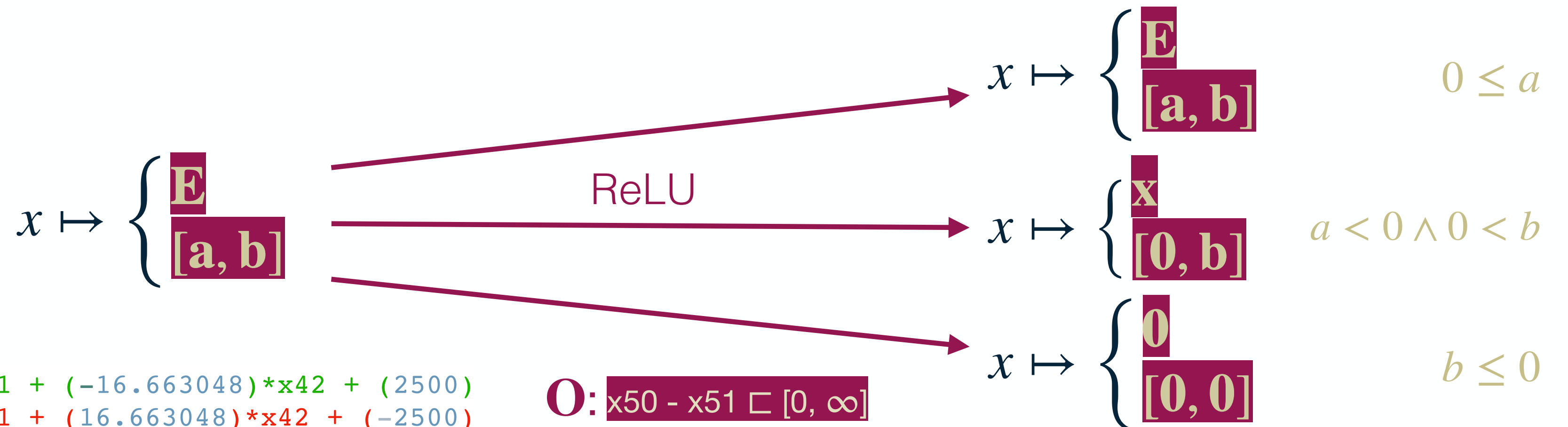
$$\mathbf{I}: x_{00}: \begin{cases} x_{00} \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} x_{01} \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} x_{02} \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} x_{03} \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} x_{04} \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} x_{05} \\ [-1, 1] \end{cases}$$

```
x10' = (0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834)
```

$$x_{10}': \begin{cases} 0.12 * x_{00} + 0.07 * x_{01} + 0.10 * x_{02} + 2.03 * x_{03} + 0.10 * x_{04} - 2.10 * x_{05} + 1.62 \\ [-2.90, 6.14] \end{cases}$$

```
x10 = ReLU(x10')
```

$$x_{10}: \begin{cases} x_{10} \\ [0, 6.14] \end{cases}$$



```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

$$\mathbf{O}: x_{50} - x_{51} \sqsubseteq [0, \infty]$$

# Safety Verification

## Abstraction #2: Symbolic Abstract Domain [Li19]

$$x \mapsto \begin{cases} E \\ [a, b] \quad a, b \in \mathbb{R} \end{cases}$$

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

$$\mathbf{I}: \begin{matrix} x00: & \begin{cases} x00 \\ [-1,1] \end{cases} & x01: & \begin{cases} x01 \\ [-1,1] \end{cases} & x02: & \begin{cases} x02 \\ [-1,1] \end{cases} & x03: & \begin{cases} x03 \\ [-1,1] \end{cases} & x04: & \begin{cases} x04 \\ [-1,1] \end{cases} & x05: & \begin{cases} x05 \\ [-1,1] \end{cases} \end{matrix}$$

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

$$x10: \begin{cases} x10 \\ [0, 6.14] \end{cases} \quad x11: \begin{cases} x11 \\ [0, 3.29] \end{cases} \quad x12: \begin{cases} x12 \\ [0, 5.02] \end{cases}$$

⋮

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

$$x40: \begin{cases} x40 \\ [0, 1054.08] \end{cases} \quad x41: \begin{cases} 0 \\ [0, 0] \end{cases} \quad x42: \begin{cases} x42 \\ [0, 191.11] \end{cases}$$

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

$$\mathbf{O}: x50 - x51: \begin{cases} -4.56 * x40 - 33.33 * x42 + 5000 \\ [-6171.35, 5000] \sqsubset [0, \infty] \end{cases}$$



# Safety Verification

## Abstraction #3: DeepPoly Abstract Domain [Singh19]

$$x \mapsto \begin{cases} [L, U] \\ [a, b] \end{cases} \quad a, b \in \mathbb{R}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$$\mathbf{I}: x_{00}: \begin{cases} [x_{00}, x_{00}] \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} [x_{01}, x_{01}] \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} [x_{02}, x_{02}] \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} [x_{03}, x_{03}] \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} [x_{04}, x_{04}] \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} [x_{05}, x_{05}] \\ [-1, 1] \end{cases}$$

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

## Abstraction #3: DeepPoly Abstract Domain [Singh19]

$$x \mapsto \begin{cases} [L, U] \\ [a, b] \end{cases} \quad a, b \in \mathbb{R}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

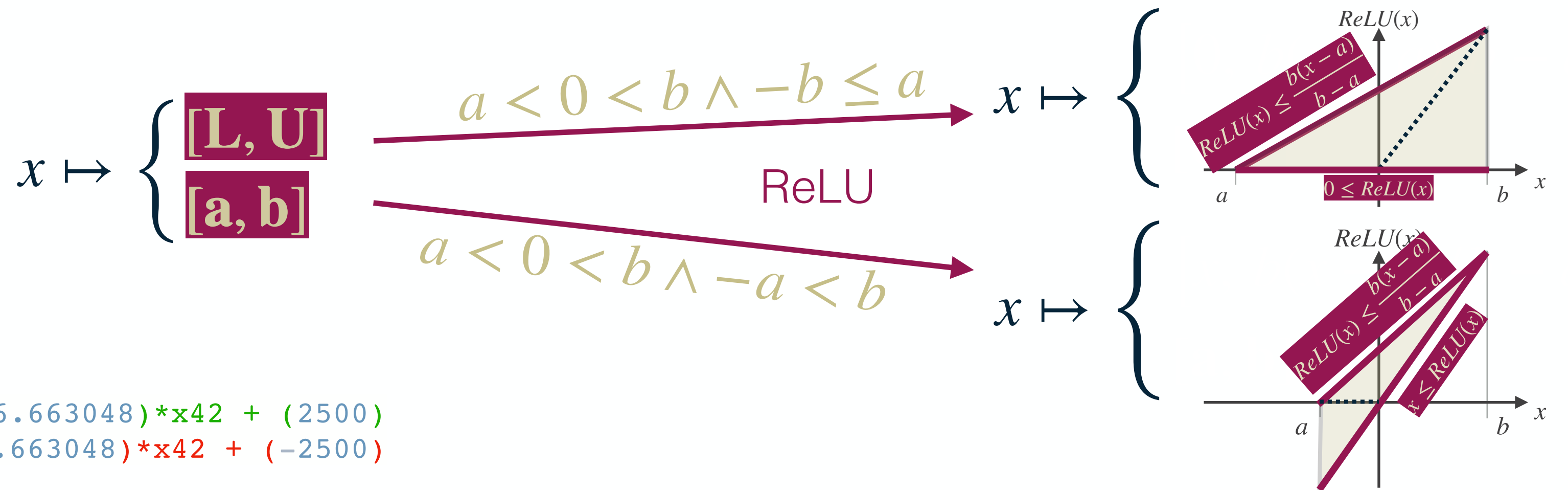
$$\mathbf{I}: x_{00}: \begin{cases} [x_{00}, x_{00}] \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} [x_{01}, x_{01}] \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} [x_{02}, x_{02}] \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} [x_{03}, x_{03}] \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} [x_{04}, x_{04}] \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} [x_{05}, x_{05}] \\ [-1, 1] \end{cases}$$

```
x10' = (0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834)
```

$$x_{10}': \begin{cases} \dots [L, U] \\ [-2.90, 6.14] \end{cases}$$

```
x10 = ReLU(x10')
```

$$x_{10}: \begin{cases} [x_{10}', 0.68 * x_{10}' + 1.97] \\ [-2.90, 6.14] \end{cases}$$



```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

## Abstraction #3: DeepPoly Abstract Domain [Singh19]

$$x \mapsto \begin{cases} [L, U] \\ [a, b] \end{cases} \quad a, b \in \mathbb{R}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

```
x10 = ReLU((0.120875)*x00 +
x11 = ReLU((0.113805)*x00 +
x12 = ReLU((0.755487)*x00 +
```

$$x_{10}: \begin{cases} [x_{10}', 0.68 * x_{10}' + 1.9] \\ [-2.90, 6.14] \end{cases}$$

⋮

```
x40 = ReLU((2.296390)*x30 +
x41 = ReLU((-0.552155)*x30
x42 = ReLU((-2.509773)*x30
```

$$x_{40}: \begin{cases} [x_{40}', 0.67 * x_{40}' + 313] \\ [-467.10, 950.38] \end{cases}$$

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

# Safety Verification

## Abstraction #2: Symbolic Abstract Domain [Li19]

$$x \mapsto \begin{cases} E \\ [a, b] \end{cases} \quad a, b \in \mathbb{R}$$

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$$I: x_{00}: \begin{cases} x_{00} \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} x_{01} \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} x_{02} \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} x_{03} \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} x_{04} \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} x_{05} \\ [-1, 1] \end{cases}$$

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

$$x_{10}: \begin{cases} x_{10} \\ [0, 6.14] \end{cases} \quad x_{11}: \begin{cases} x_{11} \\ [0, 3.29] \end{cases} \quad x_{12}: \begin{cases} x_{12} \\ [0, 5.02] \end{cases}$$

⋮

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

$$x_{40}: \begin{cases} x_{40} \\ [0, 1054.08] \end{cases} \quad x_{41}: \begin{cases} 0 \\ [0, 0] \end{cases} \quad x_{42}: \begin{cases} x_{42} \\ [0, 191.11] \end{cases}$$

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

$$O: x_{50} - x_{51}: \begin{cases} -4.56 * x_{40} - 33.33 * x_{42} + 5000 \\ [-6171.35, 5000] \sqsubset [0, \infty] \end{cases} \quad \text{🚨}$$

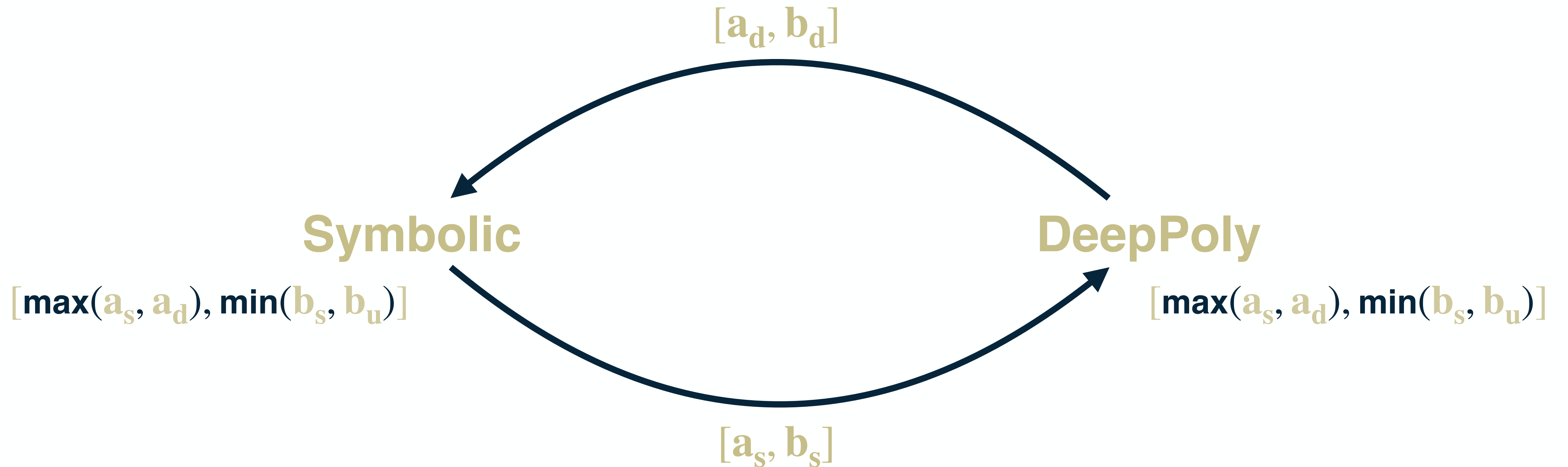
28

$$x_{10}: \begin{cases} x_{10} \\ [0, 118.63] \end{cases} \quad x_{11}: \begin{cases} x_{11} \\ [0, 118.63] \end{cases} \quad x_{12}: \begin{cases} x_{12} \\ [-142.20, 162.09] \end{cases}$$

$$O: x_{50} - x_{51}: \begin{cases} \dots \\ [-1424.80, 9072.12] \sqsubset [0, \infty] \end{cases} \quad \text{🚨}$$

# Reduced Product Domain

## Symbolic Abstract Domain & DeepPoly Abstract Domain



# Safety Verification

## Abstraction #4: Symbolic & DeepPoly Product Abstract Domain

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

$$\mathbf{I}: x_{00}: \begin{cases} x_{00} \\ [x_{00}, x_{00}] \\ [-1, 1] \end{cases} \quad x_{01}: \begin{cases} x_{01} \\ [x_{01}, x_{01}] \\ [-1, 1] \end{cases} \quad x_{02}: \begin{cases} x_{02} \\ [x_{02}, x_{02}] \\ [-1, 1] \end{cases} \quad x_{03}: \begin{cases} x_{03} \\ [x_{03}, x_{03}] \\ [-1, 1] \end{cases} \quad x_{04}: \begin{cases} x_{04} \\ [x_{04}, x_{04}] \\ [-1, 1] \end{cases} \quad x_{05}: \begin{cases} x_{05} \\ [x_{05}, x_{05}] \\ [-1, 1] \end{cases}$$

