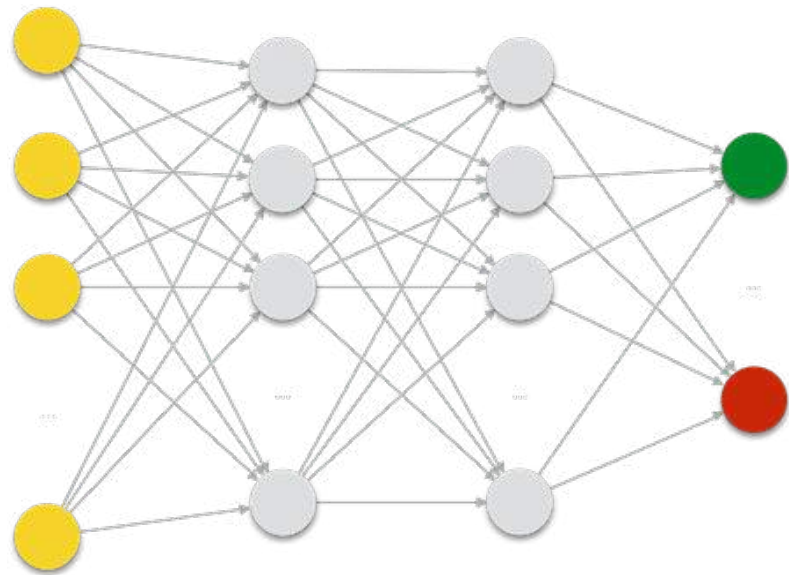


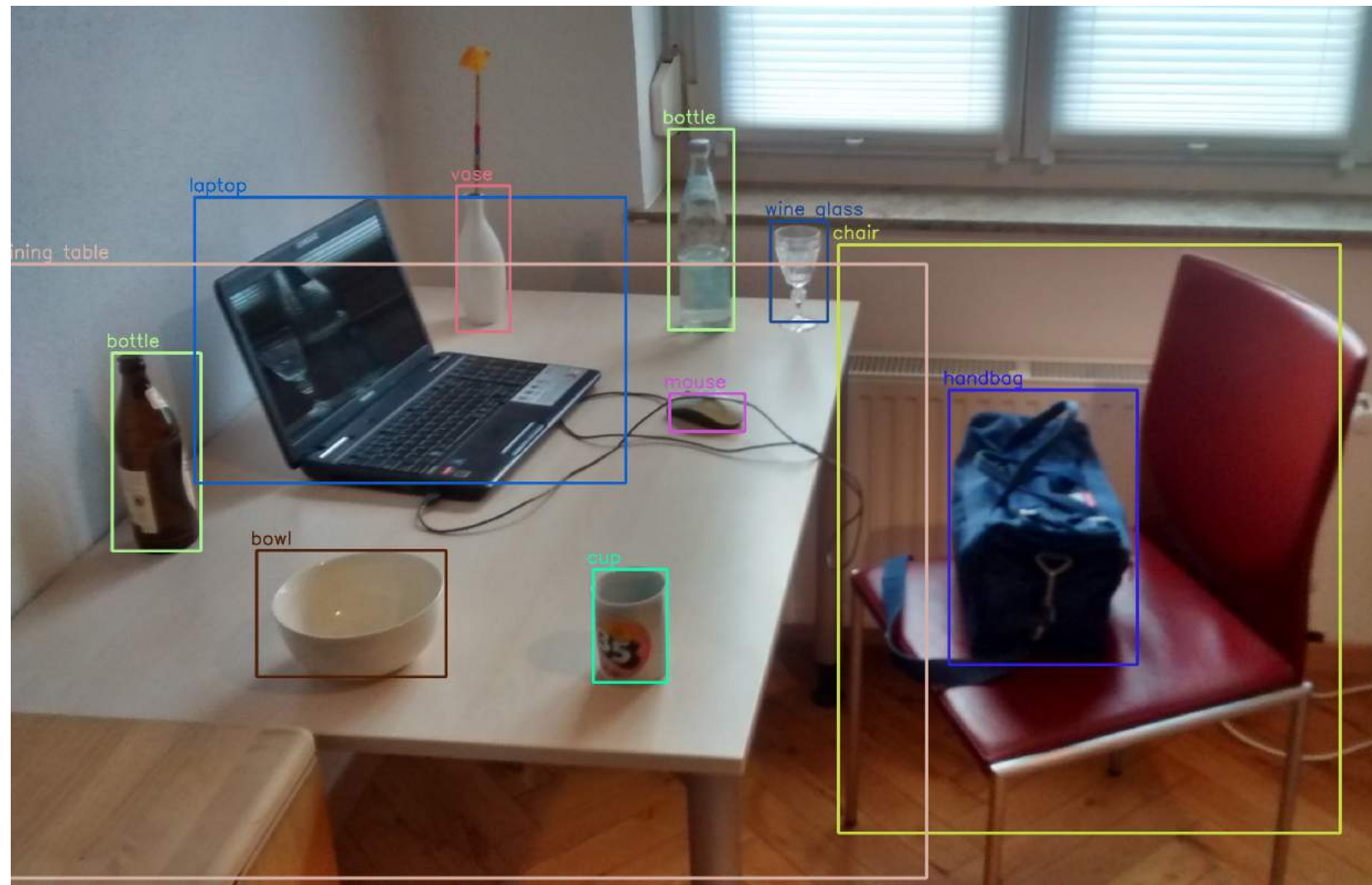
Formal Verification of Machine Learning

MPRI 2-6: Abstract Interpretation,
Application to Verification and Static Analysis



Machine Learning Revolution

Computer software able to efficiently and **autonomously perform tasks** that are difficult or even *impossible* to design using explicit programming



Examples: object recognition, image classification, speech recognition, etc.

ML in Safety-Critical Applications

Enables new functions that could not be envisioned before



Self-Driving Cars



Image-Based Taxiing, Takeoff, Landing

Aircraft Voice Control

ML in Safety-Critical Applications

Approximates complex systems and automates decision-making



Diagnosis and Drug Discovery

Deep Neural Network Compression for Aircraft Collision Avoidance Systems

Kyle D. Julian¹ and Mykel J. Kochenderfer² and Michael P. Owen³



Aircraft Collision Avoidance

Abstract—One approach to designing decision making logic for an aircraft collision avoidance system frames the problem as a Markov decision process and optimizes the system using dynamic programming. The resulting collision avoidance strategy can be represented as a numeric table. This methodology has been used in the development of the Airborne Collision Avoidance System X (ACAS X) family of collision avoidance systems for manned and unmanned aircraft, but the high dimensionality of the state space tables. To improve storage efficiency, a deep

floating point storage. A simple technique to reduce the size of the score table is to downsample the table after dynamic programming. To minimize the degradation in decision quality, states are removed in areas where the variation between values in the table are smooth. The downsampling reduces the size of the table by a factor of 180 from that produced by dynamic programming. For the rest of this paper, the downsampled ACAS Xu horizontal table is referred to as the baseline, original table.

ML in Safety-Critical Applications

¹ STAT+ ²

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

By [Casey Ross](#)³ [@caseymross](#)⁴ and Ike Swetlitz

July 25, 2018

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

Internal documents reveal that the car was at fault

By [Andrew Liptak](#) | [@AndrewLiptak](#) | Feb 25, 2017, 11:08am EST

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

[Richard Gonzales](#) November 7, 2019 10:57 PM ET



Machine Learning Pipeline

DATA SCIENCE SOFTWARE



data



data preparation



model training



model deployment



predictions



PREVIOUS LESSON

TODAY

Machine Learning Pipeline

Model Training is **Highly Non-Deterministic**



model training



model deployment



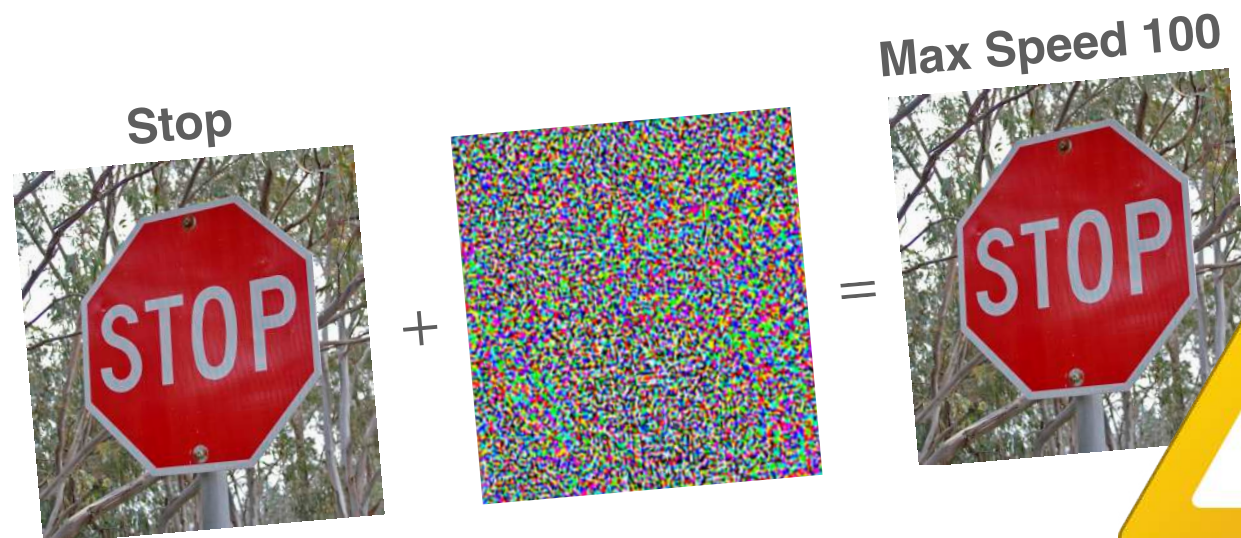
predictions



no predictability and traceability

Machine Learning Pipeline

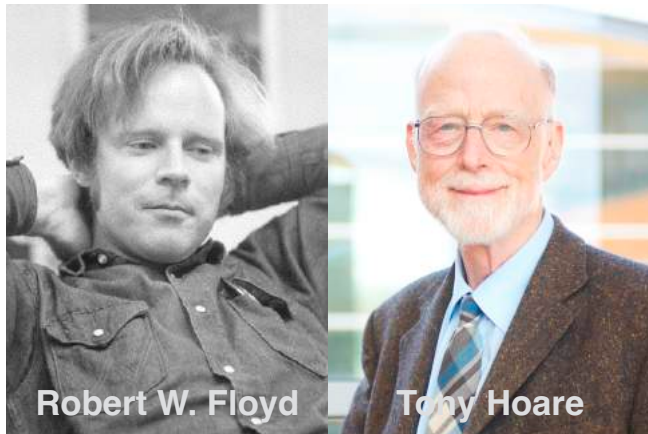
Models Only Give **Probabilistic Guarantees**



not sufficient for guaranteeing an acceptable failure rate under any circumstance

Formal Methods

Mathematical Guarantees of Safety



Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



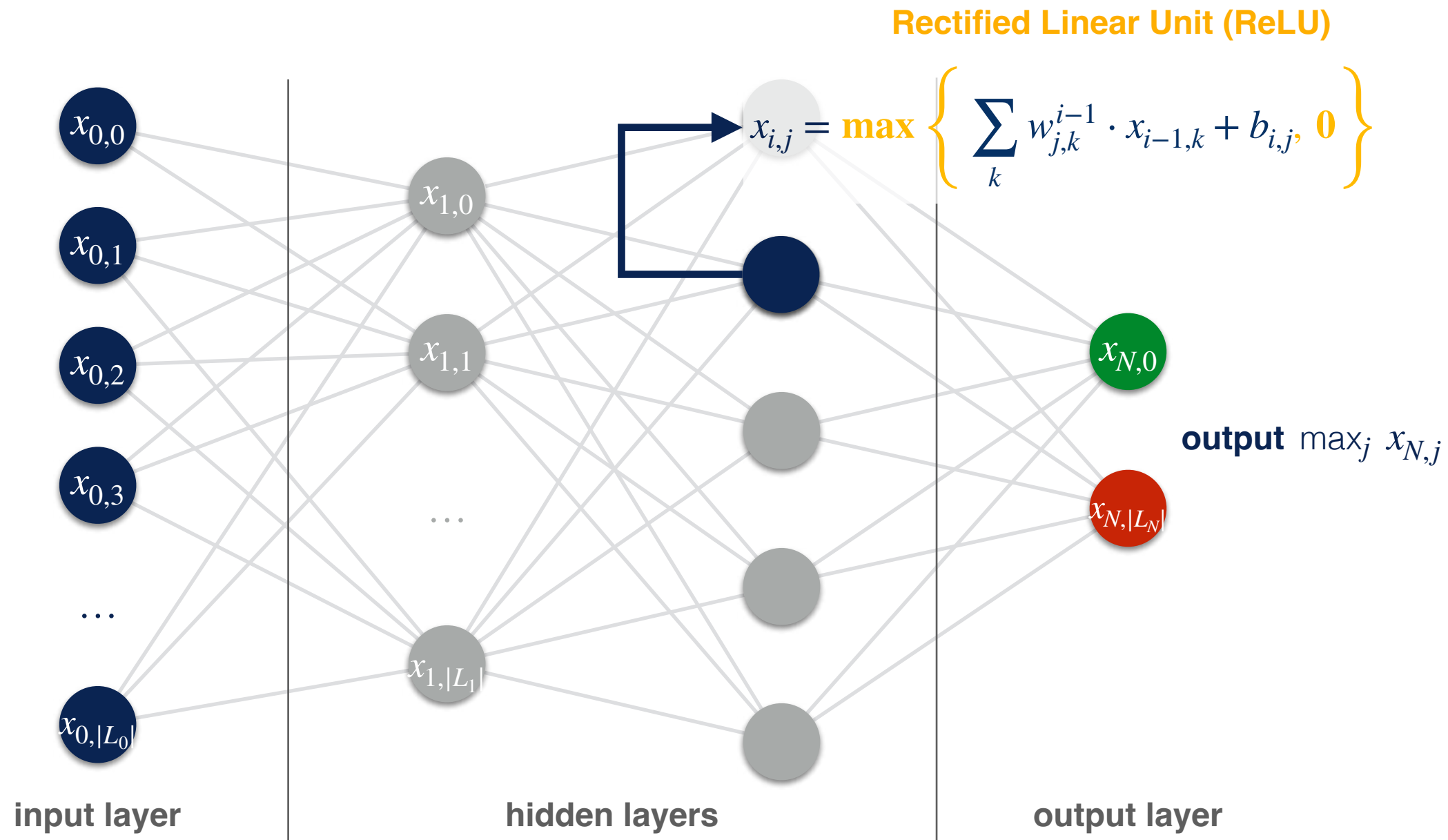
Formal Methods for Trained Models



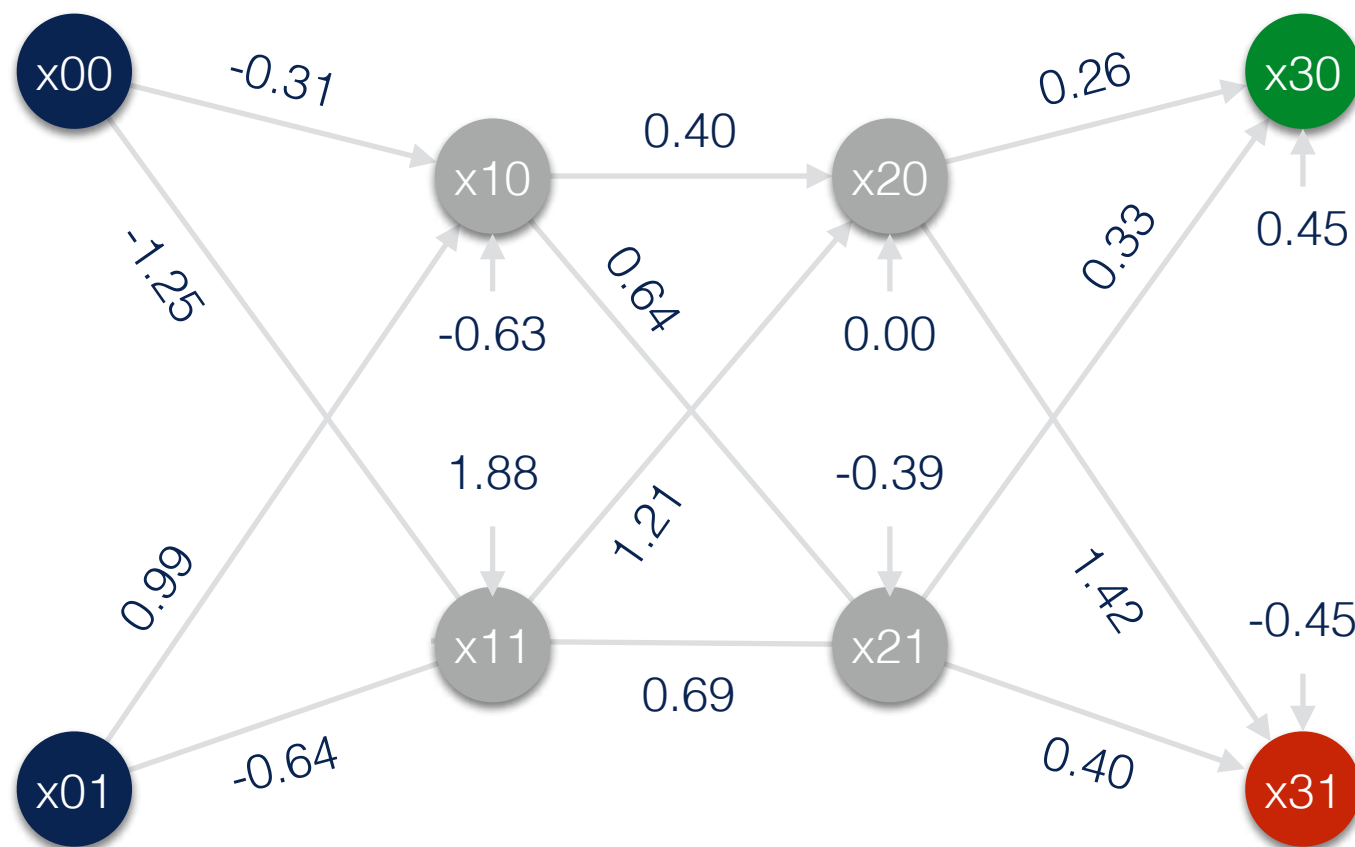
Neural Networks

Neural Networks

Feed-Forward Fully-Connected Neural Networks with ReLU Activation Functions



Feed-Forward Fully-Connected ReLU Networks as Programs



$x_{00} = \text{input}()$
 $x_{01} = \text{input}()$

$x_{10} = -0.31 * x_{00} + 0.99 * x_{01} + (-0.63)$
 $x_{11} = -1.25 * x_{00} + (-0.64) * x_{01} + 1.88$

$x_{10} = 0$ if $x_{10} < 0$ else x_{10}
 $x_{11} = 0$ if $x_{11} < 0$ else x_{11}

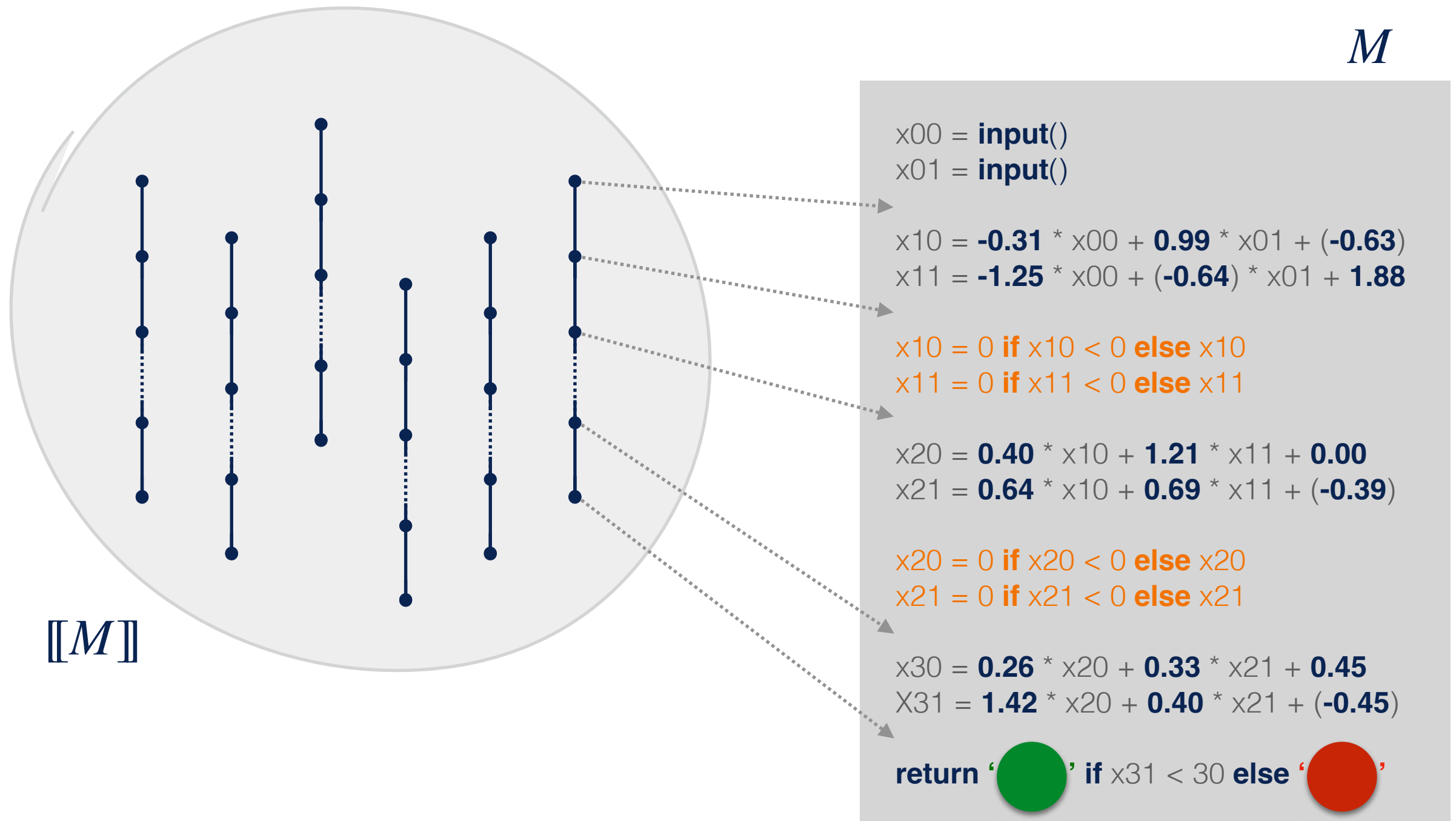
$x_{20} = 0.40 * x_{10} + 1.21 * x_{11} + 0.00$
 $x_{21} = 0.64 * x_{10} + 0.69 * x_{11} + (-0.39)$

$x_{20} = 0$ if $x_{20} < 0$ else x_{20}
 $x_{21} = 0$ if $x_{21} < 0$ else x_{21}

$x_{30} = 0.26 * x_{20} + 0.33 * x_{21} + 0.45$
 $x_{31} = 1.42 * x_{20} + 0.40 * x_{21} + (-0.45)$

return '●' if $x_{31} < 30$ else '●'

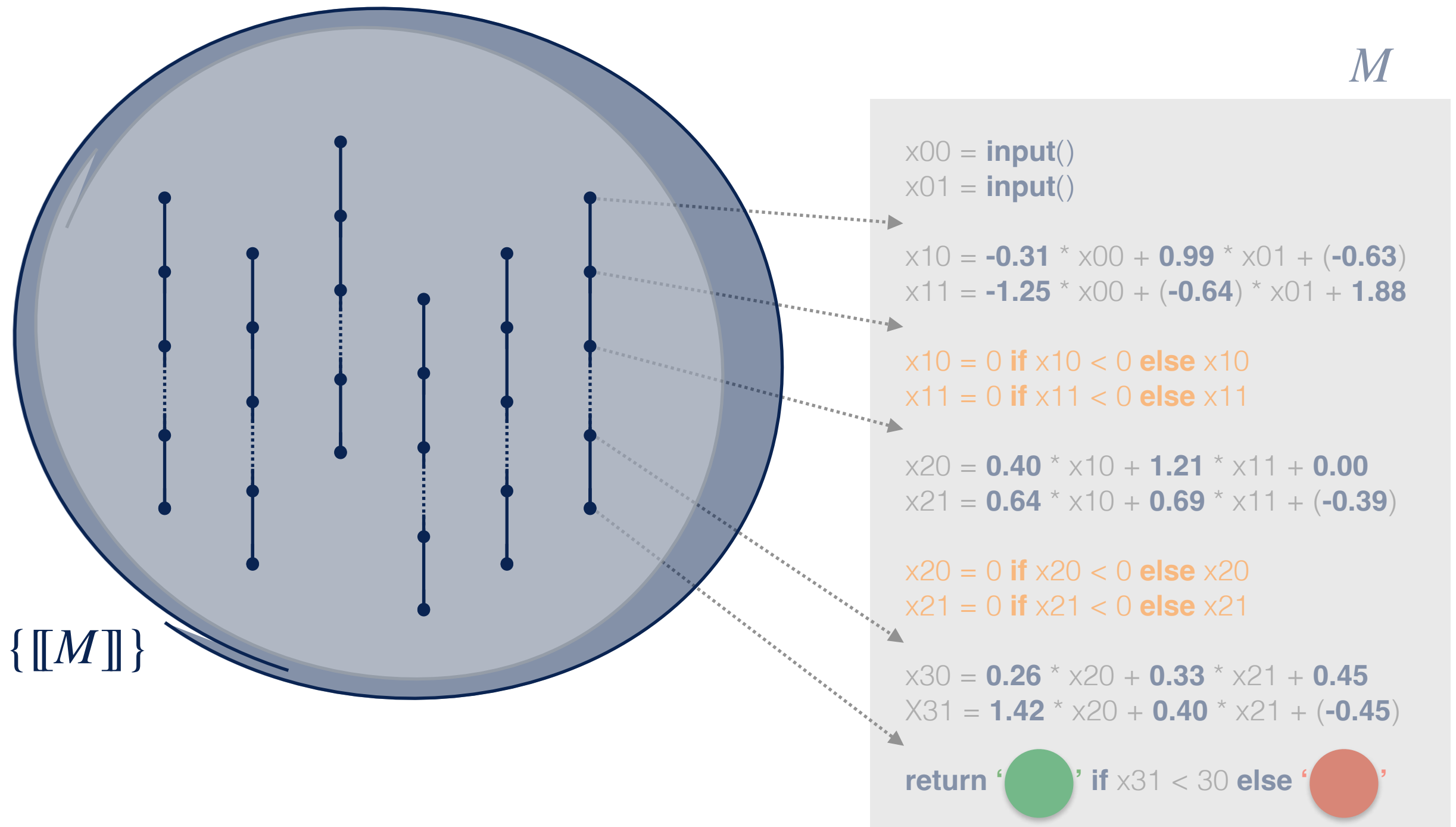
Maximal Trace Semantics





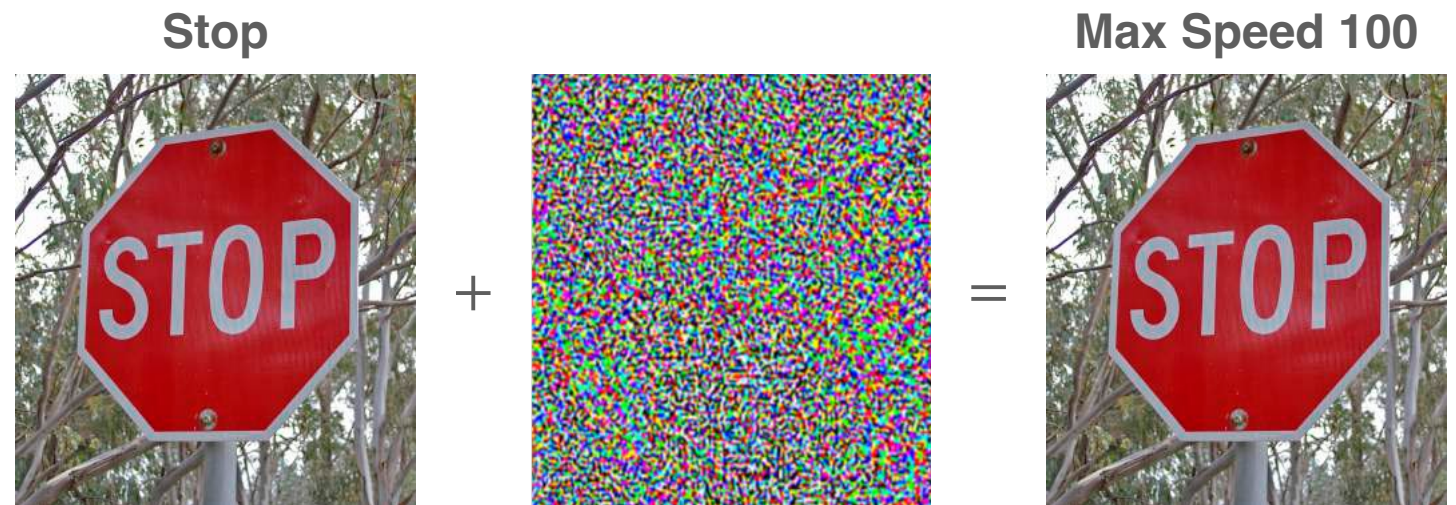
Neural Network Verification

Collecting Semantics



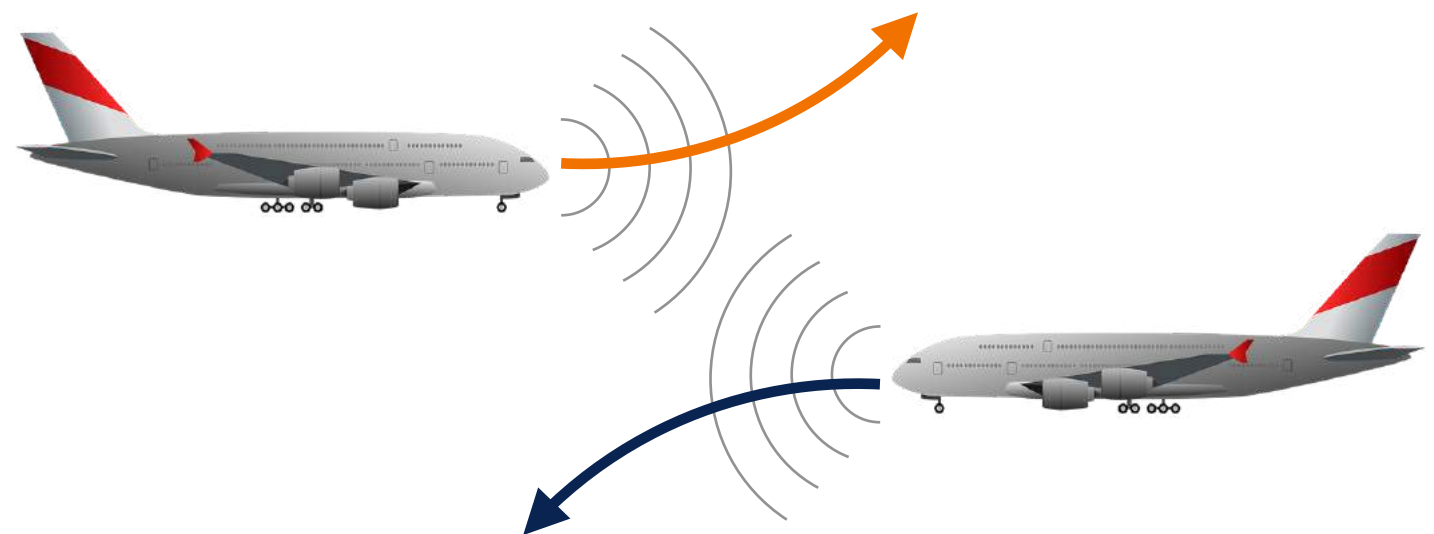
Stability

Goal G3 in [Kurd03]

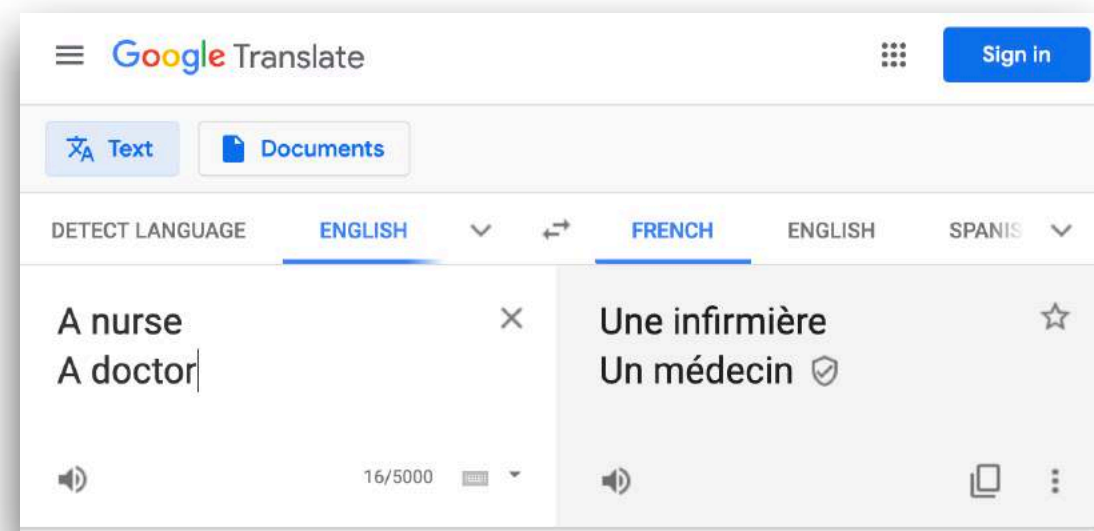


Safety

Goal G4 in [Kurd03]

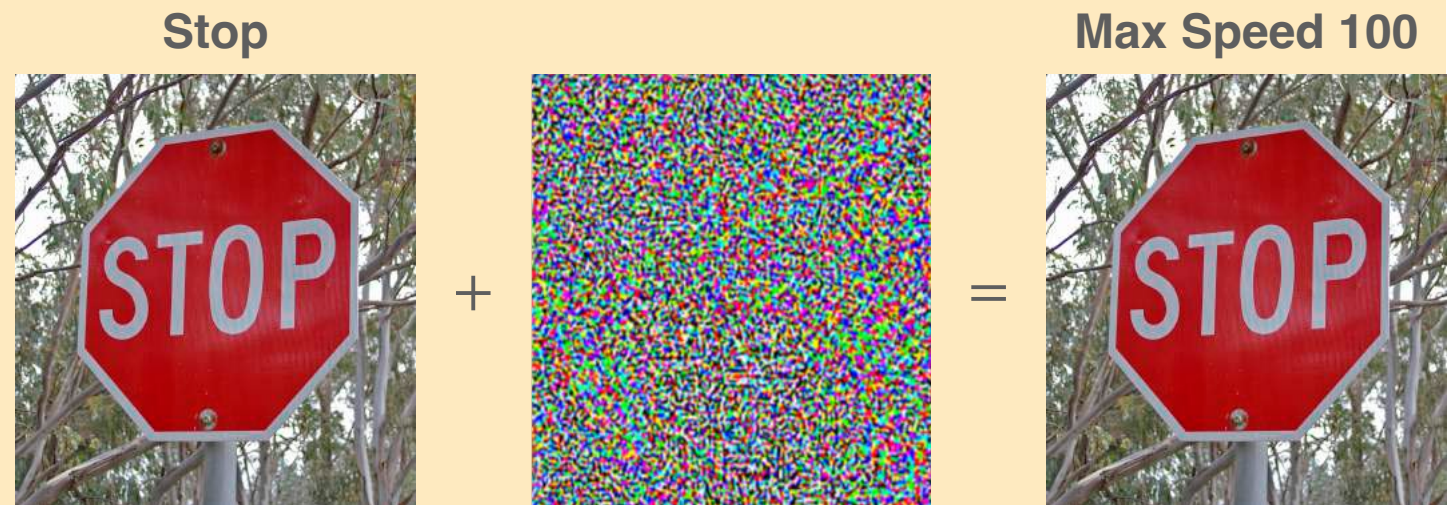


Fairness



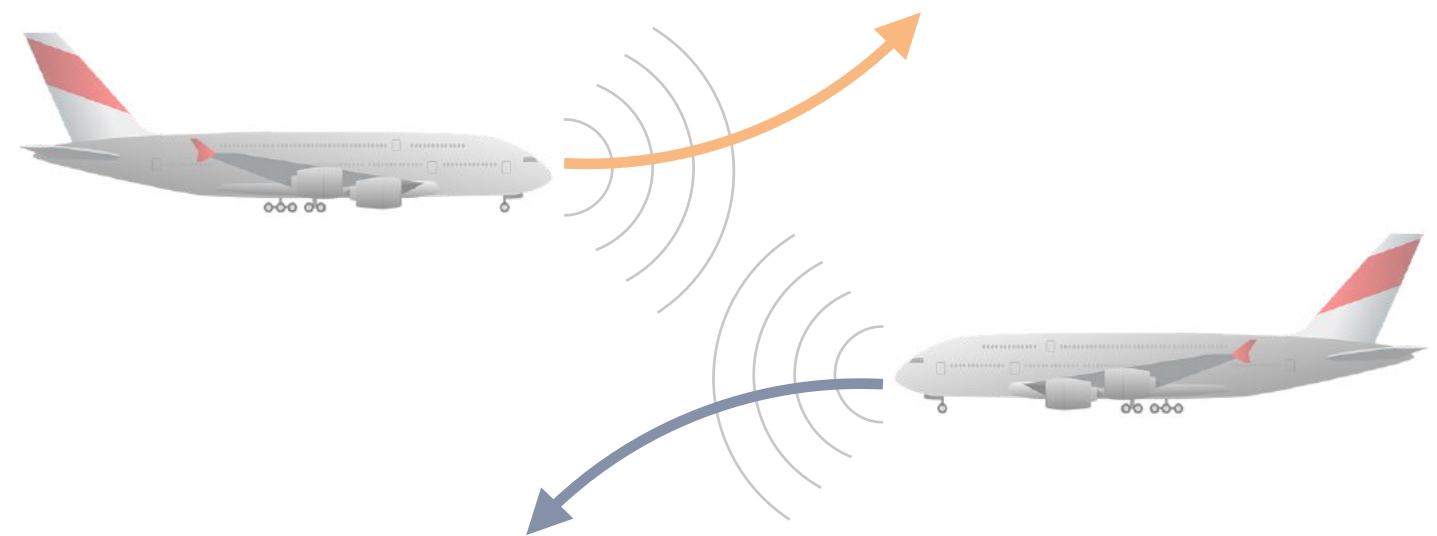
Stability

Goal G3 in [Kurd03]

