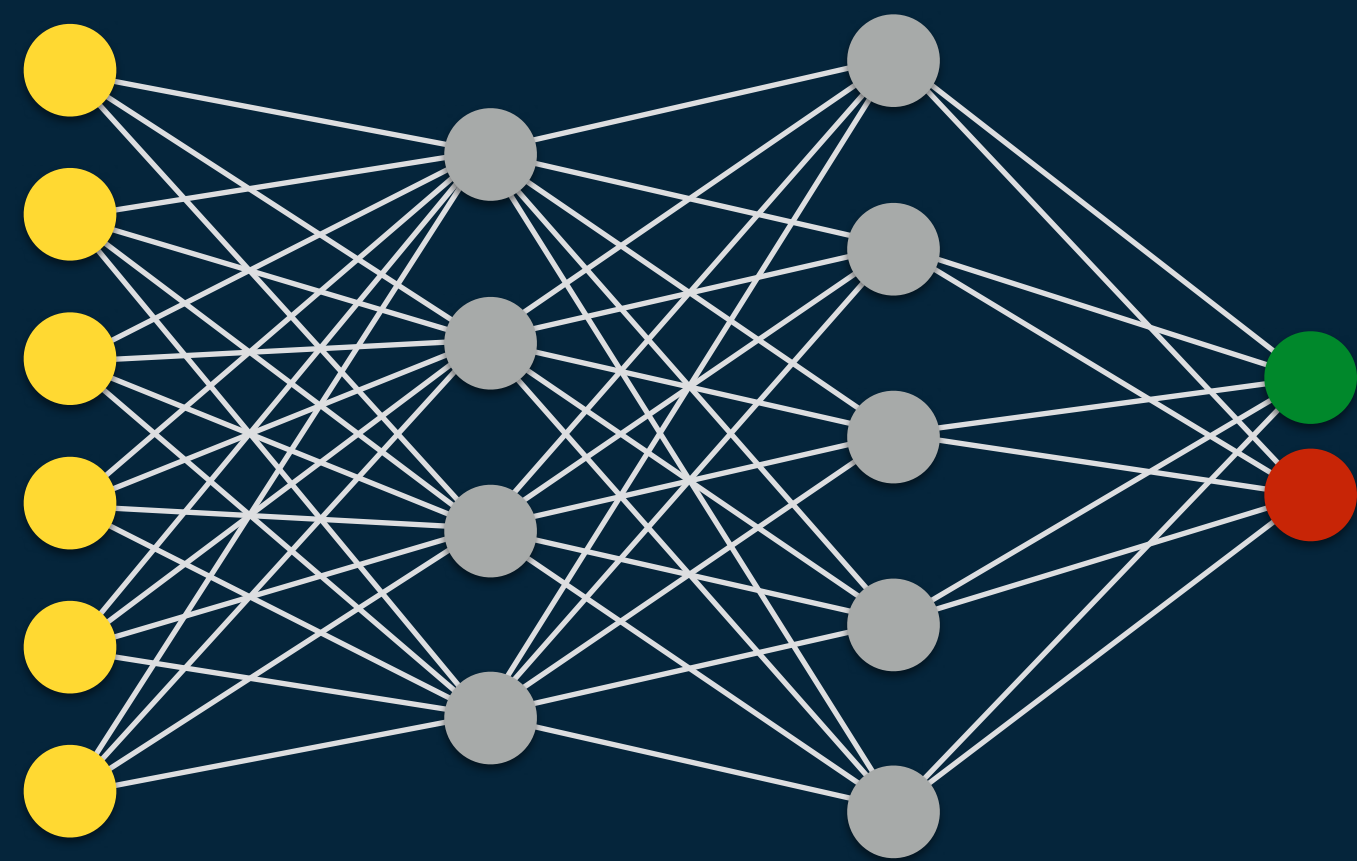


# Machine Learning Interpretability and Verification



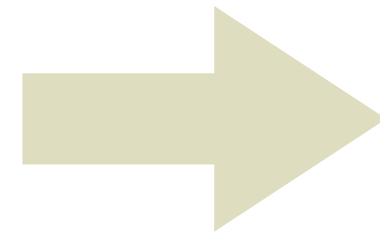
**Caterina Urban**  
Inria & École Normale Supérieure