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

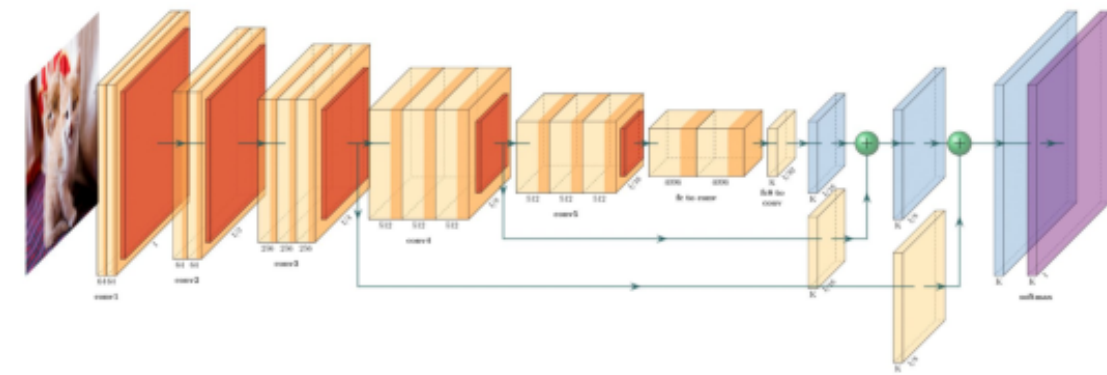
```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (2500)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

$$\mathbf{O}: x_{50} - x_{51}: \begin{cases} \vdots \\ [670.04, 5000.0] \sqsubseteq [0, \infty] \end{cases}$$



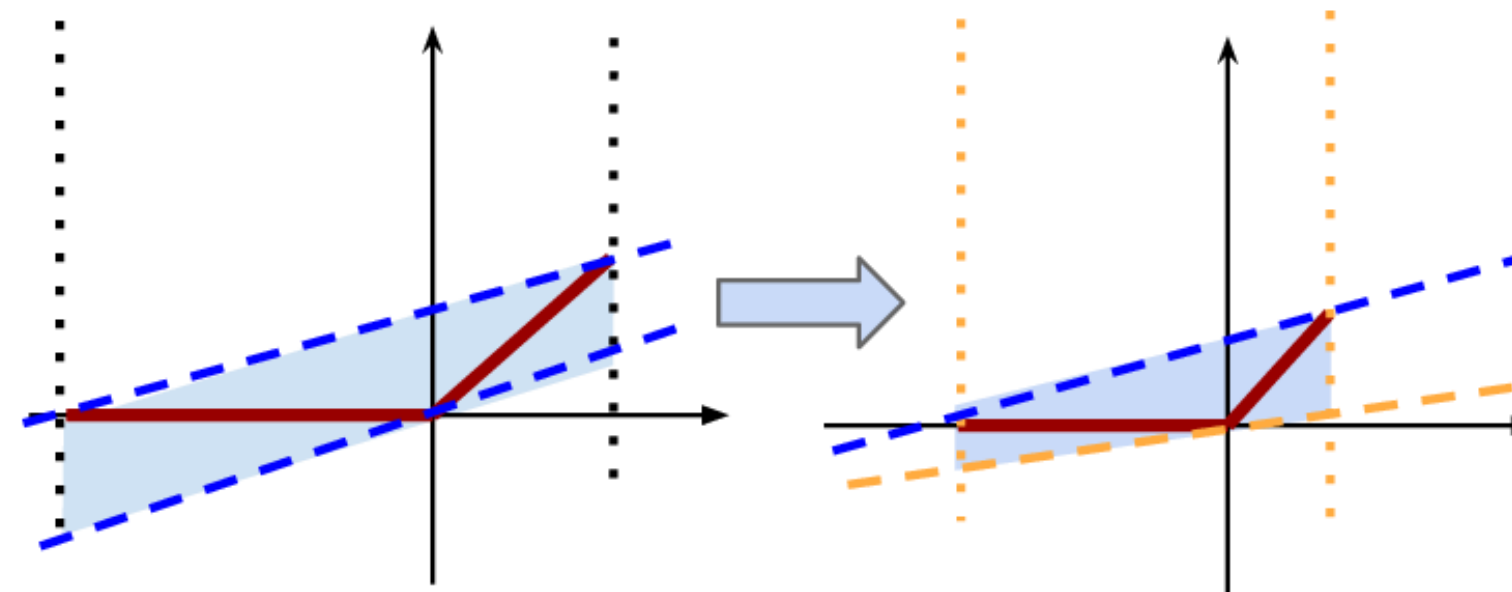
# Safety Verification

## Going Farther: $\alpha\beta$ -CROWN

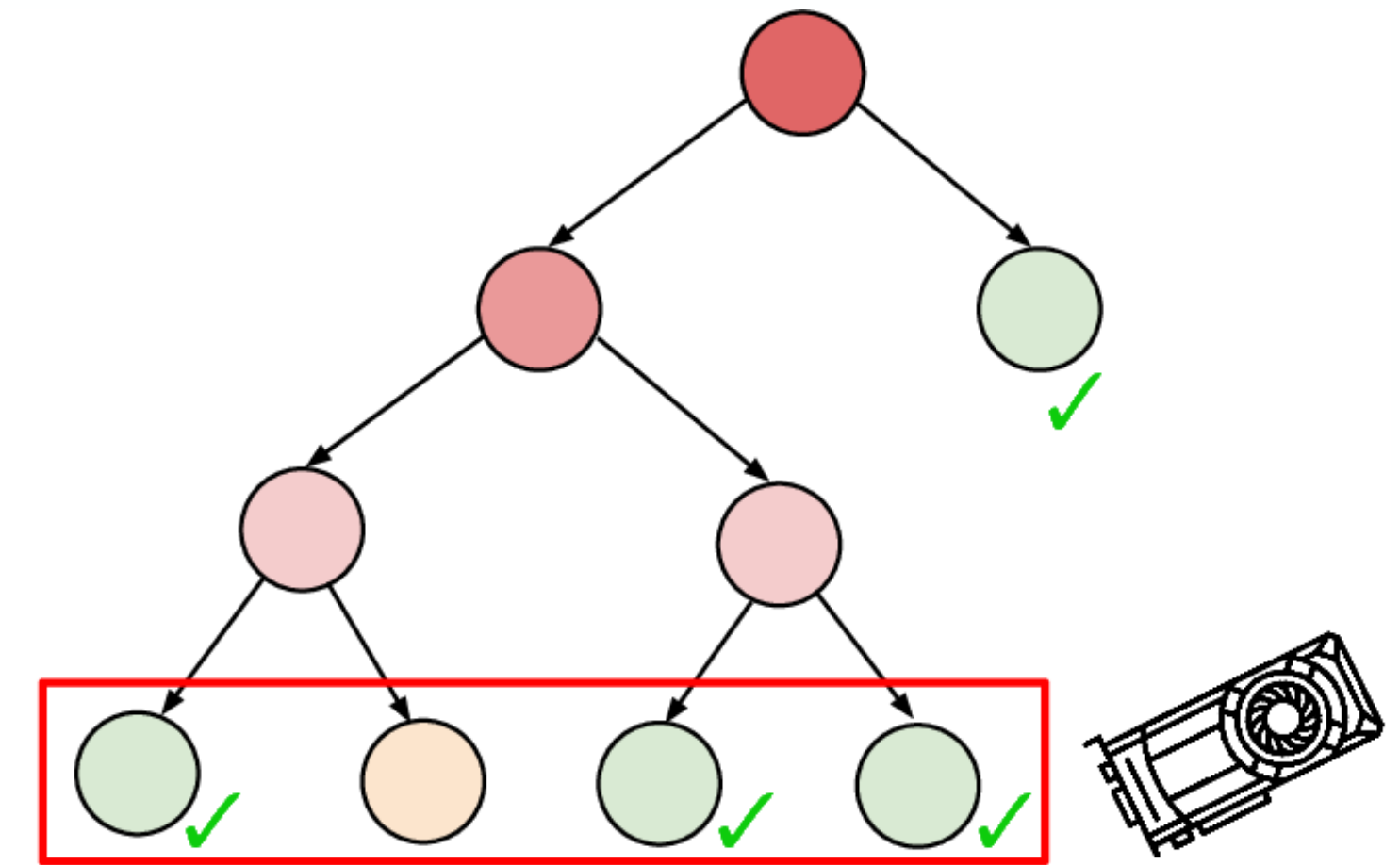


$$\min_{x \in \mathcal{C}} f(x) \geq \min_{x \in \mathcal{C}} \mathbf{a}^\top x + c$$

Efficient bound propagation (**CROWN**)



GPU optimized relaxation ( $\alpha$ -**CROWN**)



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



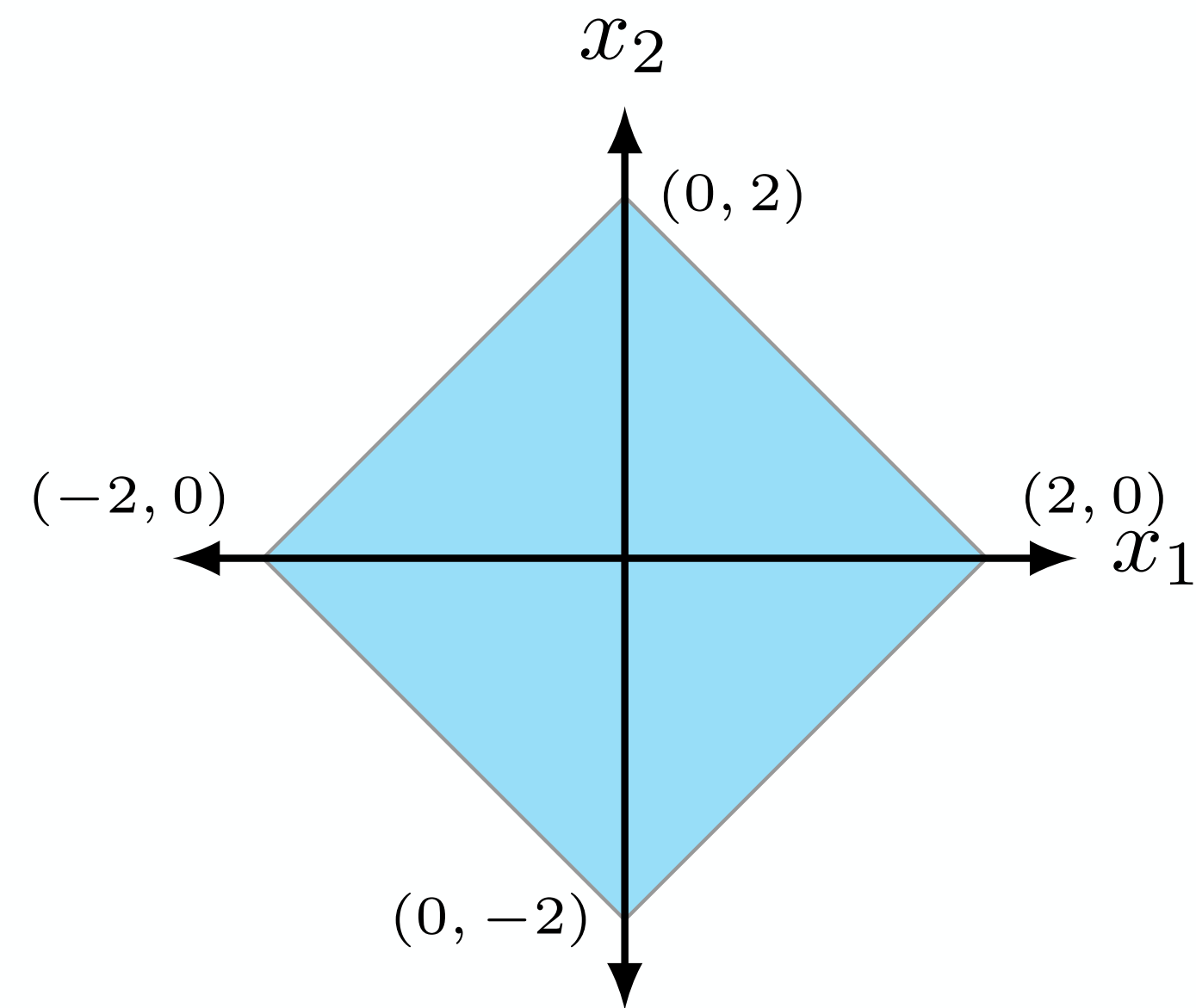
Winner of the International Verification of Neural Networks Competition since 2021

<https://github.com/Verified-Intelligence/alpha-beta-CROWN>

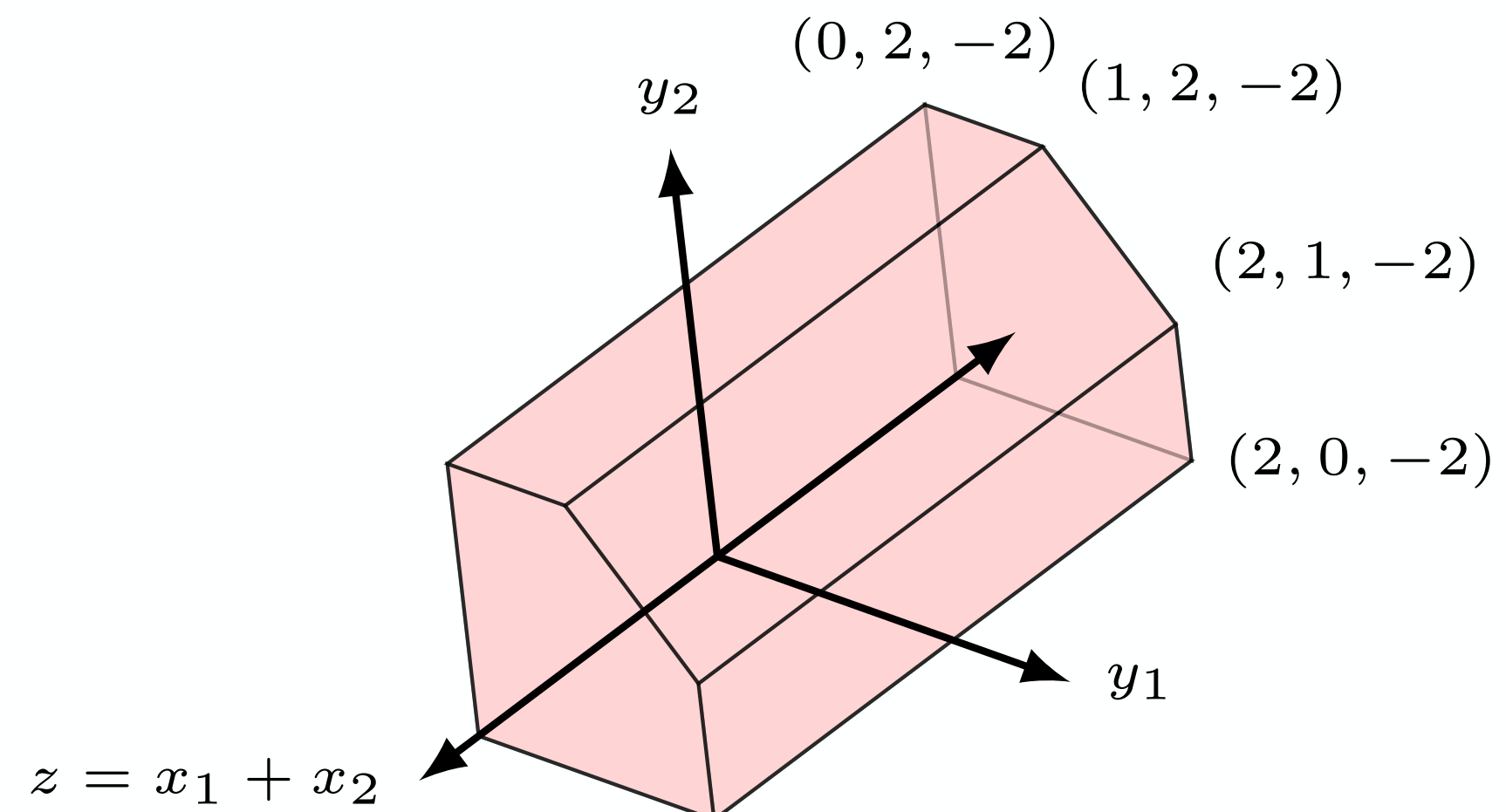


# Safety Verification

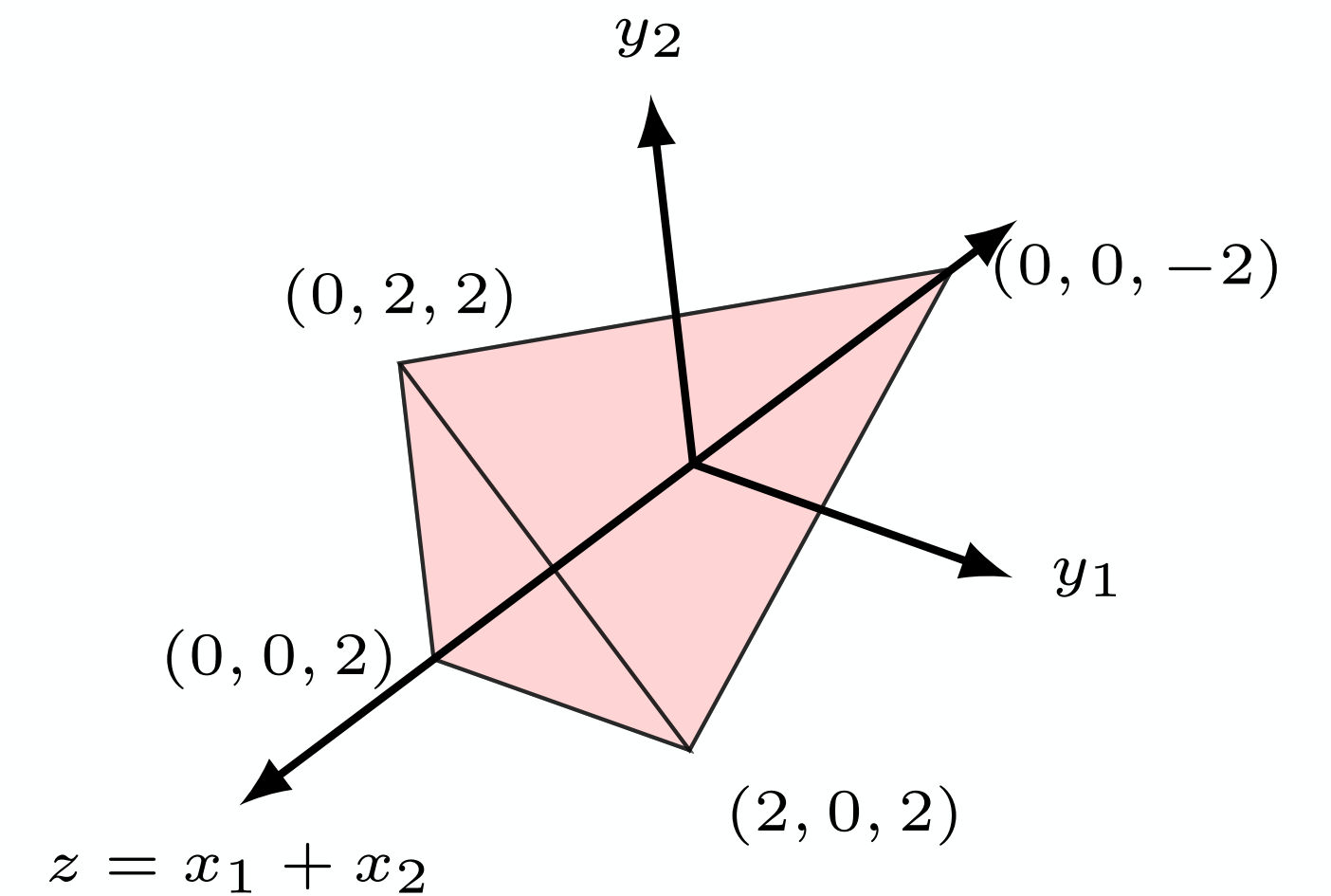
## Going Farther: Multi-Neuron Abstractions



(a) Input shape

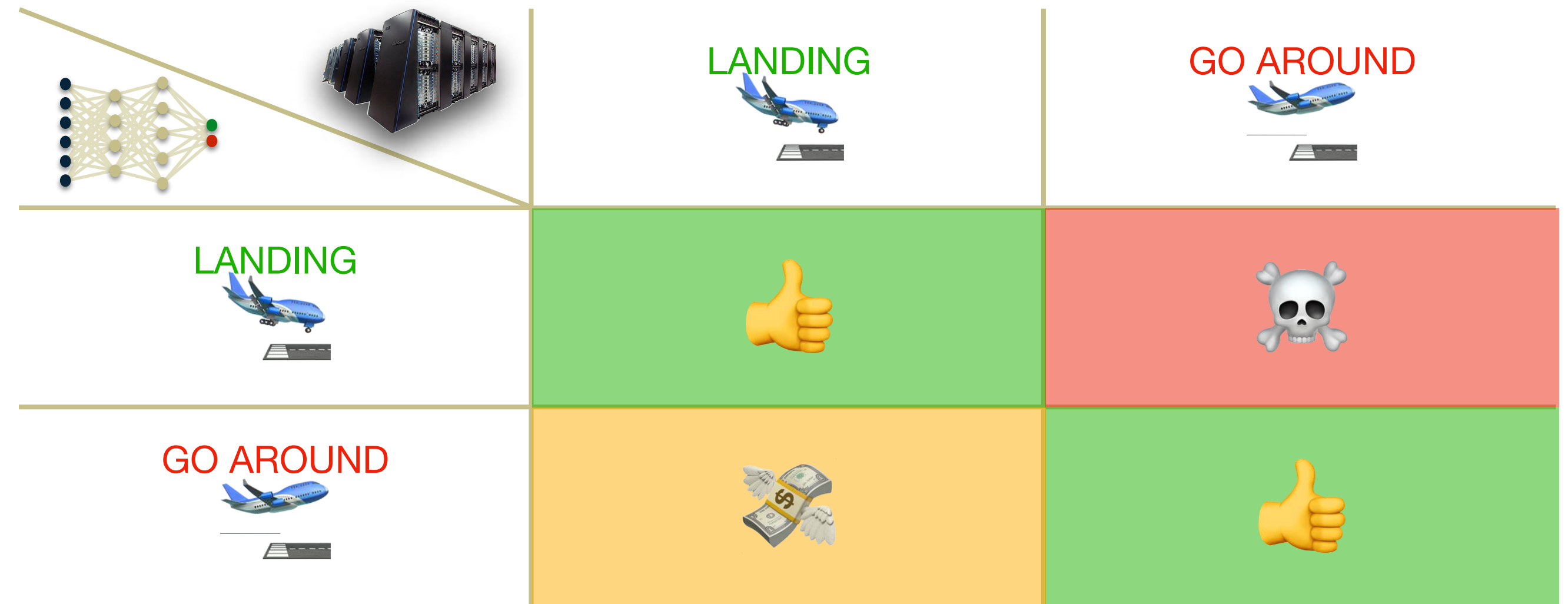


(b) 1-ReLU

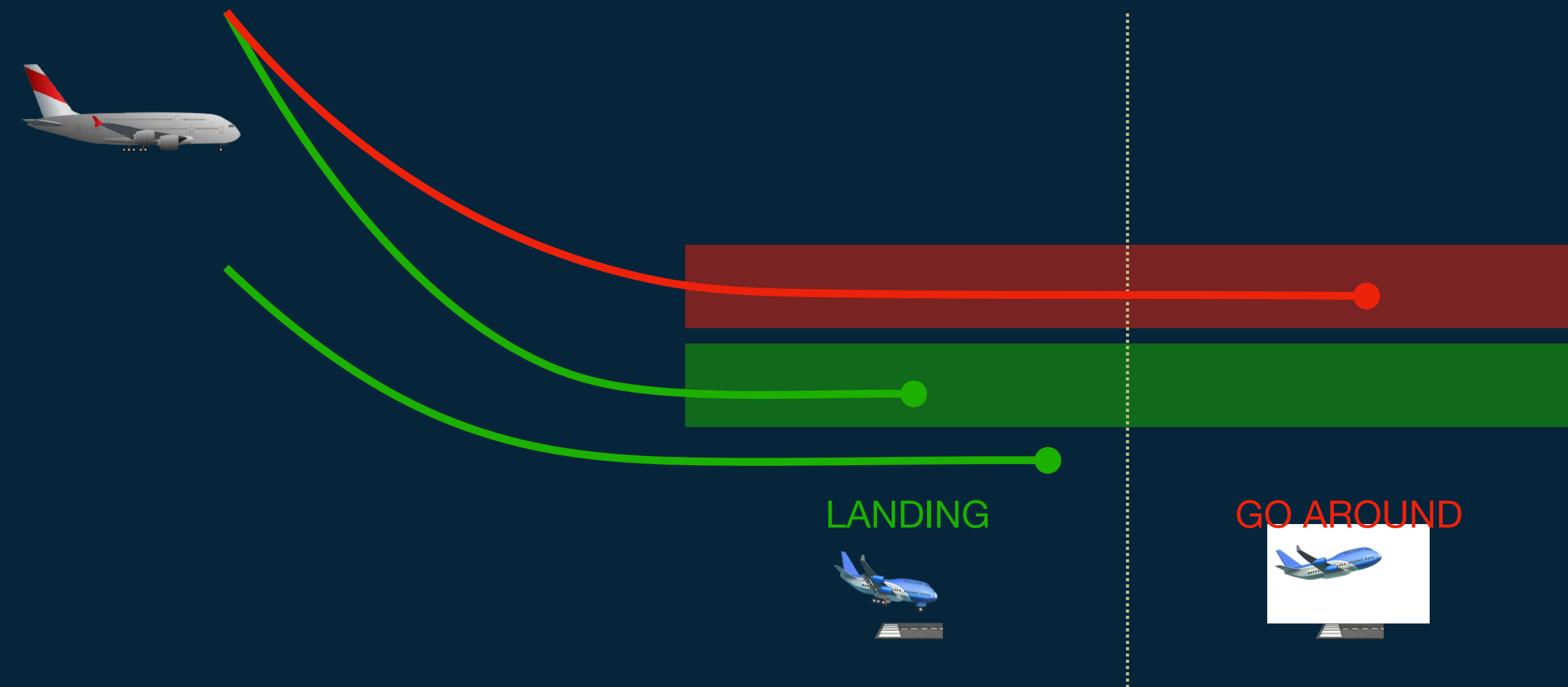


(c) 2-ReLU

# Safety

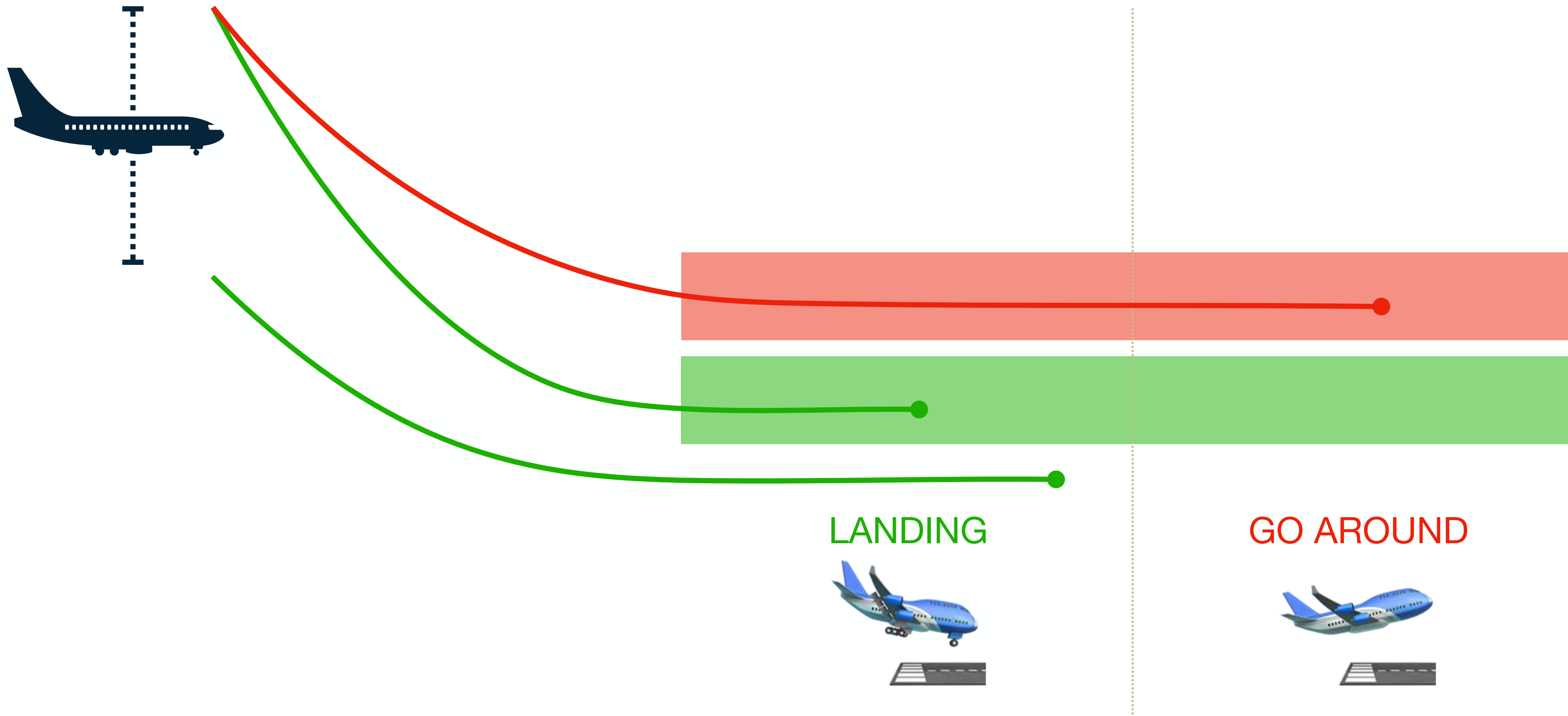


# Hypersafety



# Runway Overrun Warning

## HyperSafety of Neural Network Surrogate



# Hyperproperty Verification

## Abstract Non-Interference Properties

$\eta$ : input abstraction

$\rho$ : output abstraction

$$\mathcal{H} \stackrel{\text{def}}{=} \left\{ T \mid \forall t, t' \in T: \eta(t_0) = \eta(t'_0) \Rightarrow \rho(t_\omega) = \rho(t'_\omega) \right\}$$

$\mathcal{H}$  is the set of all forward pass traces that **satisfy** abstract non-interference with respect to  $\eta$  and  $\rho$