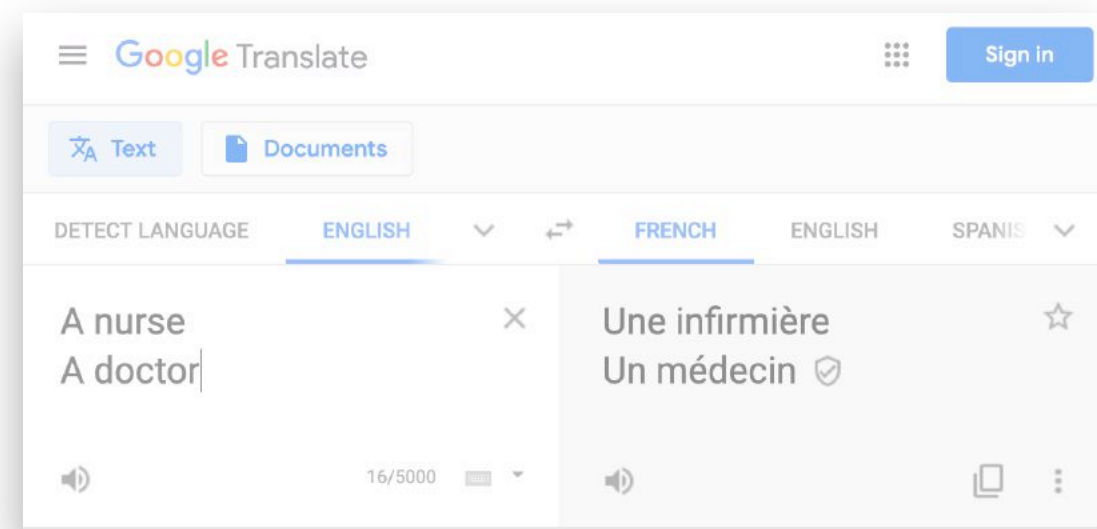


Safety

Goal G4 in [Kurd03]

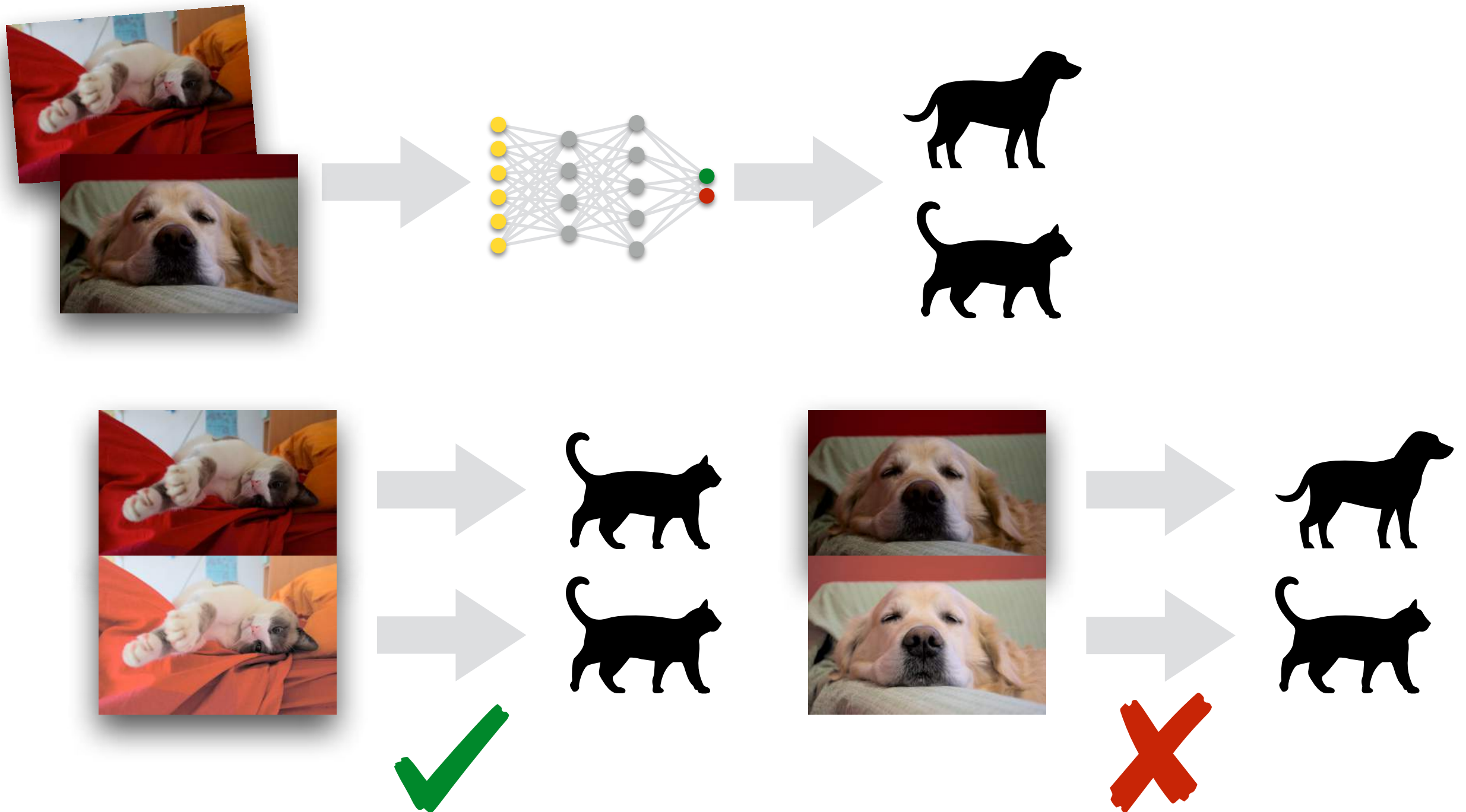


Fairness



Local Stability

The classification is **unaffected** by small input perturbations



Local Stability

Distance-Based Perturbations

$$P_{\delta,\epsilon}(\mathbf{x}) \stackrel{\text{def}}{=} \{\mathbf{x}' \in \mathcal{R}^{|L_0|} \mid \delta(\mathbf{x}, \mathbf{x}') \leq \epsilon\}$$

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

$$\mathcal{R}_x^{\delta,\epsilon} \stackrel{\text{def}}{=} \{[[M]] \in \mathcal{P}(\Sigma^*) \mid \text{STABLE}_x^{\delta,\epsilon}([M])\}$$

$\mathcal{R}_x^{\delta,\epsilon}$ is the set of all neural networks M (or, rather, their semantics $[[M]]$) that are **stable** in the neighborhood $P_{\delta,\epsilon}(\mathbf{x})$ of a given input \mathbf{x}

$$\begin{aligned} \text{STABLE}_x^{\delta,\epsilon}([M]) &\stackrel{\text{def}}{=} \forall t \in [[M]]: (\exists t' \in [[M]]: \forall 0 \leq i \leq |L_0|: t'_0(x_{0,i}) = \mathbf{x}_i) \\ &\quad \wedge (\exists \mathbf{x}' \in P_{\delta,\epsilon}(\mathbf{x}): \forall 0 \leq i \leq |L_0|: t_0(x_{0,i}) = \mathbf{x}'_i) \\ &\quad \Rightarrow \max_j t_\omega(x_{N,j}) = \max_j t'_\omega(x_{N,j}) \end{aligned}$$

Theorem

$$M \models \mathcal{R}_x^{\delta,\epsilon} \Leftrightarrow \{[[M]]\} \subseteq \mathcal{R}_x^{\delta,\epsilon}$$

Corollary

$$M \models \mathcal{R}_x^{\delta,\epsilon} \Leftrightarrow [[M]] \subseteq \bigcup \mathcal{R}_x^{\delta,\epsilon}$$

Formal Methods

Mathematical Guarantees of Safety



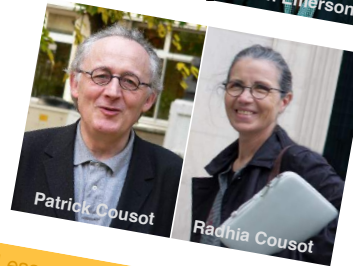
Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



Lesson 15

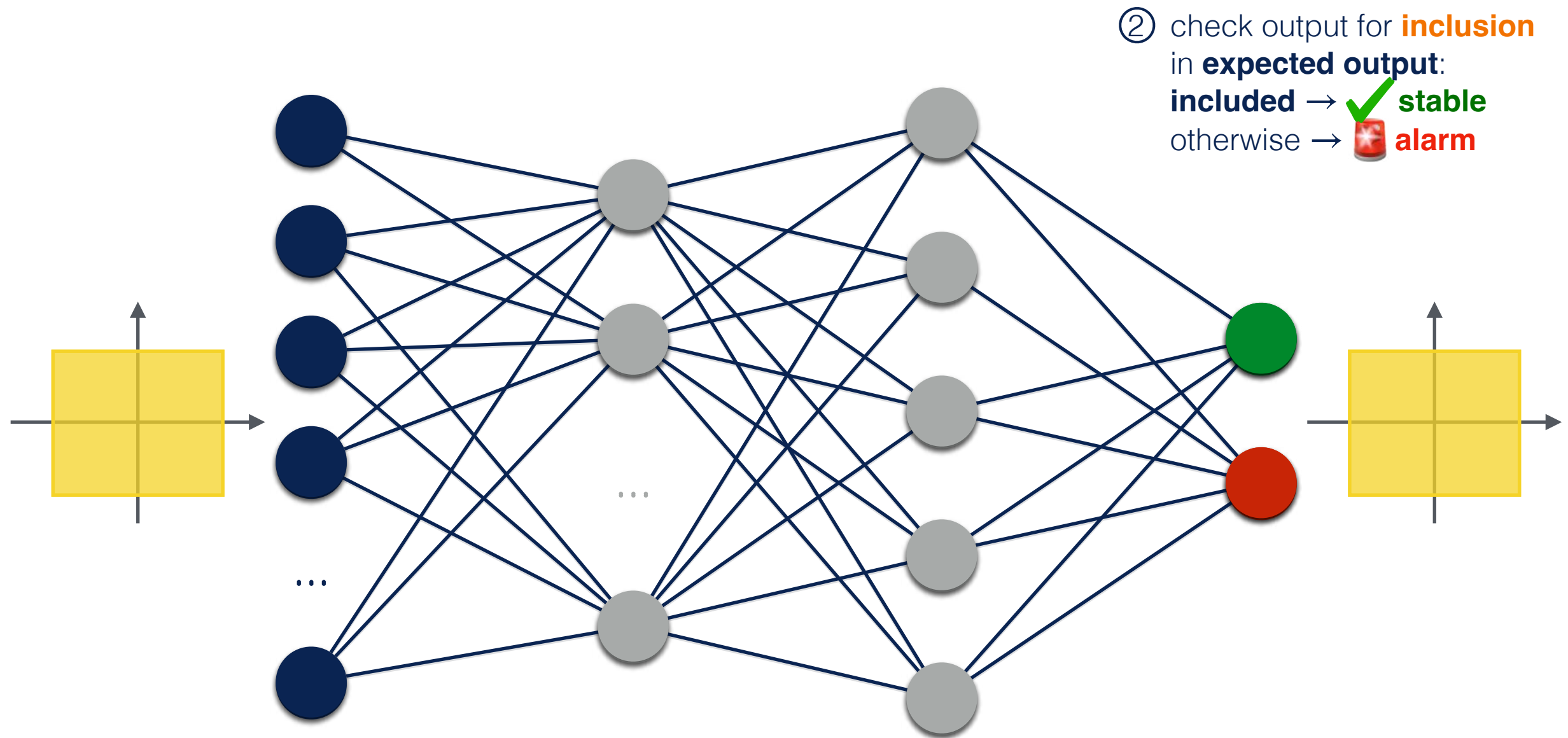
Formal Verification of Machine Learning

Caterina Urban

9

Static Analysis Methods

Forward Analysis

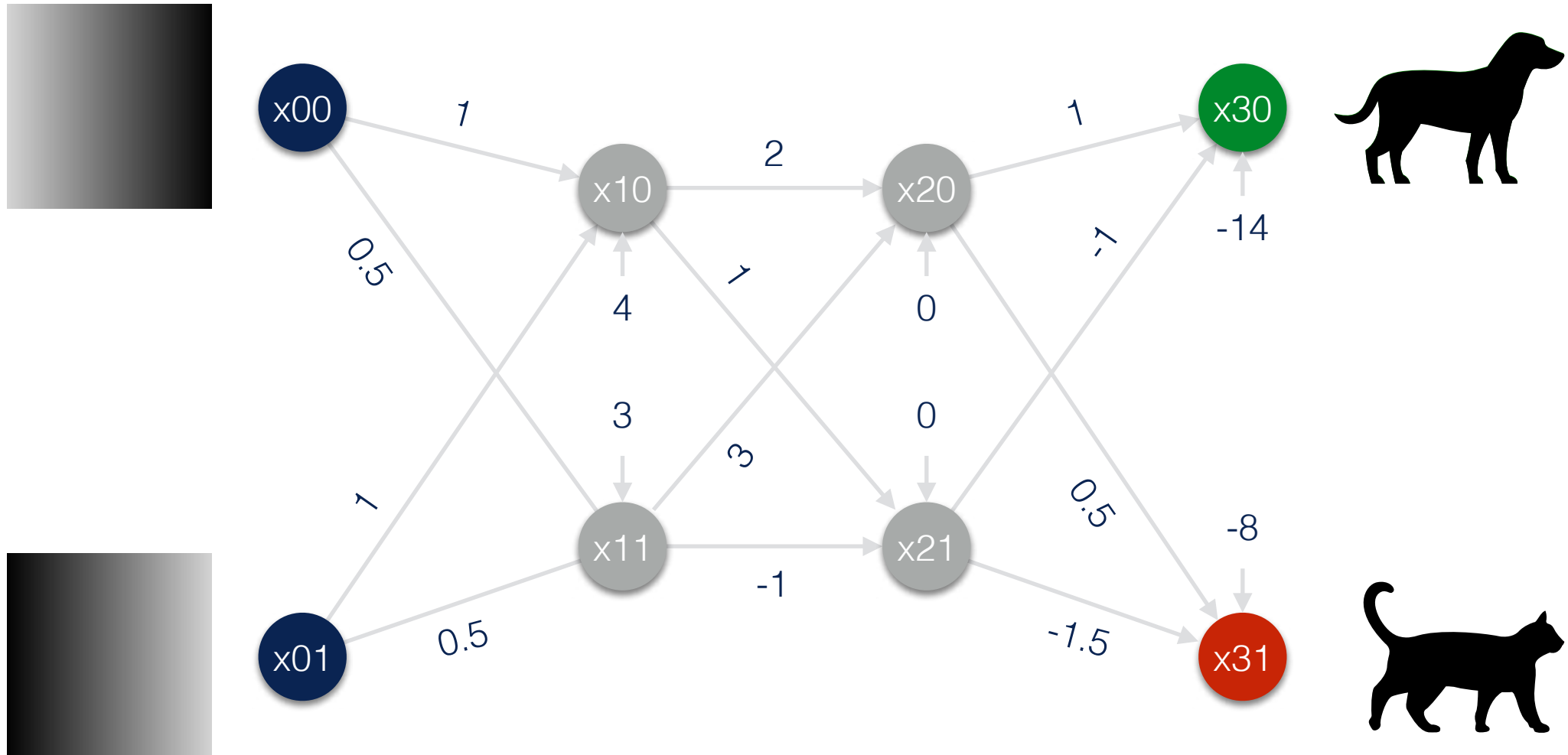


② check output for **inclusion** in **expected output**:
included → ✓ **stable**
otherwise → 🚨 **alarm**

① proceed **forwards** from **an abstraction** of all possible perturbations



Example

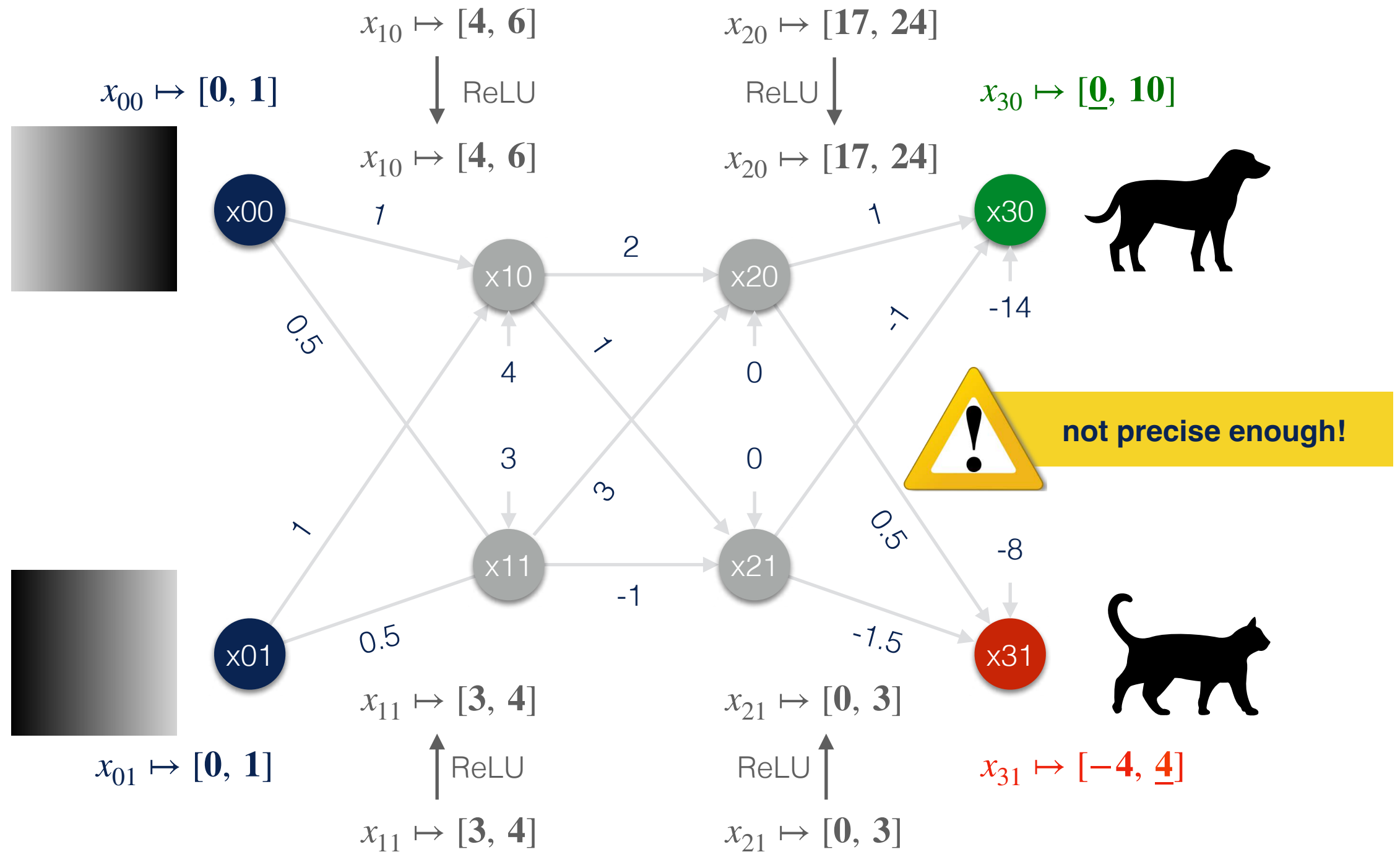


$$P(\langle 0.5, 0.75 \rangle) \stackrel{\text{def}}{=} \{ \mathbf{x} \in \mathcal{R} \times \mathcal{R} \mid 0 \leq \mathbf{x}_0 \leq 1 \wedge 0 \leq \mathbf{x}_1 \leq 1 \}$$

Interval Domain

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

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



Interval Domain

with **Symbolic Constant Propagation** [Li19]



each neuron as a **linear combination** of the inputs and the previous ReLUs

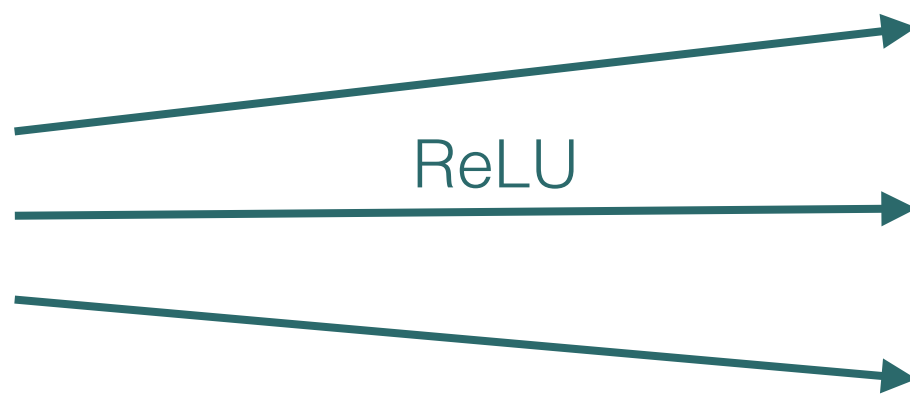
$$x_{i,j} \mapsto \begin{cases} \sum_{k=0}^{i-1} \mathbf{c}_k \cdot \mathbf{x}_k + \mathbf{c} & \mathbf{c}_k, \mathbf{c} \in \mathcal{R}^{|\mathbf{X}_k|} \\ [a, b] & a, b \in \mathcal{R} \end{cases}$$

$$\begin{aligned} x_{i-1,0} &\mapsto \mathbf{E}_{i-1,0} \\ \dots & \\ x_{i-1,j} &\mapsto \mathbf{E}_{i-1,j} \\ \dots & \end{aligned}$$

$$x_{i,j} = \sum_k w_{j,k}^{i-1} \cdot x_{i-1,k} + b_{i,j}$$

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

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



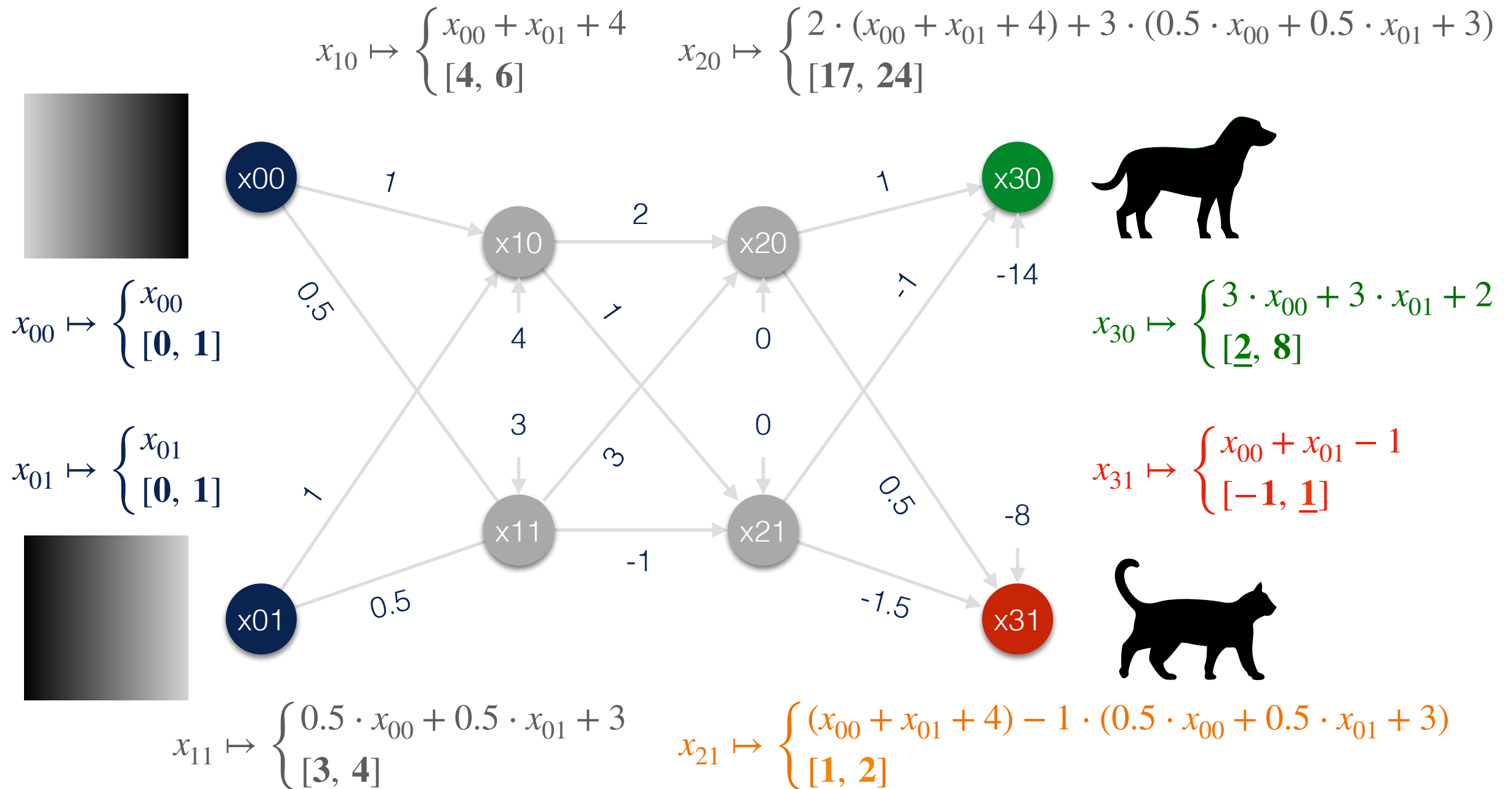
$$x_{i,j} \mapsto \begin{cases} \mathbf{E}_{i,j} \\ [a, b] \end{cases} \quad 0 \leq a$$

$$x_{i,j} \mapsto \begin{cases} \mathbf{x}_{i,j} \\ [0, b] \end{cases} \quad a < 0 \wedge 0 < b$$

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

Interval Domain

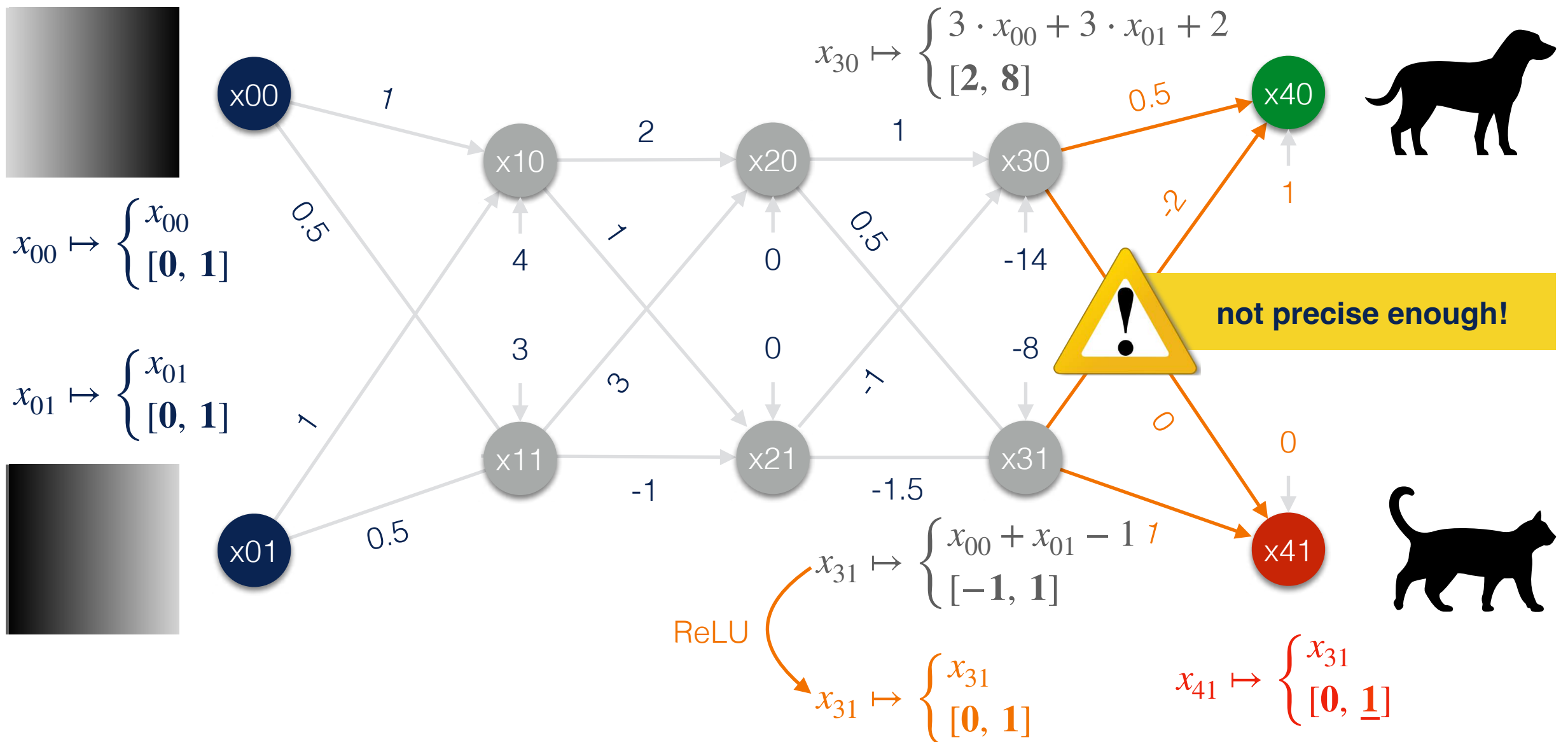
with Symbolic Constant Propagation [Li19]



Interval Domain

with **Symbolic Constant Propagation** [Li19]

$$x_{40} \mapsto \begin{cases} 1.5 \cdot x_{00} + 1.5 \cdot x_{01} + 2 \cdot x_{31} + 2 \\ [0, 5] \end{cases}$$

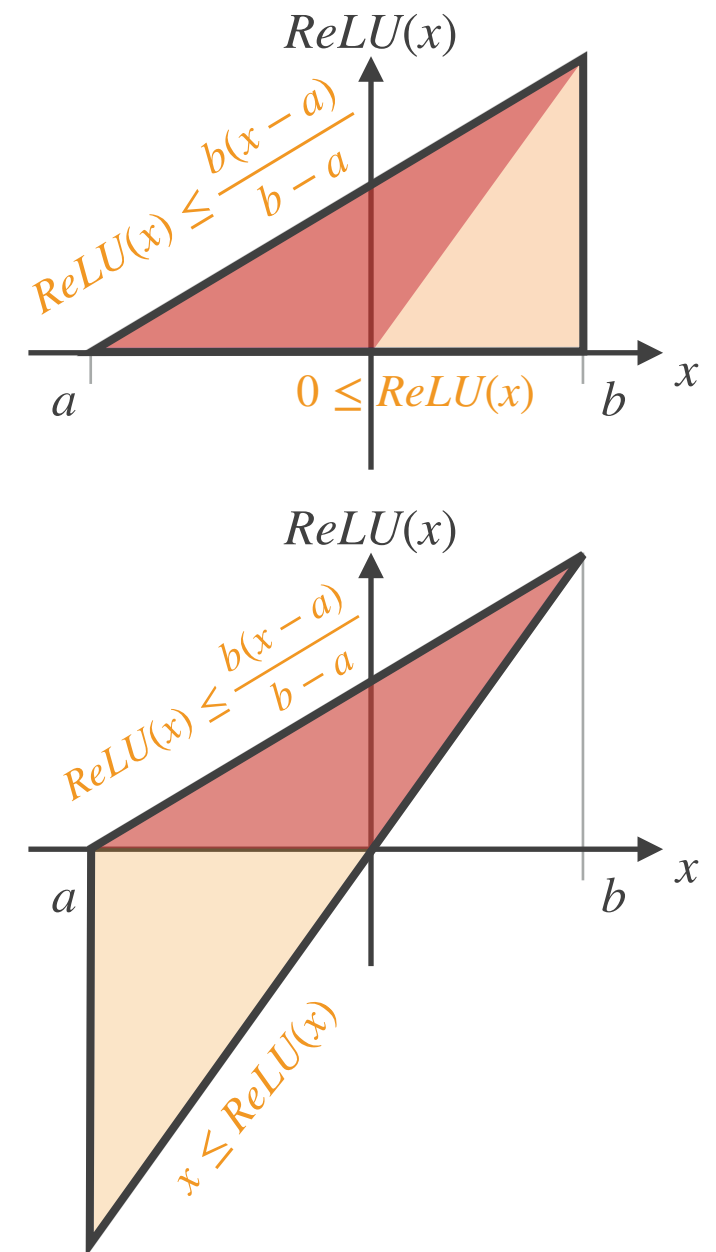
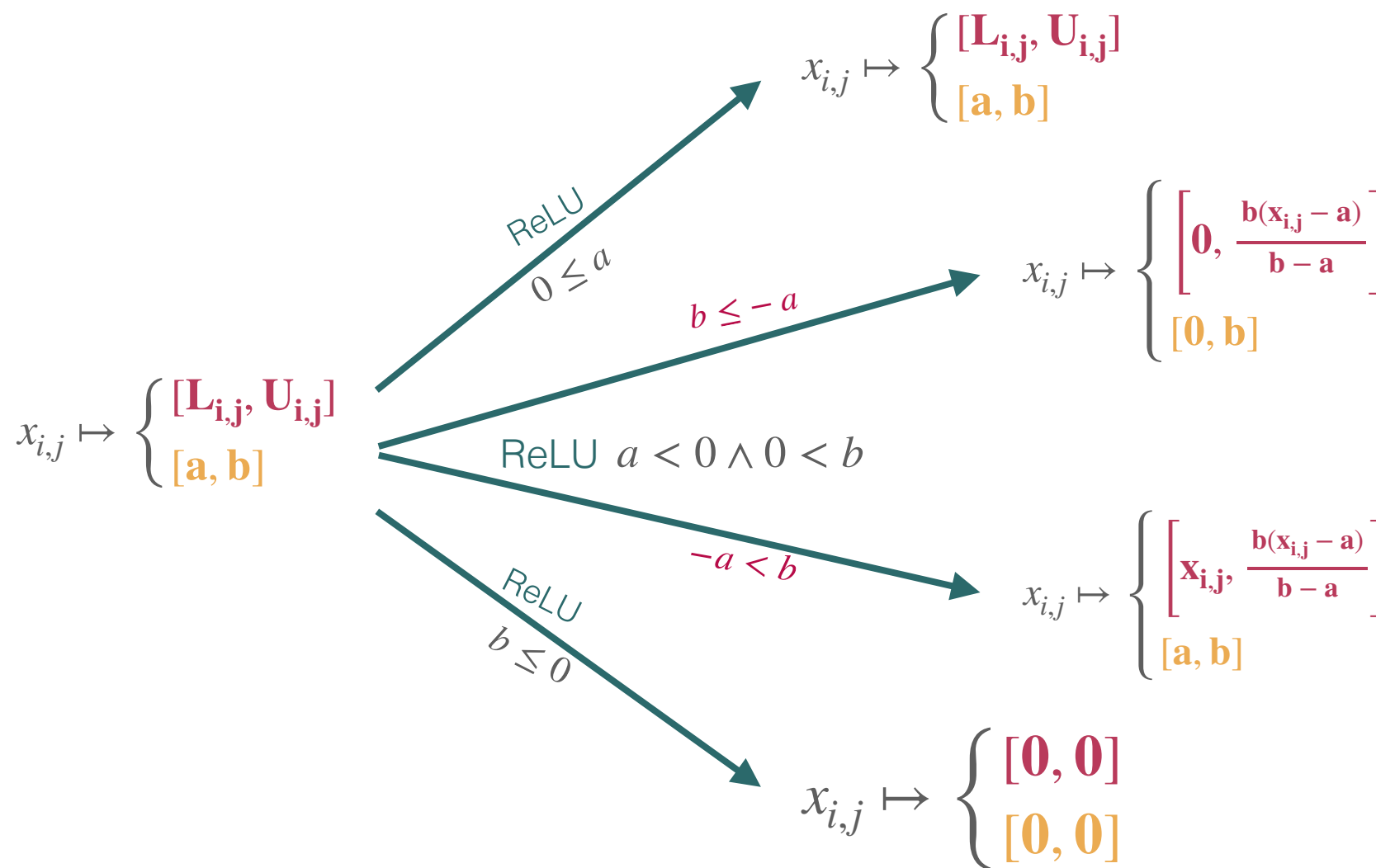