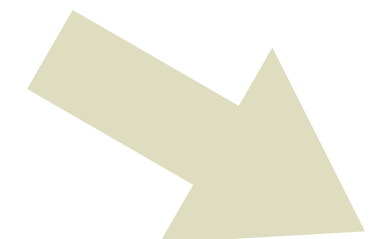
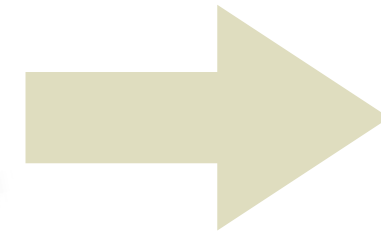
# Machine Learning in High-Stakes Systems



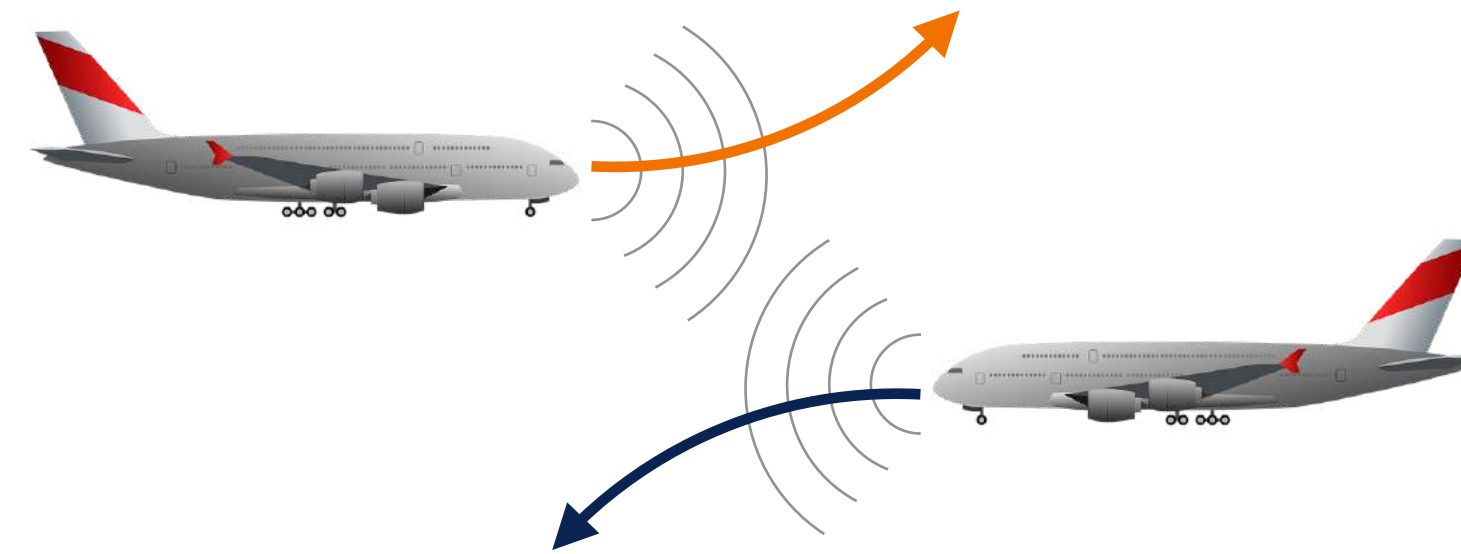
**data**



**ML software**



**perform tasks that are impossible using explicit programming**



**act as surrogate model**

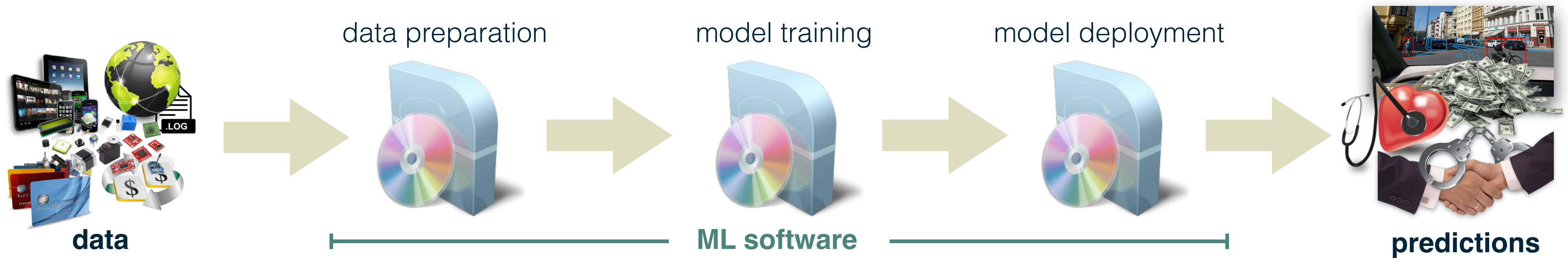