Theorem

$$M \models \mathcal{H} \Leftrightarrow \llbracket M \rrbracket \in \mathcal{H} \Leftrightarrow \{\llbracket M \rrbracket\} \subseteq \mathcal{H}$$

Corollary

$$M \models \mathcal{H} \Leftarrow \{\llbracket M \rrbracket\} \subseteq \llbracket M \rrbracket^\sharp \subseteq \mathcal{H}$$

Giacobazzi and Mastroeni. Abstract Non-Interference: A Unifying Framework for Weakening Information-Flow. In TOPS, 2018.

# Abstract Non-Interference Verification

## Example

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

$\eta$ :

ALTITUDE

$\eta(x00) = x00$
$\eta(x01) = x01$
$\eta(x02) = \top$
$\eta(x03) = x03$
$\eta(x04) = x04$
$\eta(x05) = x05$

“the risk of a runway overrun does not change when only varying the altitude at which it is measured (in the expected range) and nothing else”

$\rho$ :

$\rho(x50) = 1$ if $x50 > x51$ else $0$
$\rho(x51) = 1$ if $x51 > x50$ else $0$

# Abstract Interpretation

## 3-Step Recipe

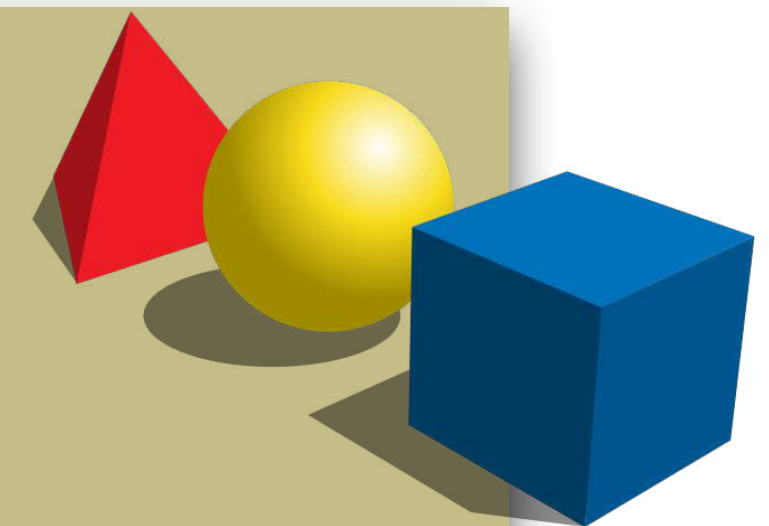
**practical tools**

targeting specific programs



**abstract semantics, abstract domains**

algorithmic approaches to decide program properties



**concrete semantics**

**mathematical models** of the program behavior



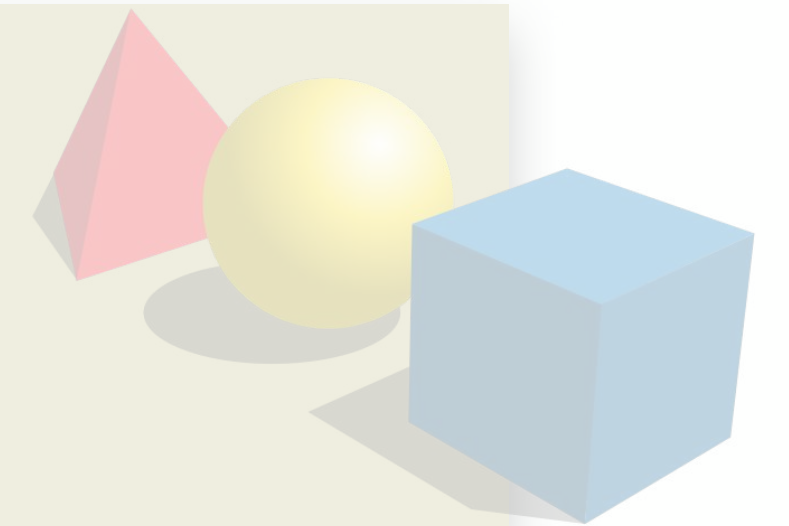
# Abstract Interpretation

## 3-Step Recipe

**practical tools**  
targeting specific programs



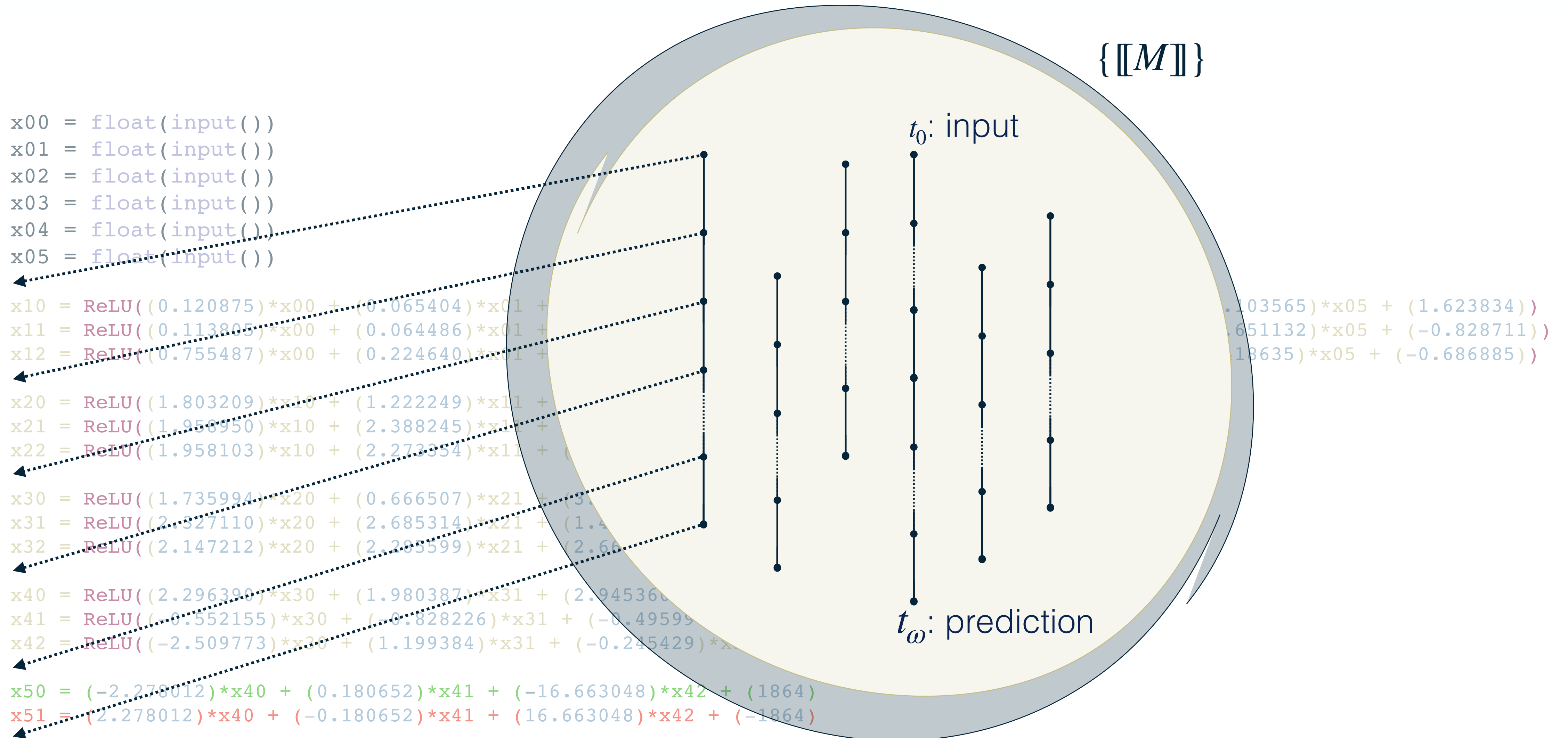
**abstract semantics, abstract domains**  
**algorithmic approaches** to decide program properties



**concrete semantics**  
**mathematical models** of the program behavior



# Collecting Semantics





# Dependency Semantics