DeepPoly [Singh19]



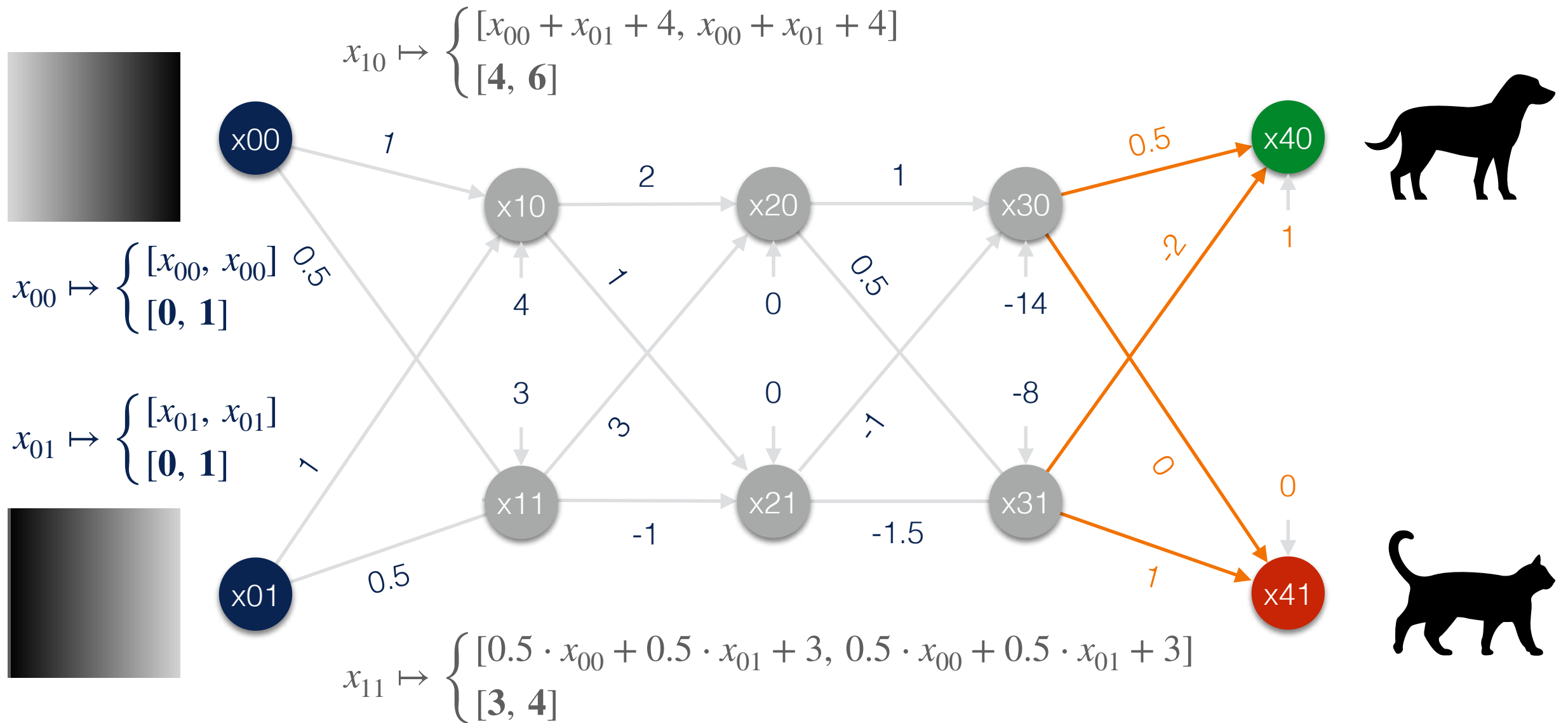
maintain symbolic lower- and upper-bounds for each neuron
+ convex ReLU approximations

$$x_{i+1,j} \mapsto \begin{cases} [\sum_k c_{i,k} \cdot x_{i,k} + c, \sum_k d_{i,k} \cdot x_{i,k} + d] & c_{i,k}, c, d_{i,k}, d \in \mathcal{R} \\ [a, b] & a, b \in \mathcal{R} \end{cases}$$

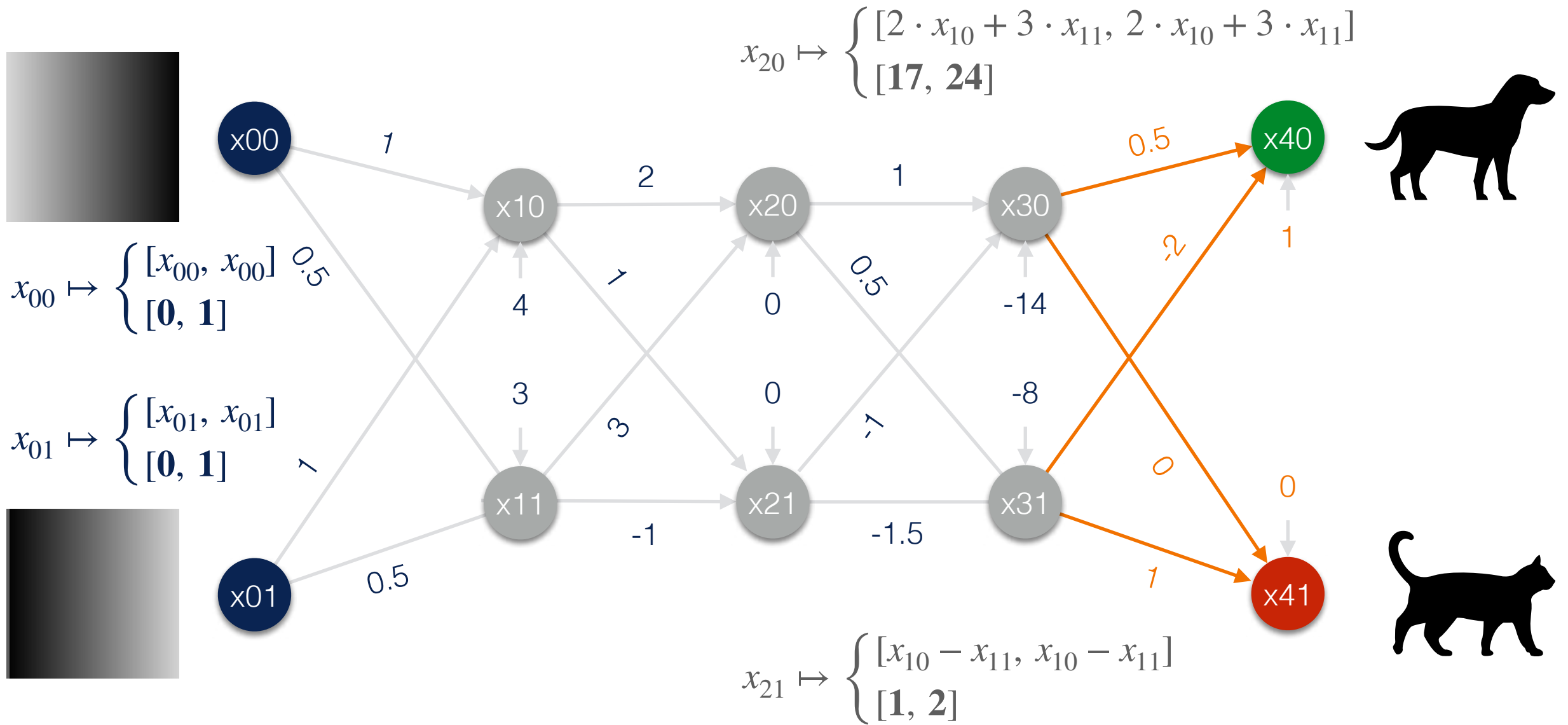


G. Singh, T. Gehr, M. Püschel, and M. Vechev - An Abstract Domain for Certifying Neural Networks (POPL 2019)

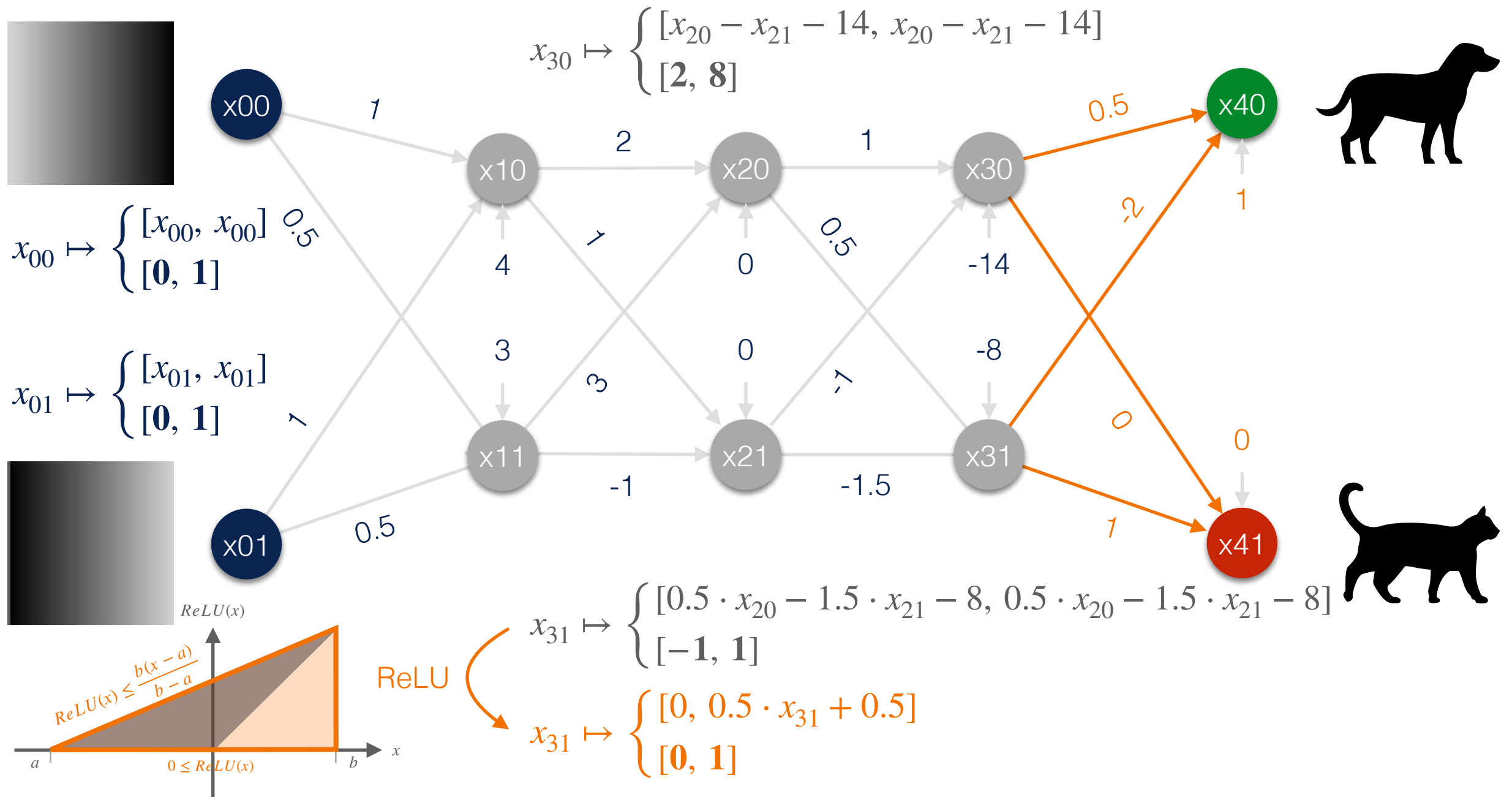
DeepPoly [Singh19]



DeepPoly [Singh19]

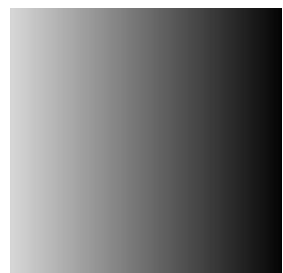


DeepPoly [Singh19]



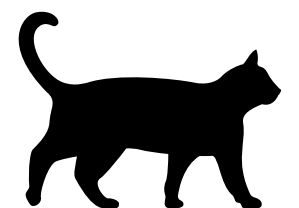
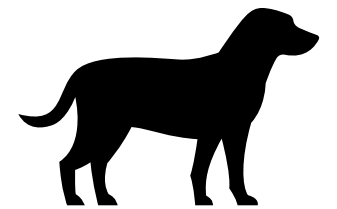
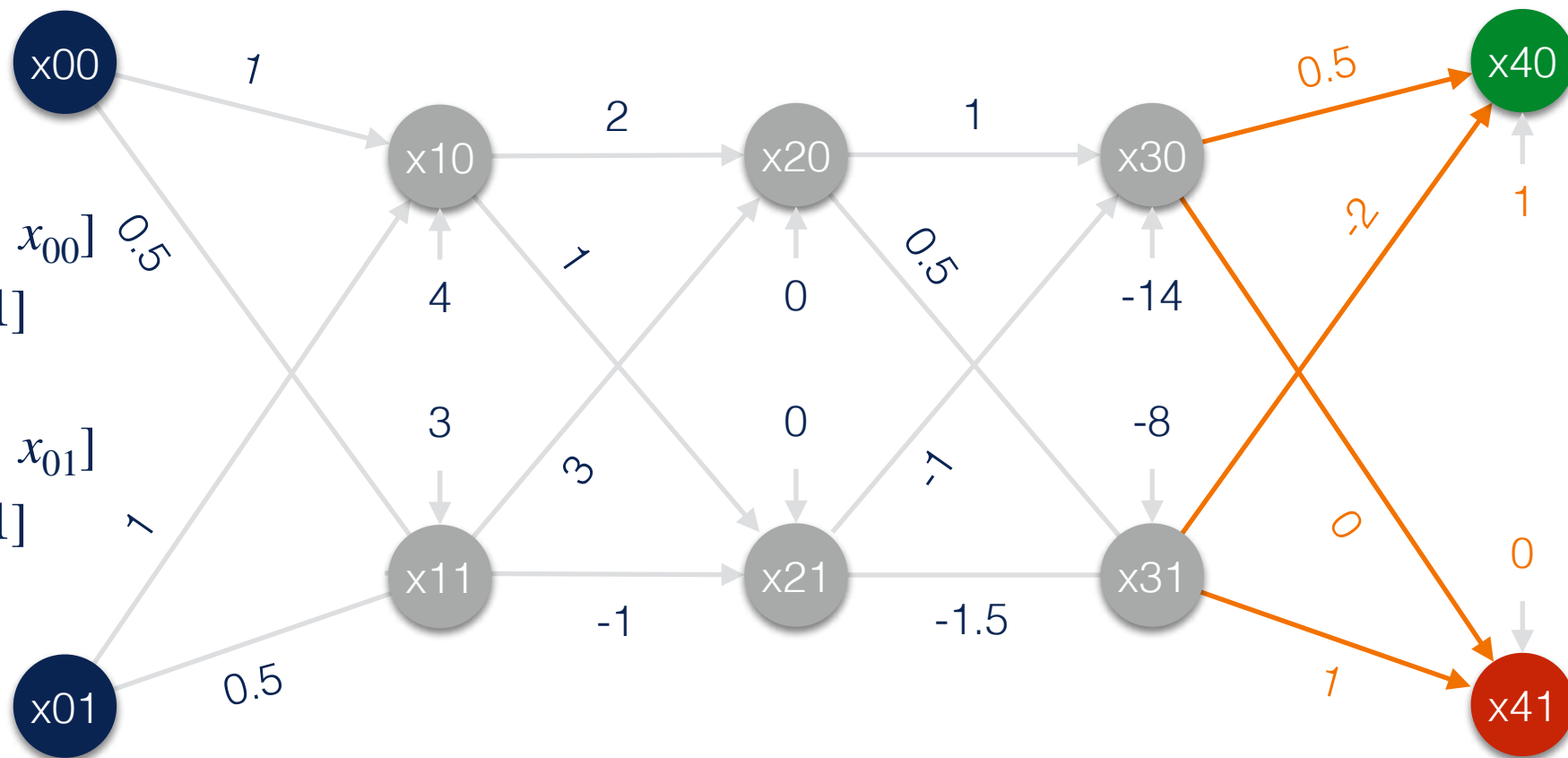
DeepPoly [Singh19]

$$x_{40} \mapsto \left\{ [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \right\}$$

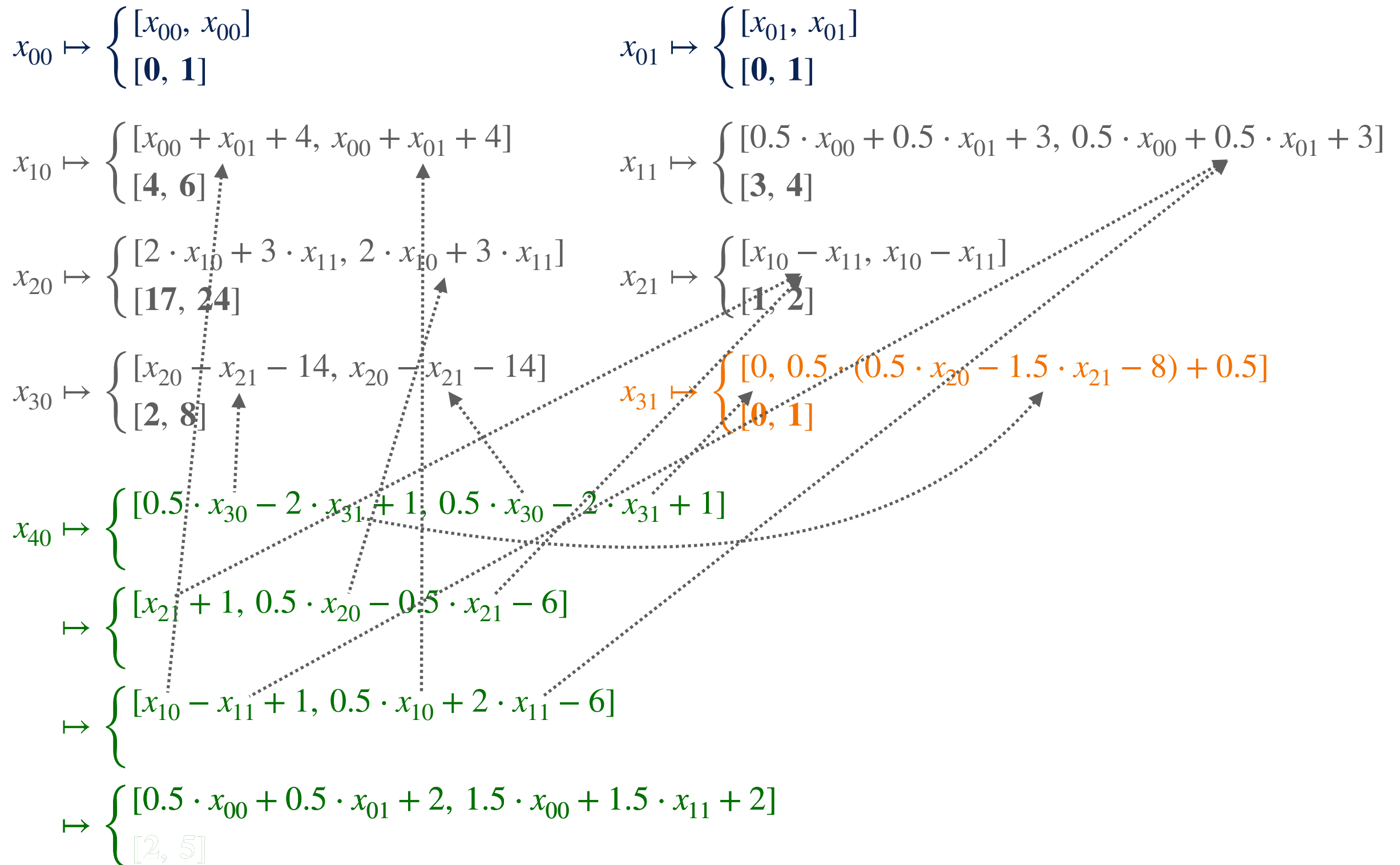


$$x_{00} \mapsto \begin{cases} [x_{00}, x_{00}] \\ [0, 1] \end{cases} \quad 0.5$$

$$x_{01} \mapsto \begin{cases} [x_{01}, x_{01}] \\ [0, 1] \end{cases} \quad 1$$

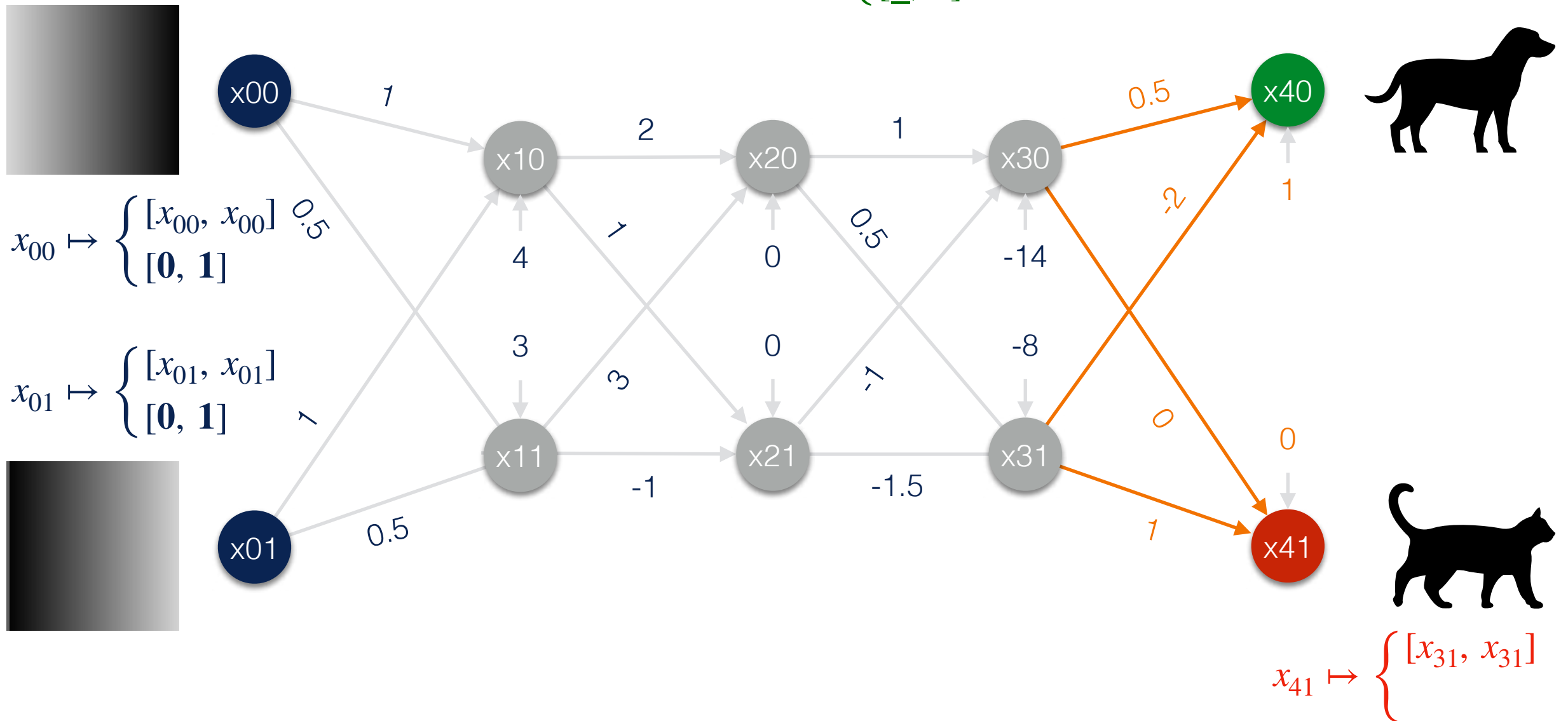


DeepPoly [Singh19]



DeepPoly [Singh19]

$$x_{40} \mapsto \begin{cases} [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \\ [\underline{2}, \underline{5}] \end{cases}$$



DeepPoly [Singh19]

$$x_{00} \mapsto \begin{cases} [x_{00}, x_{00}] \\ [\mathbf{0}, \mathbf{1}] \end{cases}$$

$$x_{01} \mapsto \begin{cases} [x_{01}, x_{01}] \\ [\mathbf{0}, \mathbf{1}] \end{cases}$$

$$x_{10} \mapsto \begin{cases} [x_{00} + x_{01} + 4, x_{00} + x_{01} + 4] \\ [\mathbf{4}, \mathbf{6}] \end{cases}$$

$$x_{11} \mapsto \begin{cases} [0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3, 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3] \\ [\mathbf{3}, \mathbf{4}] \end{cases}$$

$$x_{20} \mapsto \begin{cases} [2 \cdot x_{10} + 3 \cdot x_{11}, 2 \cdot x_{10} + 3 \cdot x_{11}] \\ [\mathbf{17}, \mathbf{24}] \end{cases}$$

$$x_{21} \mapsto \begin{cases} [x_{10} - x_{11}, x_{10} - x_{11}] \\ [\mathbf{1}, \mathbf{2}] \end{cases}$$

$$x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [\mathbf{2}, \mathbf{8}] \end{cases}$$

$$x_{31} \mapsto \begin{cases} [0, 0.5 \cdot (0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5] \\ [\mathbf{0}, \mathbf{1}] \end{cases}$$

$$x_{41} \mapsto \begin{cases} [x_{31}, x_{31}] \end{cases}$$

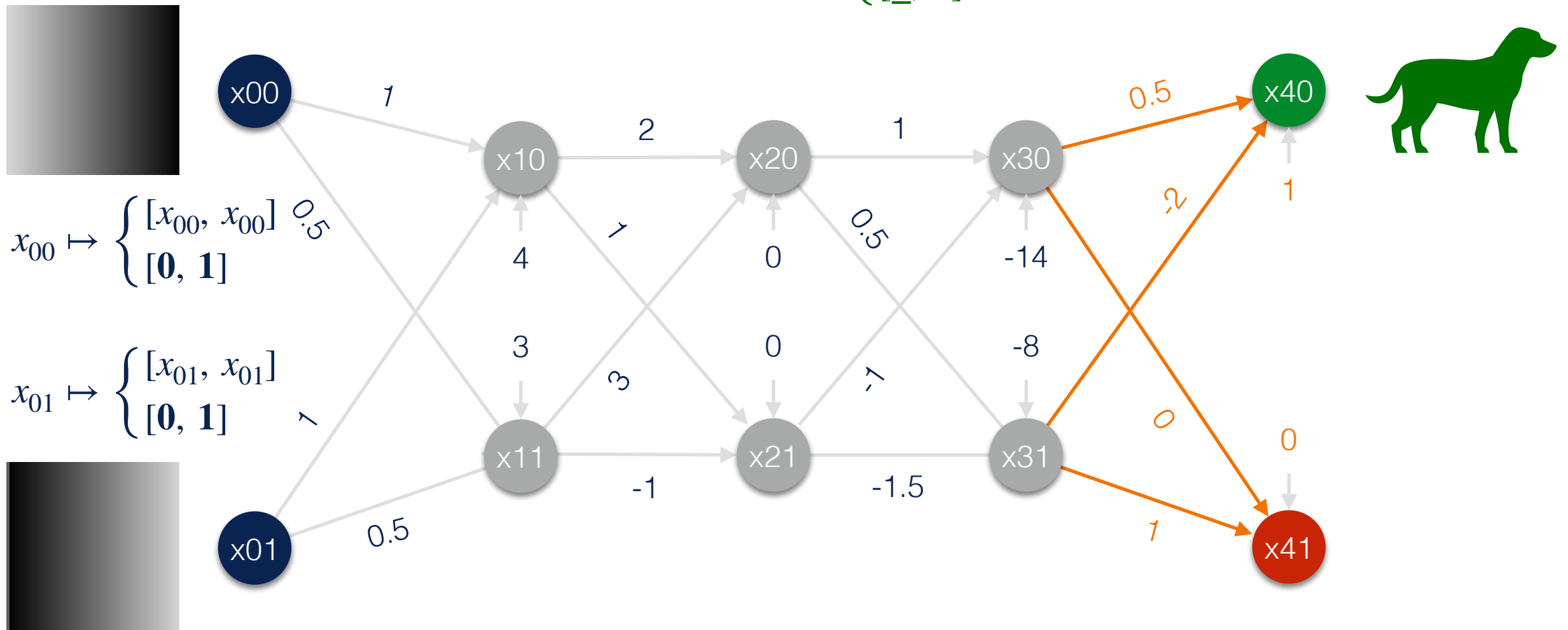
$$\mapsto \begin{cases} [0, 0.25 \cdot x_{20} - 0.75 \cdot x_{21} - 3.5] \end{cases}$$

$$\mapsto \begin{cases} [0, -0.25 \cdot x_{10} + 1.5 \cdot x_{11} - 3.5] \end{cases}$$

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

DeepPoly [Singh19]

$$x_{40} \mapsto \begin{cases} [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \\ [\underline{2}, \underline{5}] \end{cases}$$

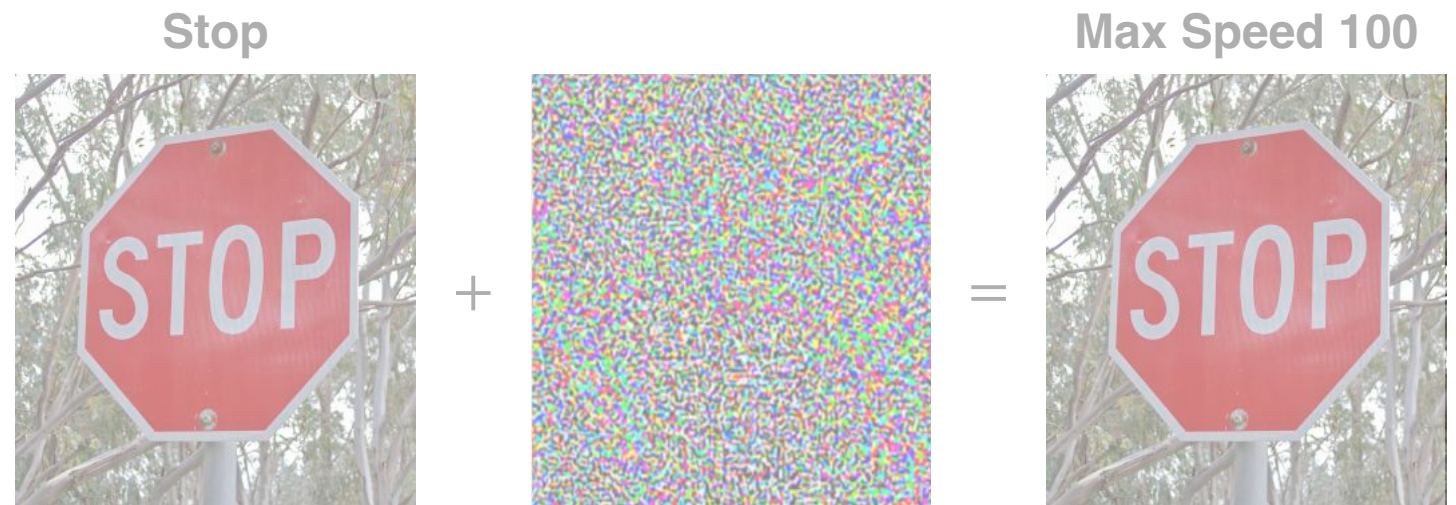


Other Static Analysis Methods

- **T. Gehr, M. Mirman, D. Drachler-Cohen, P. Tsankov, S. Chaudhuri, and M. Vechev.** *AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation.* In S&P, 2018.
the first use of abstract interpretation for verifying neural networks
- **G. Singh, T. Gehr, M. Mirman, M. Püschel, and M. Vechev.** *Fast and Effective Robustness Certification.* In NeurIPS, 2018.
a custom zonotope domain for certifying neural networks
- **G. Singh, R. Ganvir, M. Püschel, and M. Vechev.** *Beyond the Single Neuron Convex Barrier for Neural Network Certification.* In NeurIPS, 2019.
a framework to jointly approximate k ReLU activations
- **M. N. Müller, G. Makarchuk, G. Singh, M. Püschel, and M. Vechev.** *PRIMA: General and Precise Neural Network Certification via Scalable Convex Hull Approximations.* In POPL, 2022.
a multi-neuron abstraction via a convex-hull approximation algorithm

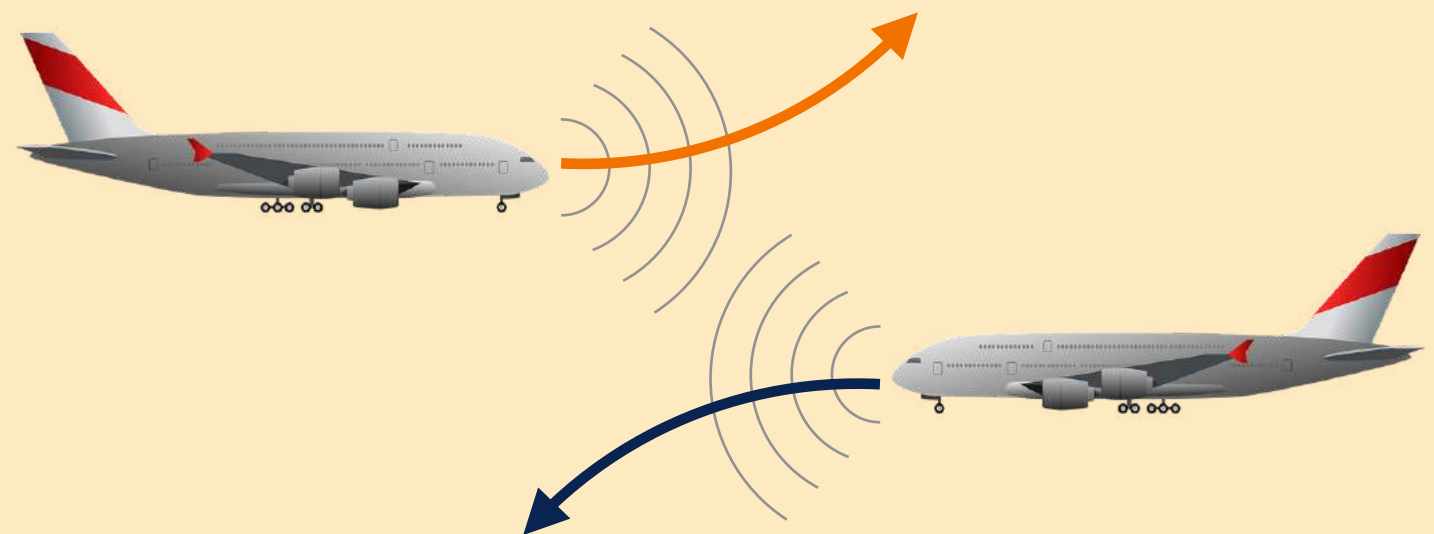
Stability

Goal G3 in [Kurd03]