**automate decision-making**

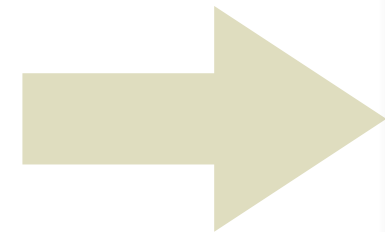


# Machine Learning Development Process

## Machine Learning Pipeline



# Models Only Give Probabilistic Guarantees



model deployment

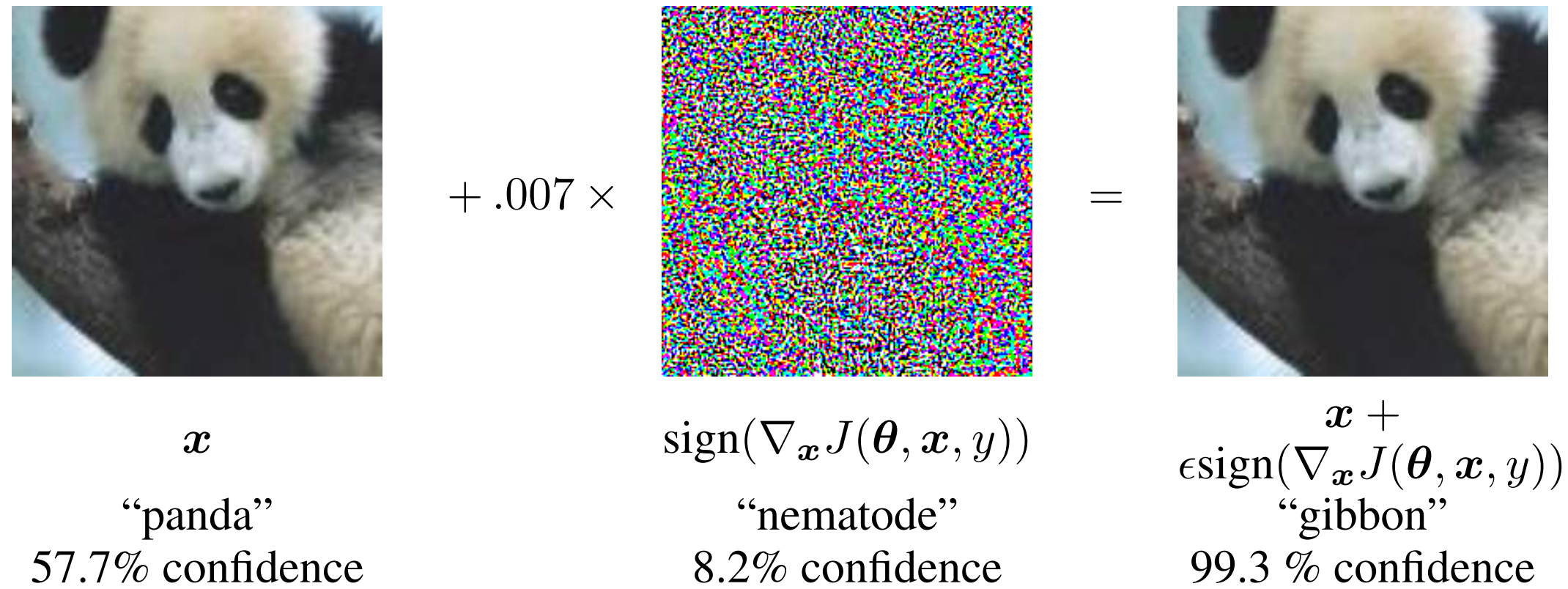


not sufficient for guaranteeing an **acceptable failure rate** under all circumstances