```

x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())

```

```

x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.064486)*x02 + (0.064486)*x03 + (0.064486)*x04 + (0.064486)*x05)
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.064486)*x02 + (0.064486)*x03 + (0.064486)*x04 + (0.064486)*x05)
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.224640)*x02 + (0.224640)*x03 + (0.224640)*x04 + (0.224640)*x05)

```

```

x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (1.222249)*x12)
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.388245)*x12)
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (2.273354)*x12)

```

```

x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (0.666507)*x22)
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (2.685314)*x22)
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.285599)*x22)

```

```

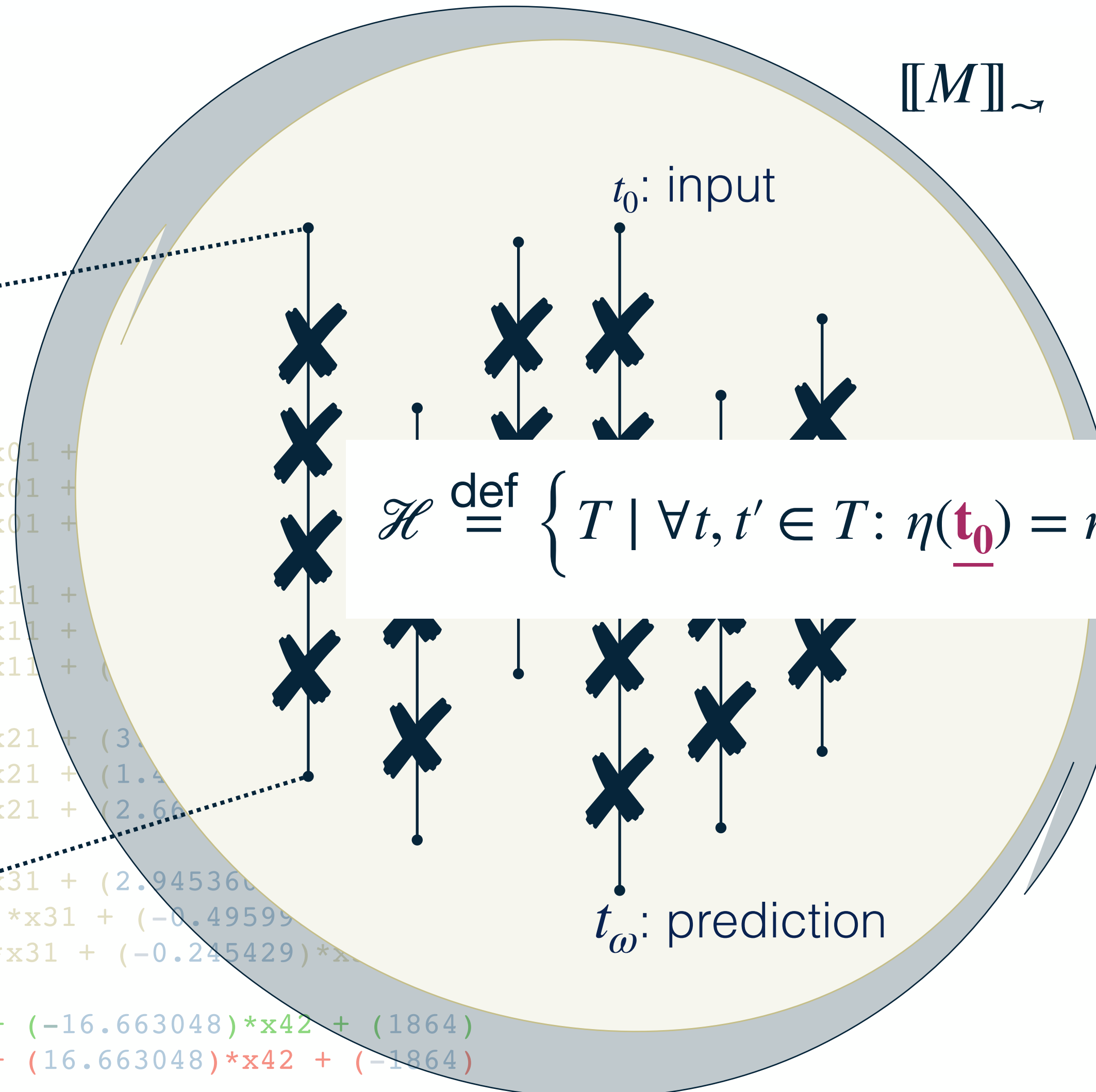
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (1.980387)*x32)
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.828226)*x32)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (1.199384)*x32)

```

```

x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)

```



# Parallel Semantics

```

x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())

```

```

x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.065404)*x02 + (0.065404)*x03 + (0.065404)*x04 + (0.065404)*x05)
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.064486)*x02 + (0.064486)*x03 + (0.064486)*x04 + (0.064486)*x05)
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.224640)*x02 + (0.224640)*x03 + (0.224640)*x04 + (0.224640)*x05)

```

```

x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (1.222249)*x12 + (1.222249)*x13 + (1.222249)*x14 + (1.222249)*x15)
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.388245)*x12 + (2.388245)*x13 + (2.388245)*x14 + (2.388245)*x15)
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (2.273354)*x12 + (2.273354)*x13 + (2.273354)*x14 + (2.273354)*x15)

```

```

x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (0.666507)*x22 + (0.666507)*x23 + (0.666507)*x24 + (0.666507)*x25)
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (2.685314)*x22 + (2.685314)*x23 + (2.685314)*x24 + (2.685314)*x25)
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.285599)*x22 + (2.285599)*x23 + (2.285599)*x24 + (2.285599)*x25)

```

```

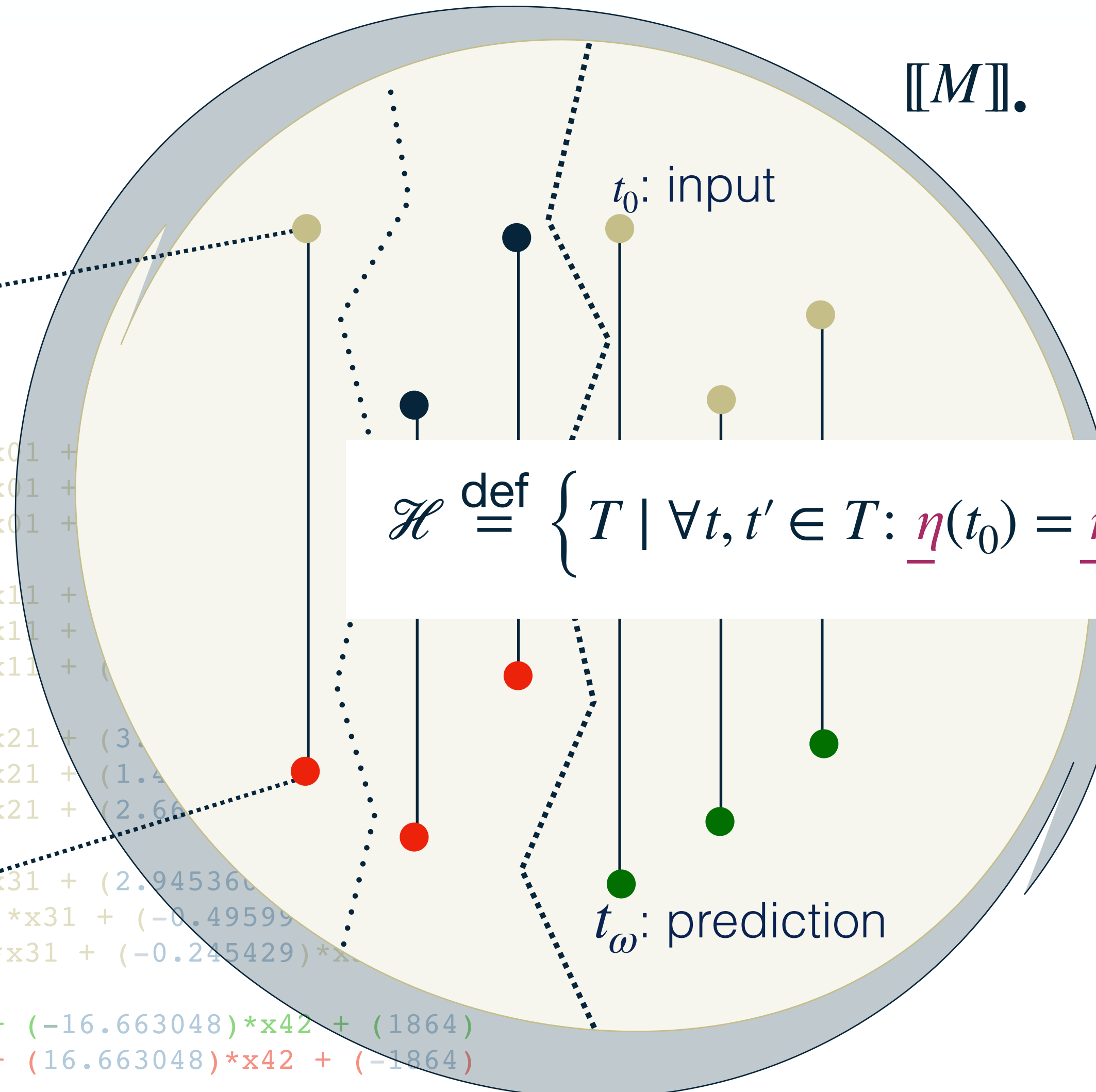
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (1.980387)*x32 + (1.980387)*x33 + (1.980387)*x34 + (1.980387)*x35)
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.828226)*x32 + (-0.828226)*x33 + (-0.828226)*x34 + (-0.828226)*x35)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (1.199384)*x32 + (1.199384)*x33 + (1.199384)*x34 + (1.199384)*x35)

```

```

x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)

```



$$\mathcal{H} \stackrel{\text{def}}{=} \left\{ T \mid \forall t, t' \in T: \underline{\eta}(t_0) = \underline{\eta}(t'_0) \Rightarrow \underline{\rho}(t_\omega) = \underline{\rho}(t'_\omega) \right\}$$

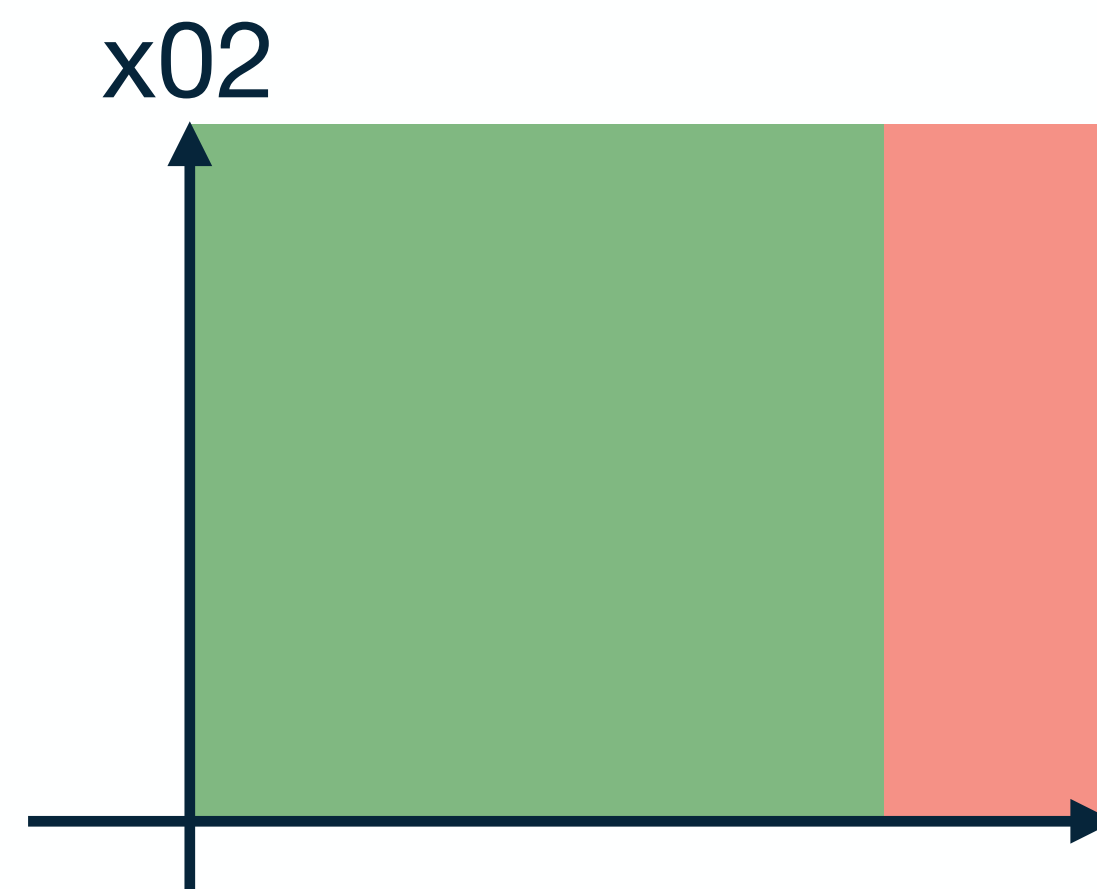
# Hyperproperty Verification

## Abstract Non-Interference Properties

$$\mathcal{H} \stackrel{\text{def}}{=} \left\{ T \mid \forall t, t' \in T: \eta(t_0) = \eta(t'_0) \Rightarrow \rho(t_\omega) = \rho(t'_\omega) \right\}$$

Lemma

$$M \models \mathcal{H} \Leftrightarrow \forall A, B \in \llbracket M \rrbracket.: \rho(A_\omega) \neq \rho(B_\omega) \Rightarrow \eta(A_0) \neq \eta(B_0)$$



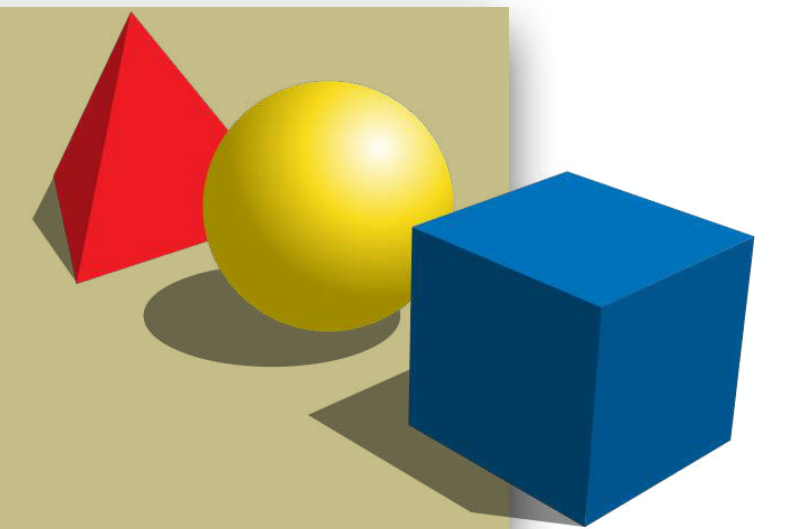
# Abstract Interpretation

## 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



# Hyperproperty Verification [Urban20]

## Naïve Static Backward Analysis

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

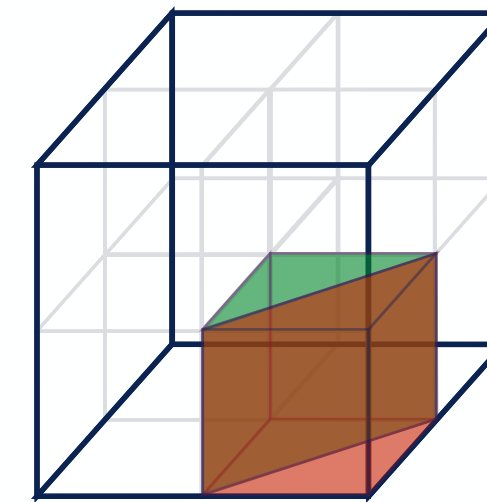
```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