Safety

Goal G4 in [Kurd03]



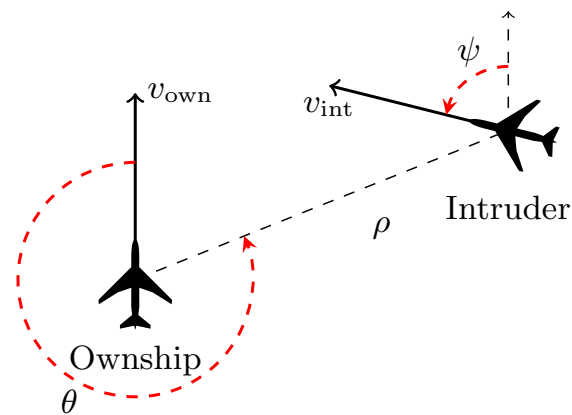
Fairness



ACAS Xu [Julian16][Katz17]

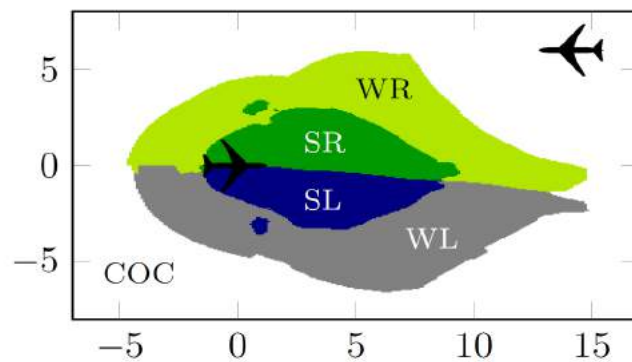
Airborne Collision Avoidance System for Unmanned Aircraft

implemented using **45 feed-forward fully-connected ReLU networks**



5 input sensor measurements

- ρ : distance from ownship to intruder
- θ : angle to intruder relative to ownship heading direction
- ψ : heading angle to intruder relative to ownship heading direction
- v_{own} : speed of ownship
- v_{int} : speed of intruder

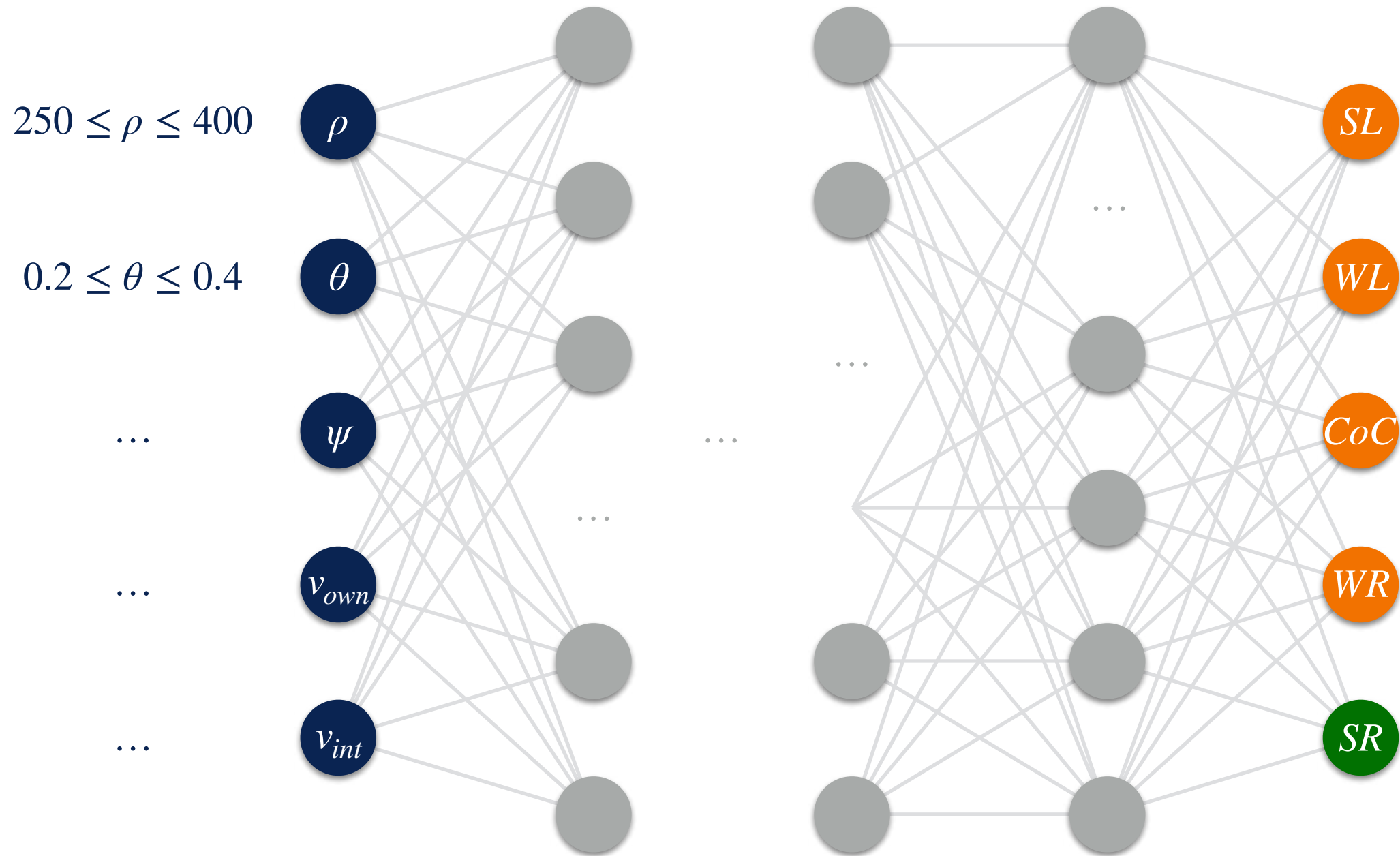


5 output horizontal advisories

- Strong Left
- Weak Left
- Clear of Conflict
- Weak Right
- Strong Right

ACAS Xu Properties [Katz17]

Example: “if intruder is **near** and approaching **from the left**, go **Strong Right**”



Safety

Input-Output Properties

I: input specification

O: output specification

$$\mathcal{S}_{\mathbf{O}}^{\mathbf{I}} \stackrel{\text{def}}{=} \{ \llbracket M \rrbracket \in \mathcal{P}(\Sigma^*) \mid \text{SAFE}_{\mathbf{O}}^{\mathbf{I}}(\llbracket M \rrbracket) \}$$

$\mathcal{S}_{\mathbf{O}}^{\mathbf{I}}$ is the set of all neural networks M (or, rather, their semantics $\llbracket M \rrbracket$) that **satisfy** the input and output specification **I** and **O**

$$\text{SAFE}_{\mathbf{O}}^{\mathbf{I}}(\llbracket M \rrbracket) \stackrel{\text{def}}{=} \forall t \in \llbracket M \rrbracket : t_0 \models \mathbf{I} \Rightarrow t_\omega \models \mathbf{O}$$

Theorem

$$M \models \mathcal{S}_{\mathbf{O}}^{\mathbf{I}} \Leftrightarrow \{ \llbracket M \rrbracket \} \subseteq \mathcal{S}_{\mathbf{O}}^{\mathbf{I}}$$

Corollary

$$M \models \mathcal{S}_{\mathbf{O}}^{\mathbf{I}} \Leftrightarrow \llbracket M \rrbracket \subseteq \bigcup \mathcal{S}_{\mathbf{O}}^{\mathbf{I}}$$

Formal Methods

Mathematical Guarantees of Safety



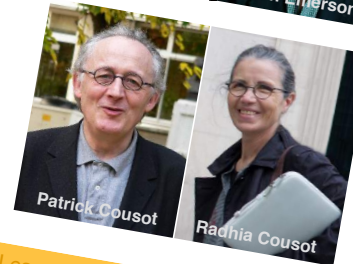
Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



Lesson 15

Formal Verification of Machine Learning

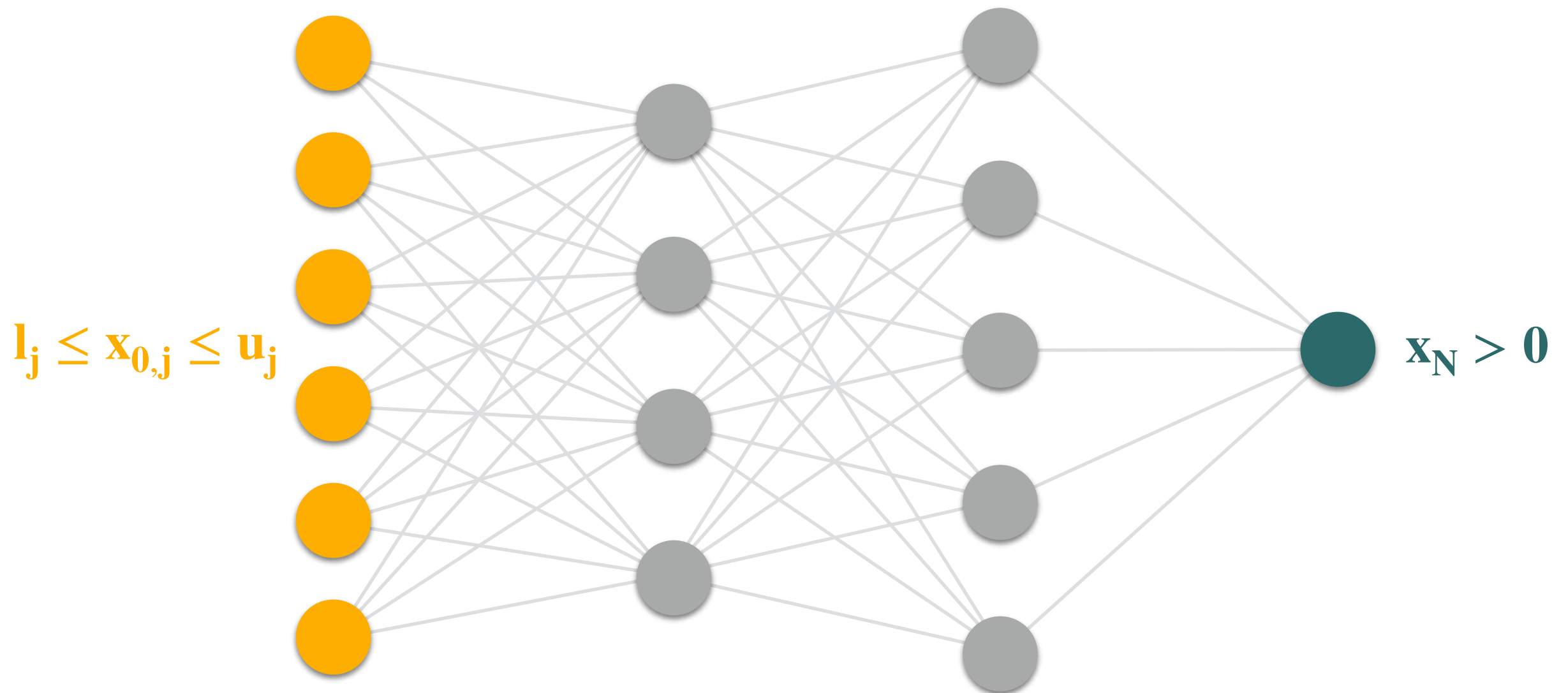
Caterina Urban

9

Model Checking Methods

Safety

Example



SMT-Based Methods

Verification Reduced to Constraint Satisfiability

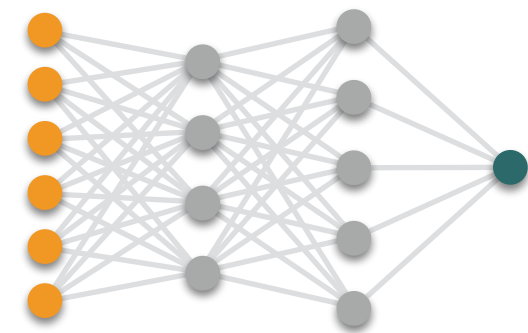
$$\mathbf{l}_j \leq \mathbf{x}_{0,j} \leq \mathbf{u}_j \quad j \in \{0, \dots, |\mathbf{X}_0|\}$$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \quad i \in \{0, \dots, n-1\}$$

$$x_{i,j} = \max\{0, \hat{x}_{i,j}\} \quad \begin{array}{l} i \in \{1, \dots, n-1\}, \\ j \in \{0, \dots, |\mathbf{X}_i|\} \end{array}$$

$$\mathbf{x}_N \leq \mathbf{0}$$

input specification



(negation of)
output specification

satisfiable → ~~X~~ counterexample
otherwise → ✓ safe

Planet



use **approximations** to reduce the solution search space

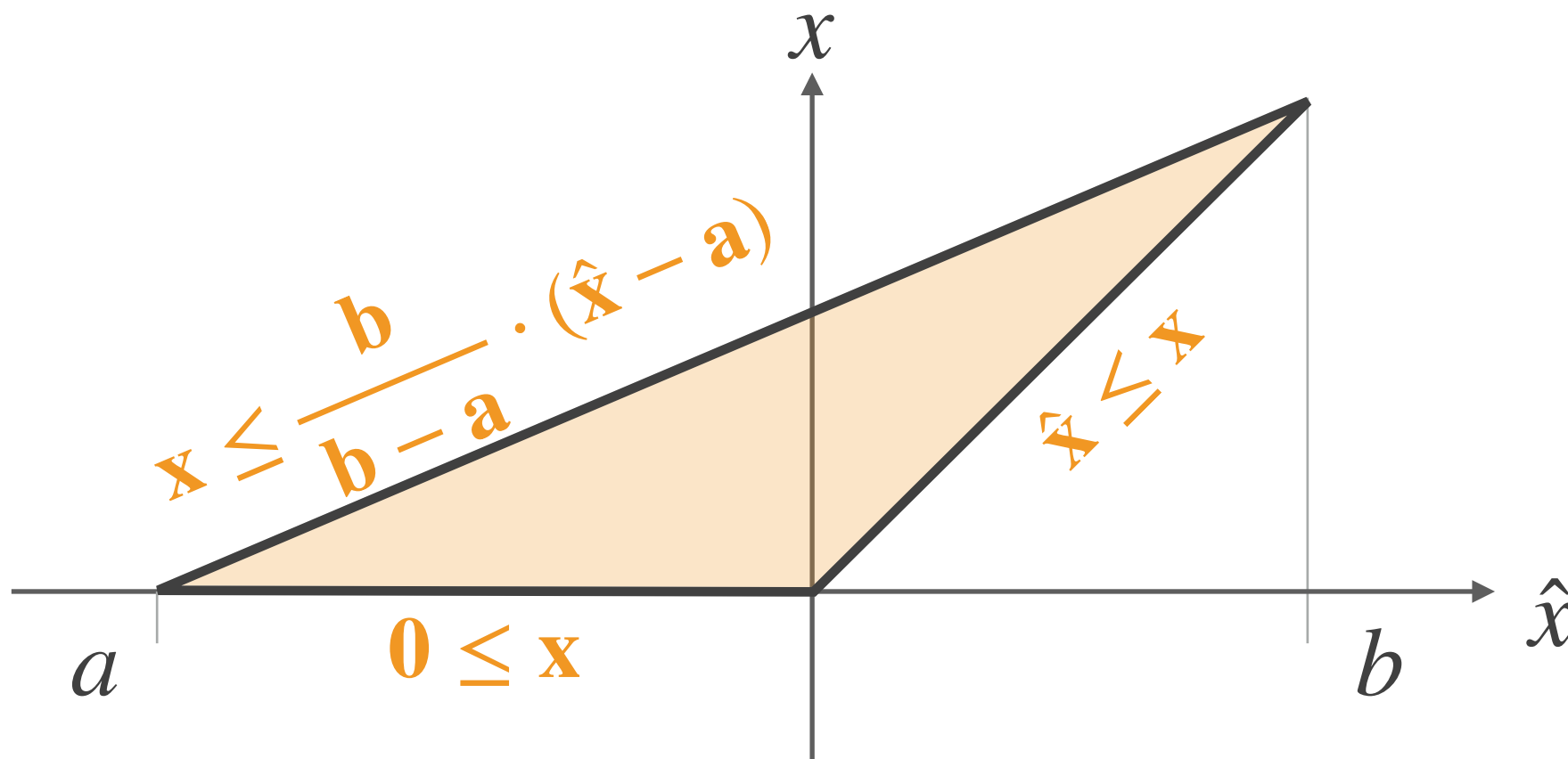
$$x_{i,j} = \max\{0, \hat{x}_{i,j}\}$$



$$0 \leq x_{i,j}$$

$$\hat{x}_{i,j} \leq x_{i,j}$$

$$x_{i,j} \leq \frac{b_{i,j}}{b_{i,j} - a_{i,j}} \cdot (\hat{x}_{i,j} - a_{i,j})$$

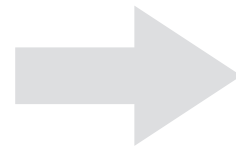


Reluplex



based on the **simplex algorithm** extended to support ReLUs

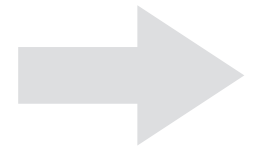
Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}_{ij}
x_{ij}	v_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	v_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	\hat{v}'_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	0
...	...
x_N	v_N



Reluplex



based on the extended

Follow-up Work

G. Katz et al. - The Marabou Framework for Verification and Analysis of Deep Neural Networks (CAV 2019)

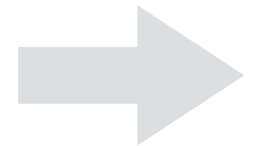
Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}_{ij}
x_{ij}	v_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	v_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	\hat{v}'_{ij}
...	...
x_N	v_N



Variable	Value
x_{00}	v_{00}
...	...
\hat{x}_{ij}	\hat{v}'_{ij}
x_{ij}	0
...	...
x_N	v_N



G. Katz et al. - Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks (CAV 2017)

Other SMT-Based Methods

- **L. Pulina and A. Tacchella.** *An Abstraction-Refinement Approach to Verification of Artificial Neural Networks.* In CAV, 2010.
the first formal verification method for neural networks
- **O. Bastani, Y. Ioannou, L. Lampropoulos, D. Vytiniotis, A. Nori, and A. Criminisi.** *Measuring Neural Net Robustness with Constraints.* In NeurIPS, 2016.
an approach for finding the nearest adversarial example according to the L_∞ distance
- **X. Huang, M. Kwiatkowska, S. Wang, and M. Wu.** *Safety Verification of Deep Neural Networks.* In CAV, 2017.
an approach for proving local robustness to adversarial perturbations
- **N. Narodytska, S. Kasiviswanathan, L. Ryzhyk, M. Sagiv, and T. Walsh.** *Verifying Properties of Binarized Deep Neural Networks.* In AAAI, 2018.
C. H. Cheng, G. Nührenberg, C. H. Huang, and H. Ruess. *Verification of Binarized Neural Networks via Inter-Neuron Factoring.* In VSTTE, 2018.
approaches focusing on binarized neural networks

MILP-Based Methods

Verification Reduced to Mixed Integer Linear Program

$$l_j \leq \mathbf{x}_{0,j} \leq u_j$$

$$j \in \{0, \dots, |\mathbf{X}_0|\}$$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j}$$

$$i \in \{0, \dots, n-1\}$$

$$x_{i,j} = \delta_{i,j} \cdot \hat{x}_{i,j}$$

$$\delta_{i,j} \in \{0, 1\}$$

$$\delta_{i,j} = 1 \Rightarrow \hat{x}_{i,j} \geq 0$$

$$i \in \{1, \dots, n-1\}$$

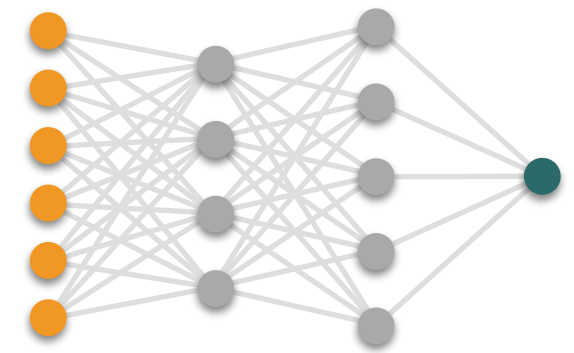
$$\delta_{i,j} = 0 \Rightarrow \hat{x}_{i,j} < 0$$

$$j \in \{0, \dots, |\mathbf{X}_i|\}$$

min \mathbf{x}_N

min $\mathbf{x}_N \leq \mathbf{0} \rightarrow \text{X counterexample}$
otherwise $\rightarrow \text{✓ safe}$

input specification



objective function

MILP-Based Methods

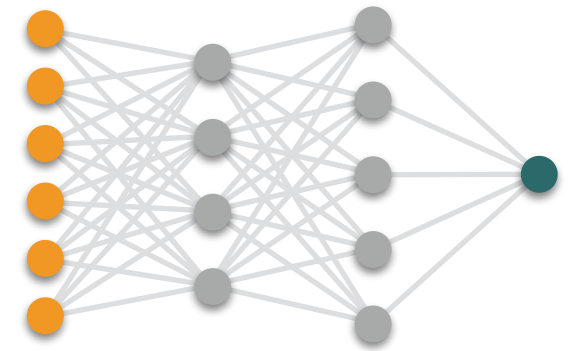
Bounded Encoding with **Symmetric Bounds**

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \quad i \in \{0, \dots, n-1\}$$

$$0 \leq x_{i,j} \leq \mathbf{M}_{i,j} \cdot \delta_{i,j} \quad \delta_{i,j} \in \{0, 1\}$$

$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{M}_{i,j} \cdot (1 - \delta_{i,j}) \quad i \in \{1, \dots, n-1\}$$

$$\mathbf{M}_{i,j} = \max\{-l_i, u_i\} \quad j \in \{0, \dots, |\mathbf{X}_i|\}$$



Sherlock

Output Range Analysis



use **local search** to speed up the MILP solver

$$l_j \leq \mathbf{x}_{0,j} \leq u_j$$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j}$$

$$0 \leq x_{i,j} \leq \mathbf{M}_{i,j} \cdot \delta_{i,j}$$

$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{M}_{i,j} \cdot (1 - \delta_{i,j})$$

$$\mathbf{M}_{i,j} = \max\{-l_j, u_j\}$$

$$\mathbf{x}_N < \mathbf{L}$$

sample random input \mathbf{X}
and evaluate output \mathbf{L}



Sherlock

Output Range Analysis



use **local search** to speed up the MILP solver

$$l_j \leq \mathbf{x}_{0,j} \leq u_j$$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j}$$

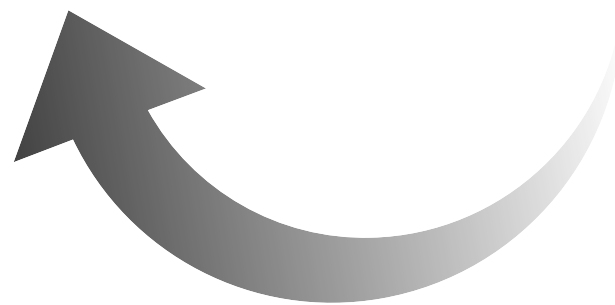
$$0 \leq x_{i,j} \leq \mathbf{M}_{i,j} \cdot \delta_{i,j}$$

$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{M}_{i,j} \cdot (1 - \delta_{i,j})$$

$$\mathbf{M}_{i,j} = \max\{-l_j, u_j\}$$

$$\mathbf{x}_N < \hat{\mathbf{L}}$$

find another input $\hat{\mathbf{X}}$
such that $\hat{\mathbf{L}} \leq \mathbf{x}_N$



Sherlock

Output Range Analysis



use **local search** to speed up the MILP solver

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

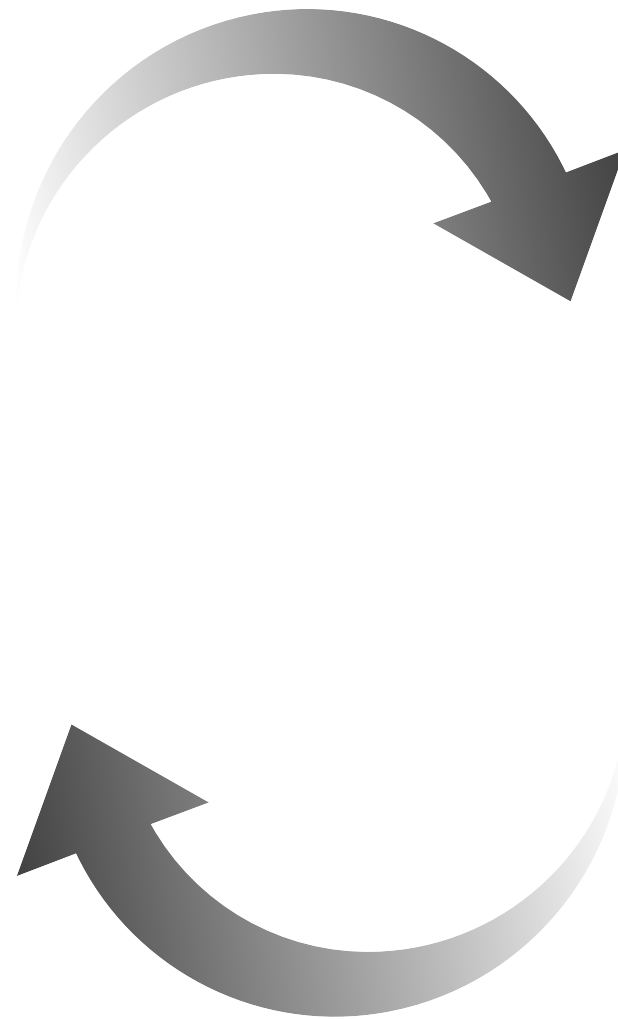
$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j}$$

$$0 \leq x_{i,j} \leq \mathbf{M}_{i,j} \cdot \delta_{i,j}$$

$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{M}_{i,j} \cdot (1 - \delta_{i,j})$$

$$\mathbf{M}_{i,j} = \max\{-l_j, u_j\}$$

$$\mathbf{x}_N < \hat{\mathbf{L}}$$



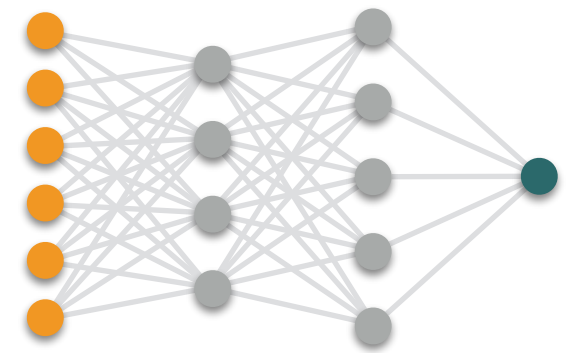
find another input $\hat{\mathbf{X}}$
such that $\hat{\mathbf{L}} \leq \mathbf{x}_N$

MILP-Based Methods

Bounded Encoding with **Asymmetric Bounds**

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \quad i \in \{0, \dots, n-1\}$$

$$0 \leq x_{i,j} \leq \mathbf{u}_{i,j} \cdot \delta_{i,j} \quad \delta_{i,j} \in \{0, 1\}$$
$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{l}_{i,j} \cdot (1 - \delta_{i,j}) \quad i \in \{1, \dots, n-1\}$$
$$j \in \{0, \dots, |\mathbf{X}_i|\}$$



MIPVerify

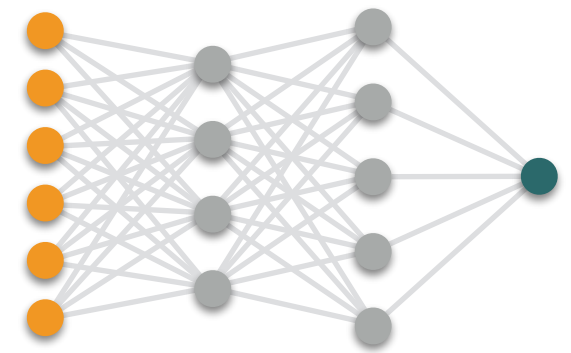
Finding Nearest Adversarial Example

$$\min_{\mathbf{X}'} \mathbf{d}(\mathbf{X}, \mathbf{X}')$$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \quad i \in \{0, \dots, n-1\}$$

$$0 \leq x_{i,j} \leq \mathbf{u}_{i,j} \cdot \delta_{i,j} \quad \delta_{i,j} \in \{0, 1\}$$
$$\hat{x}_{i,j} \leq x_{i,j} \leq \hat{x}_{i,j} - \mathbf{l}_{i,j} \cdot (1 - \delta_{i,j}) \quad i \in \{1, \dots, n-1\}$$
$$j \in \{0, \dots, |\mathbf{X}_i|\}$$

$$\mathbf{x}_N \neq \mathbf{0}$$



Other MILP-Based Methods

- **R. Bunel, I. Turkaslan, P. H. S. Torr, P. Kohli, and M. P. Kumar.** *A Unified View of Piecewise Linear Neural Network Verification.* In NeurIPS, 2018.
a unifying verification framework for piecewise-linear ReLU neural networks
- **C.-H. Cheng, G. Nührenberg, and H. Rues.** *Maximum Resilience of Artificial Neural Networks.* In ATVA, 2017.
an approach for finding a lower bound on robustness to adversarial perturbations
- **M. Fischetti and J. Jo.** *Deep Neural Networks and Mixed Integer Linear Optimization.* 2018.
an approach for feature visualization and building adversarial examples

Formal Methods

Mathematical Guarantees of Safety



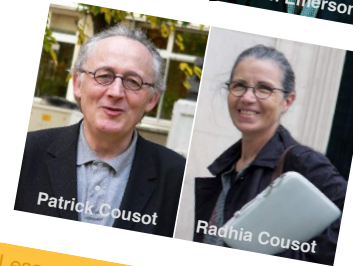
Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



Lesson 15

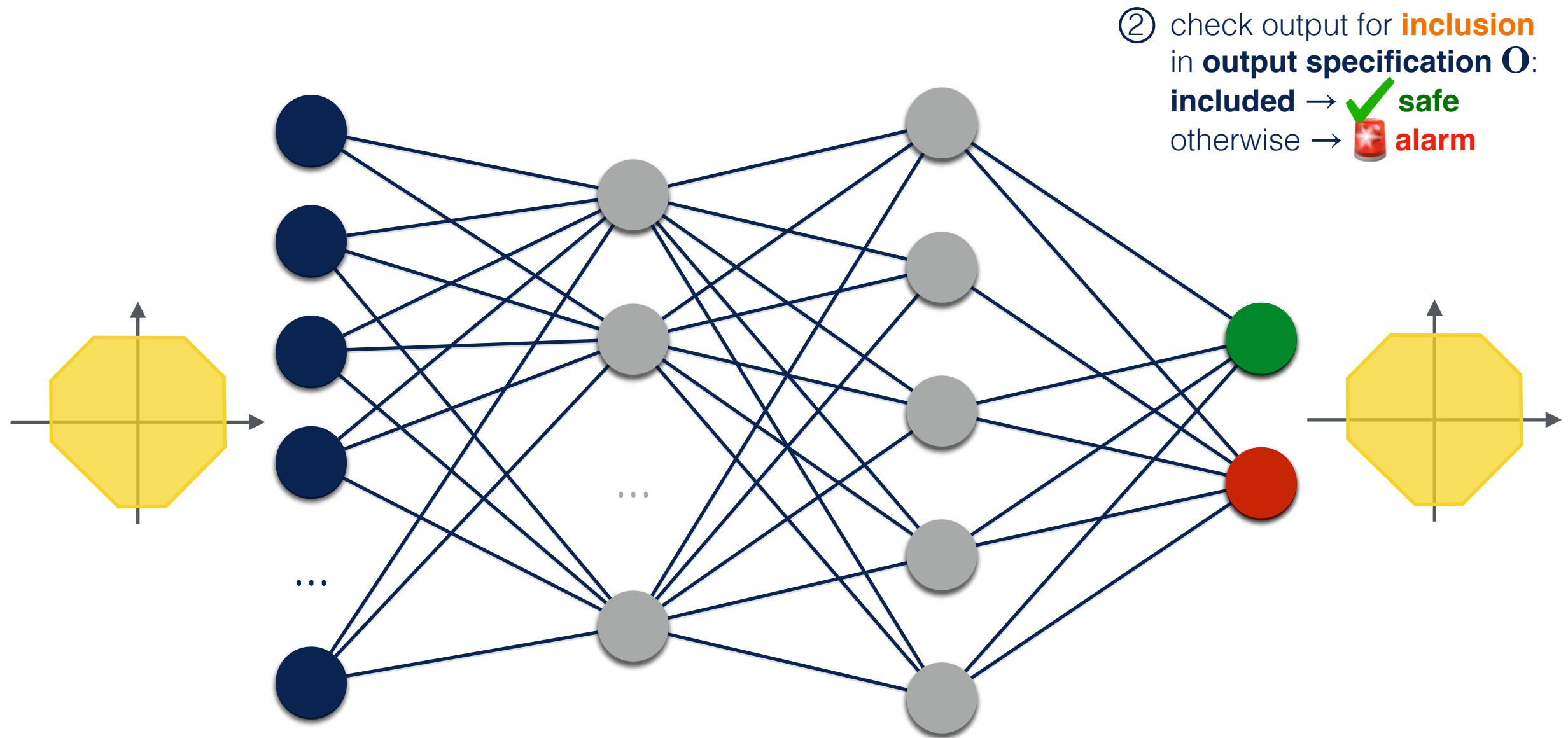
Formal Verification of Machine Learning

Caterina Urban

9

Static Analysis Methods

Forward Analysis

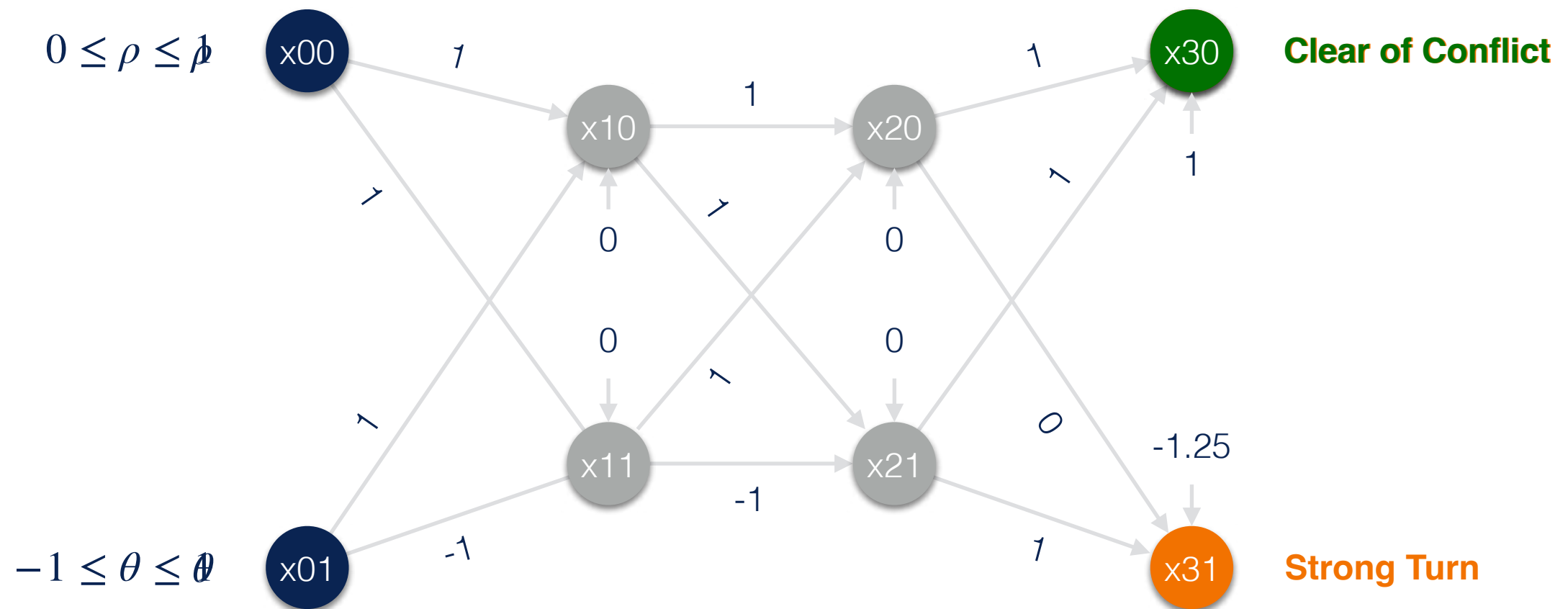


② check output for **inclusion** in **output specification O**:
included → ✓ **safe**
otherwise → 🚨 **alarm**