# Models Only Give Probabilistic Guarantees

## Adversarial Examples



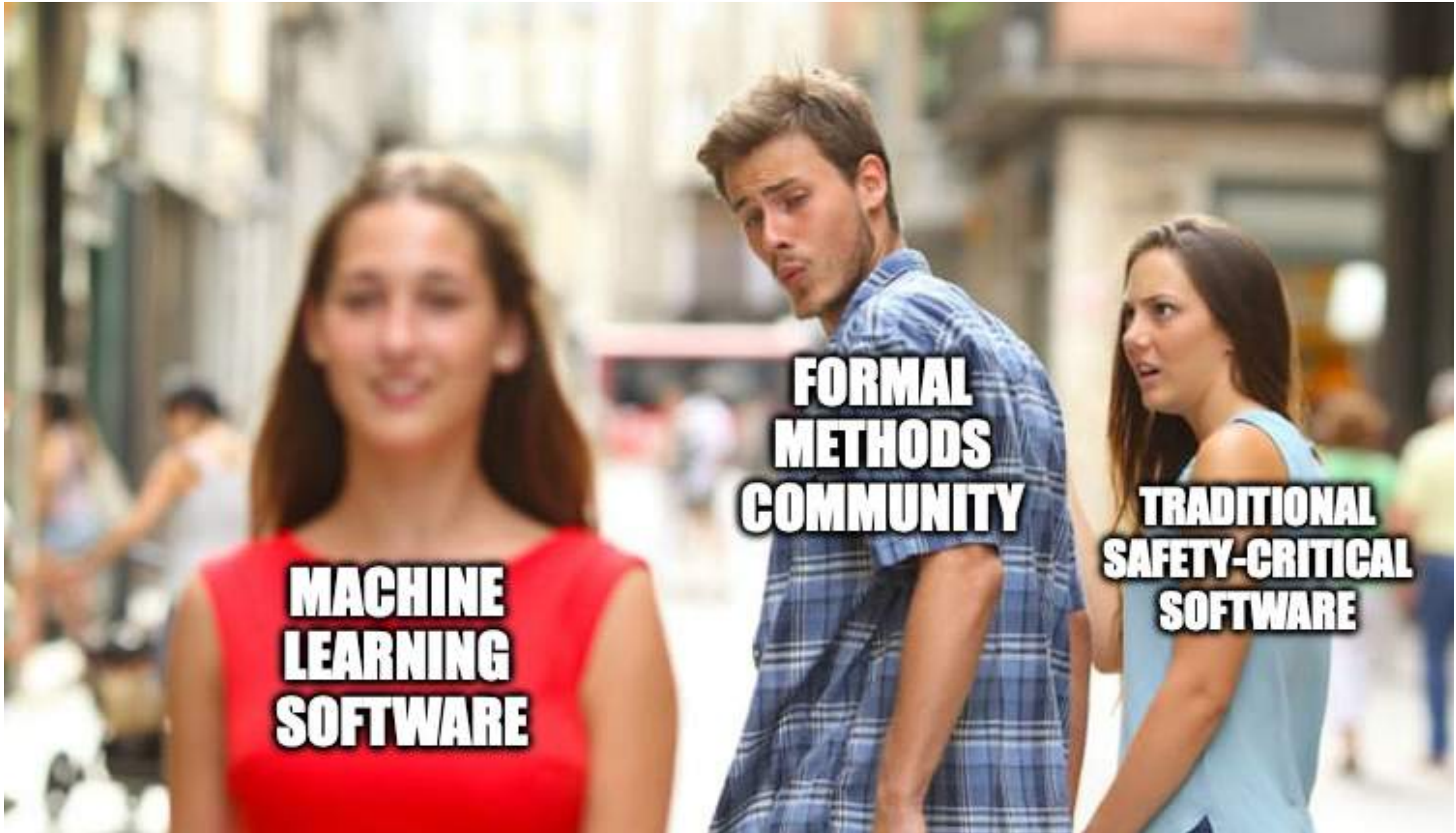
Published as a conference paper at ICLR 2015