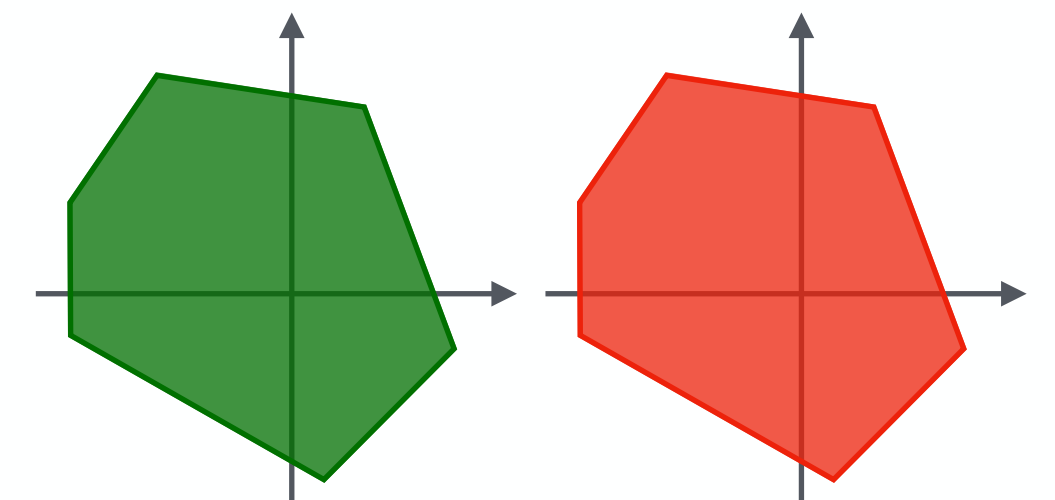
```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```



- ① check for **disjunction** in corresponding **input partitions**:  
**disjoint** → ✓ **safe**  
otherwise → 🚨 **alarm**

- ① start from an **abstraction** for each possible classification outcome

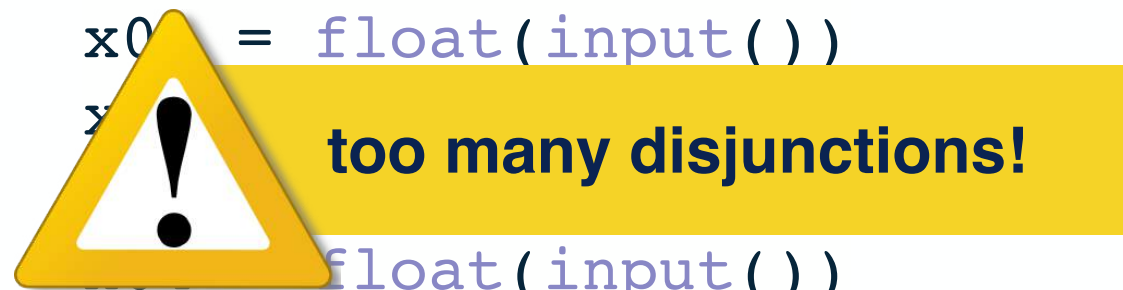


# Naïve Static Backward Analysis

## Disjunctive Polyhedra Abstract Domain

```
x00 = float(input())
```

```
x01 = float(input())
```



```
x02 = float(input())
```

```
4096 x05 = float(input())
```

```
4096 x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
```

```
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
```

```
512 x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
512 x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
```

```
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
```

```
64 x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
64 x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
```

```
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
```

```
8 x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

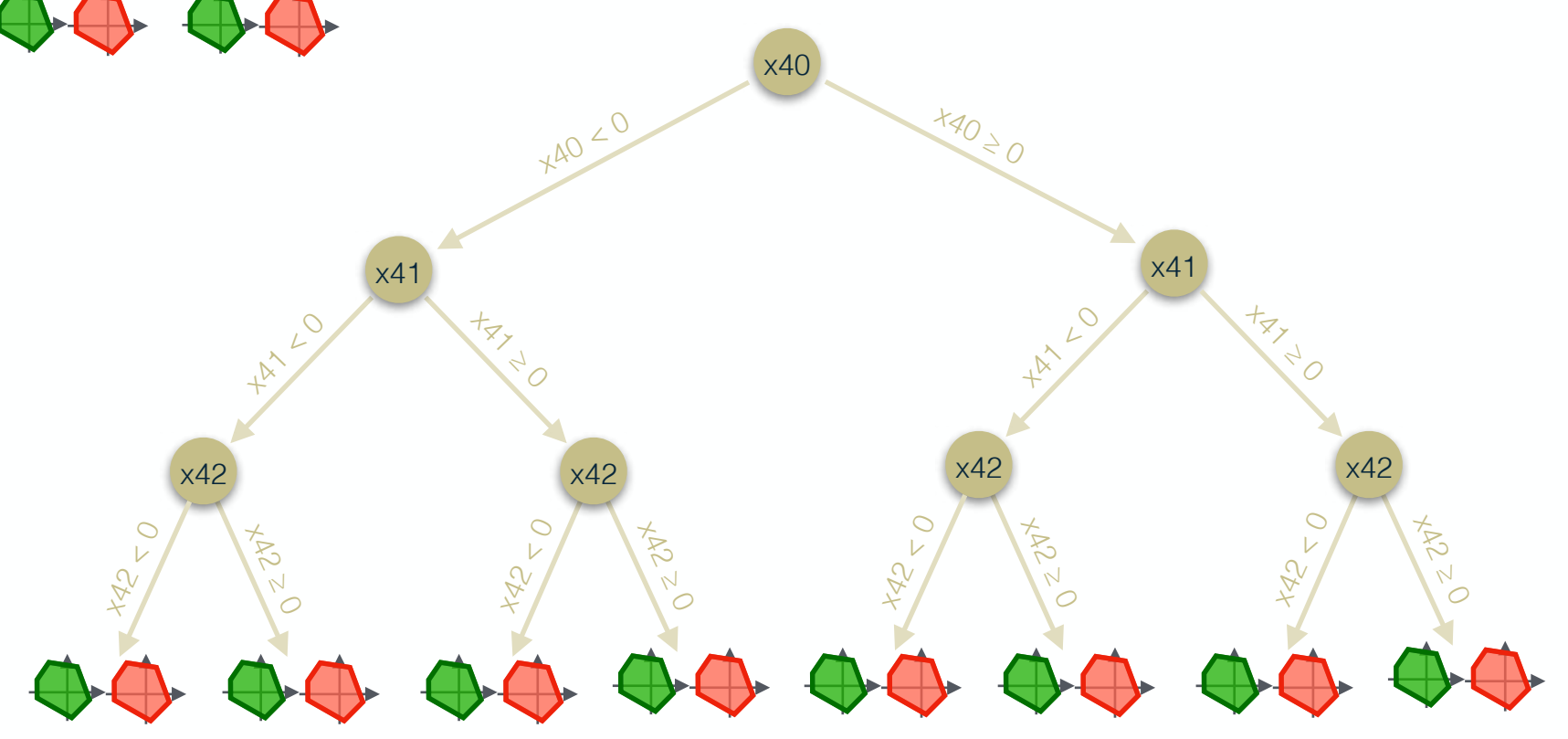
```
8 x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
```

```
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
```

```
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
```

```
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```



$(-4.556024) * x_{40} + (0.361304) * x_{41} + (-33.326096) * x_{42} + (3728) > 0$   
 $(4.556024) * x_{40} + (33.326096) * x_{42} - 3728 > 0$

# Hyperproperty Verification [Urban20]

## Static Forward Analysis

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

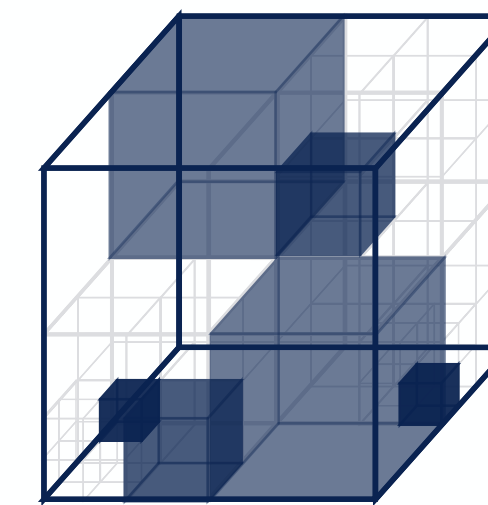
```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

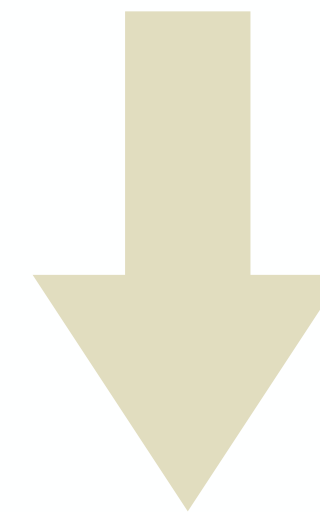
```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

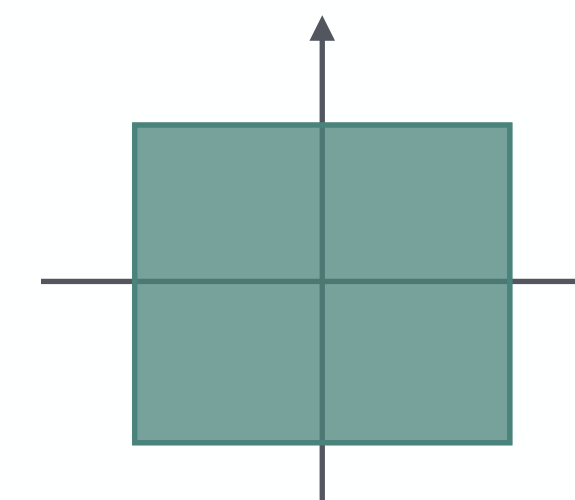
```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```



① start from a **partition** of the input space



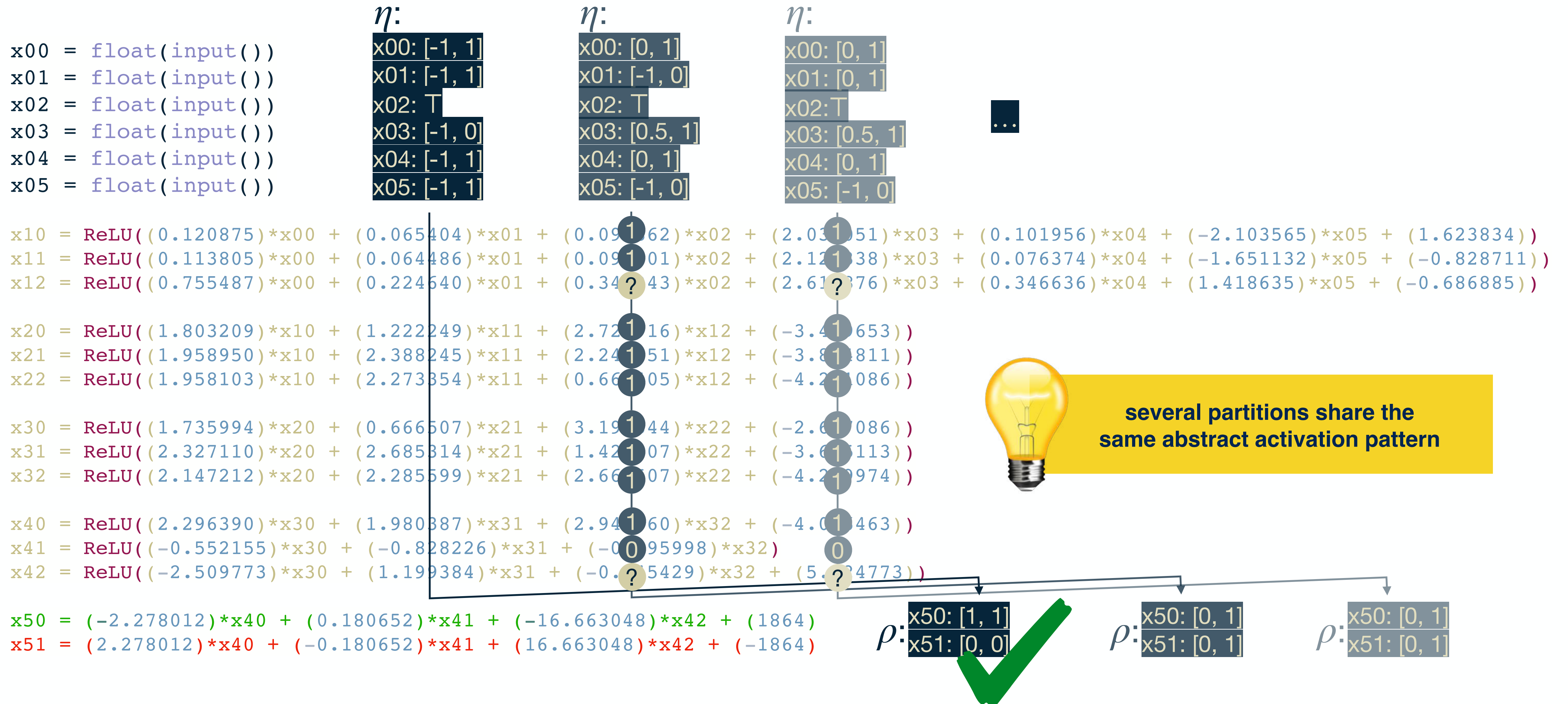
② proceed **forwards in parallel** from all partitions



③ check output for:  
- **unique classification outcome** → ✓ **safe**  
- **abstract activation pattern**

# Static Forward Analysis

## Symbolic & DeepPoly Product Abstract Domain





# Hyperproperty Verification [Urban20]

## Static Backward Analysis

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

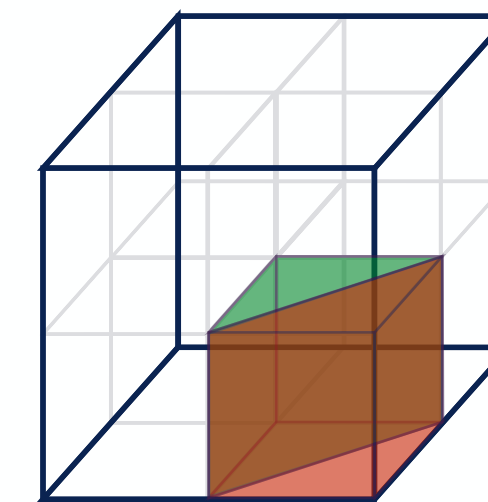
```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

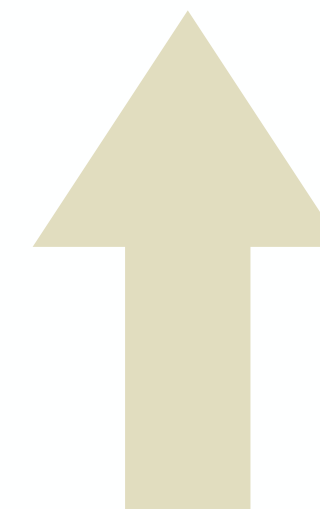
```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

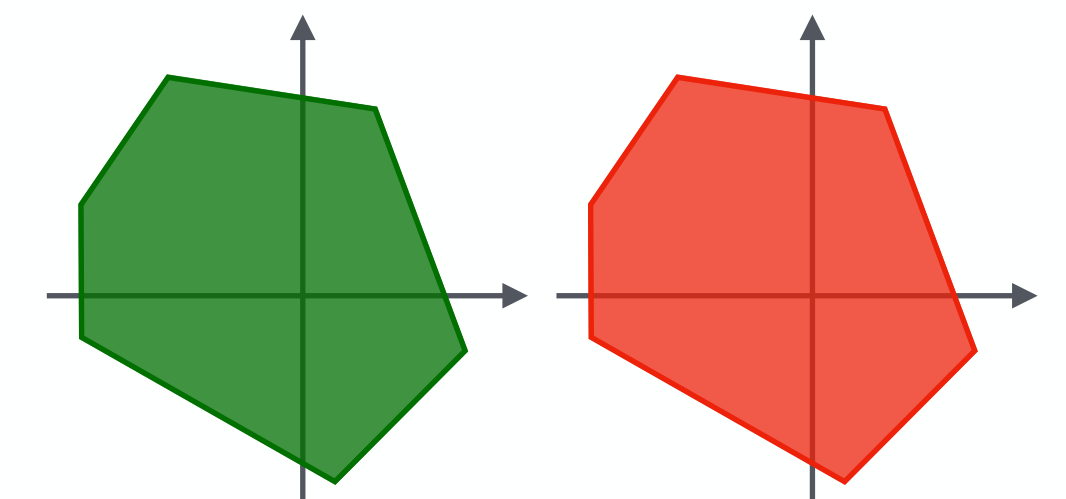


- ① check for **disjunction** in corresponding **input partitions**:  
**disjoint** → ✓ **safe**  
otherwise → 🚨 **alarm**