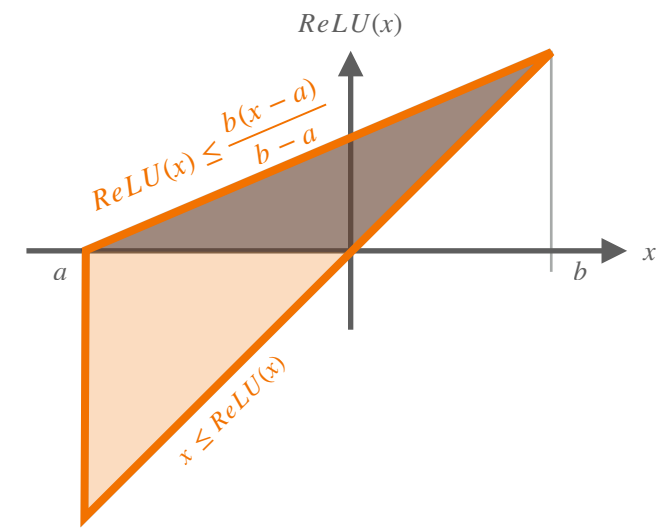
① proceed **forwards** from **an abstraction** of the input specification **I**



Example



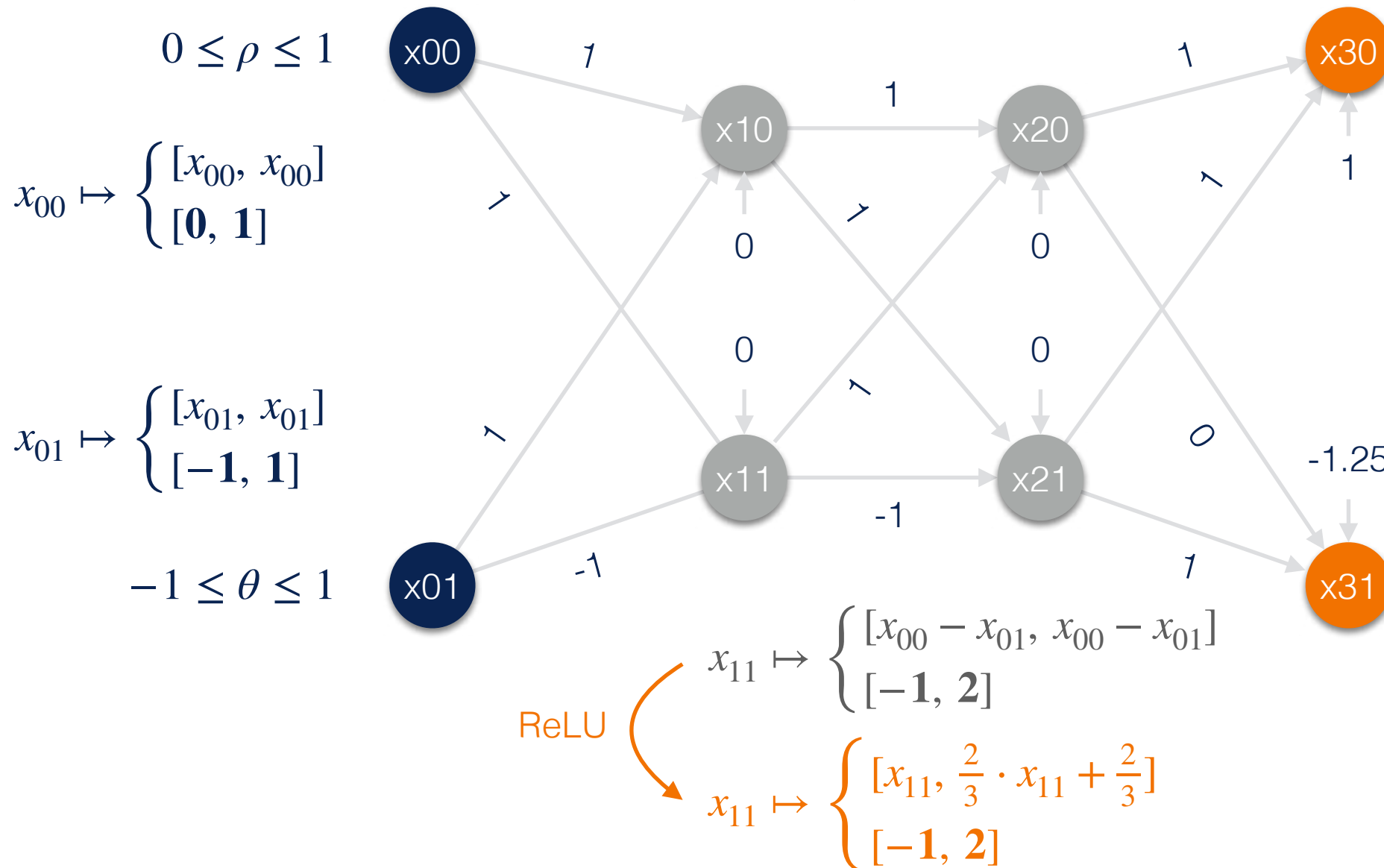
DeepPoly Domain [Singh19]



ReLU

$$x_{10} \mapsto \begin{cases} [x_{10}, \frac{2}{3} \cdot x_{10} + \frac{2}{3}] \\ [-1, 2] \end{cases}$$

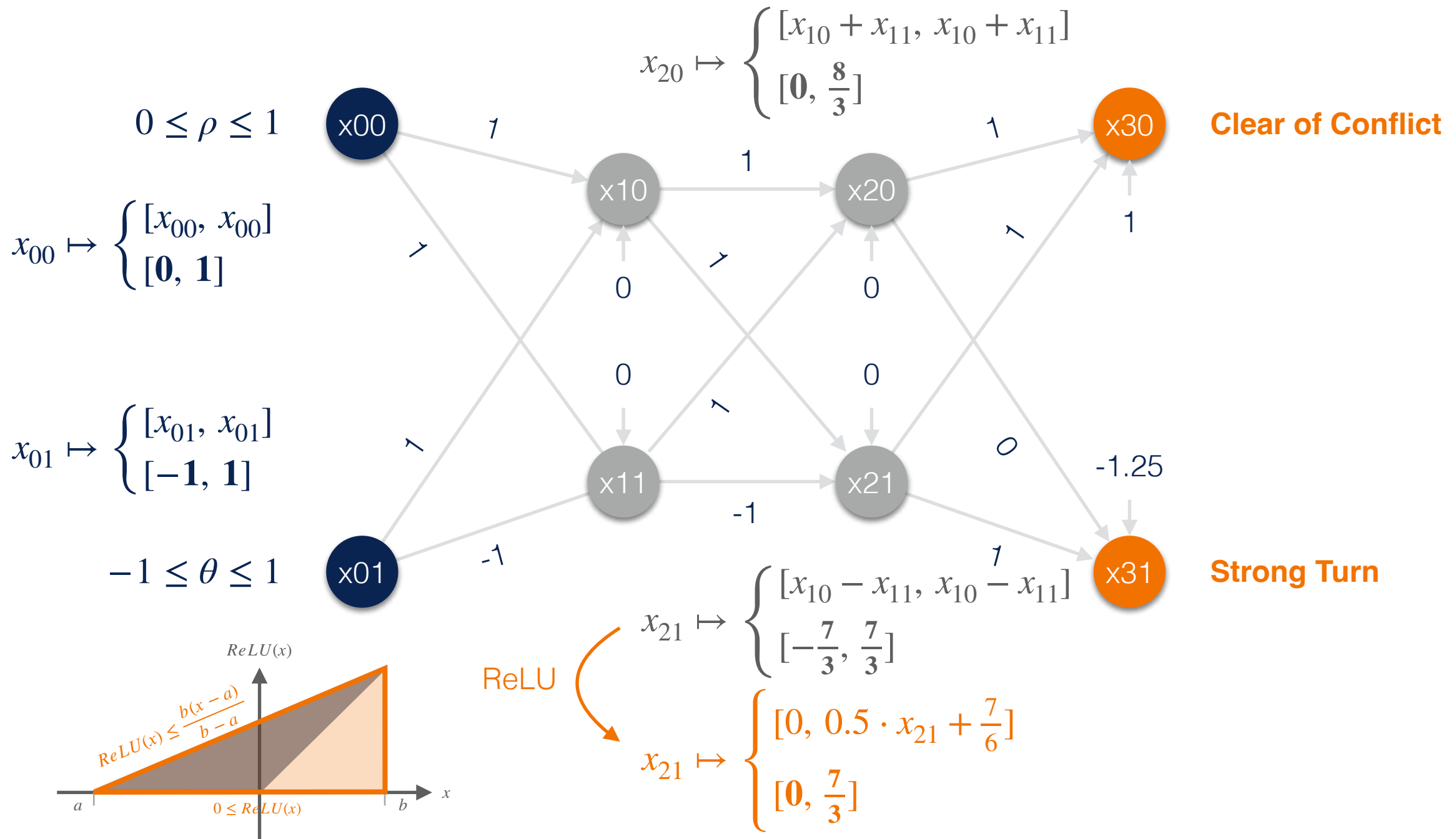
$$x_{10} \mapsto \begin{cases} [x_{00} + x_{01}, x_{00} + x_{01}] \\ [-1, 2] \end{cases}$$



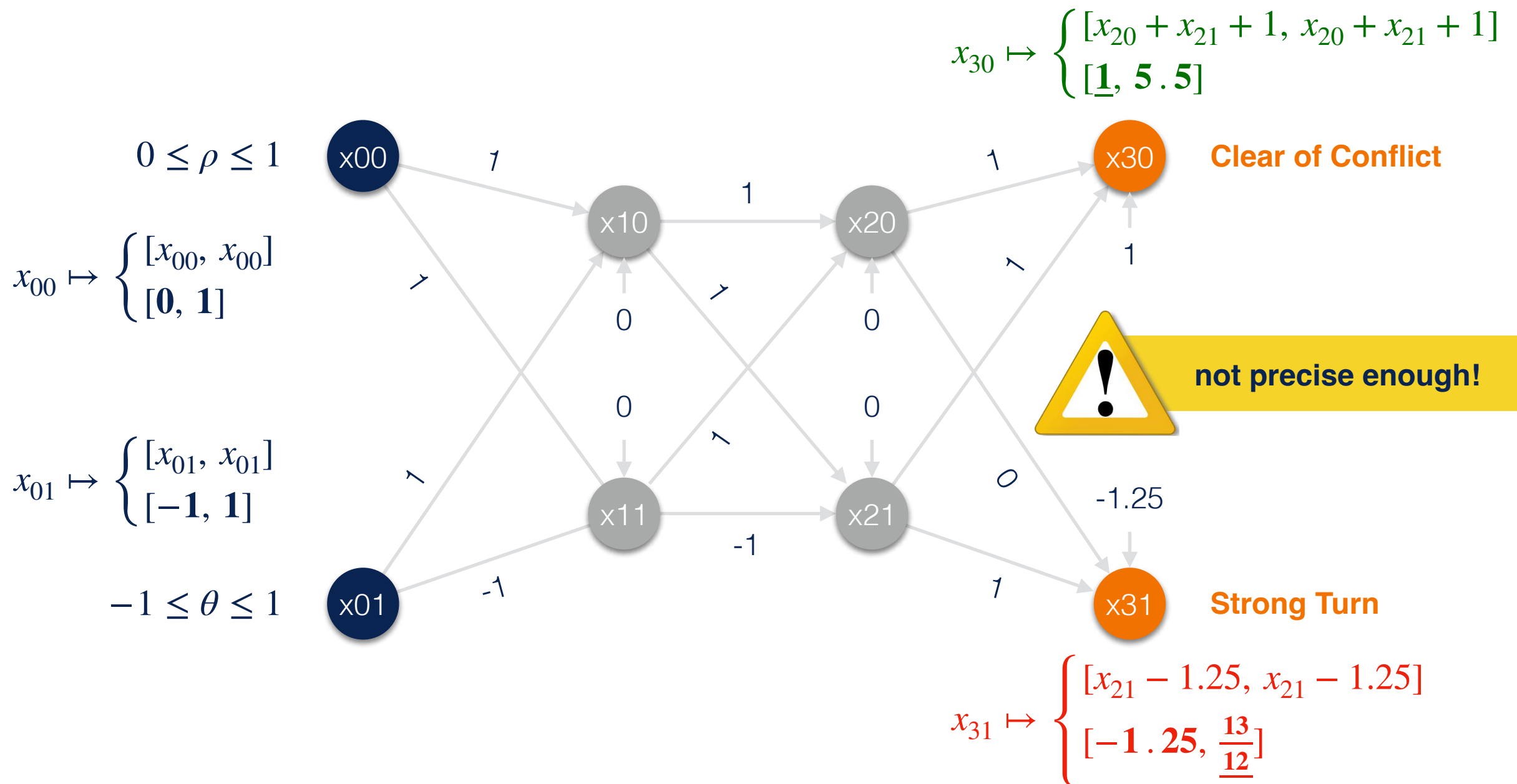
Clear of Conflict

Strong Turn

DeepPoly Domain [Singh19]

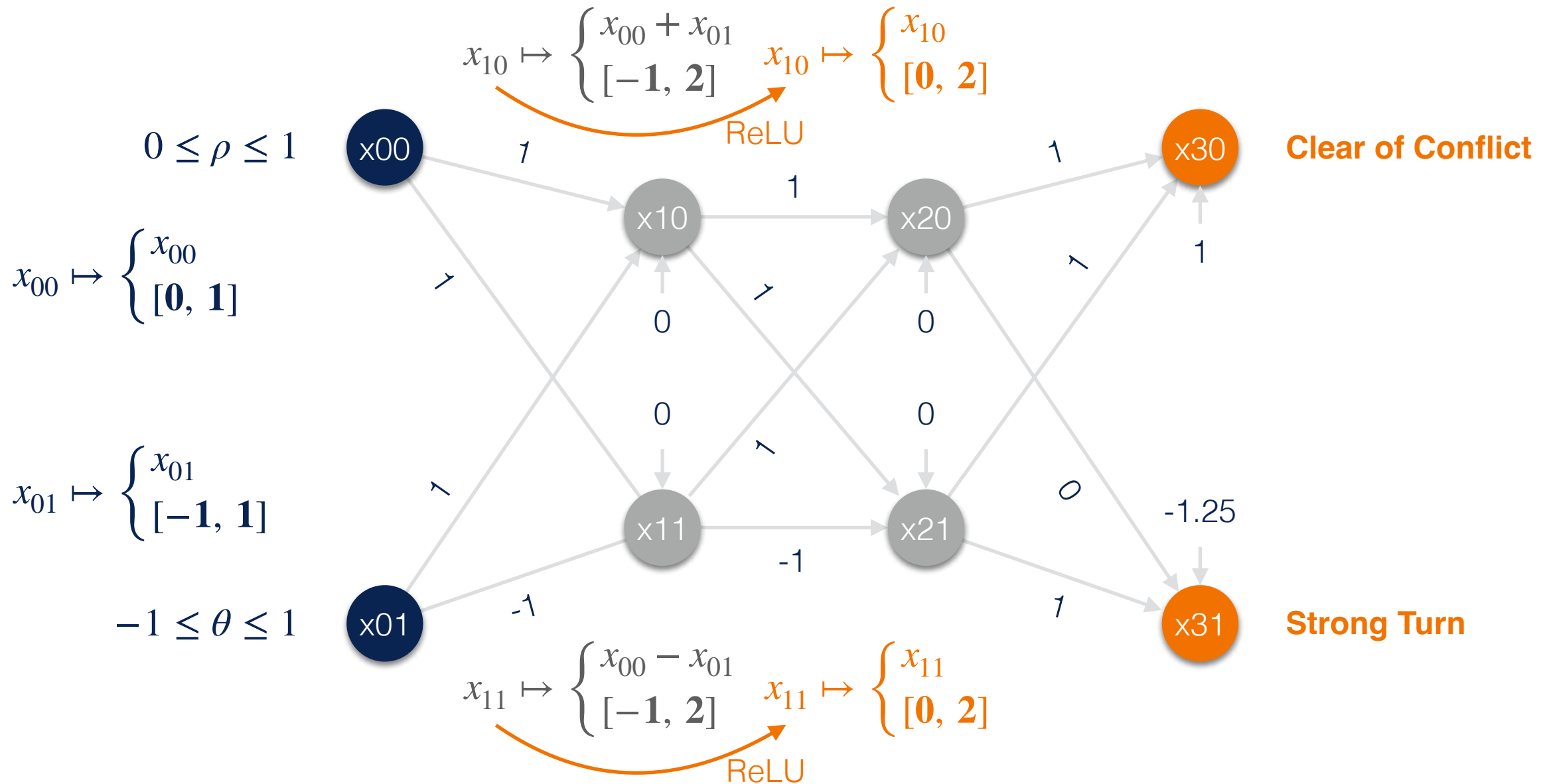


DeepPoly Domain [Singh19]



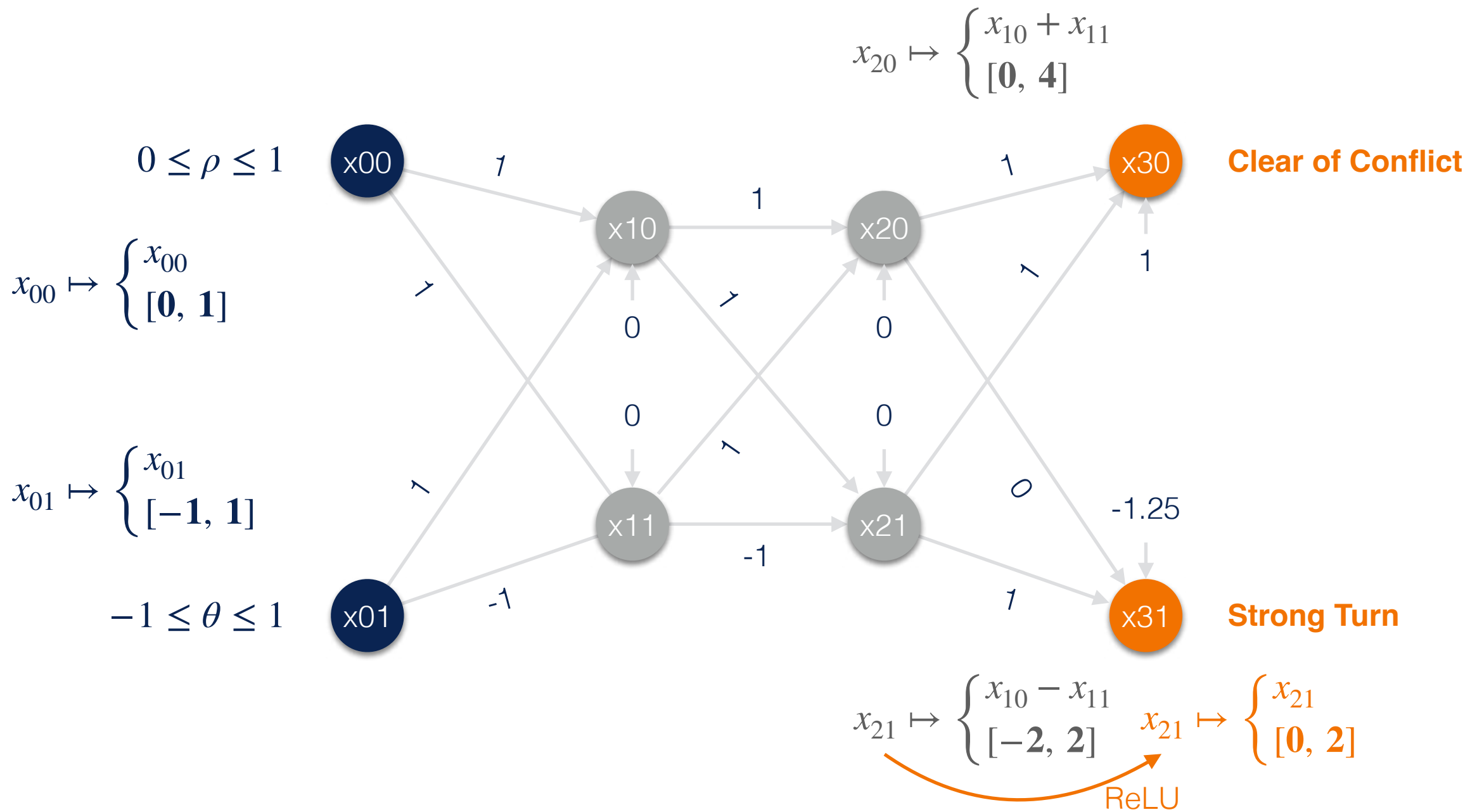
Interval Domain

with **Symbolic Constant Propagation** [Li19]



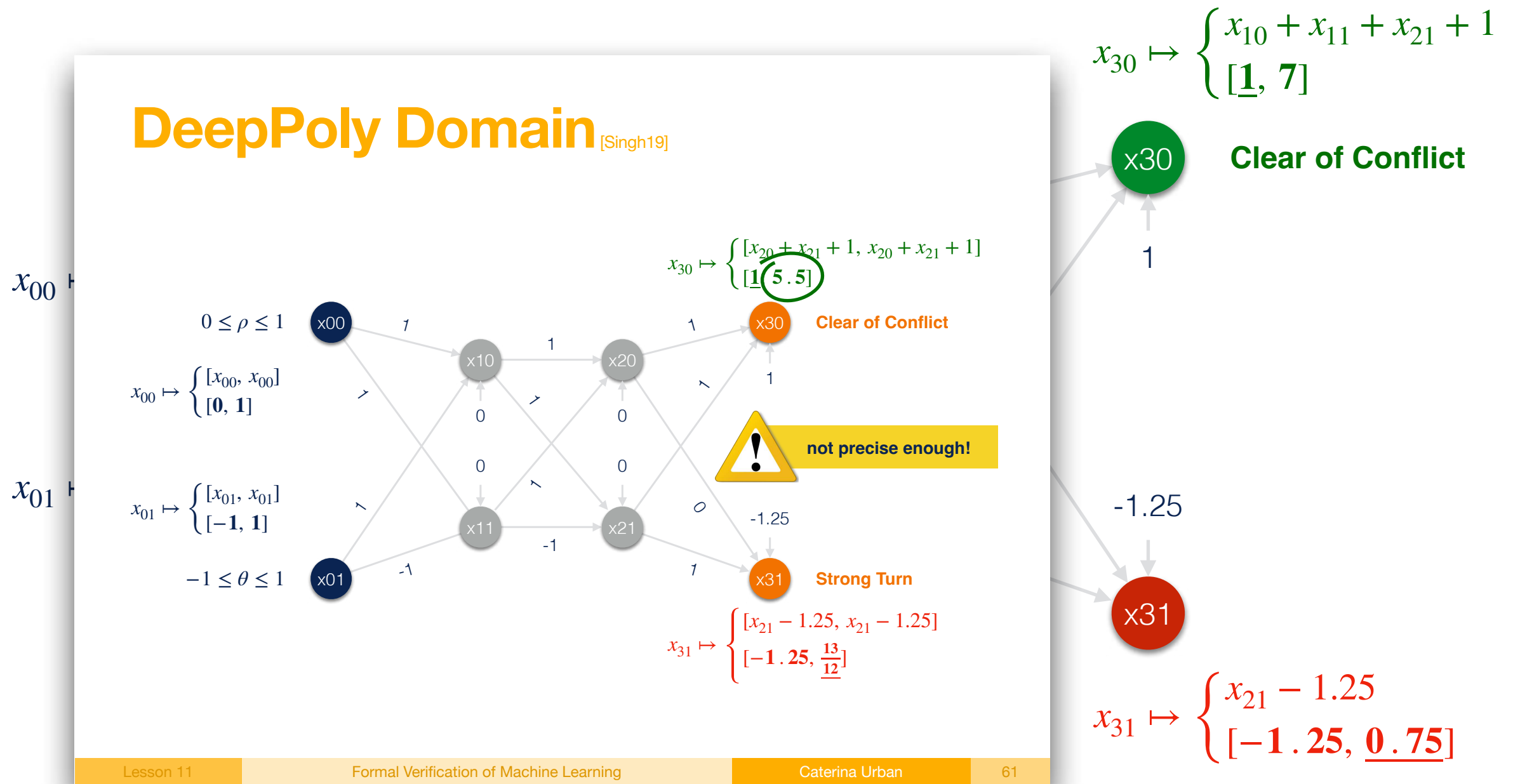
Interval Domain

with **Symbolic Constant Propagation** [Li19]

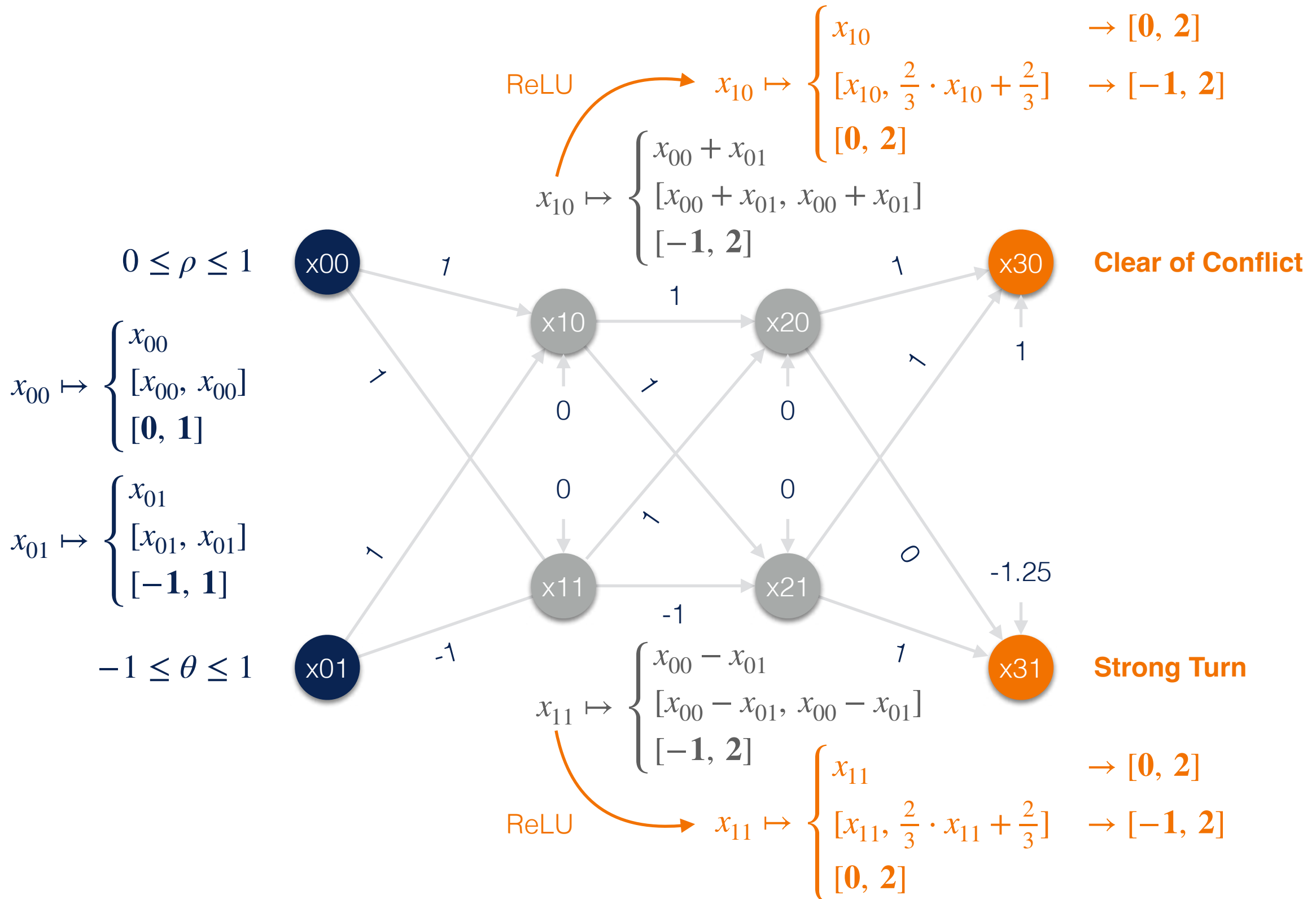


Interval Domain

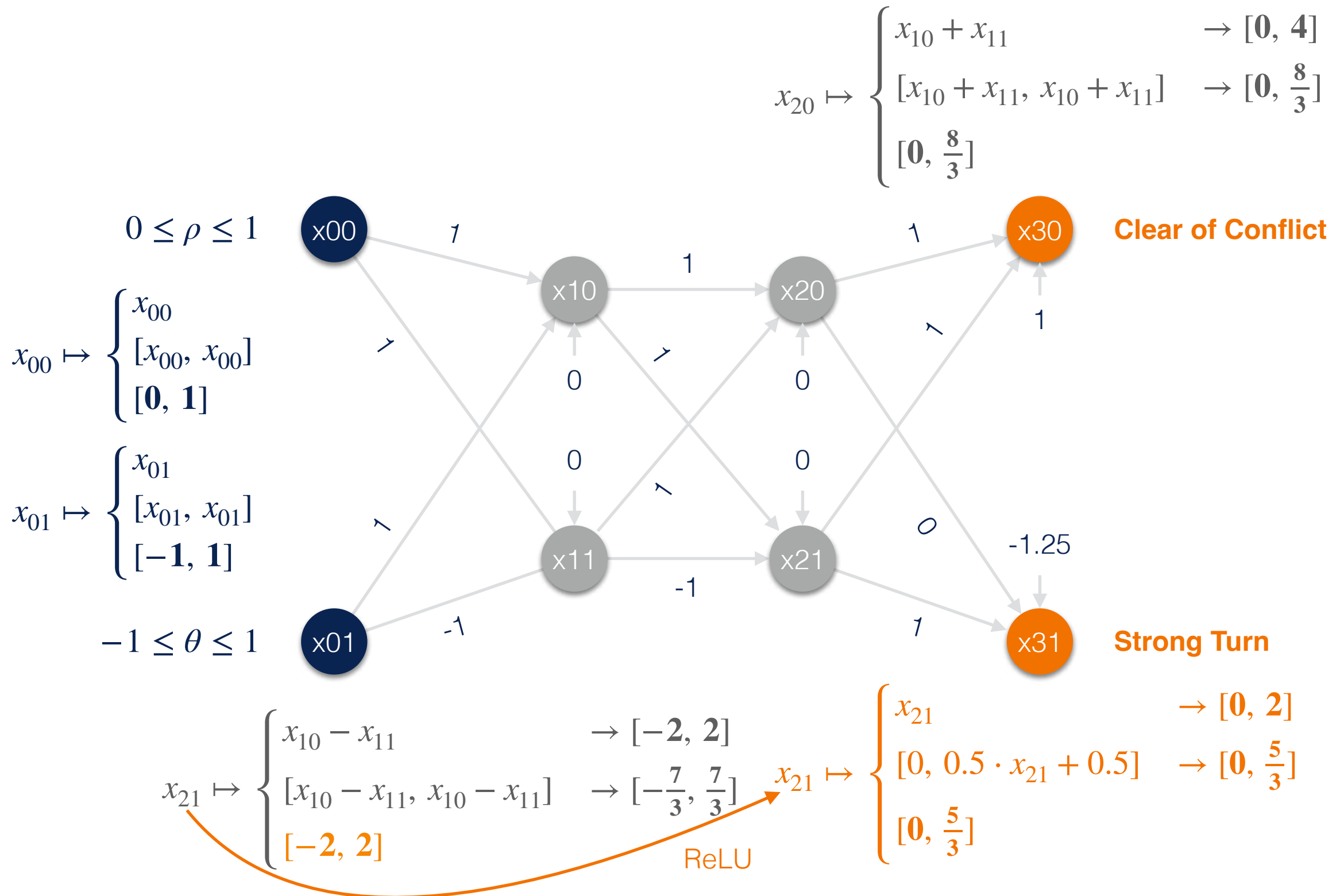
with Symbolic Constant Propagation [Li19]



Product Domain [Mazzucato21]

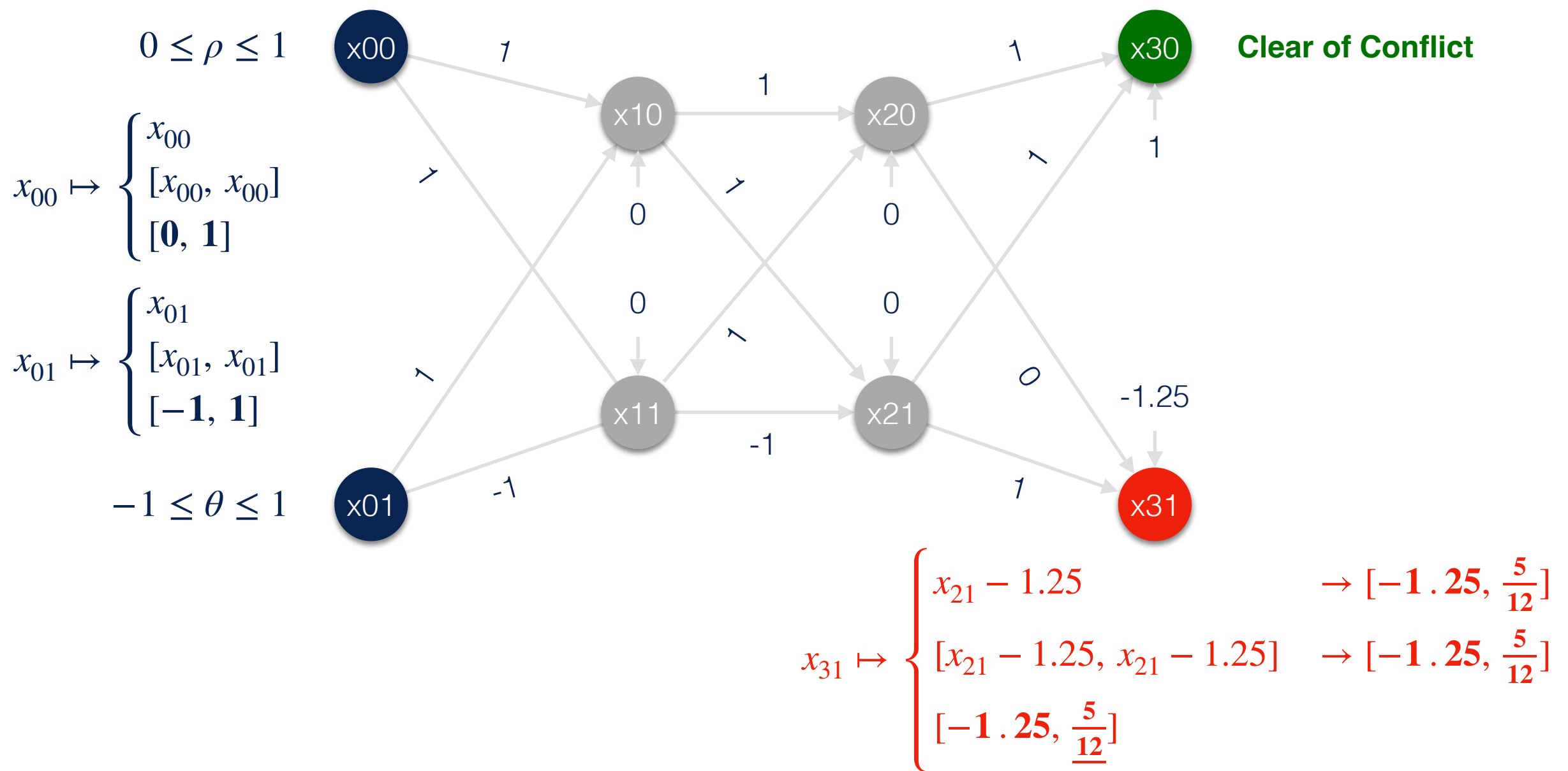


Product Domain [Mazzucato21]