---

### EXPLAINING AND HARNESSING ADVERSARIAL EXAMPLES

**Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy**  
Google Inc., Mountain View, CA  
{goodfellow, shlens, szegedy}@google.com





**MACHINE  
LEARNING  
SOFTWARE**

**FORMAL  
METHODS  
COMMUNITY**

**TRADITIONAL  
SAFETY-CRITICAL  
SOFTWARE**

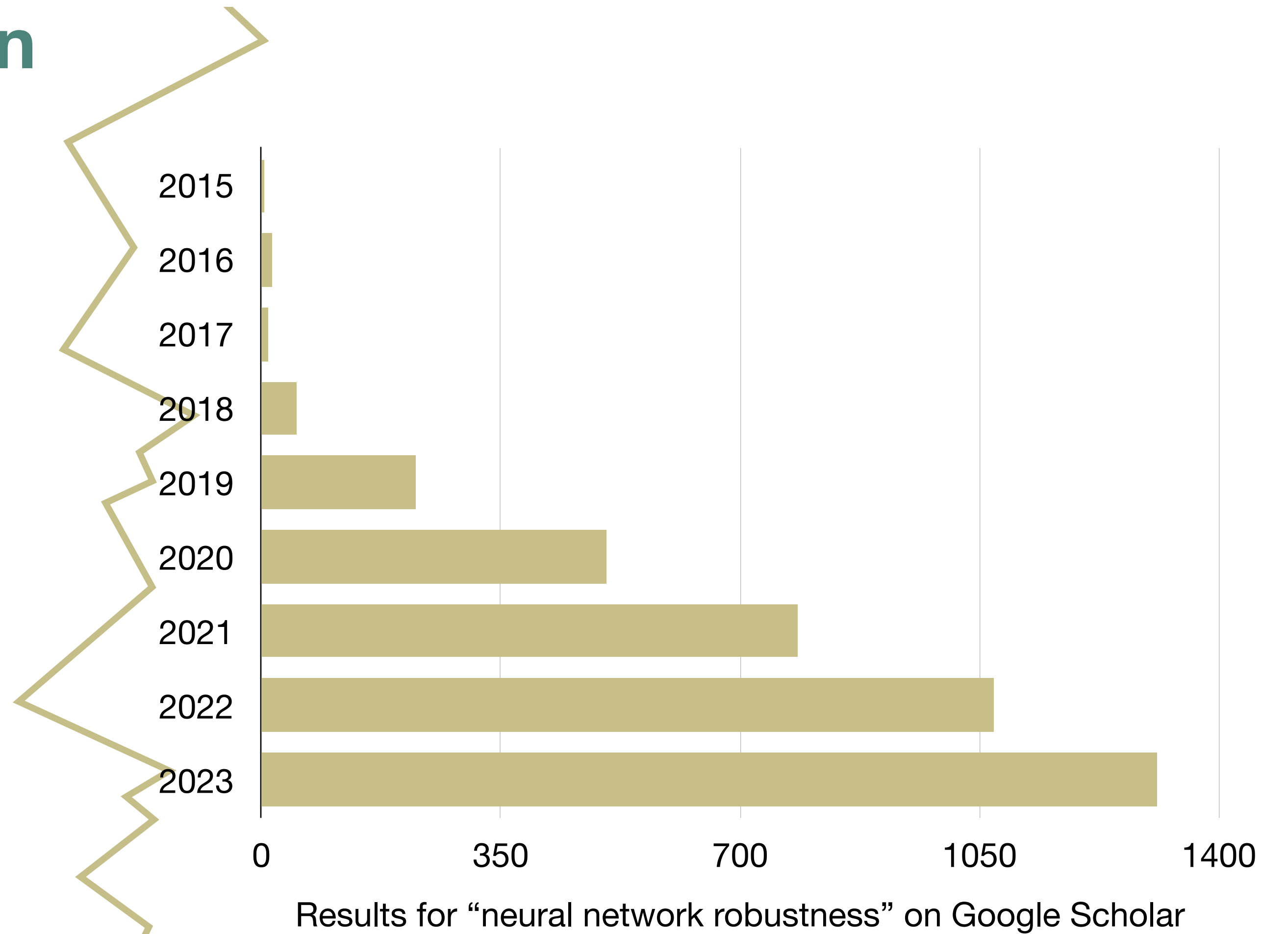


# Models Only Give Probabilistic Guarantees

## Local Robustness Verification



Machine Learning Community



Formal Methods Community

# Machine Learning Development Process

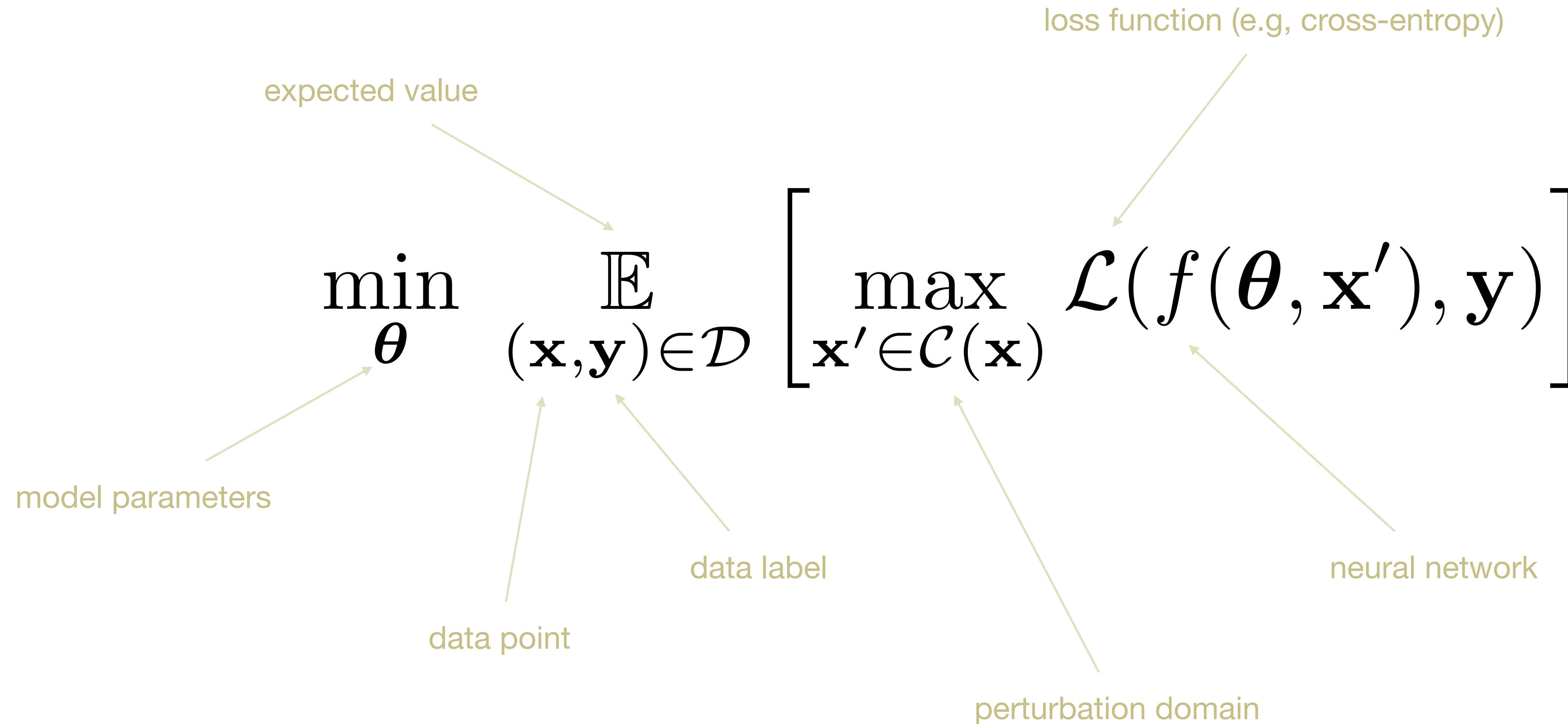