- ② proceed **backwards** in parallel **for each abstract activation pattern**

- ① start from an **abstraction** for each possible classification outcome



# Static Backward Analysis

## Symbolic & DeepPoly Product Abstract Domain

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$\eta$ :

```
x00: [0, 1]
x01: [-1, 0]
x02: T
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

$\eta$ :

```
x00: [0, 1]
x01: [0, 1]
x02: T
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

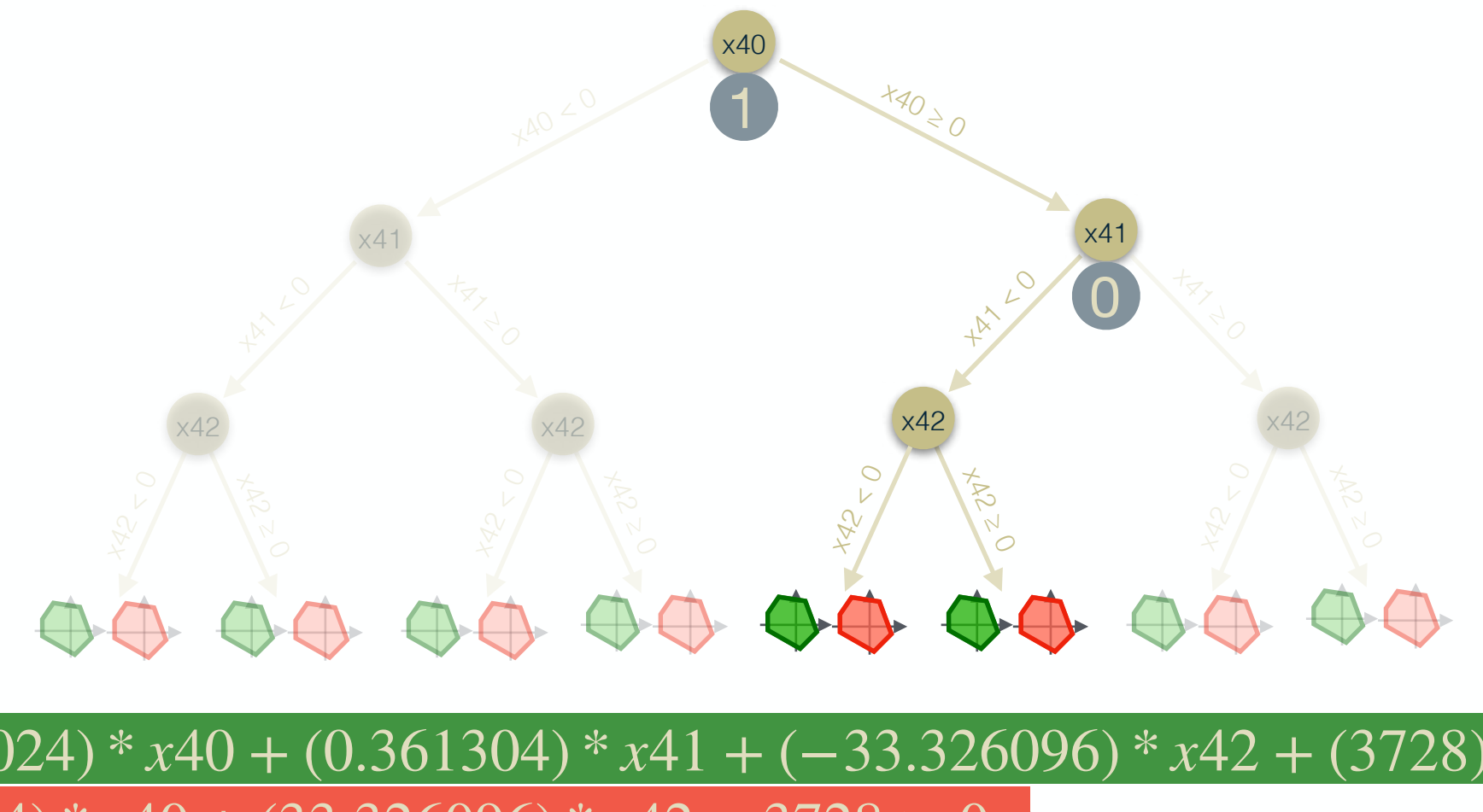
```
1 x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
1 x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
? x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
1 x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
1 x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
1 x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
1 x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))
1 x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))
1 x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
1 x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
0 x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
? x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```



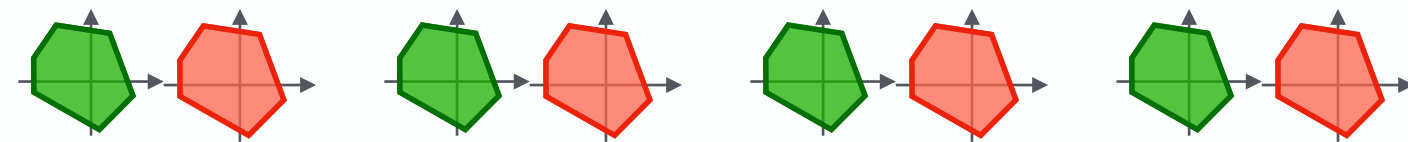
$(-4.556024) * x40 + (0.361304) * x41 + (-33.326096) * x42 + (3728) > 0$

$(4.556024) * x40 + (33.326096) * x42 - 3728 > 0$

# Static Backward Analysis

## Symbolic & DeepPoly Product Abstract Domain

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```



$\eta$ :

```
x00: [0, 1]
x01: [-1, 0]
x02: T
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

$\eta$ :

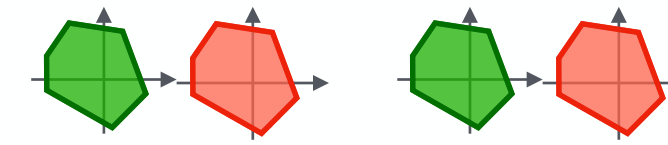
```
x00: [0, 1]
x01: [0, 1]
x02: T
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

counterexample

x00: 1	x00: 1
x01: 1	x01: 1
x02: -1	x02: 1
x03: 1	x03: 1
x04: 1	x04: 1
x05: -1	x05: -1

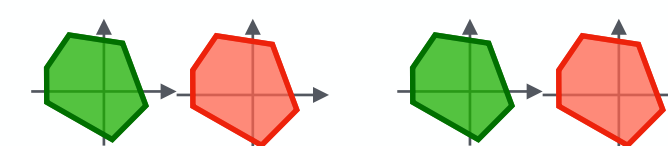
```
1 x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
1 x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
? x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
1 x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
1 x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
1 x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```



⋮

```
1 x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
0 x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
? x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```



```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

$(-4.556024) * x40 + (0.361304) * x41 + (-33.326096) * x42 + (3728) > 0$

$(4.556024) * x40 + (33.326096) * x42 - 3728 > 0$

# Abstract Interpretation

## 3-Step Recipe

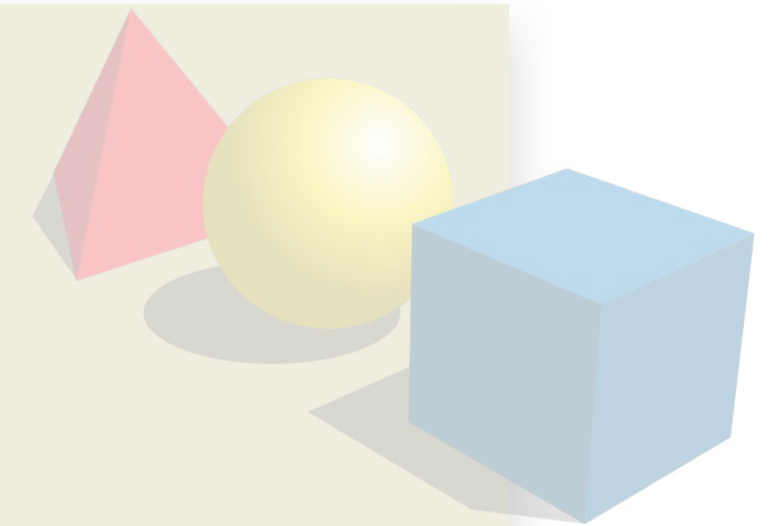
### practical tools

targeting specific programs



### abstract semantics, abstract domains

algorithmic approaches to decide program properties



### concrete semantics

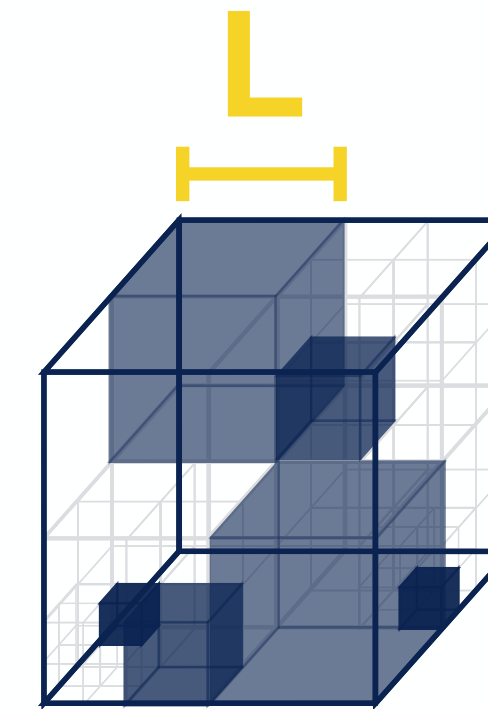
mathematical models of the program behavior