Product Domain [Mazzucato21]

$$x_{30} \mapsto \begin{cases} x_{10} + x_{11} + x_{21} + 1 & \rightarrow [1, \frac{20}{3}] \\ [x_{20} + x_{21} + 1, x_{20} + x_{21} + 1] & \rightarrow [1, 4.5] \\ \underline{[1, 4.5]} \end{cases}$$



Formal Methods

Mathematical Guarantees of Safety



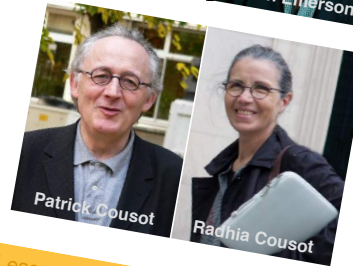
Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



Lesson 15

Formal Verification of Machine Learning

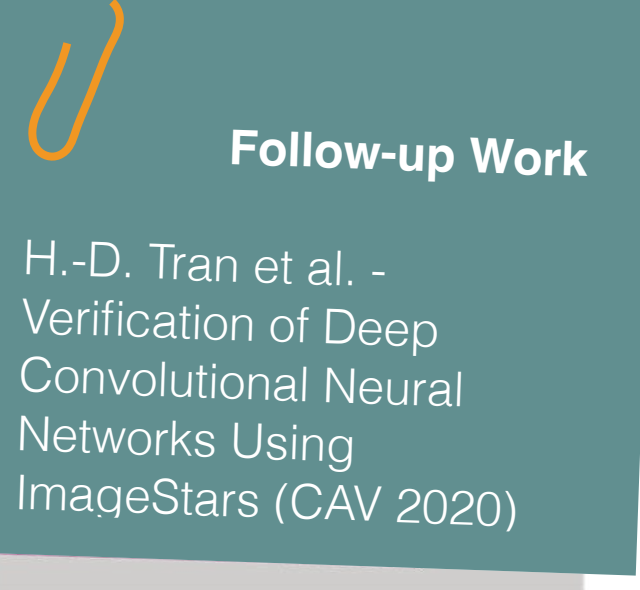
Caterina Urban

9

Other Complete Methods

Star Sets

Exact Static Analysis Method

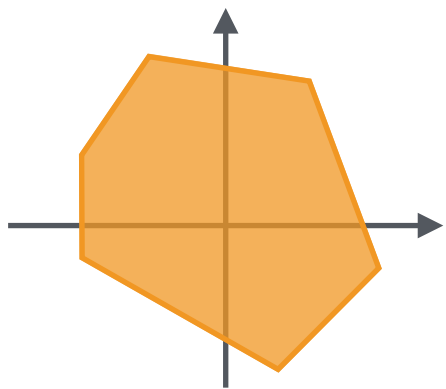


$$\Theta \stackrel{\text{def}}{=} \langle c, V, P \rangle$$

$c \in \mathcal{R}^n$: center

$V = \{v_1, \dots, v_m\}$: basis vectors in \mathcal{R}^n

$P: \mathcal{R}^m \rightarrow \{ \perp, \top \}$: predicate



$$\llbracket \Theta \rrbracket = \left\{ x \mid x = c + \sum_{i=1}^m \alpha_i v_i \text{ such that } P(\alpha_1, \dots, \alpha_m) = \top \right\}$$

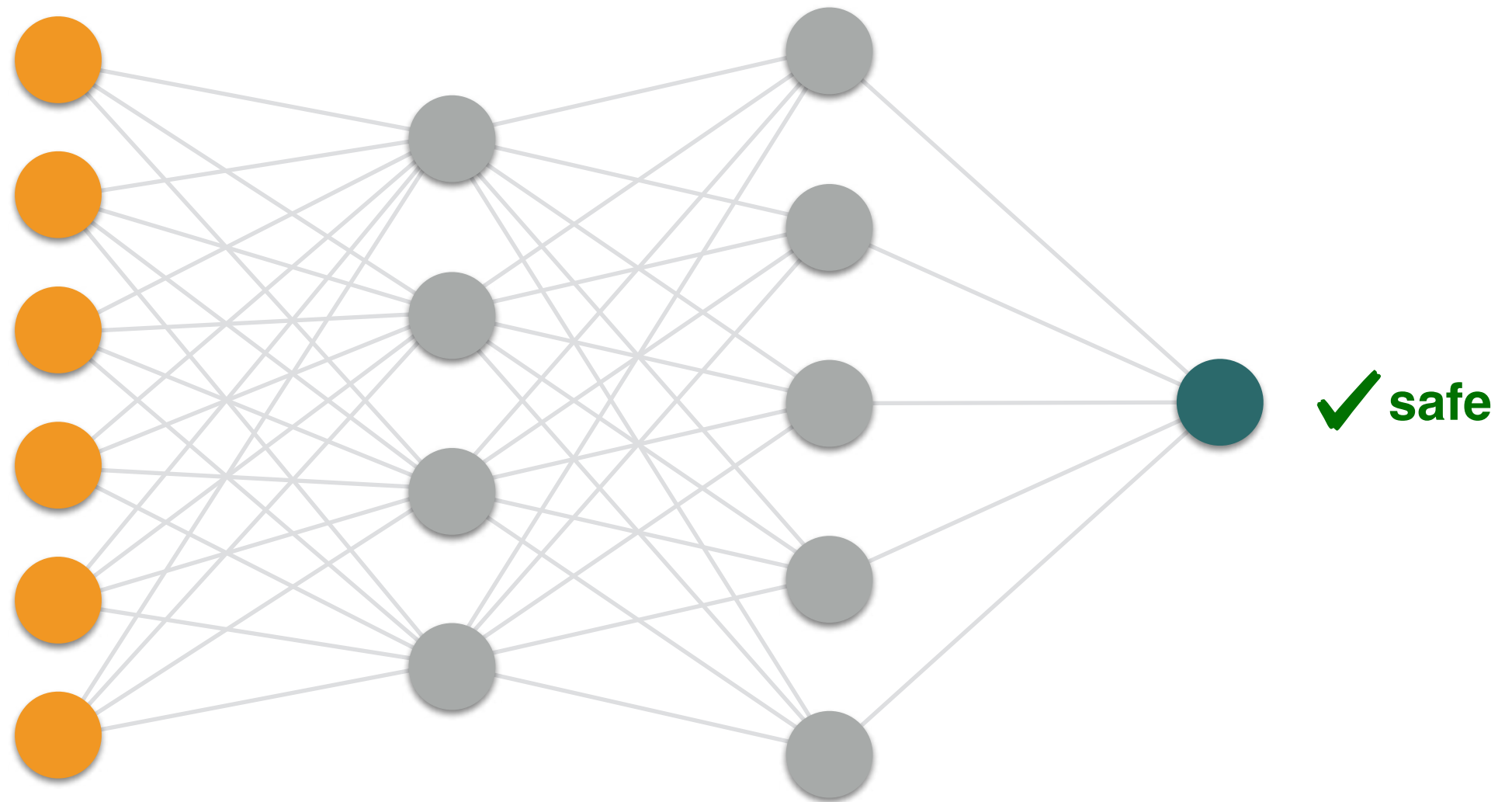
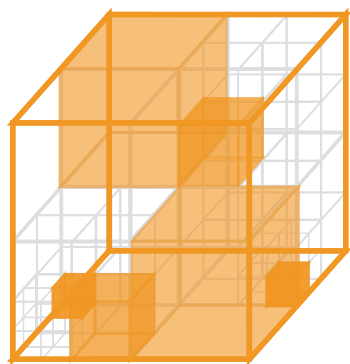
- fast and cheap **affine mapping operations** \rightarrow neural network layers
- inexpensive **intersections with half-spaces** \rightarrow ReLU activations

ReluVal



use symbolic propagation
+ **iterative** input **refinement**

Asymptotically Complete Method



S. Wang et al. - Formal Security Analysis of Neural Networks Using Symbolic Intervals (USENIX Security 2018)

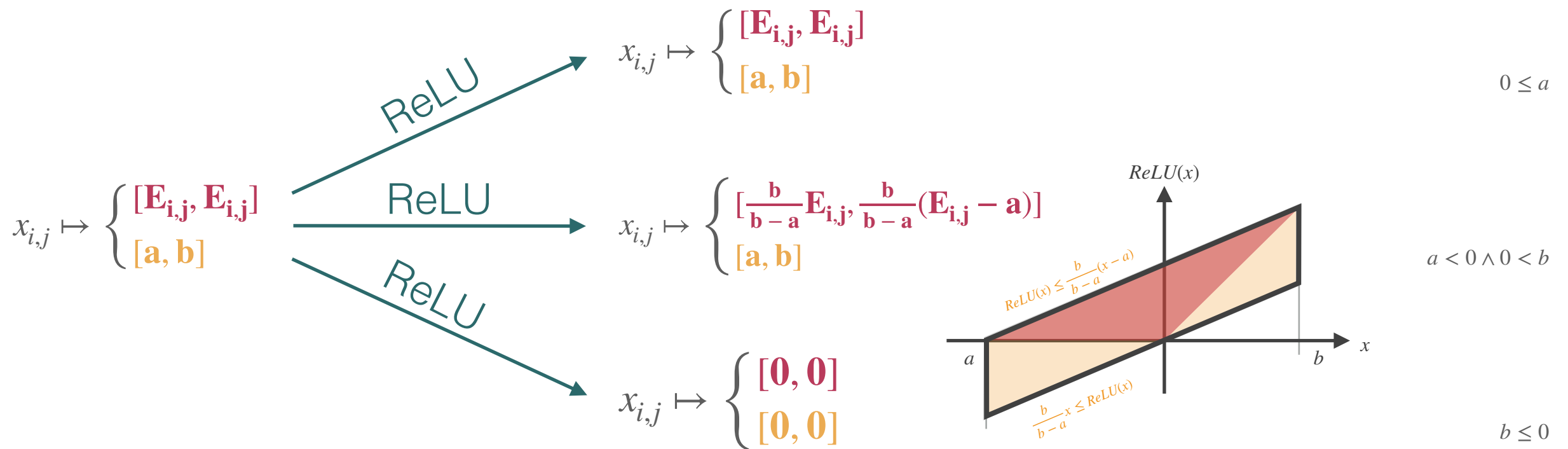
Neurify

Asymptotically Complete Method



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

$$x_{i,j} \mapsto \begin{cases} [\sum_k c_{0,k} \cdot x_{0,k} + c, \sum_k d_{0,k} \cdot x_{0,k} + d] & c_{0,k}, c, d_{0,k}, d \in \mathcal{R} \\ [a, b] & a, b \in \mathcal{R} \end{cases}$$



Further Complete Methods

- **W. Ruan, X. Huang, and M. Kwiatkowska.** *Reachability Analysis of Deep Neural Networks with Provable Guarantees.* In IJCAI, 2018.
a global optimization-based approach for verifying Lipschitz continuous neural networks
- **G. Singh, T. Gehr, M. Püschel, and M. Vechev.** *Boosting Robustness Certification of Neural Networks.* In ICLR, 2019.
an approach combining abstract interpretation and (mixed integer) linear programming

Formal Methods

Mathematical Guarantees of Safety



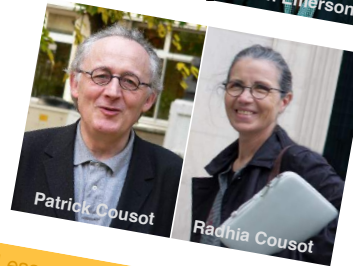
Deductive Verification

- extremely **expressive**
- **relies on the user** to guide the proof



Model Checking

- analysis of a **model** of the software
- **sound and complete** with respect to the model



Static Analysis

- analysis of the software at some level of **abstraction**
- fully **automatic** and **sound** by construction
- generally **not complete**



Lesson 15

Formal Verification of Machine Learning

Caterina Urban

9

Other Incomplete Methods

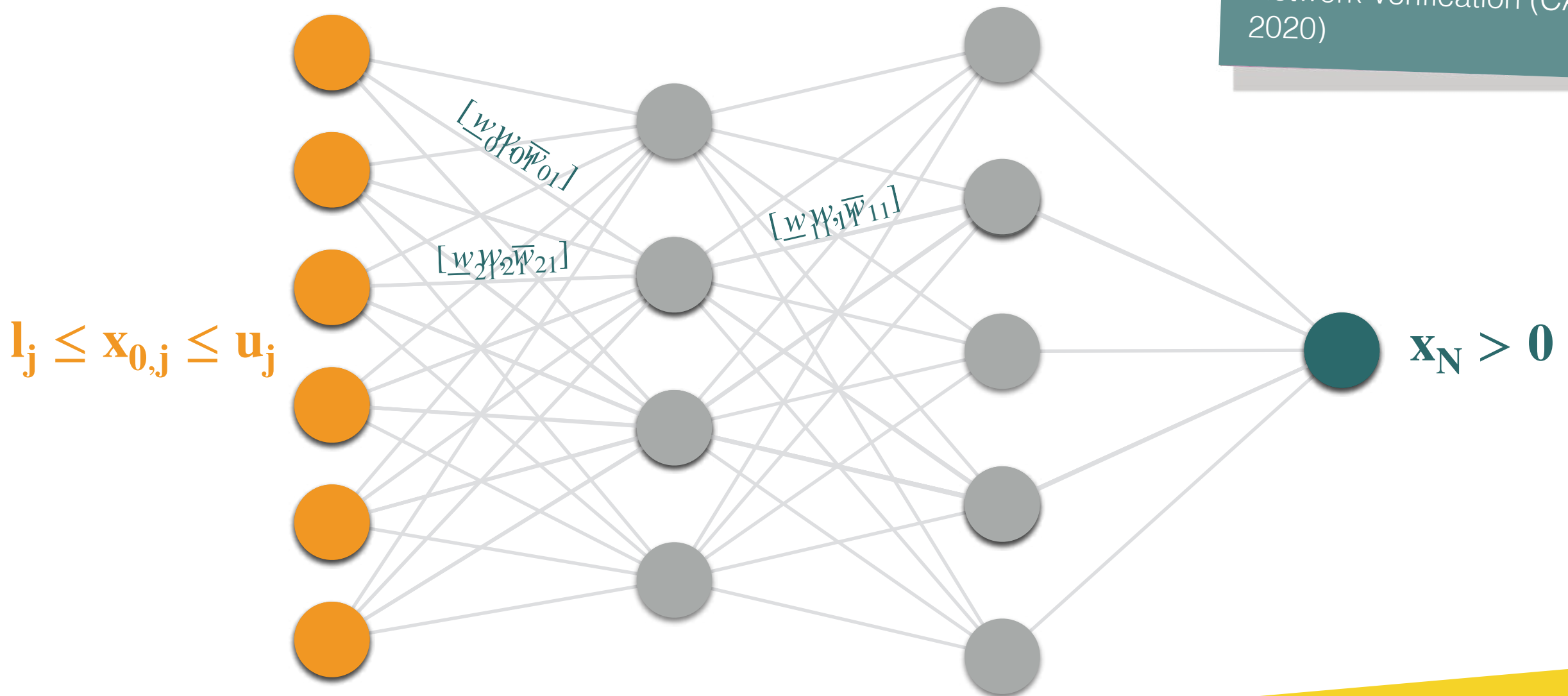
Interval Neural Networks

Abstraction-Based Method



Related Work

Y. Y. Elboher et al. - An Abstraction-Based Framework for Neural Network Verification (CAV 2020)



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

P. Prabhakar and Z. R. Afza - Abstraction based Output Range Analysis for Neural Networks (NeurIPS 2019)

Further Incomplete Methods

- **W. Xiang, H.-D. Tran, and T. T. Johnson.** *Output Reachable Set Estimation and Verification for Multi-Layer Neural Networks.* 2018.
an approach combining **simulation** and **linear programming**
- **K. Dvijotham, R. Stanforth, S. Gowal, T. Mann, and P. Kohli.** *A Dual Approach to Scalable Verification of Deep Networks.* In UAI, 2018.
an approach based on **duality** for verifying **neural networks**

Further Incomplete Methods

- **E. Wong and Z. Kolter.** *Provable Defenses Against Adversarial Examples via the Convex Outer Adversarial Polytope.* In ICML, 2018.
 - **A. Raghunathan, J. Steinhardt, and P. Liang.** *Certified Defenses against Adversarial Examples.* In ICML, 2018.
 - **T.-W. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, L. Daniel, D. Boning, and I. Dhillon.** *Towards Fast Computation of Certified Robustness for ReLU Networks.* In ICML, 2018.
 - **H. Zhang, T.-W. Weng, P.-Y. Chen, C.-J. Hsieh, and L. Daniel.** *Efficient Neural Network Robustness Certification with General Activation Functions.* In NeurIPS, 2018.
- approaches for finding a lower bound on robustness to adversarial perturbations**

Further Incomplete Methods

- **A. Boopathy, T.-W. Weng, P.-Y. Chen, S. Liu, and L. Daniel.** *CNN-Cert: An Efficient Framework for Certifying Robustness of Convolutional Neural Networks.* In AAI, 2019.
approach focusing on convolutional neural networks
- **C.-Y. Ko, Z. Lyu, T.-W. Weng, L. Daniel, N. Wong, and D. Lin.** *POPQORN: Quantifying Robustness of Recurrent Neural Networks.* In ICML, 2019.
H. Zhang, M. Shinn, A. Gupta, A. Gurfinkel, N. Le, and N. Narodytska. *Verification of Recurrent Neural Networks for Cognitive Tasks via Reachability Analysis.* In ECAI, 2020.
approaches focusing on recurrent neural networks
- **D. Gopinath, H. Converse, C. S. Pasareanu, and A. Taly.** *Property Inference for Deep Neural Networks.* In ASE, 2019.
an approach for inferring safety properties of neural networks

Complete Methods

Advantages

sound and complete

Disadvantages

soundness not typically guaranteed
with respect to **floating-point arithmetic**

do not scale to large models

often **limited** to certain
model **architectures**

suffer from **false positives**

Disadvantages

able to scale to large models

sound often also with respect to
floating-point arithmetic

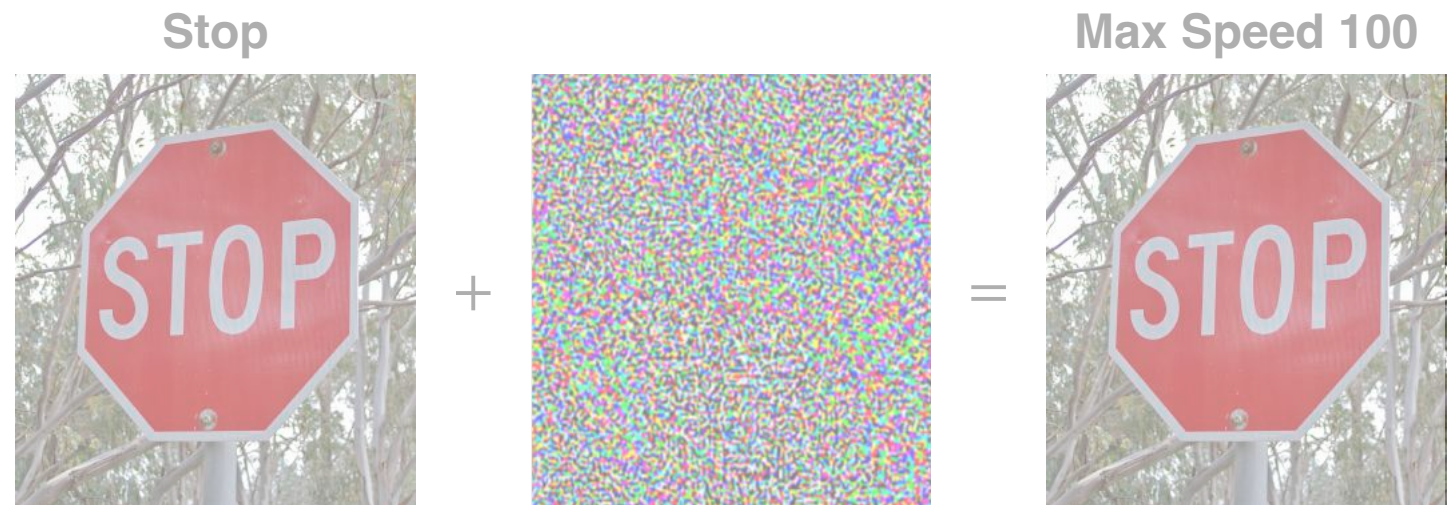
less limited to certain
model architectures

Advantages

Incomplete Methods

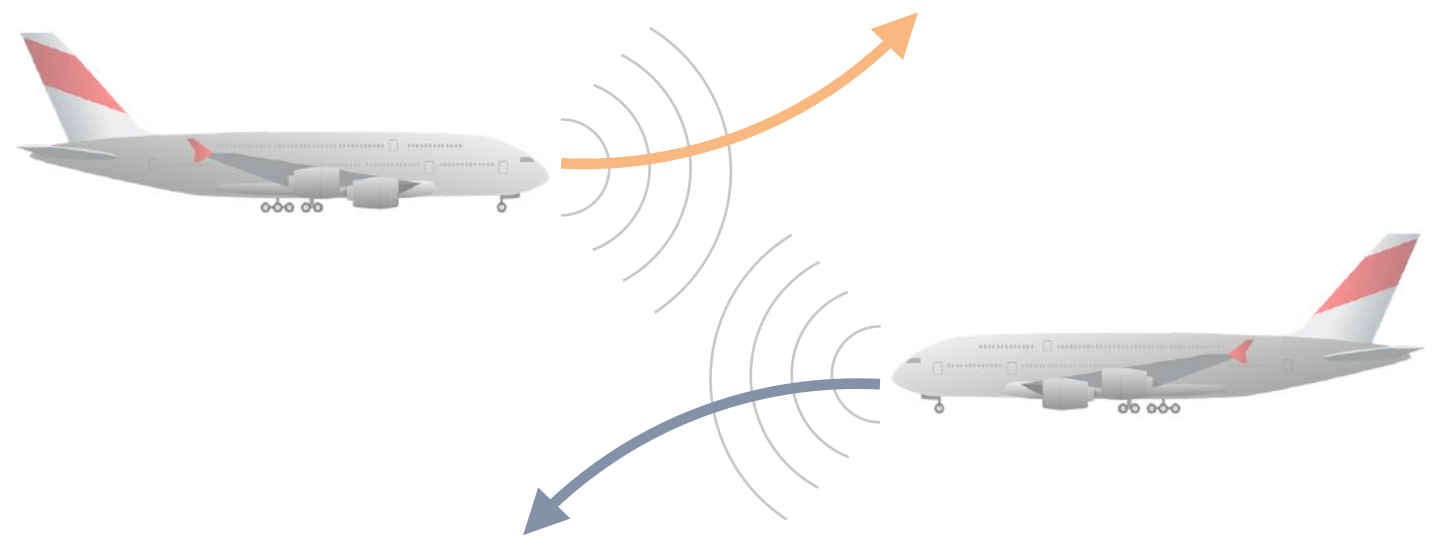
Stability

Goal G3 in [Kurd03]

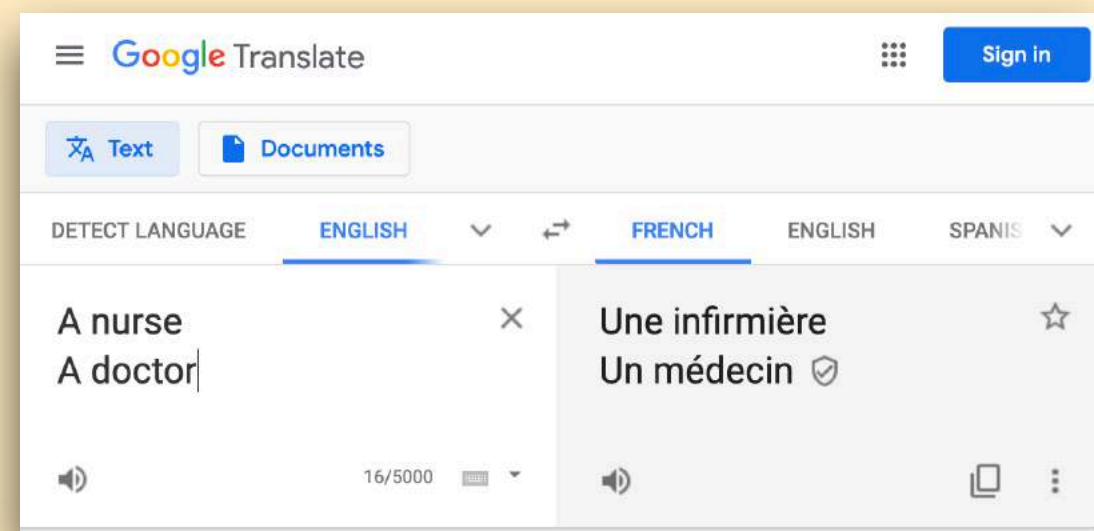


Safety

Goal G4 in [Kurd03]



Fairness



ML Impacts Our Society



WIRED

In 2019, predictive algorithms will start to make better decisions

will start to

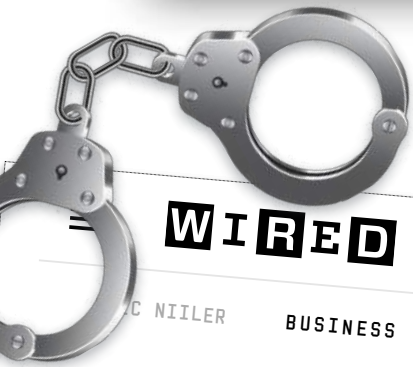
Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

**D CHECKS ARE
FOR A HOME**

By Colin Lecher | @colinlecher | Feb 1, 2019, 8:00am EST



WIRED

BUSINESS

MORE ▾ SIGN IN

03.25.2019 07:00 AM

Can AI Be a Fair Judge in Court? Estonia Thinks So

Estonia plans to use an artificial intelligence program to handle small-claims cases, part of a push to make government services smarter.

BUSINESS NEWS OCTOBER 10, 2018 / 5:12 AM / A YEAR AGO

Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin



Translation tutorial: 21 fairness definitions and their politics

Arvind Narayanan
@random_walker



Tutorial: 21 fairness definitions and their politics

19,759 views • Mar 1, 2018

196 6 SHARE SAVE



Arvind Narayanan
226 subscribers

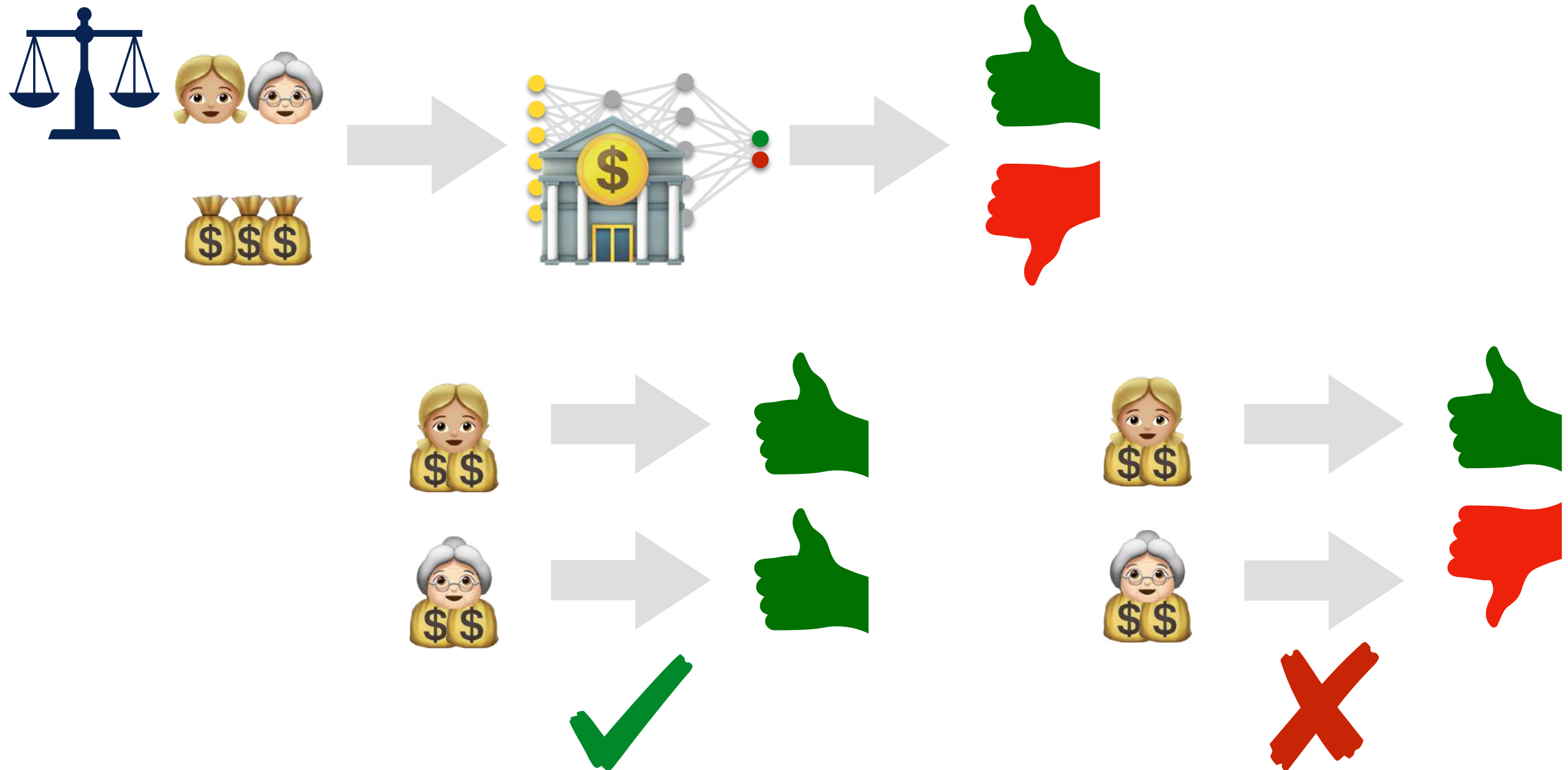
SUBSCRIBE

Computer scientists and statisticians have devised numerous mathematical criteria to define what it means for a classifier or a model to be fair. The proliferation of these definitions represents an attempt to make technical sense of

SHOW MORE

Dependency Fairness [Galhotra17]

The classification is **independent of** the values of the **sensitive inputs**



Dependency Fairness

$$\mathcal{F}_i \stackrel{\text{def}}{=} \{ \llbracket M \rrbracket \in \mathcal{P}(\Sigma^*) \mid \text{UNUSED}_i(\llbracket M \rrbracket) \}$$

\mathcal{F}_i is the set of all neural networks M (or, rather, their semantics $\llbracket M \rrbracket$) that **do not use** the value of the sensitive input node $x_{0,i}$ for classification

$$\begin{aligned} \text{UNUSED}_i(\llbracket M \rrbracket) \stackrel{\text{def}}{=} & \forall t \in \llbracket M \rrbracket, v \in \mathcal{R} : t_0(x_{0,i}) \neq v \Rightarrow \exists t' \in \llbracket M \rrbracket : \\ & (\forall 0 \leq j \leq |L_0| : j \neq i \Rightarrow t_0(x_{0,j}) = t'_0(x_{0,j})) \\ & \wedge t'_0(x_{0,i}) = v \\ & \wedge \max_j t_\omega(x_{N,j}) = \max_j t'_\omega(x_{N,j}) \end{aligned}$$

Intuitively: **any possible classification outcome** is possible from any value of the sensitive input node $x_{0,i}$

Input Data (Non-)Usage

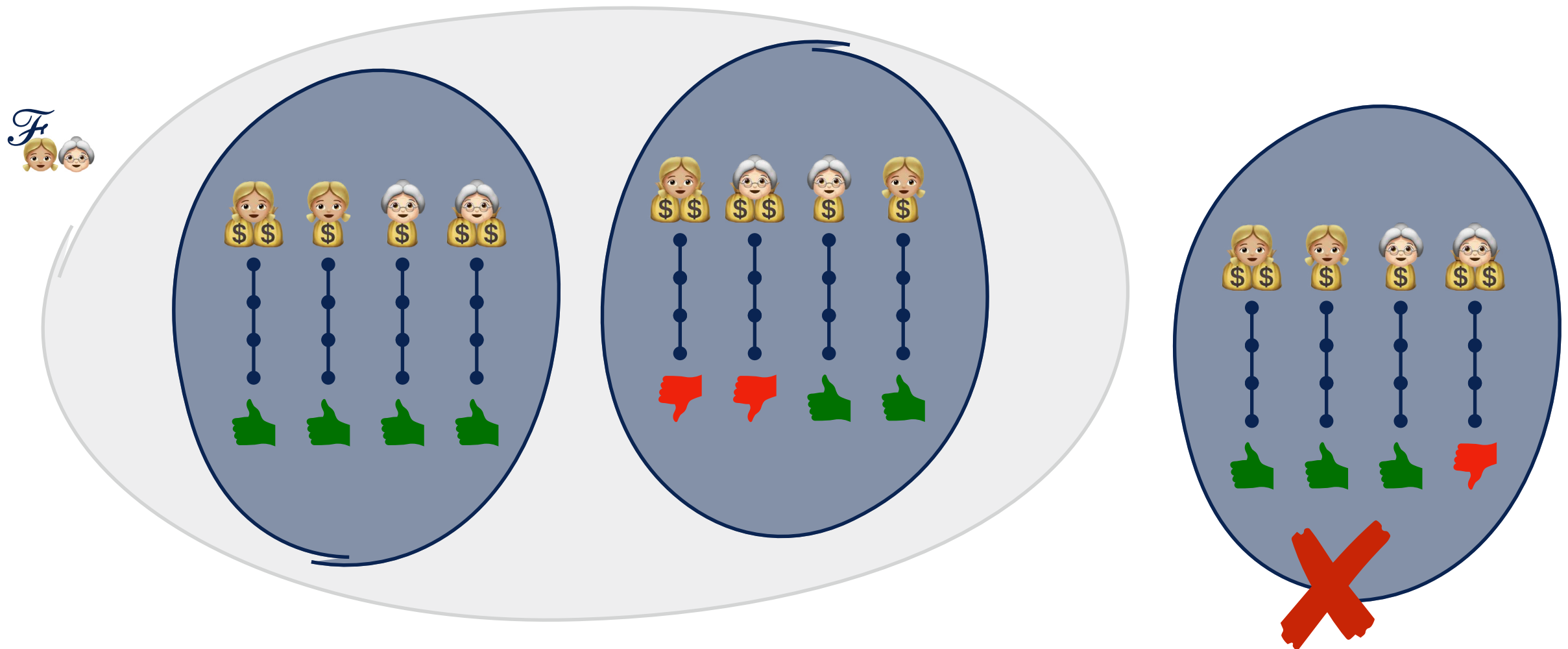
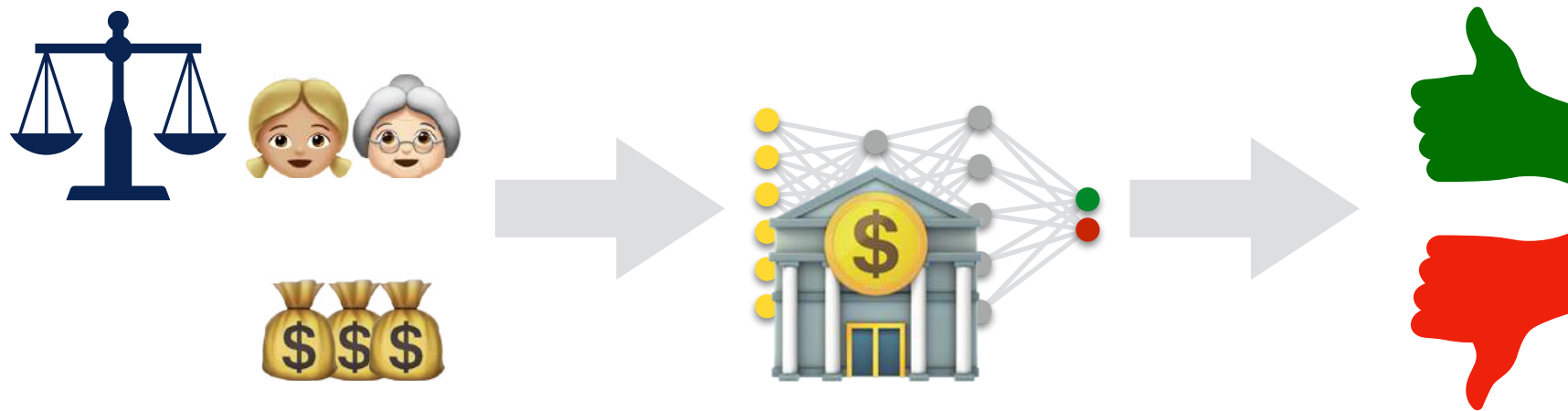
$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\mathcal{N}_J is the set of all programs P (or, rather, their semantics $\llbracket P \rrbracket$) that **do not use** the value of the input variables in $J \subseteq I_P$

$$\begin{aligned} \text{UNUSED}_J(\llbracket P \rrbracket) \stackrel{\text{def}}{=} & \forall t \in \llbracket P \rrbracket, v \in \mathcal{V} : t_0(J) \neq v \Rightarrow \exists t' \in \llbracket P \rrbracket : \\ & (\forall 0 \leq i \leq |I_P| : i \notin J \Rightarrow t_0(i) = t'_0(i)) \\ & \wedge t'_0(J) = v \\ & \wedge t_\omega = t'_\omega \end{aligned}$$