## Machine Learning Pipeline





# Robust Training

## Minimizing the Worst-Case Loss for Each Input

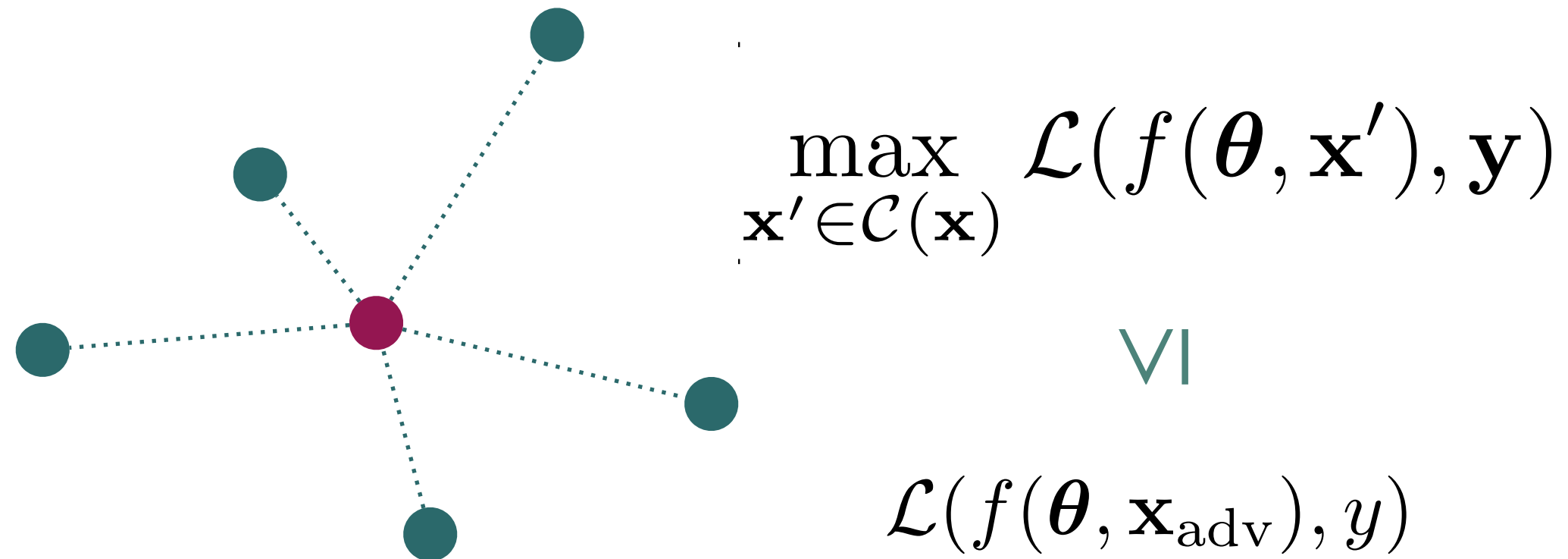


# Robust Training

Minimizing the Worst-Case Loss for Each Input

## Adversarial Training

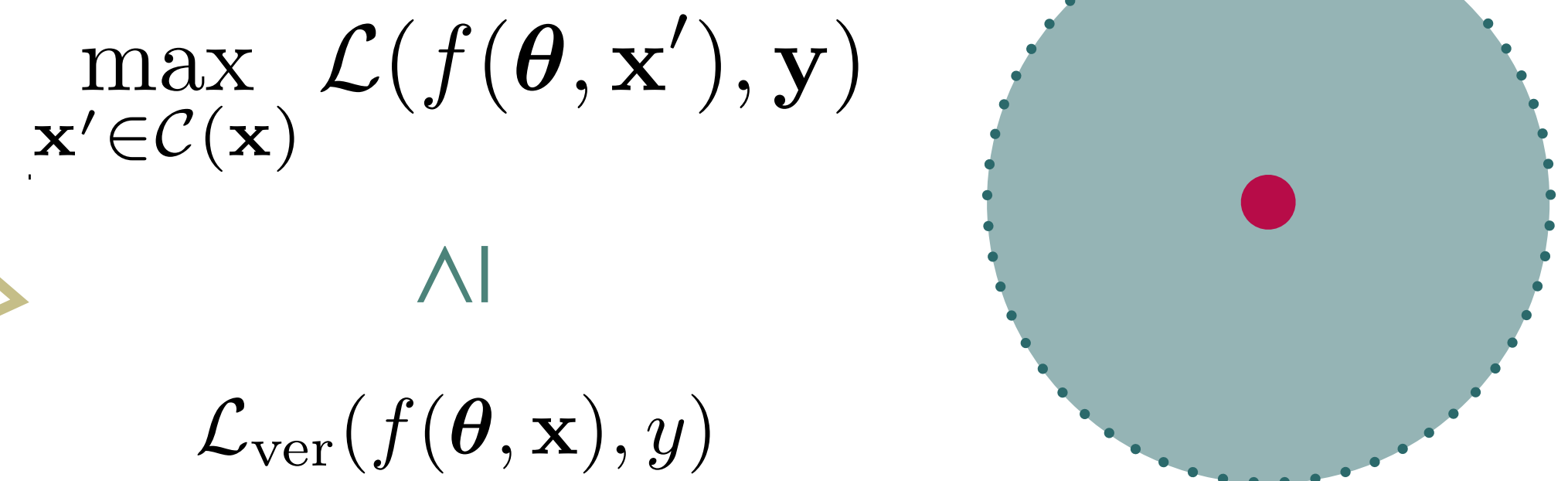
Minimizing a Lower-Bound on the Worst-Case Loss



Machine Learning Community

## Certified Training

Minimizing an Upper-Bound on the Worst-Case Loss



Formal Methods Community



# Robust Training

## Certified Training



Machine Learning Community

Table 7: Comparison of the standard (Acc.), adversarial (Adv. Acc), and certified (Cert. Acc.) accuracy