# Hyperproperty Verification [Urban20]

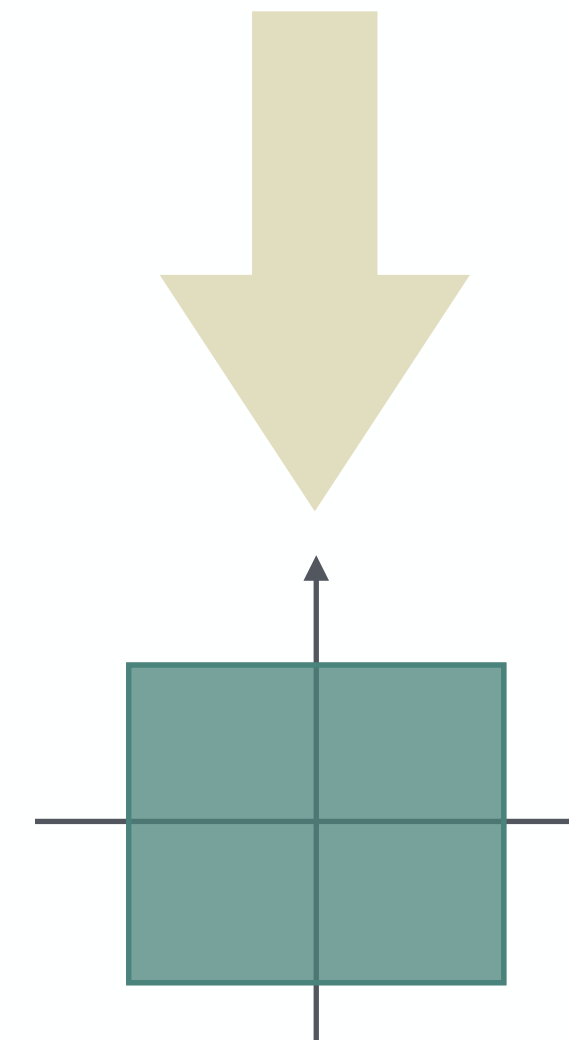
## Static Forward Analysis

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```



① **iteratively** partition the input space

```
① x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
① x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
? x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))  
  
? x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
? x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
? x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))  
  
? x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
① x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
① x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))  
  
① x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
① x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
① x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))  
  
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```



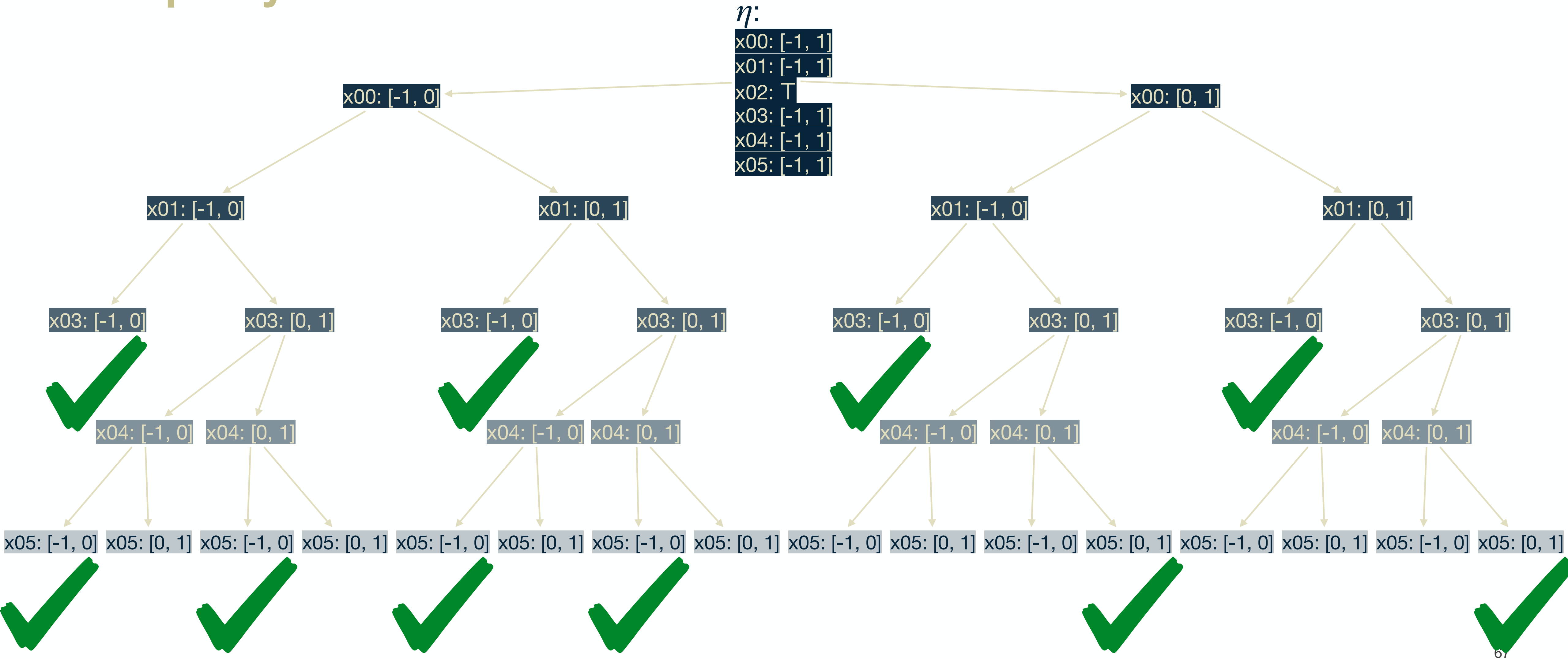
② proceed **forwards in parallel** from all partitions

③ check output for:  
- **unique classification outcome** → ✓ **safe**  
- **abstract activation pattern**



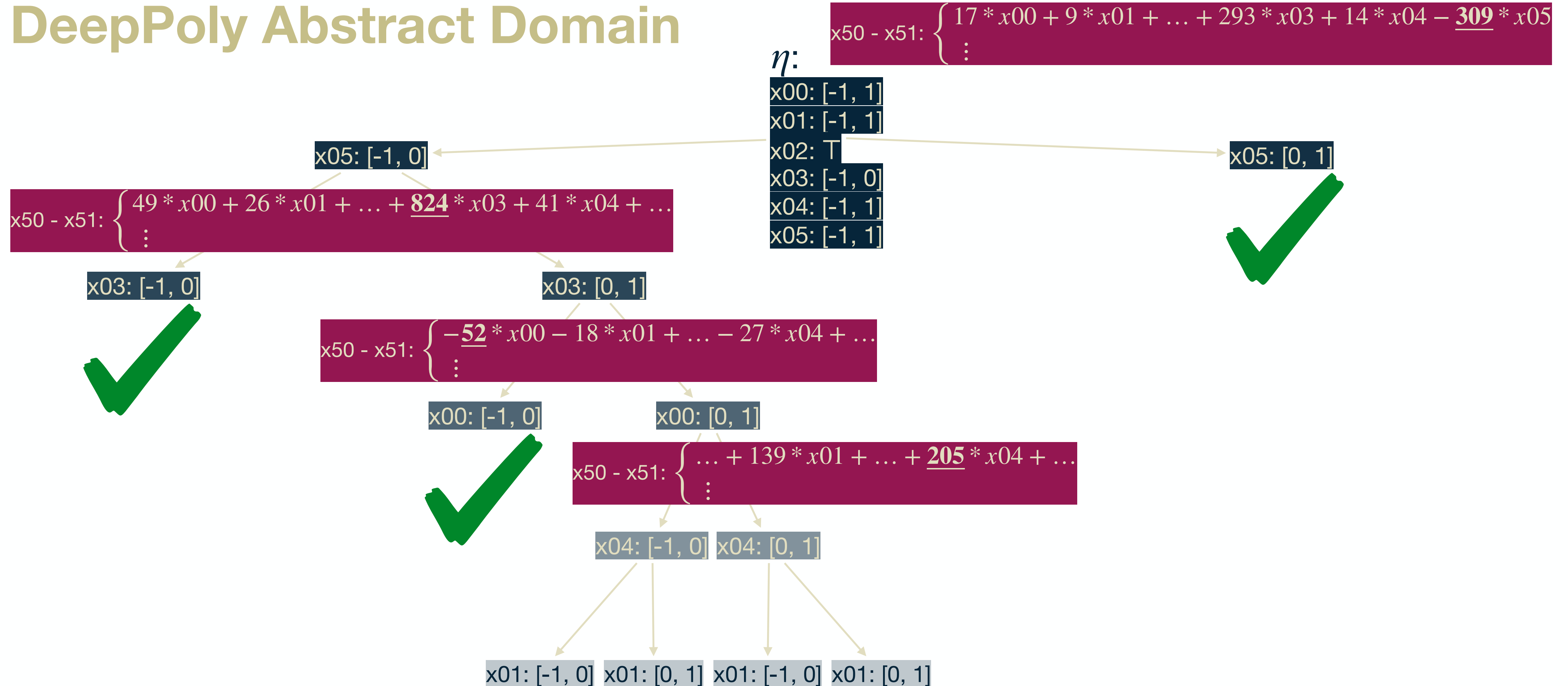
# Partitioning Strategies: Interval Range

## DeepPoly Abstract Domain



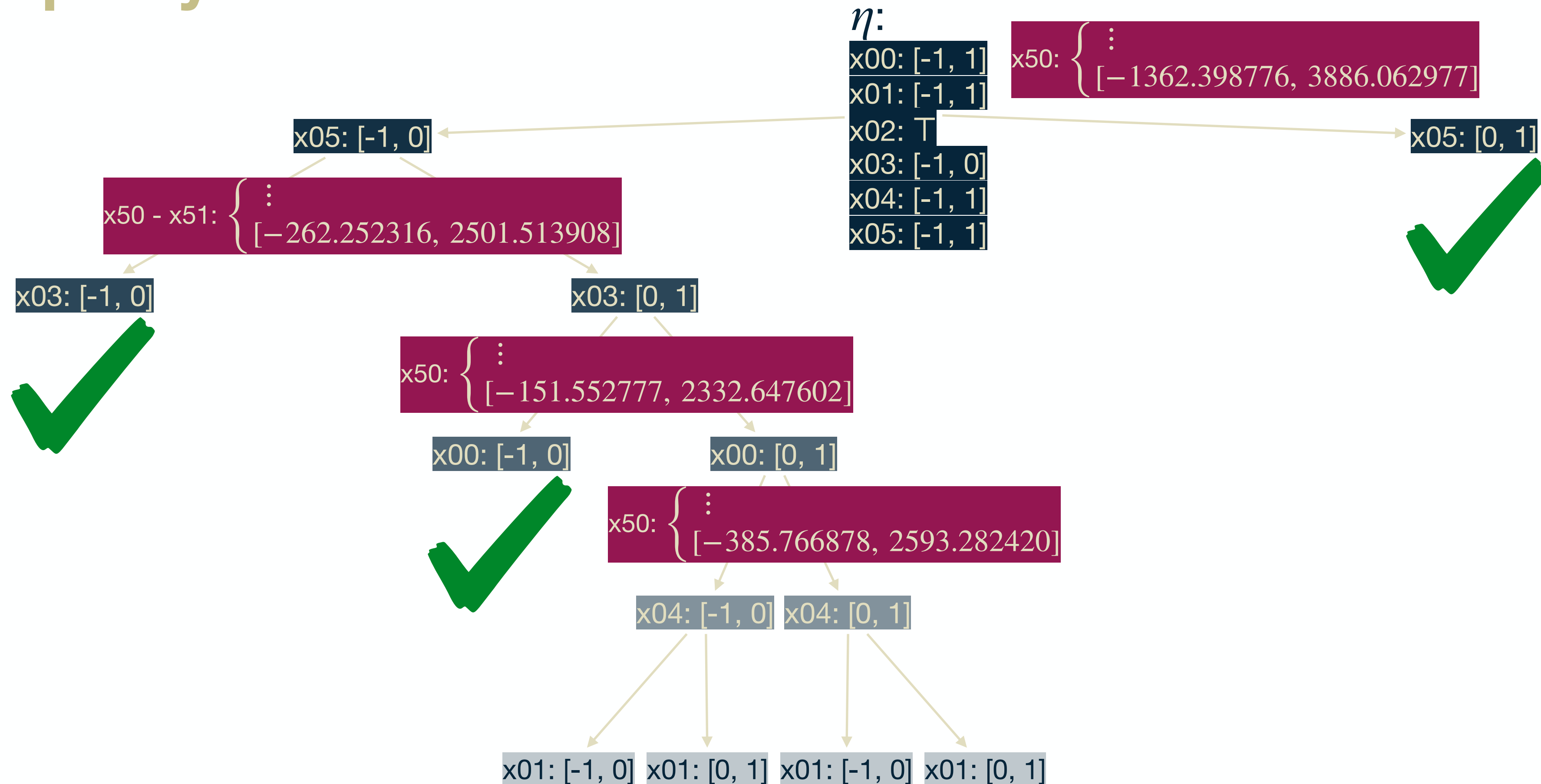
# Partitioning Strategies: ReCIPH

## DeepPoly Abstract Domain



# Input Refinement $\nRightarrow$ Output Refinement

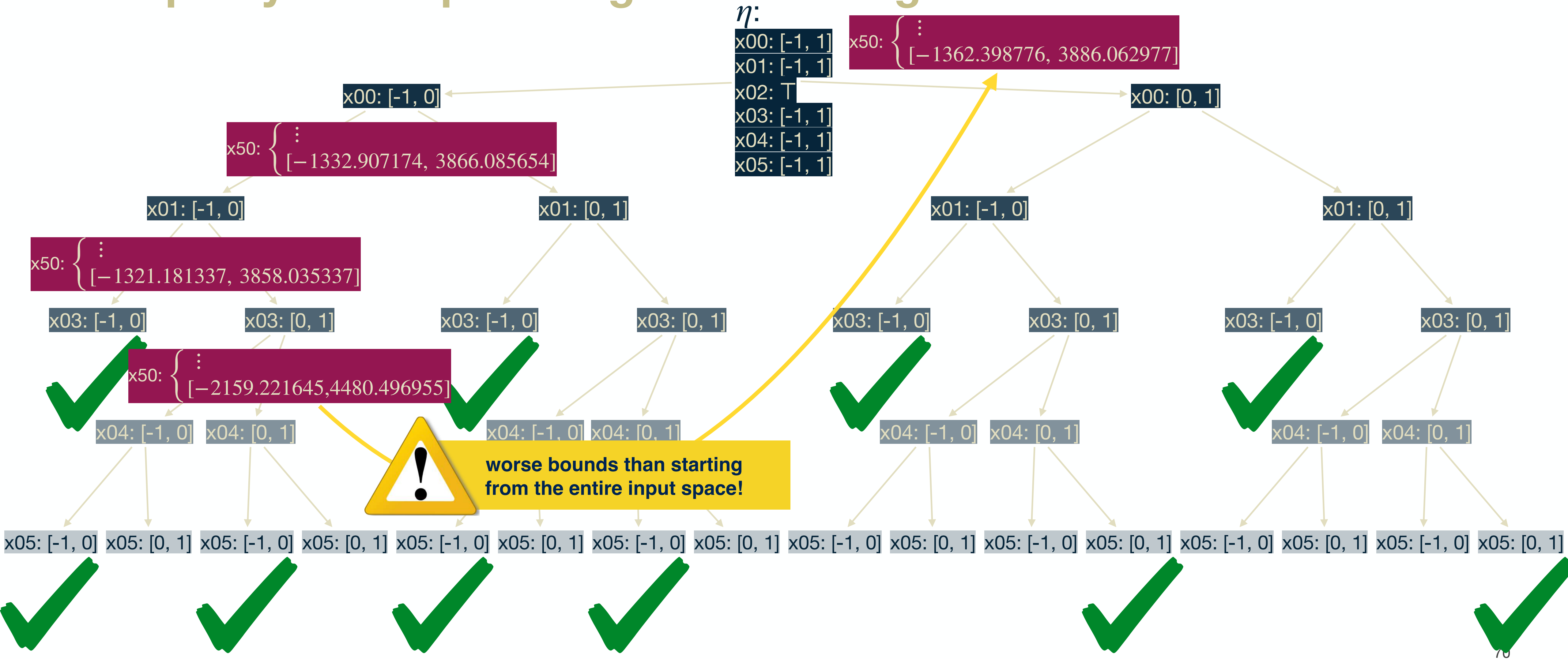
## DeepPoly Abstract Domain





# Input Refinement $\nRightarrow$ Output Refinement

## DeepPoly with Input Range Partitioning



# Scalability-vs-Precision Tradeoff

## Analyzed Input Space Percentage

L	U	Boxes	Symbolic	DeepPoly		Product	
				Input Range Partitioning	ReCIPH	Input Range Partitioning	ReCIPH
1	2	46,9 %	46,9 %	68,8 %	87,5 %	90,6 %	90,6 %
	6	46,9 %	46,9 %	68,8 %	87,5 %	90,6 %	90,6 %
0.5	2	76,9 %	89,2 %	100,0 %	100,0 %	100,0 %	100,0 %
	6	84,4 %	89,9 %	100,0 %	100,0 %	100,0 %	100,0 %

## Execution Time

L	U	Boxes	Symbolic	DeepPoly		Product	
				Input Range Partitioning	ReCIPH	Input Range Partitioning	ReCIPH
1	2	0,08s	0,14s	0,26s	0,11s	0,26s	0,12s
	6	0,16s	0,31s	0,51s	0,20s	0,35s	0,20s
0.5	2	8,88s	5,76s	2,60s	1,61s	2,10s	1,61s
	6	64,67s	40,90s	2,65s	1,63s	2,10s	1,62s

# Neural Network Verification

# Neural Network Explainability

# Abductive Explanations (AXp) [Marques-Silva21]

## Subset-Minimal Set of Input Features Sufficient for Ensuring Prediction

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

**X:**

```
x00: 1
x01: 1
x02: -1
x03: 1
x04: 1
x05: -1
```

**AXp = { x02 }**

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01)
x11 = ReLU((0.113805)*x00 + (0.064486)*x01)
x12 = ReLU((0.755487)*x00 + (0.224640)*x01)
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11)
x21 = ReLU((1.958950)*x10 + (2.388245)*x11)
x22 = ReLU((1.958103)*x10 + (2.273354)*x11)
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21)
x31 = ReLU((2.327110)*x20 + (2.685314)*x21)
x32 = ReLU((2.147212)*x20 + (2.285599)*x21)
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31)
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31)
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

### Static Backward Analysis

#### Symbolic & DeepPoly Product Abstract Domain

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

$\eta$ :

```
x00: [0, 1]
x01: [-1, 0]
x02: 1
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

$\eta$ :

```
x00: [0, 1]
x01: [0, 1]
x02: 1
x03: [0.5, 1]
x04: [0, 1]
x05: [-1, 0]
```

counterexample

```
x00: 1
x01: 1
x02: -1
x03: 1
x04: 1
x05: -1
```

```
x00: 1
x01: 1
x02: 1
x03: 1
x04: 1
x05: -1
```

```

1 x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))
1 x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))
? x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))

1 x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))
1 x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))
1 x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))

:

1 x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))
0 x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)
? x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))

x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)

```

$(-4.556024) * x40 + (0.361304) * x41 + (-33.326096) * x42 + (3728) > 0$

$(4.556024) * x40 + (33.326096) * x42 - 3728 > 0$

```
834))
8711))
885))
```

# Over-Approximating One AXp

Drop (i.e., Free) Input Features While AXp Condition Holds

SAME PREDICTION

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

```
x:
x00: [-1, 1]
x01: [-1, 1]
x02: -1
x03: [-1, 1]
x04: [-1, 1]
x05: [-1, 1]
```

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (1.222249)*x12 + (-2.103565)*x05 + (1.623834))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.388245)*x12 + (-2.103565)*x05 + (1.623834))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (2.273354)*x12 + (-2.103565)*x05 + (1.623834))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (0.666507)*x22 + (-2.103565)*x05 + (1.623834))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (2.685314)*x22 + (-2.103565)*x05 + (1.623834))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.285599)*x22 + (-2.103565)*x05 + (1.623834))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (1.980387)*x32 + (-2.103565)*x05 + (1.623834))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.828226)*x32 + (-2.103565)*x05 + (1.623834))
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (1.199384)*x32 + (-2.103565)*x05 + (1.623834))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-2500)
```

{ x00, x01, x02, x03, x04, x05 } → x51

Free x00: { x01, x02, x03, x04, x05 } → x51

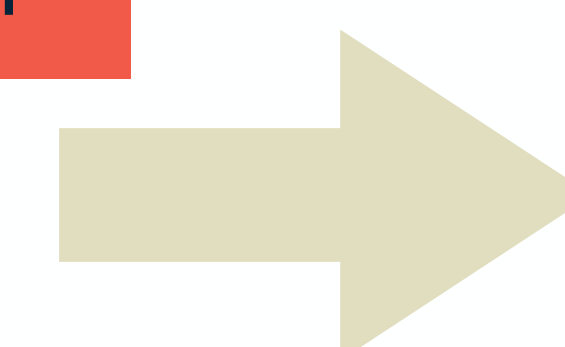
Free x01: { x02, x03, x04, x05 } → x51

Free x02: { x03, x04, x05 } → ~~x51~~

Free x03: { x02, x04, x05 } → x51

Free x04: { x02, x05 } → x51

Free x05: { x02 } → x51



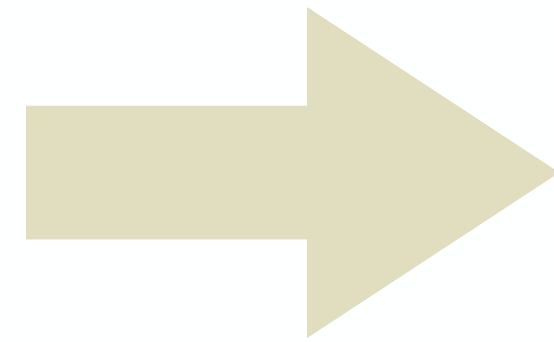
AXp = { x02 }

# Over-Approximating One AXp

Drop (i.e., Free) Input Features While AXp Condition Holds

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

**X:**  
x00: 1  
x01: 1  
x02: -1  
x03: 1  
x04: 1  
x05: -1



**BOXES**

$\text{oAXp} = \{ x02, x03, x05 \}$

**SYMBOLIC**

$\text{oAXp} = \{ x00, x02, x03 \}$   
 $\text{oAXp} = \{ x02, x03, x05 \}$

**DEEPPOLY**

$\text{oAXp} = \{ x02, x03 \}$   
 $\text{oAXp} = \{ x02, x05 \}$

**= PRODUCT**

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

# Contrastive Explanations (CXp) [Marques-Silva21]

Subset-Minimal Set of Input Features Sufficient for Changing Prediction

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

**X:**  
x00: 1  
x01: 1  
x02: -1  
x03: 1  
x04: 1  
x05: -1

**CXp = { x02 }**

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01)
x11 = ReLU((0.113805)*x00 + (0.064486)*x01)
x12 = ReLU((0.755487)*x00 + (0.224640)*x01)
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11)
x21 = ReLU((1.958950)*x10 + (2.388245)*x11)
x22 = ReLU((1.958103)*x10 + (2.273354)*x11)
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21)
x31 = ReLU((2.327110)*x20 + (2.685314)*x21)
x32 = ReLU((2.147212)*x20 + (2.285599)*x21)
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31)
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31)
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31)
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

## Static Backward Analysis

### Symbolic & DeepPoly Product Abstract Domain

x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())

$\eta$ :

x00:	[0, 1]
x01:	[-1, 0]
x02:	T
x03:	[0.5, 1]
x04:	[0, 1]
x05:	[-1, 0]

✓

$\eta$ :

x00:	[0, 1]
x01:	[0, 1]
x02:	T
x03:	[0.5, 1]
x04:	[0, 1]
x05:	[-1, 0]

✗

counterexample

x00:	1	x00:	1
x01:	1	x01:	1
x02:	-1	x02:	1
x03:	1	x03:	1
x04:	1	x04:	1
x05:	-1	x05:	-1

① x10 = ReLU((0.120875)\*x00 + (0.065404)\*x01 + (0.097862)\*x02 + (2.030051)\*x03 + (0.101956)\*x04 + (-2.103565)\*x05 + (1.623834))

① x11 = ReLU((0.113805)\*x00 + (0.064486)\*x01 + (0.090701)\*x02 + (2.123338)\*x03 + (0.076374)\*x04 + (-1.651132)\*x05 + (-0.828711))

? x12 = ReLU((0.755487)\*x00 + (0.224640)\*x01 + (0.344943)\*x02 + (2.619876)\*x03 + (0.346636)\*x04 + (1.418635)\*x05 + (-0.686885))

① x20 = ReLU((1.803209)\*x10 + (1.222249)\*x11 + (2.725716)\*x12 + (-3.489653))

① x21 = ReLU((1.958950)\*x10 + (2.388245)\*x11 + (2.245851)\*x12 + (-3.834811))

① x22 = ReLU((1.958103)\*x10 + (2.273354)\*x11 + (0.662405)\*x12 + (-4.211086))

⋮

① x40 = ReLU((2.296390)\*x30 + (1.980387)\*x31 + (2.945360)\*x32 + (-4.096463))

① x41 = ReLU((-0.552155)\*x30 + (-0.828226)\*x31 + (-0.495998)\*x32)

? x42 = ReLU((-2.509773)\*x30 + (1.199384)\*x31 + (-0.245429)\*x32 + (5.024773))

x50 = (-2.278012)\*x40 + (0.180652)\*x41 + (-16.663048)\*x42 + (1864)     (-4.556024)\*x40 + (0.361304)\*x41 + (-33.326096)\*x42 + (3728) > 0

x51 = (2.278012)\*x40 + (-0.180652)\*x41 + (16.663048)\*x42 + (-1864)     (4.556024)\*x40 + (33.326096)\*x42 - 3728 > 0

834 ) )  
8711 ) )  
885 ) )

# Under-Approximating One CXp

Drop (i.e., Fix) Input Features While CXp Condition Holds

DIFFERENT PREDICTIONS