Intuitively: **any possible program outcome** is possible from any value of the input variable $x_{0,i}$

Dependency Fairness



Dependency Fairness

$$\mathcal{F}_i \stackrel{\text{def}}{=} \{ \llbracket M \rrbracket \in \mathcal{P}(\Sigma^*) \mid \text{UNUSED}_i(\llbracket M \rrbracket) \}$$

\mathcal{F}_i is the set of all neural networks M (or, rather, their semantics $\llbracket M \rrbracket$) that **do not use** the value of the sensitive input node $x_{0,i}$ for classification

$$\begin{aligned} \text{UNUSED}_i(\llbracket M \rrbracket) \stackrel{\text{def}}{=} & \forall t \in \llbracket M \rrbracket, v \in \mathcal{R} : t_0(x_{0,i}) \neq v \Rightarrow \exists t' \in \llbracket M \rrbracket : \\ & (\forall 0 \leq j \leq |L_0| : j \neq i \Rightarrow t_0(x_{0,j}) = t'_0(x_{0,j})) \\ & \wedge t'_0(x_{0,i}) = v \\ & \wedge \max_j t_\omega(x_{N,j}) = \max_j t'_\omega(x_{N,j}) \end{aligned}$$

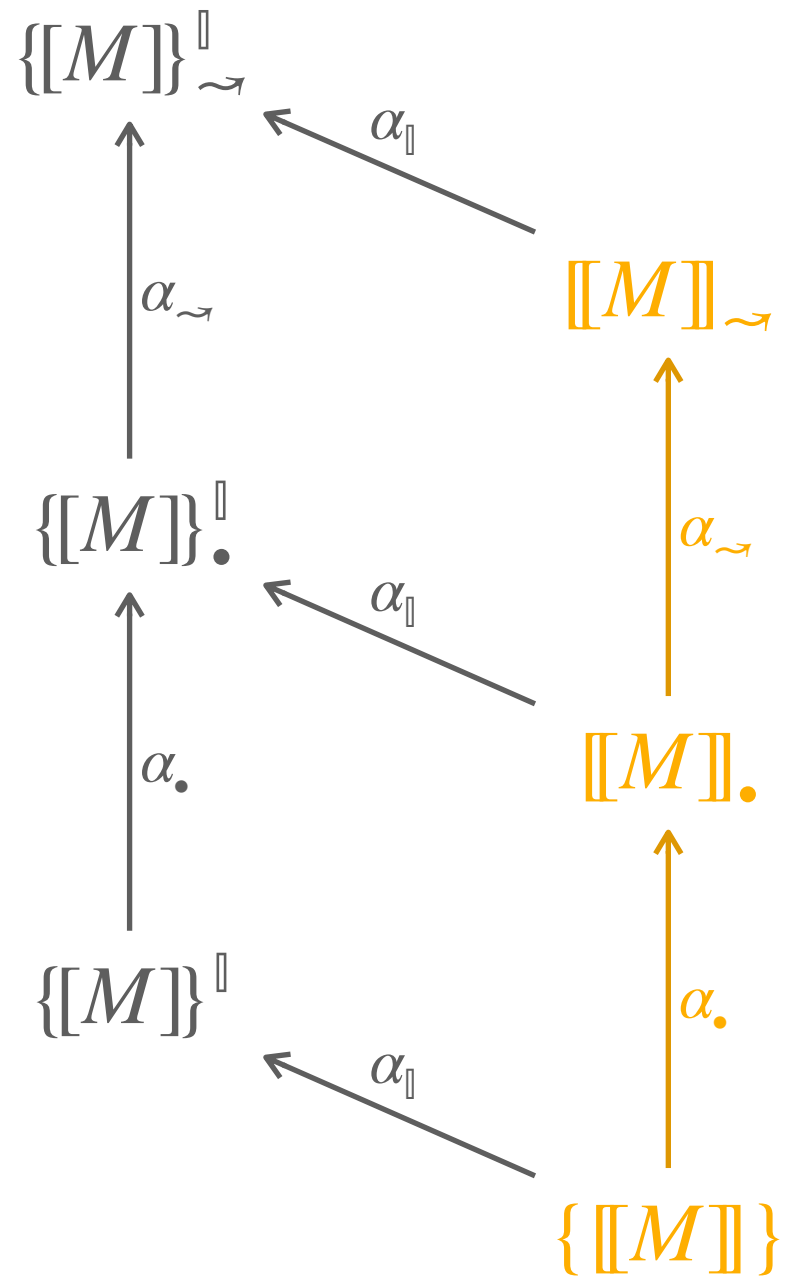
Intuitively: **any possible classification outcome** is possible from **any value** of the sensitive input node $x_{0,i}$

Theorem

$$M \models \mathcal{F}_i \Leftrightarrow \{ \llbracket M \rrbracket \} \subseteq \mathcal{F}_i$$

Hierarchy of Semantics

parallel semantics

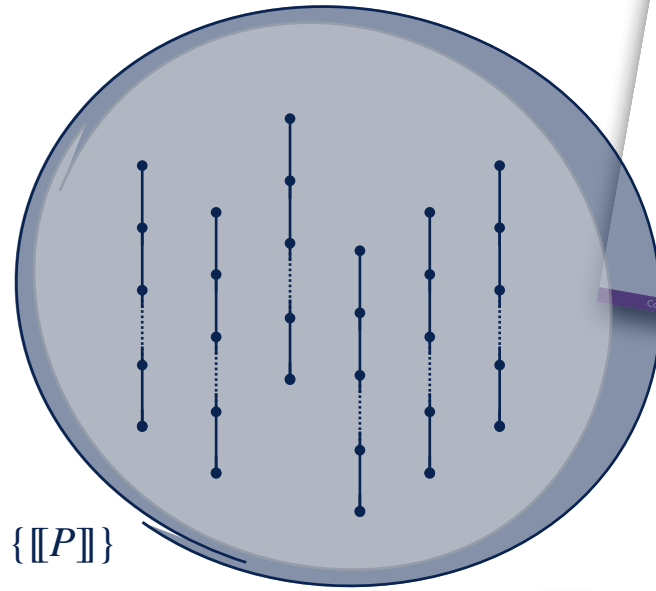


dependency semantics

outcome semantics

collecting semantics

Collecting Semantics



General collecting semantics

The collecting semantics $Col : Prog \rightarrow \mathcal{P}(\mathcal{P}(\Sigma^*))$ is the **strongest property** of a program

Hence: $Col(prog) \stackrel{def}{=} \{\{prog\}\}$

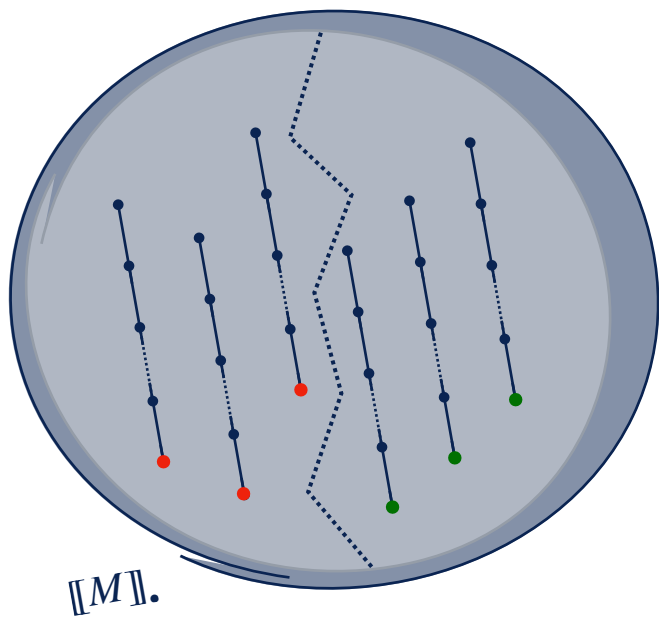
Benefits: uniformity of semantics and properties, \subseteq information order

- given a program $prog$ and a property $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$ the **verification problem** is an inclusion check:

$$Col(prog) \subseteq P$$
- generally, the collecting semantics **cannot be computed**, we settle for a weaker property S^\sharp that
 - is sound: $Col(prog) \subseteq S^\sharp$
 - implies the desired property: $S^\sharp \subseteq P$

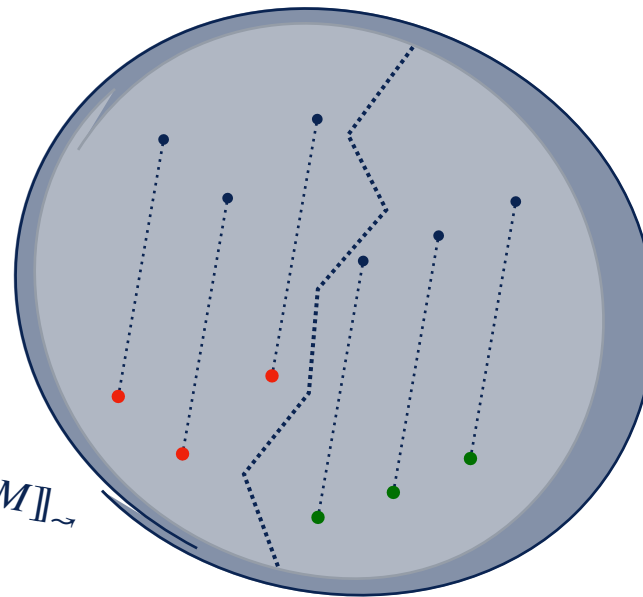
Course 2 Program Semantics and Properties Antoine Mott p. 24 / 98

Outcome Semantics



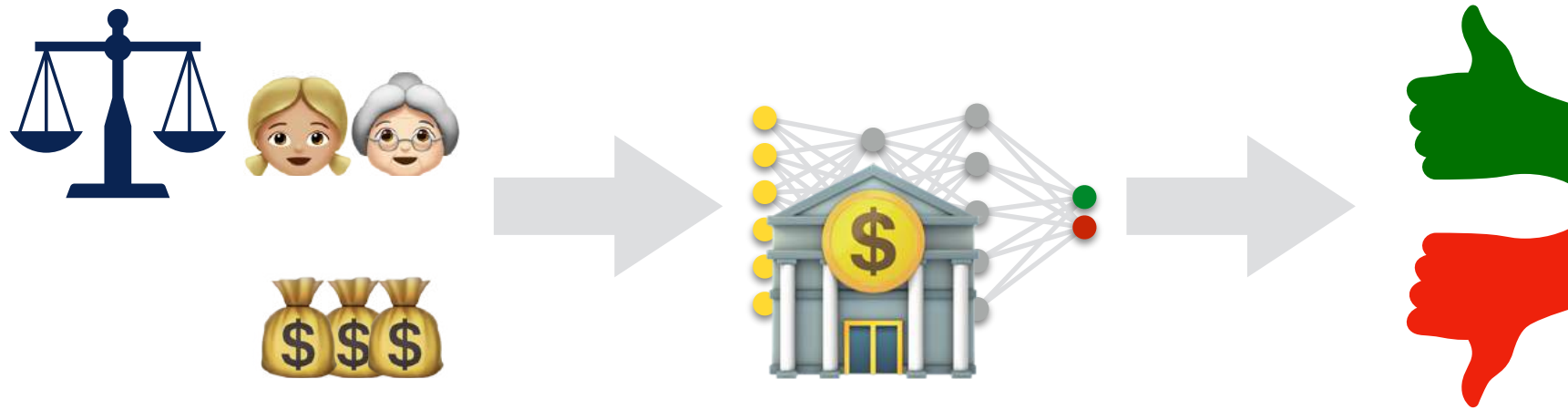
partitioning a set of traces that satisfies input data (non-)usage **with respect to the program outcome** yields sets of traces that also satisfy input data (non-)usage

Dependency Semantics

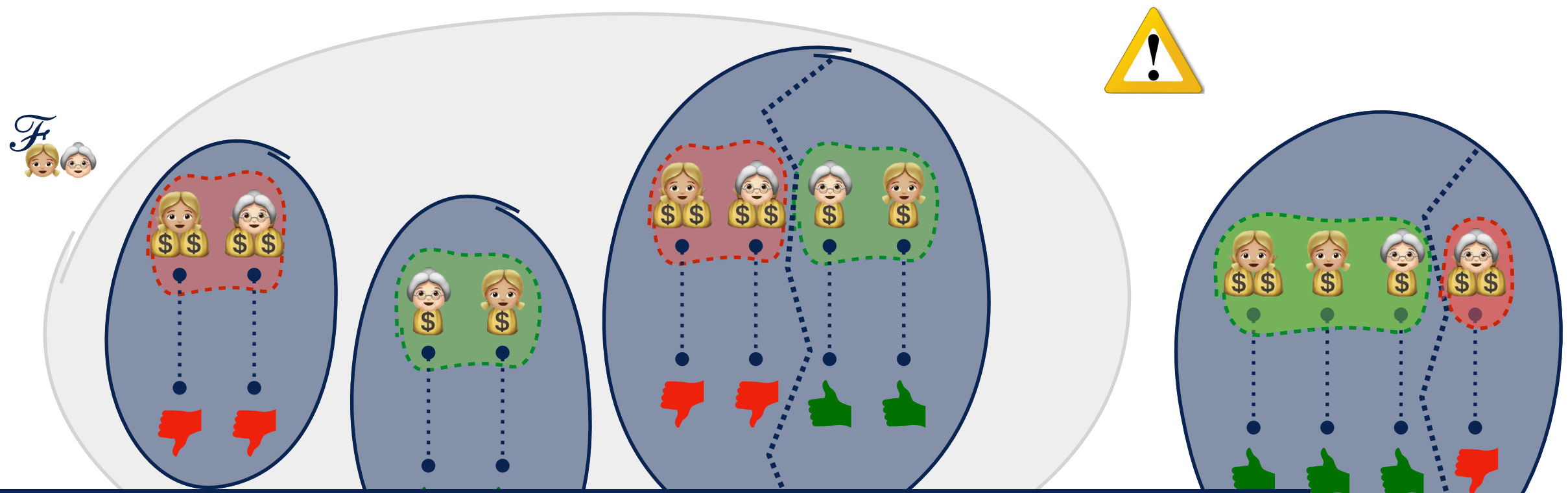


to reason about input data (non-)usage **we do not need to consider all intermediate computations** between the initial and final states of a trace (if any)

Dependency Semantics



partitioning with respect to the outcome classification induces a partition of the space of values of the input nodes *used for classification*



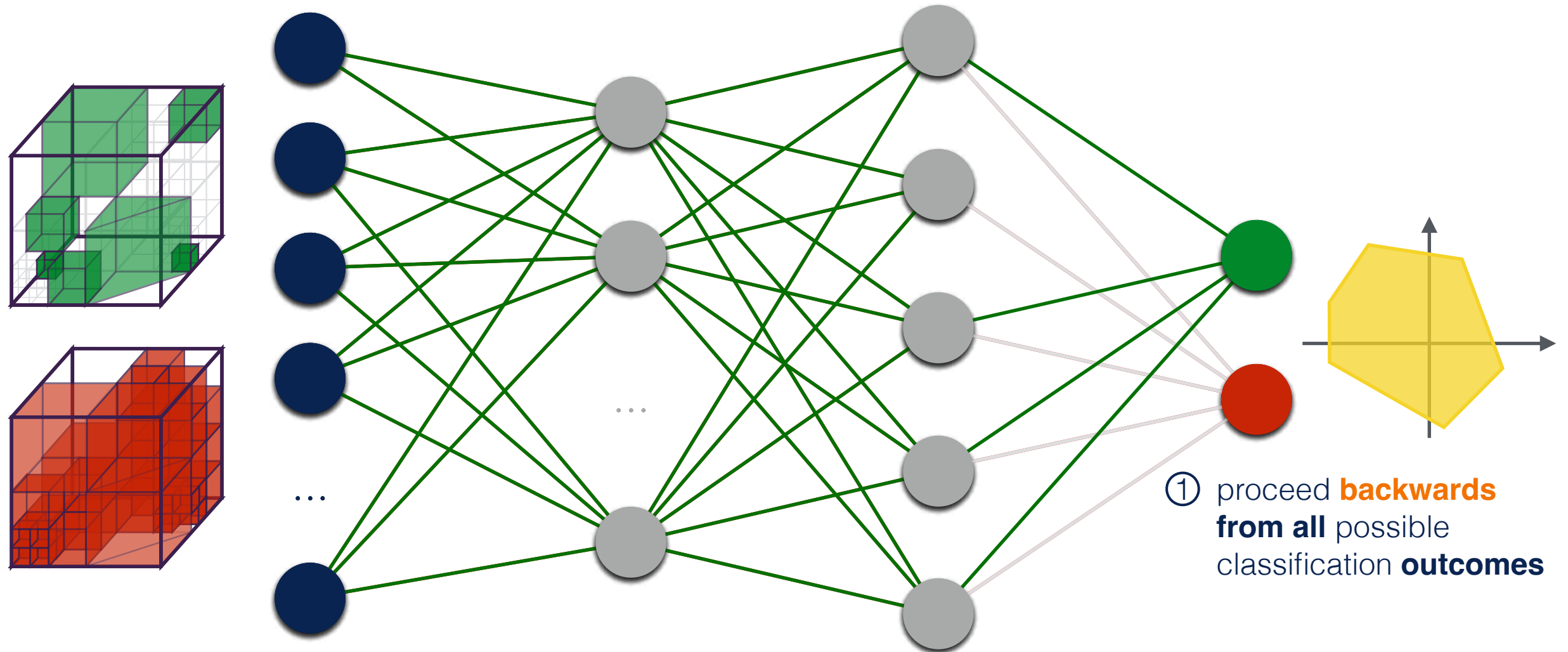
Lemma

$$M \models \mathcal{F}_i \Leftrightarrow \forall A, B \in \llbracket M \rrbracket_{\sim} : (A_{\omega} \neq B_{\omega} \Rightarrow A_0|_{\neq i} \cap B_0|_{\neq i} = \emptyset)$$

Naïve Abstraction

Naïve Backward Analysis

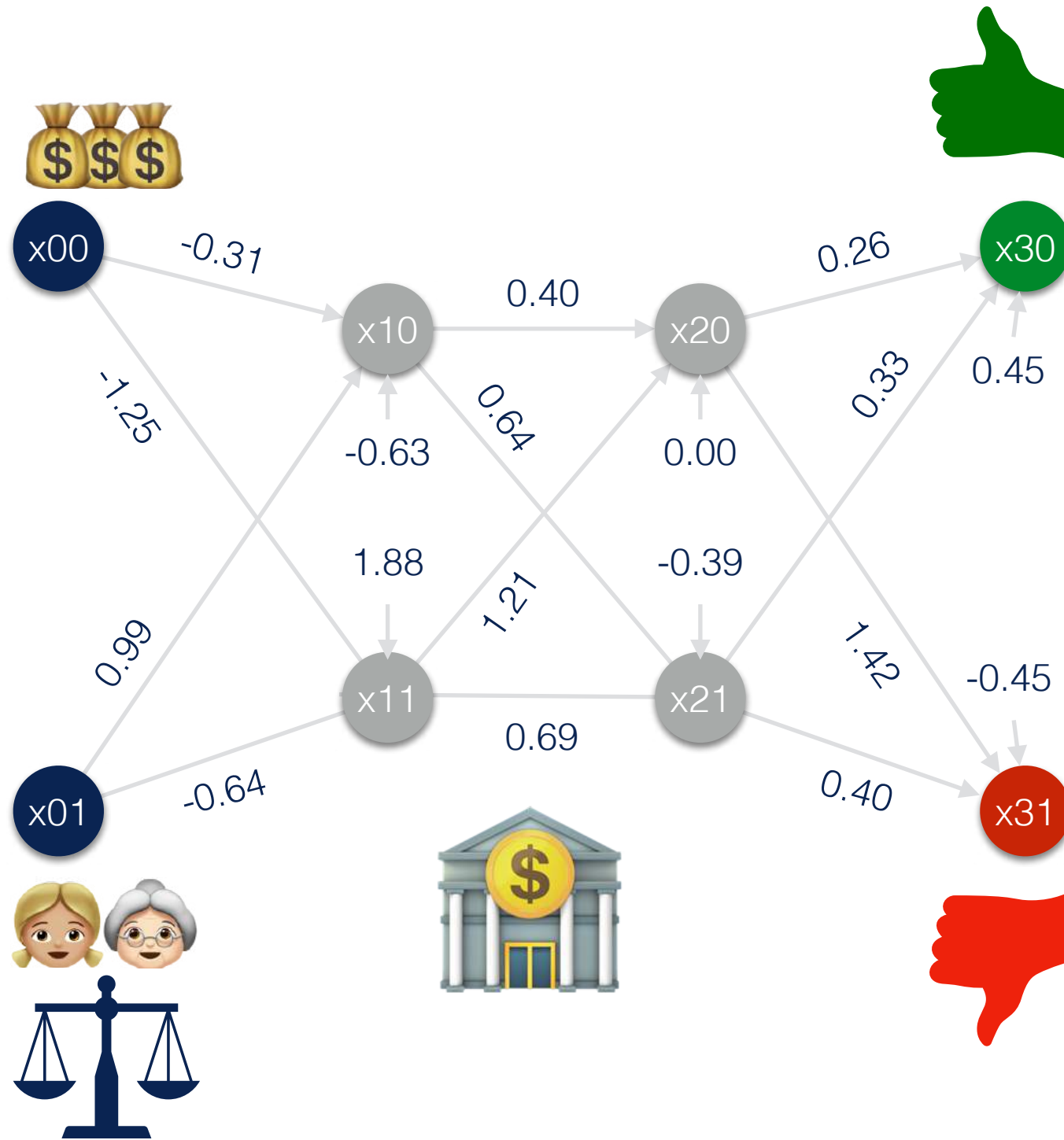
② **forget** the values of the **sensitive input** nodes



① proceed **backwards** from **all** possible classification **outcomes**

③ check for **intersection**:
empty → ✓ **fair**
otherwise → 🚨 **alarm**

Naïve Backward Analysis



```

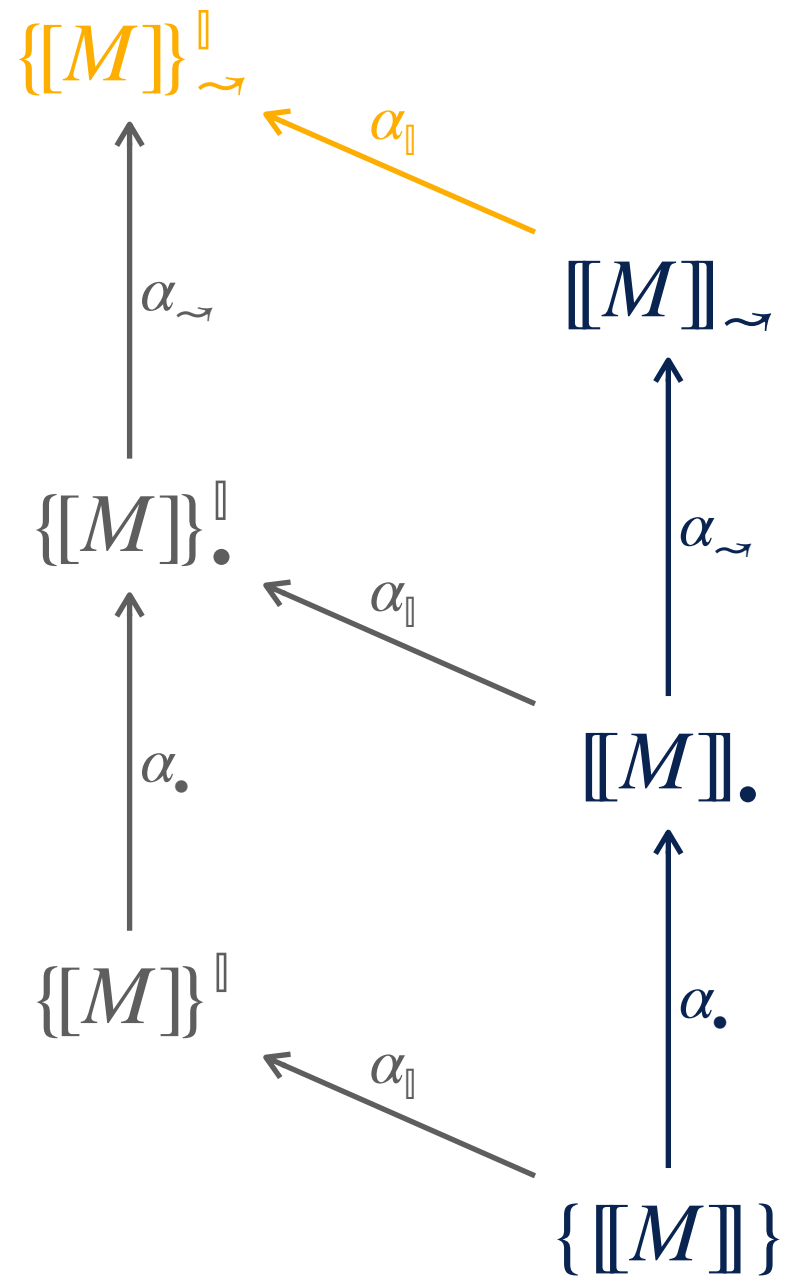
x00 = input()
x01 = input()
x10 = -0.31 * x00 + 0.99 * x01 + (-0.63)
x11 = -1.25 * x00 + (-0.64) * x01 + 1.88
x10 = 0 if x10 < 0 else x10
x11 = 0 if x11 < 0 else x11
x20 = 0.40 * x10 + 1.21 * x11 + 0.00
x21 = 0.64 * x10 + 0.69 * x11 + (-0.39)
x20 = 0 if x20 < 0 else x20
x21 = 0 if x21 < 0 else x21
1.16 * x20 + 0.07 * x21 ≤ 0.90
1.16 * x20 + 0.07 * x21 ≥ 0.90
x30 = 0.26 * x20 + 0.33 * x21 + 0.45
x31 = 1.42 * x20 + 0.40 * x21 + (-0.45)
x30 ≥ x31
x31 ≥ x30
return 'thumbs up' if x31 < 30 else 'thumbs down'
    
```

too many disjunctions!

Back to the Semantics...

Hierarchy of Semantics

parallel semantics

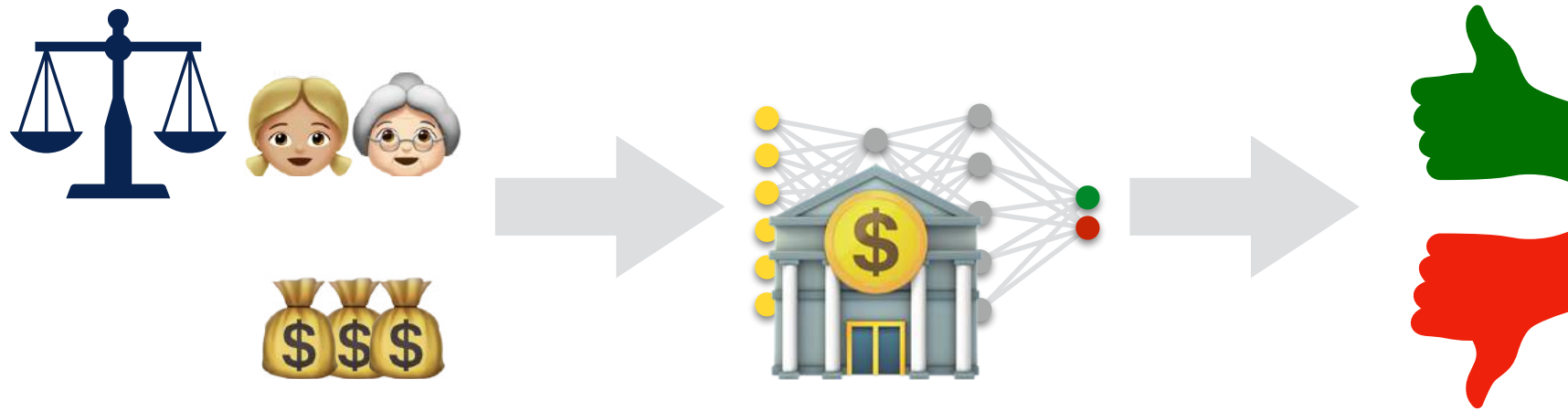


dependency semantics

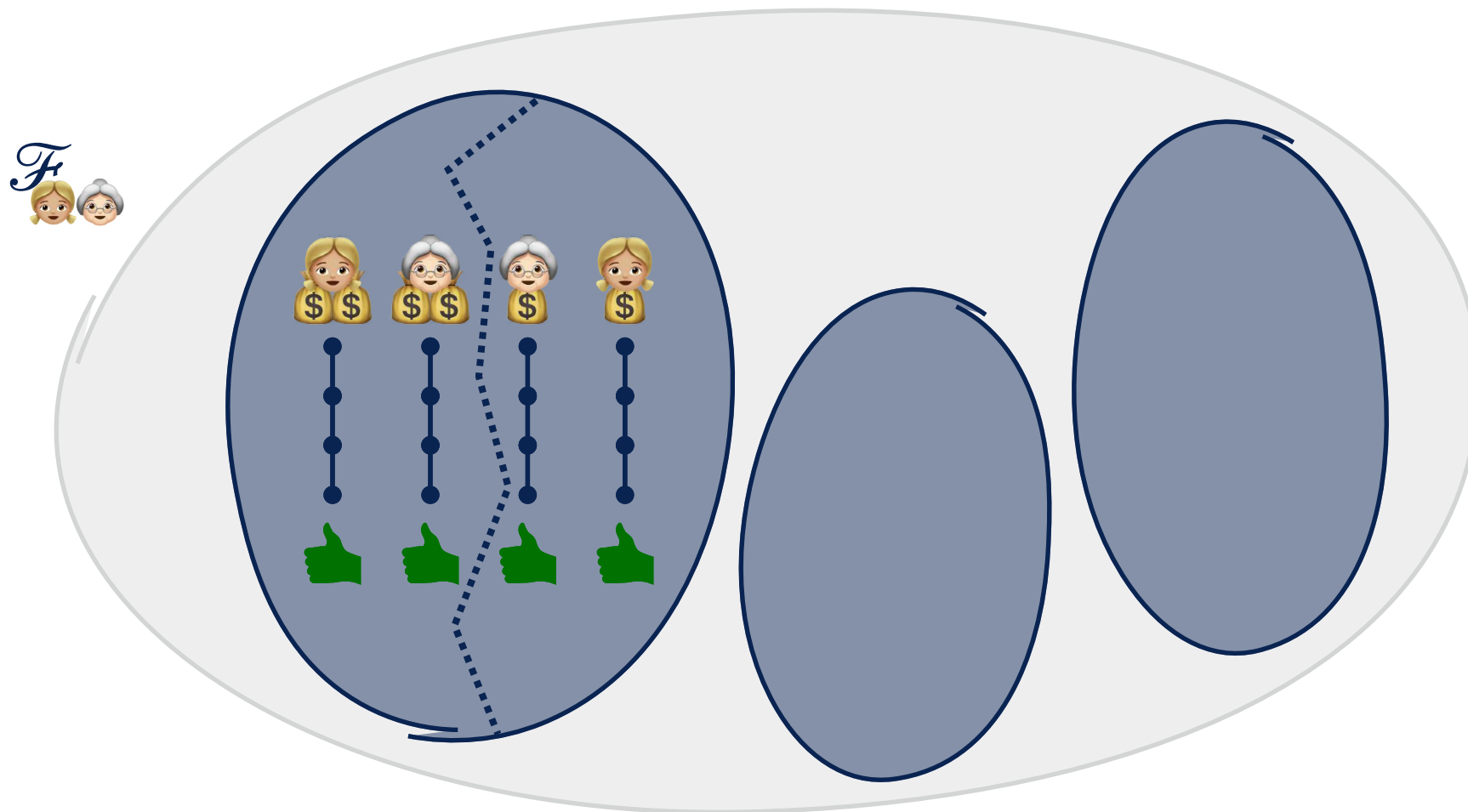
outcome semantics

collecting semantics

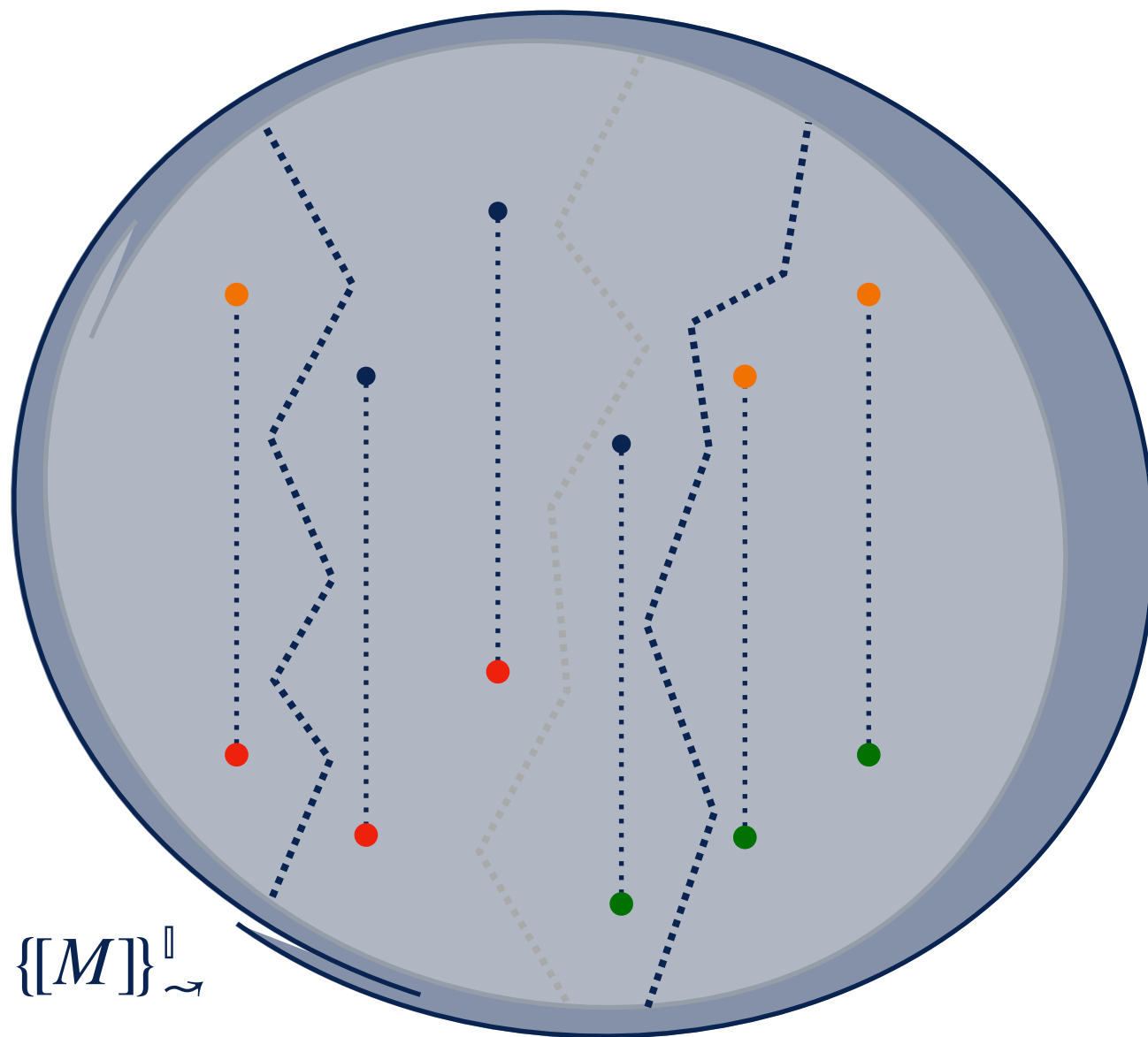
Parallel Semantics




 **partitioning** a set of traces that satisfies dependency fairness **with respect to the non-sensitive inputs** yields sets of traces that also satisfy dependency fairness

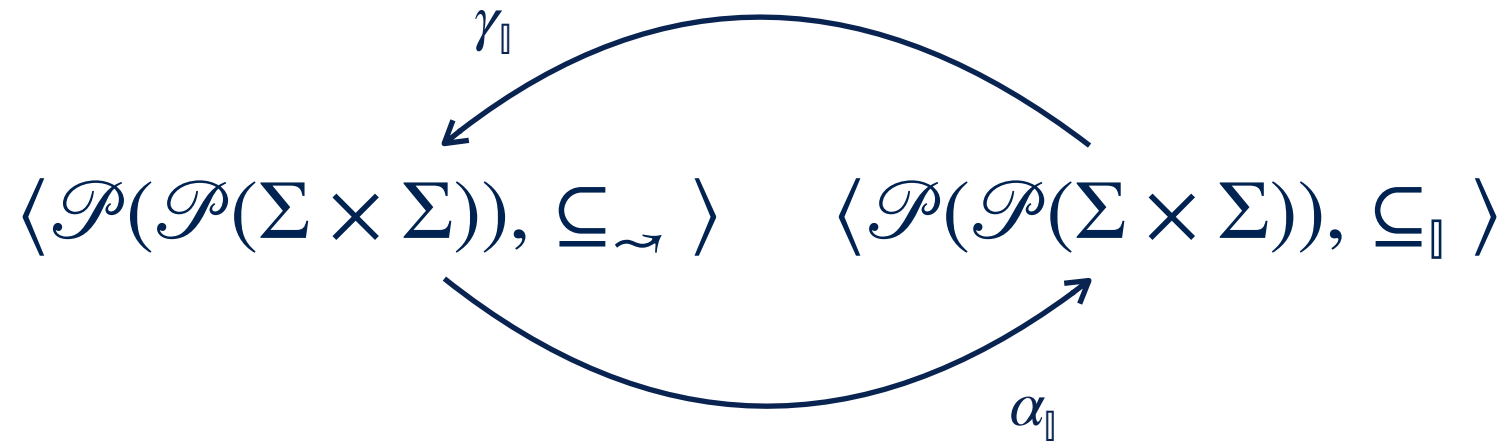


Parallel Semantics



💡 **partitioning** a set of traces that satisfies dependency fairness **with respect to the non-sensitive inputs** yields sets of traces that also satisfy dependency fairness

Parallel Semantics



$$\alpha_{\parallel}(S) \stackrel{\text{def}}{=} \{ \{ \langle t_0, t_\omega \rangle \in R \mid t_0 \in I \} \mid R \in S \wedge I \in \mathbb{I} \} \quad \text{parallel abstraction}$$

$$\begin{aligned} \llbracket M \rrbracket_{\sim}^{\mathbb{I}} &\stackrel{\text{def}}{=} \alpha_{\parallel}(\llbracket M \rrbracket_{\sim}) \\ &= \{ \{ \langle t_0, t_\omega \rangle \in \Sigma \times \Sigma \mid t \in \llbracket M \rrbracket \wedge t_0 \in I \wedge t_\omega \in O \} \mid I \in \mathbb{I} \wedge O \in \mathbb{O} \} \end{aligned}$$

Theorem

$$M \models \mathcal{F}_i \Leftrightarrow \gamma_{\sim}(\llbracket M \rrbracket_{\sim}^{\mathbb{I}}) \subseteq \mathcal{F}_i$$

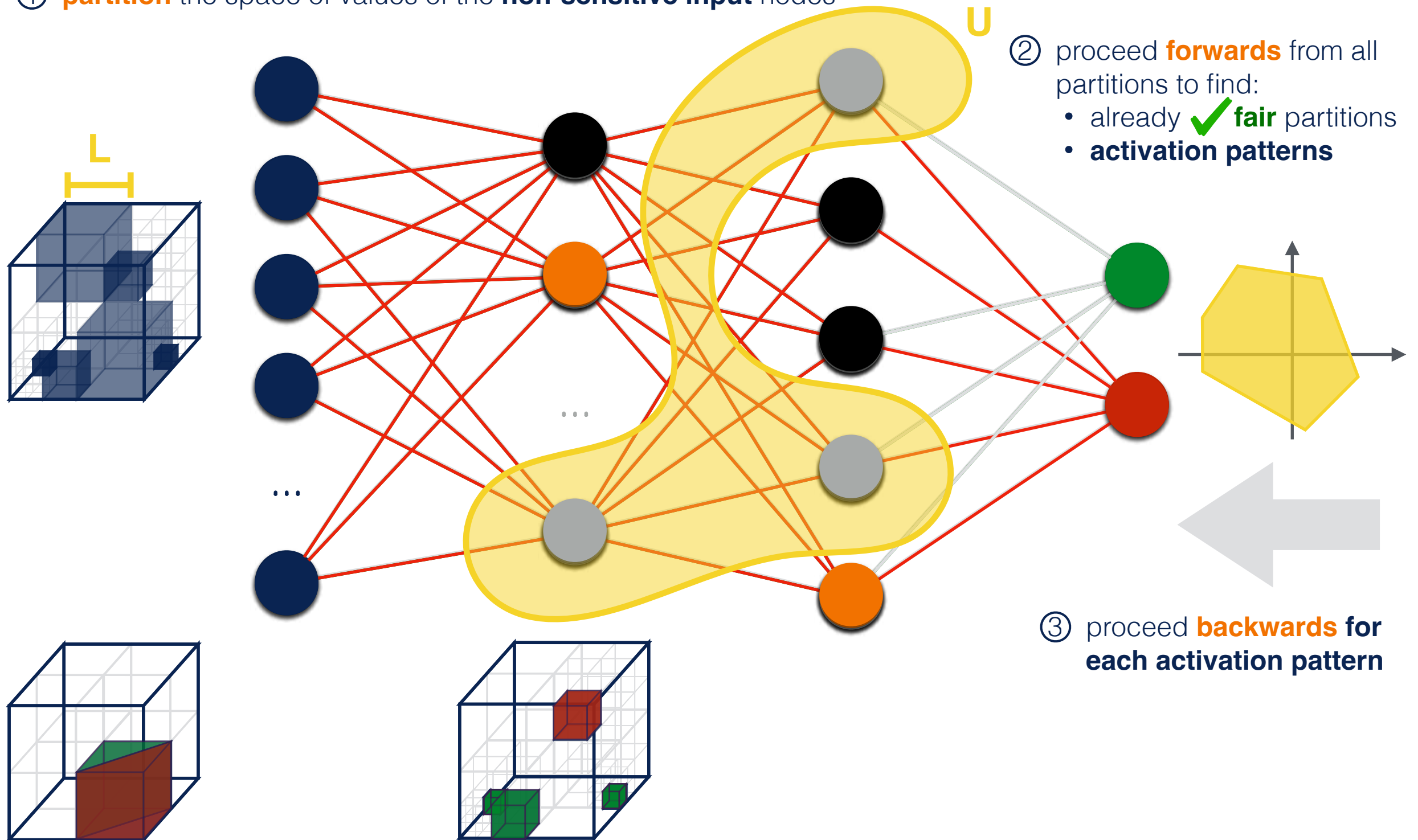
Lemma

$$M \models \mathcal{F}_i \Leftrightarrow \forall I \in \mathbb{I}: \forall A, B \in \llbracket M \rrbracket_{\sim}^{\mathbb{I}}: (A_\omega^I \neq B_\omega^I \Rightarrow A_0^I|_{\neq i} \cap B_0^I|_{\neq i} = \emptyset)$$

Better Abstraction

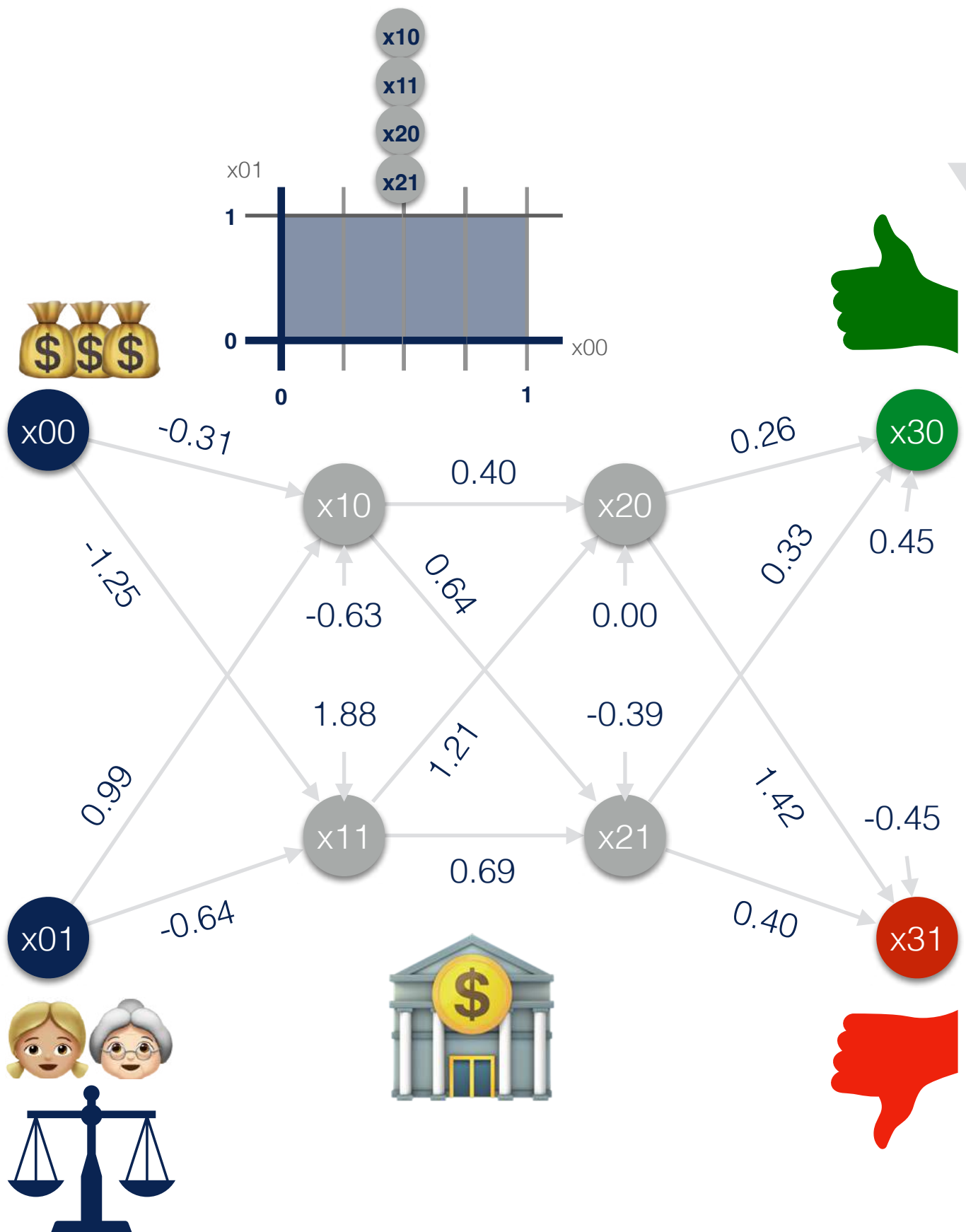
Forward and Backward Analysis

① **partition** the space of values of the **non-sensitive input** nodes



$$L = 0.25$$

$$U = 2$$



```

x00 = input()
x01 = input()

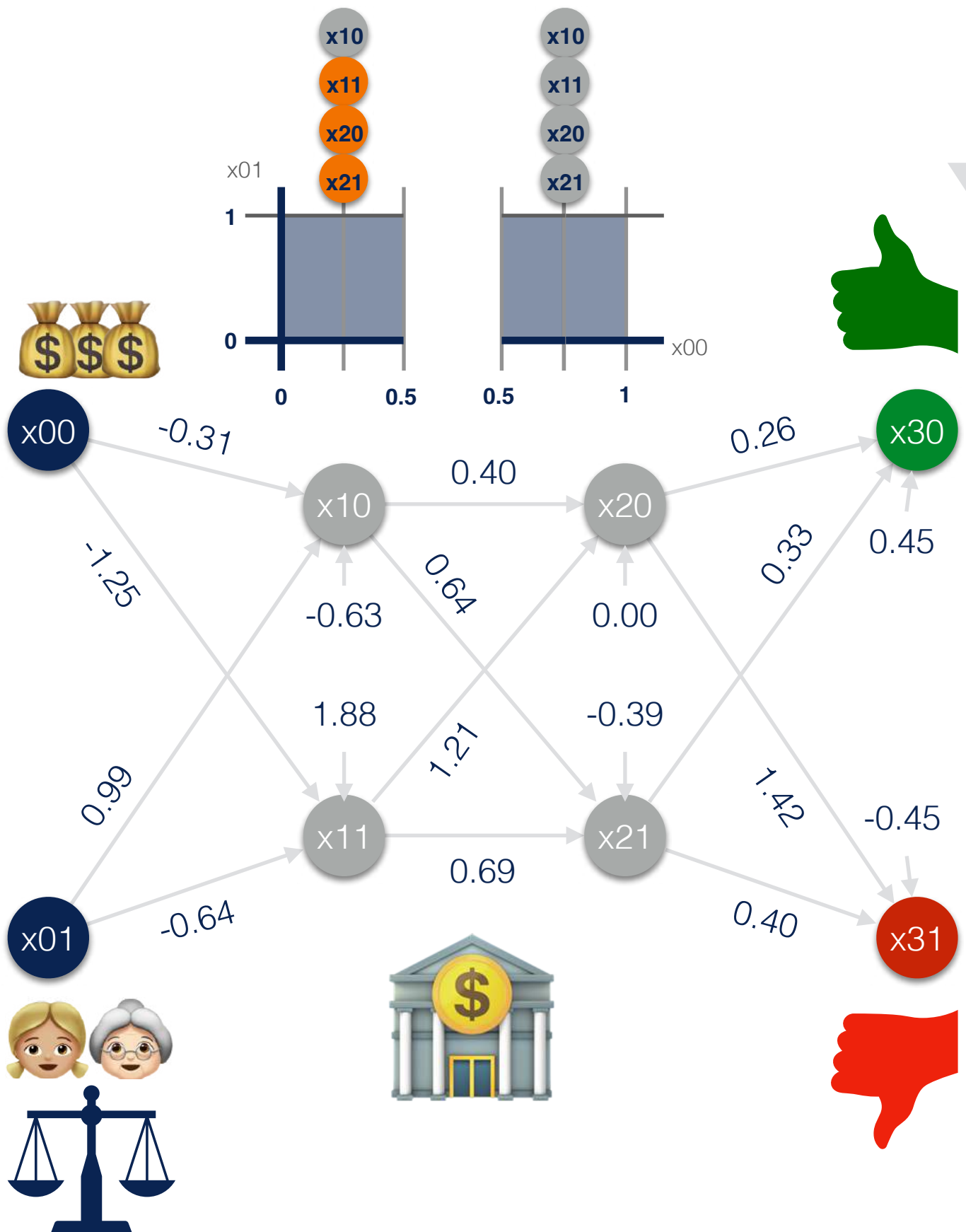
x10 = -0.31 * x00 + 0.99 * x01 + (-0.63)
x11 = -1.25 * x00 + (-0.64) * x01 + 1.88

x20 = 0.40 * x10 + 1.21 * x11 + 0.00
x21 = 0.64 * x10 + 0.69 * x11 + (-0.39)

x30 = 0.26 * x20 + 0.33 * x21 + 0.45
x31 = 1.42 * x20 + 0.40 * x21 + (-0.45)

return '👍' if x31 < 30 else '👎'

```



$L = 0.25$
 $U = 2$

```

x00 = input()
x01 = input()

x10 = -0.31 * x00 + 0.99 * x01 + (-0.63)
x11 = -1.25 * x00 + (-0.64) * x01 + 1.88

x20 = 0 if x20 < 0 else x20
x21 = 0 if x21 < 0 else x21

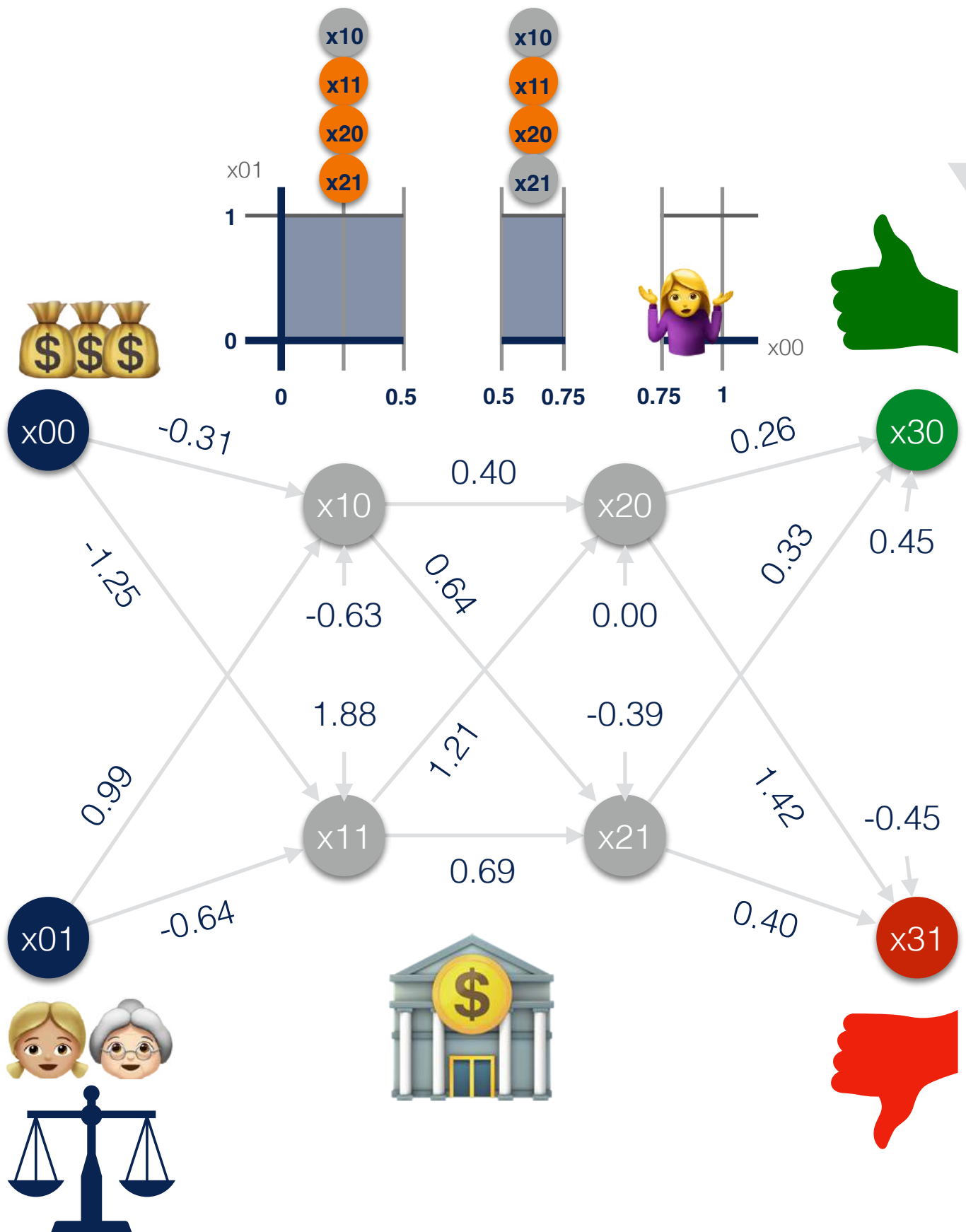
x20 = 0.40 * x10 + 1.21 * x11 + 0.00
x21 = 0.64 * x10 + 0.69 * x11 + (-0.39)

x20 = 0 if x20 < 0 else x20
x21 = 0 if x21 < 0 else x21

x30 = 0.26 * x20 + 0.33 * x21 + 0.45
x31 = 1.42 * x20 + 0.40 * x21 + (-0.45)

return '👍' if x31 < 30 else '👎'

```

$L = 0.25$
 $U = 2$

```

x00 = input()
x01 = input()

x10 = -0.31 * x00 + 0.99 * x01 + (-0.63)
x11 = -1.25 * x00 + (-0.64) * x01 + 1.88

x20 = 0 if x10 < 0 else x10
x21 = 0 if x11 < 0 else x11

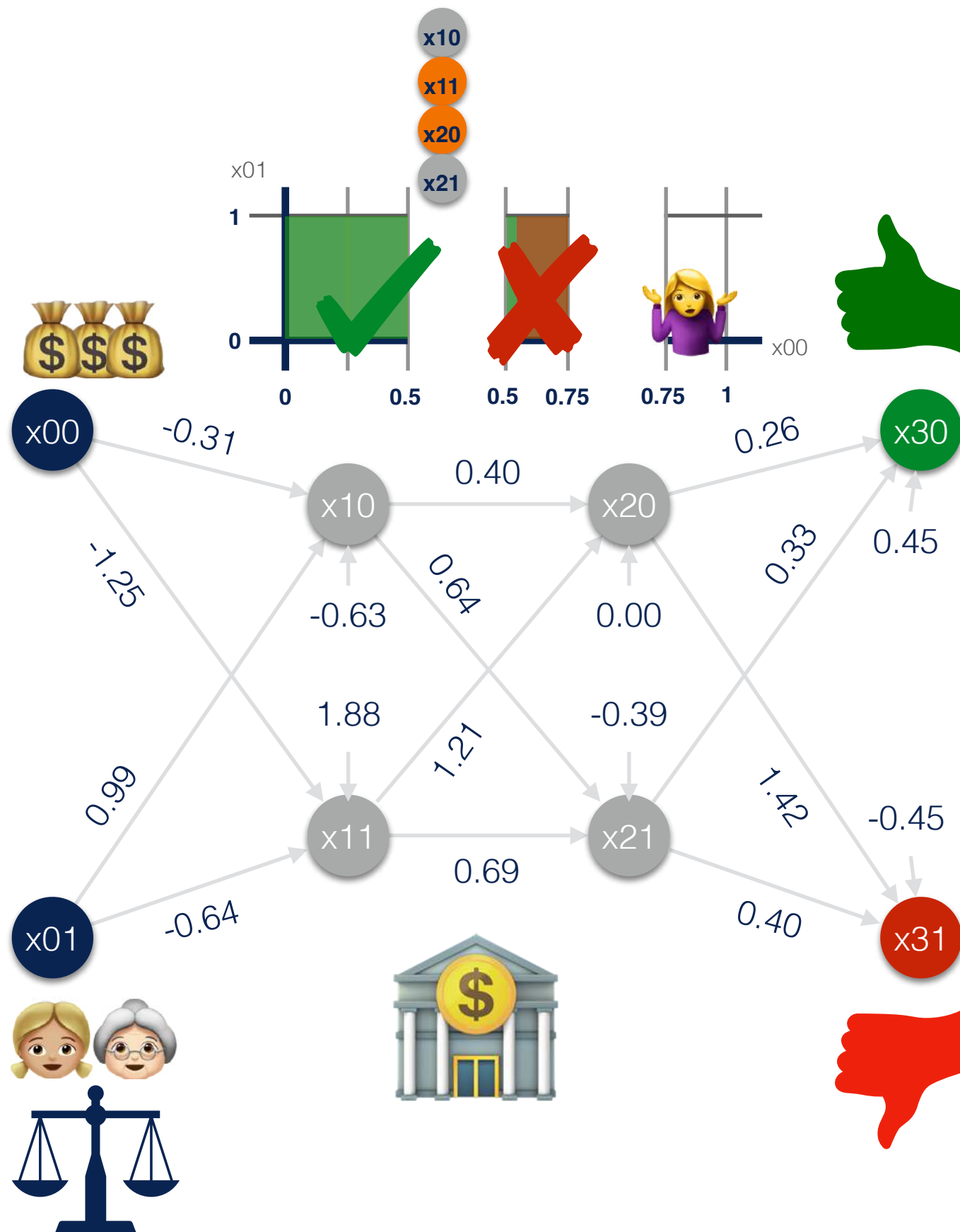
x20 = 0.40 * x10 + 1.21 * x11 + 0.00
x21 = 0.64 * x10 + 0.69 * x11 + (-0.39)

x20 = 0 if x20 < 0 else x20
x21 = 0 if x21 < 0 else x21

x30 = 0.26 * x20 + 0.33 * x21 + 0.45
x31 = 1.42 * x20 + 0.40 * x21 + (-0.45)

return '👍' if x31 < 30 else '👎'

```



$$L = 0.25$$

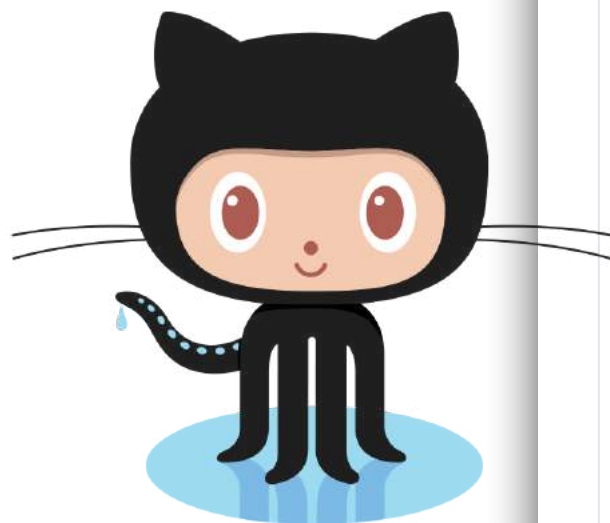
$$U = 2$$

```

x00 = input()
x01 = input()
x10 = -0.31 * x00 + 0.99 * x01 + (-0.63)
x11 = -1.25 * x00 + (-0.64) * x01 + 1.88
x10 = 0 if x10 < 0 else x10
x11 = 0 if x11 < 0 else x11
x20 = 0.40 * x10 + 1.21 * x11 + 0.00
x21 = 0.64 * x10 + 0.69 * x11 + (-0.39)
x20 = 0 if x20 < 0 else x20
x21 = 0 if x21 < 0 else x21
1.16 * x20 + 0.07 * x21 ≤ 0.90
1.16 * x20 + 0.07 * x21 ≥ 0.90
x30 = 0.26 * x20 + 0.33 * x21 + 0.45
x31 = 1.42 * x20 + 0.40 * x21 + (-0.45)
x30 ≥ x31
x31 ≥ x30
return 'thumbs up' if x31 < 30 else 'thumbs down'

```

Libra



caterinaurban / **Libra**

<> Code Issues Pull requests Actions Projects Security Insights


master 2 branches 0 tags Go to file Code

caterinaurban README 9f830db on Aug 8 53 commits

src	RQ5 and RQ6 reproducibility	4 months ago
.gitignore	RQ1 reproducibility	4 months ago
LICENSE	Initial prototype	2 years ago
README.md	RQ5 and RQ6 reproducibility	4 months ago
README.pdf	README	4 months ago
icon.png	icon	4 months ago
libra.png	icon	4 months ago
requirements.txt	some documentation	4 months ago
setup.py	some documentation	4 months ago

README.md

Libra



Nowadays, machine-learned software plays an increasingly important role in critical decision-making in our social, economic, and civic lives.

About
No description or website provided.

#abstract-interpretation
#static-analysis
#machine-learning
#neural-networks #fairness

Readme
MPL-2.0 License

Releases
No releases published

Packages
No packages published

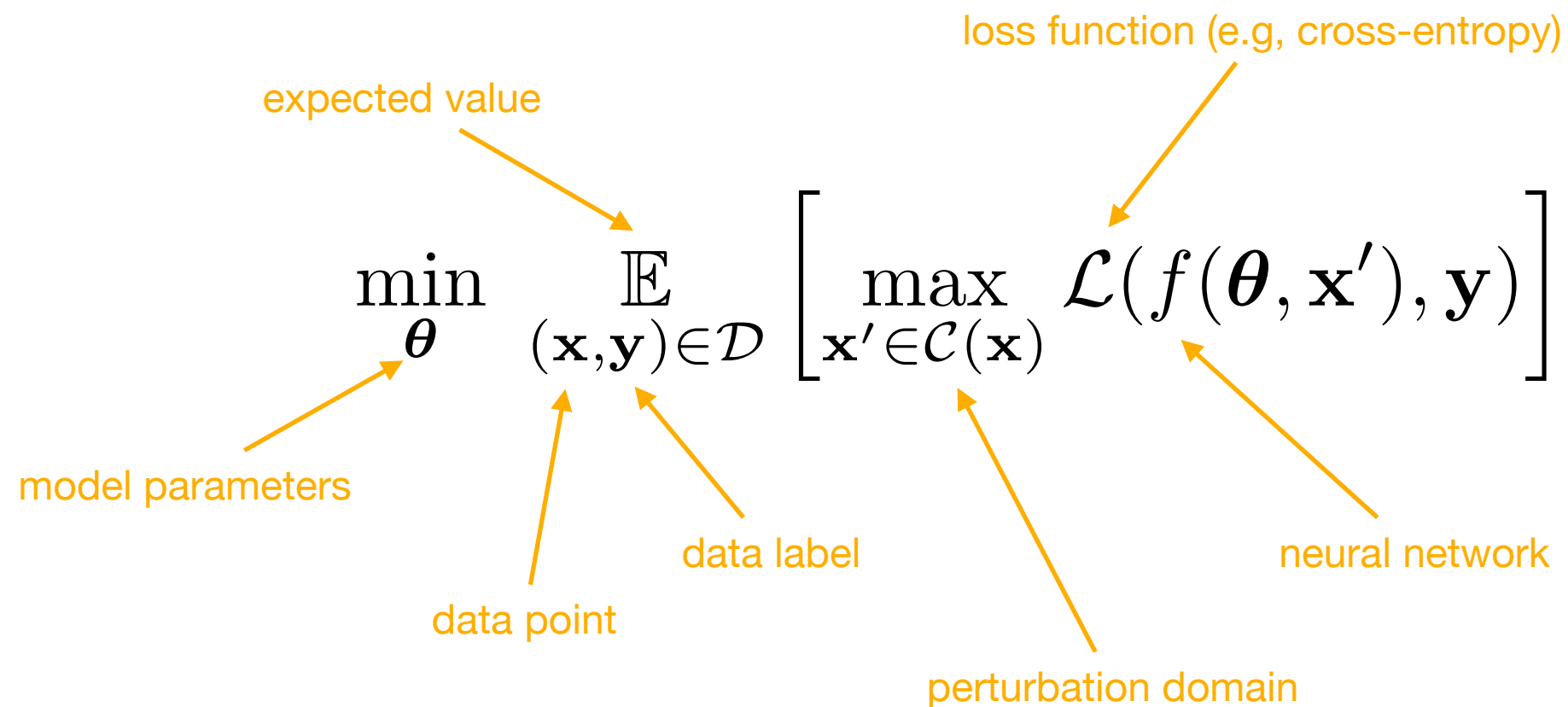
Languages

- Python 98.7%
- Shell 1.3%

Formal Methods for Model Training

Robust Training

Minimizing the Worst-Case Loss for Each Input



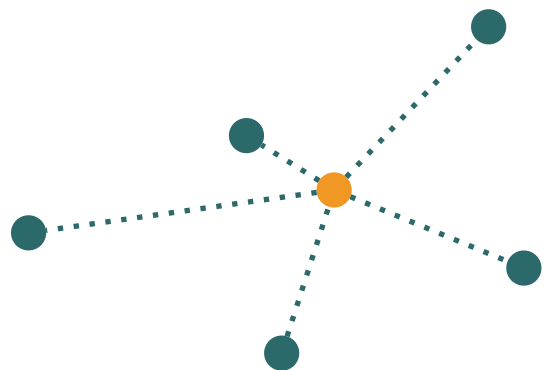
Robust Training

Minimizing the Worst-Case Loss for Each Input

Adversarial Training

Minimizing a Lower Bound on the Worst-Case Loss for Each Input

$$\max_{\mathbf{x}' \in \mathcal{C}(\mathbf{x})} \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}'), \mathbf{y})$$



VI

$$\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), \mathbf{y})$$



generate adversarial inputs
and use them as training data

Robust Training

Minimizing the Worst-Case Loss for Each Input

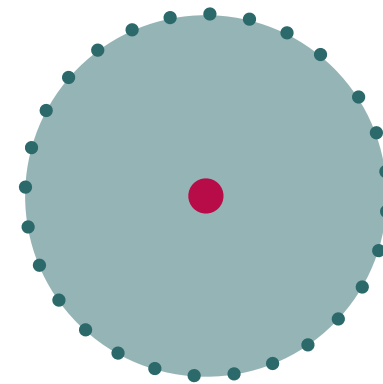
$$\max_{\mathbf{x}' \in \mathcal{C}(\mathbf{x})} \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}'), \mathbf{y})$$

\wedge

$$\mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), \mathbf{y})$$

Certified Training

Minimizing an **Upper Bound on the Worst-Case Loss** for Each Input



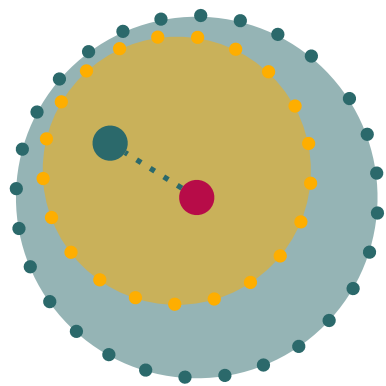
use upper bound as regularizer to encourage robustness

Robust Training

Minimizing the Worst-Case Loss for Each Input

Hybrid Training

Minimizing an **Approximation of the Worst-Case Loss** that **Contains and Adversarial Example** for Each Input



use approximation as regularizer to encourage robustness

M. N. Müller et al. - Certified Training: Small Boxes Are All You Need (ICLR 2023)

99xv1 [cs.LG] 23 May 2023

Expressive Losses for Verified Robustness via Convex Combinations

Alessandro De Palma
Imperial College London¹
adepalma@ic.ac.uk

Rudy Bunel
Google DeepMind

Krishnamurthy (Dj) Dvijotham
Google DeepMind

M. Pawan Kumar
Google DeepMind

Robert Stanforth
Google DeepMind

Alessio Lomuscio
Imperial College London

Abstract

In order to train networks for verified adversarial robustness, previous work typically over-approximates the worst-case loss over (subsets of) perturbation regions or induces verifiability on top of adversarial training. The key to state-of-the-art performance lies in the expressivity of the employed loss function, which should be able to match the tightness of the verifiers to be employed post-training. We formalize a definition of expressivity, and show that it can be satisfied via simple convex combinations between adversarial attacks and IBP bounds. We then show that the resulting algorithms, named CC-IBP and MTL-IBP, yield state-of-the-art results across a variety of settings in spite of their conceptual simplicity. In particular, for ℓ_∞ perturbations of radius $1/255$ on TinyImageNet and downscaled ImageNet, MTL-IBP improves on the best standard and verified accuracies from the literature by from 1.98% to 3.92% points while only relying on single-step adversarial attacks.

Introduction

... successes [24, 34, 50], serious concerns over the trust-
... the vulnerability of neural networks to
... critical domains. As a result,
... neural net-

Bibliography

[Kurd03] **Zeshan Kurd, Tim Kelly.** Establishing Safety Criteria for Artificial Neural Networks. In KES, pages 63-169, 2003.

[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.

[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.

[Mazzucato21] **Denis Mazzucato and Caterina Urban.** Reduced Products of Abstract Domains for Fairness Certification of Neural Networks. In SAS, 2021.

[Julian16] **Kyle D. Julian, Jessica Lopez, Jeffrey S. Brush, Michael P. Owen, Mykel J. Kochenderfer.** Policy Compression for Aircraft Collision Avoidance Systems. In DASC, pages 1–10, 2016.

Bibliography

[Katz17] **Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, Mykel J. Kochenderfer.** Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In CAV, pages 97–117, 2017.

[Galhotra17] **Sainyam Galhotra, Yuriy Brun, and Alexandra Meliou.** Fairness Testing: Testing Software for Discrimination. In FSE, pages 498–510, 2017.

[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.

[Urban21] **Caterina Urban and Antoine Miné.** A Review of Formal Methods applied to Machine Learning. <https://arxiv.org/abs/2104.02466>, 2021.

[Müller23] **Mark Niklas Müller, Franziska Eckert, Marc Fischer, Martin Vechev.** Certified Training: Small Boxes Are All You Need. In ICLR, 2023.