```
x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())
```

X:   
x00: 1   
x01: 1   
x02: [-1, 1]   
x03: 1   
x04: 1   
x05: -1

{ x00, x01, x02, x03, x04, x05 } →



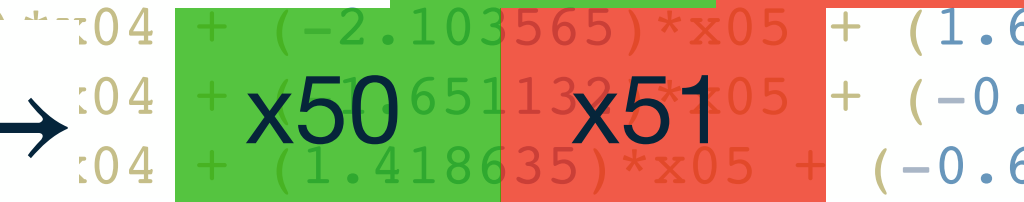
Fix x00:

{ x01, x02, x03, x04, x05 } →



Fix x01:

{ x02, x03, x04, x05 } →



Fix x02:

{ x03, x04, x05 } →



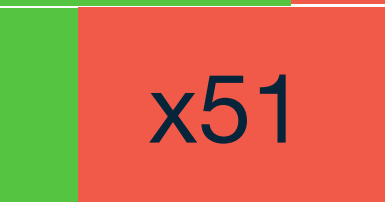
Fix x03:

{ x02, x04, x05 } →



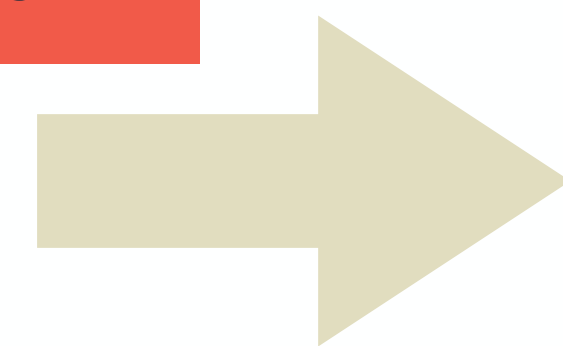
Fix x04:

{ x02, x05 } →



Fix x05:

{ x02 } →



CXp = { x02 }

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.007862)*x02 + (0.000051)*x03 + (0.101056)*x04 + (-2.103565)*x05 + (1.623834))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (0.007862)*x12 + (0.000051)*x13 + (0.101056)*x14 + (-2.103565)*x15 + (1.623834))
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (0.007862)*x12 + (0.000051)*x13 + (0.101056)*x14 + (-2.103565)*x15 + (1.623834))
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.007862)*x12 + (0.000051)*x13 + (0.101056)*x14 + (-2.103565)*x15 + (1.623834))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (0.007862)*x22 + (0.000051)*x23 + (0.101056)*x24 + (-2.103565)*x25 + (1.623834))
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (0.007862)*x22 + (0.000051)*x23 + (0.101056)*x24 + (-2.103565)*x25 + (1.623834))
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (0.007862)*x22 + (0.000051)*x23 + (0.101056)*x24 + (-2.103565)*x25 + (1.623834))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (0.007862)*x32 + (0.000051)*x33 + (0.101056)*x34 + (-2.103565)*x35 + (1.623834))
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (0.007862)*x32 + (0.000051)*x33 + (0.101056)*x34 + (-2.103565)*x35 + (1.623834))
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (0.007862)*x32 + (0.000051)*x33 + (0.101056)*x34 + (-2.103565)*x35 + (1.623834))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

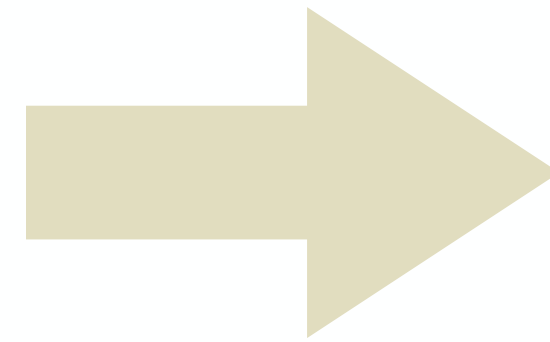


# Under-Approximating One CXp

Drop (i.e., Fix) Input Features While CXp Condition Holds

```
x00 = float(input())  
x01 = float(input())  
x02 = float(input())  
x03 = float(input())  
x04 = float(input())  
x05 = float(input())
```

**X:**  
x00: 1  
x01: 1  
x02: -1  
x03: 1  
x04: 1  
x05: -1



**BOXES**

```
WAXp = { x02 }  
WAXp = { x03 }  
WAXp = { x05 }
```

**SYMBOLIC**

```
WAXp = { x02 }  
WAXp = { x03 }  
WAXp = { x00, x05 }
```

**DEEPPOLY**

```
WAXp = { x02 }  
WAXp = { x03, x05 }  
= PRODUCT
```

```
x10 = ReLU((0.120875)*x00 + (0.065404)*x01 + (0.097862)*x02 + (2.030051)*x03 + (0.101956)*x04 + (-2.103565)*x05 + (1.623834))  
x11 = ReLU((0.113805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.651132)*x05 + (-0.828711))  
x12 = ReLU((0.755487)*x00 + (0.224640)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885))
```

```
x20 = ReLU((1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.489653))  
x21 = ReLU((1.958950)*x10 + (2.388245)*x11 + (2.245851)*x12 + (-3.834811))  
x22 = ReLU((1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086))
```

```
x30 = ReLU((1.735994)*x20 + (0.666507)*x21 + (3.192344)*x22 + (-2.627086))  
x31 = ReLU((2.327110)*x20 + (2.685314)*x21 + (1.424807)*x22 + (-3.695113))  
x32 = ReLU((2.147212)*x20 + (2.285599)*x21 + (2.665507)*x22 + (-4.299974))
```

```
x40 = ReLU((2.296390)*x30 + (1.980387)*x31 + (2.945360)*x32 + (-4.096463))  
x41 = ReLU((-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32)  
x42 = ReLU((-2.509773)*x30 + (1.199384)*x31 + (-0.245429)*x32 + (5.024773))
```

```
x50 = (-2.278012)*x40 + (0.180652)*x41 + (-16.663048)*x42 + (1864)  
x51 = (2.278012)*x40 + (-0.180652)*x41 + (16.663048)*x42 + (-1864)
```

# Verification and Explainability

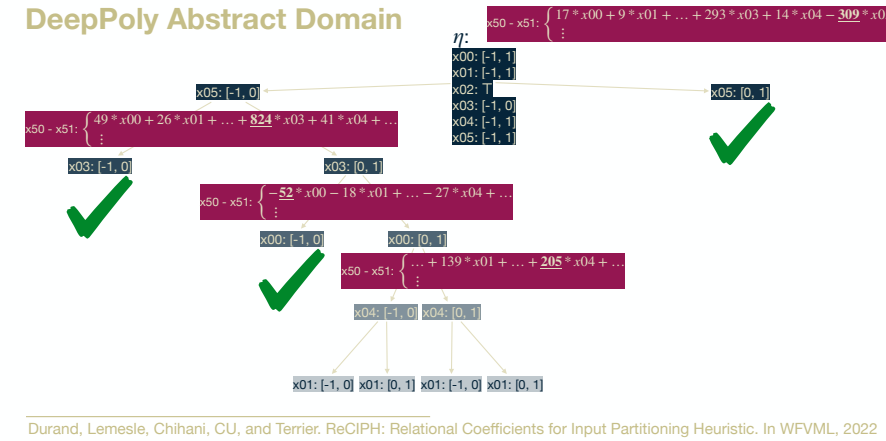
## Safety-Critical Neural Networks



practical tools  
targeting specific programs

### Partitioning Strategies: ReCIPH

DeepPoly Abstract Domain



### Scalability-vs-Precision Tradeoff

Analyzed Input Space Percentage

L	U	Boxes	Symbolic	DeepPoly		Product	
				Input Range Partitioning	ReCIPH	Input Range Partitioning	ReCIPH
1	2	46.9 %	46.9 %	68.8 %	87.5 %	90.6 %	90.6 %
0.5	6	84.4 %	89.2 %	100.0 %	100.0 %	100.0 %	100.0 %

Execution Time

L	U	Boxes	Symbolic	DeepPoly		Product	
				Input Range Partitioning	ReCIPH	Input Range Partitioning	ReCIPH
1	2	0.08s	0.14s	0.26s	0.11s	0.26s	0.12s
0.5	6	0.16s	0.31s	0.51s	0.20s	0.35s	0.20s
2	2	8.88s	5.76s	2.60s	1.61s	2.10s	1.61s
6	6	64.67s	40.90s	2.65s	1.63s	2.10s	1.62s



algorithmic approaches  
to decide program properties

### Hyperproperty Verification

Static Forward Analysis

```

x00 = float(input())
x01 = float(input())
x02 = float(input())
x03 = float(input())
x04 = float(input())
x05 = float(input())

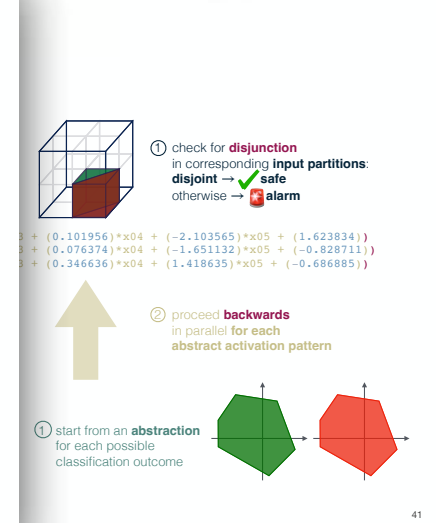
x10 = ReLU(0.12075)*x00 + (0.06484)*x01 + (0.09782)*x02 + (2.30351)*x03 + (0.01108)*x04 + (-2.10356)*x05 + (1.62384)
x11 = ReLU(0.111805)*x00 + (0.064486)*x01 + (0.090701)*x02 + (2.123338)*x03 + (0.076374)*x04 + (-1.65132)*x05 + (-0.828711)
x12 = ReLU(0.755487)*x00 + (0.224440)*x01 + (0.344943)*x02 + (2.619876)*x03 + (0.346636)*x04 + (1.418635)*x05 + (-0.686885)

x20 = ReLU(1.803209)*x10 + (1.222249)*x11 + (2.725716)*x12 + (-3.499651)
x21 = ReLU(1.958950)*x10 + (2.388245)*x11 + (2.248593)*x12 + (-0.834831)
x22 = ReLU(1.958103)*x10 + (2.273354)*x11 + (0.662405)*x12 + (-4.211086)

x30 = ReLU(1.735994)*x20 + (0.666507)*x21 + (-3.192244)*x22 + (-2.627086)
x31 = ReLU(2.327110)*x20 + (2.485314)*x21 + (1.424807)*x22 + (-3.695112)
x32 = ReLU(2.187212)*x20 + (2.285989)*x21 + (2.665093)*x22 + (-4.239973)

x40 = ReLU(2.296390)*x30 + (1.960387)*x31 + (2.945360)*x32 + (-4.056443)
x41 = ReLU(-0.552155)*x30 + (-0.828226)*x31 + (-0.495998)*x32 + (-0.828711)
x42 = ReLU(-2.509773)*x30 + (1.193984)*x31 + (-0.245429)*x32 + (5.024773)

x50 = (-2.278012)*x40 + (0.186652)*x41 + (-0.603048)*x42 + (1.864)
x51 = (2.78022)*x40 + (-0.186652)*x41 + (0.603048)*x42 + (-1.864)
    
```



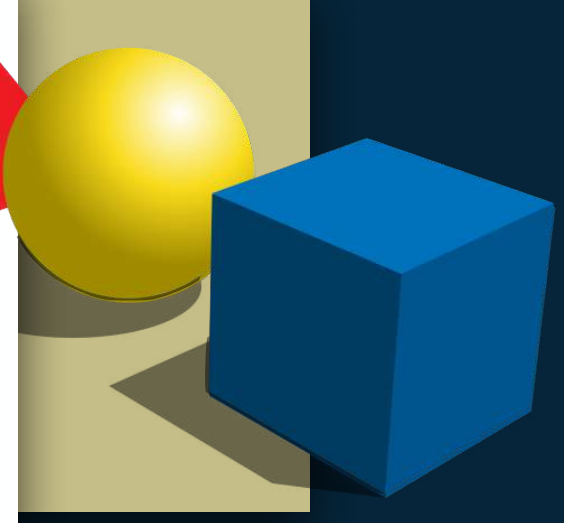
### Reduced Product Domain

Symbolic Abstract Domain & DeepPoly Abstract Domain

Symbolic:  $\max(a_s, a_d), \min(b_s, b_d)$

DeepPoly:  $\max(a_s, a_d), \min(b_s, b_d)$

Symbolic and DeepPoly domains are connected by bidirectional arrows, with a central box labeled  $[a_s, b_s]$ .



mathematical models  
of the program behavior

### Parallel Semantics

Diagram showing a neural network structure with input  $t_0$  and prediction  $t_{out}$ . The semantics is defined as:

$$\mathcal{H}_p^{\parallel} \stackrel{\text{def}}{=} \{ T \mid \forall t, t' \in T: \eta(t_0) = \eta(t'_0) \Rightarrow \rho(t_{out}) = \rho(t'_{out}) \}$$

### Hyperproperty Verification

Abstract Non-Interference Properties

$$\mathcal{H}_p^{\parallel} \stackrel{\text{def}}{=} \{ T \mid \forall t, t' \in T: \eta(t_0) = \eta(t'_0) \Rightarrow \rho(t_{out}) = \rho(t'_{out}) \}$$

Lemma:  $M \models \mathcal{H}_p^{\parallel} \Leftrightarrow \forall I \in \mathbb{I}: \forall A, B \in \mathbb{I}[M]: \rho(A'_0) \cap \rho(B'_0) = \perp \Rightarrow \rho(A'_0) \cap \rho(B'_0) = \perp$

Diagram showing two boxes labeled x02, one with a warning sign and one with a thumbs up.



THANKS!

# References

[Li19] **Jianlin Li, Jiangchao Liu, Pengfei Yang, Liqian Chen, Xiaowei Huang, and Lijun Zhang.** Analyzing Deep Neural Networks with Symbolic Propagation: Towards Higher Precision and Faster Verification. In SAS, page 296–319, 2019.

**symbolic abstraction**

[Singh19] **Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin T. Vechev.** An Abstract Domain for Certifying Neural Networks. In POPL, pages 41:1 - 41:30, 2019.

**deepoly abstraction**

[Urban20] **Caterina Urban, Maria Christakis, Valentin Wüstholtz, and Fuyuan Zhang.** Perfectly Parallel Fairness Certification of Neural Networks. In OOPSLA, pages 185:1–185:30, 2020.

**hypersafety verification**

[Marques-Silva21] **João Marques-Silva, Thomas Gerspacher, Martin C. Cooper, Alexey Ignatiev, and Nina Narodytska.** Explanations for Monotonic Classifiers. In ICML, pages 7469-7479, 2021.

[Wu23] **Min Wu, Haoze Wu, Clark W. Barrett.** VeriX: Towards Verified Explainability of Deep Neural Networks. In NeurIPS, 2023.

**logic-based explanations**