for different certified training methods on the full CIFAR-10 test set. We use MN-BAB (Ferrari et al., 2022) to compute all certified and adversarial accuracies.

$\epsilon_\infty$	Training Method	Source	Acc. [%]	Adv. Acc. [%]	Cert. Acc. [%]
2/255	COLT	Balunovic & Vechev (2020)	78.42	<b>66.17</b>	61.02
	CROWN-IBP	Zhang et al. (2020) <sup>†</sup>	71.27	59.58	58.19
	IBP	Shi et al. (2021)	-	-	-
	SABR	this work	<b>79.52</b>	65.76	<b>62.57</b>
8/255	COLT	Balunovic & Vechev (2020)	51.69	31.81	27.60
	CROWN-IBP	Zhang et al. (2020) <sup>†</sup>	45.41	33.33	33.18
	IBP	Shi et al. (2021)	48.94	35.43	<b>35.30</b>
	SABR	this work	<b>52.00</b>	<b>35.70</b>	35.25

**ROBUSTBENCH** Leaderboards Paper FAQ Contribute Model Zoo 🚀

Leaderboard: CIFAR-10,  $\ell_\infty = 8/255$ , untargeted attack

Show 15 entries Search: Papers, architectures, v

Rank	Method	Standard accuracy	AutoAttack robust accuracy	Best known robust accuracy	AA eval. potentially unreliable	Extra data	Architecture	Venue
	Robust Principles: Architectural Design Principles for Adversarially Robust CNNs <small>It uses additional 50M synthetic images in training.</small>	93.27%	71.07%	71.07%	×	×	RaWideResNet-70-16	BMVC 2023

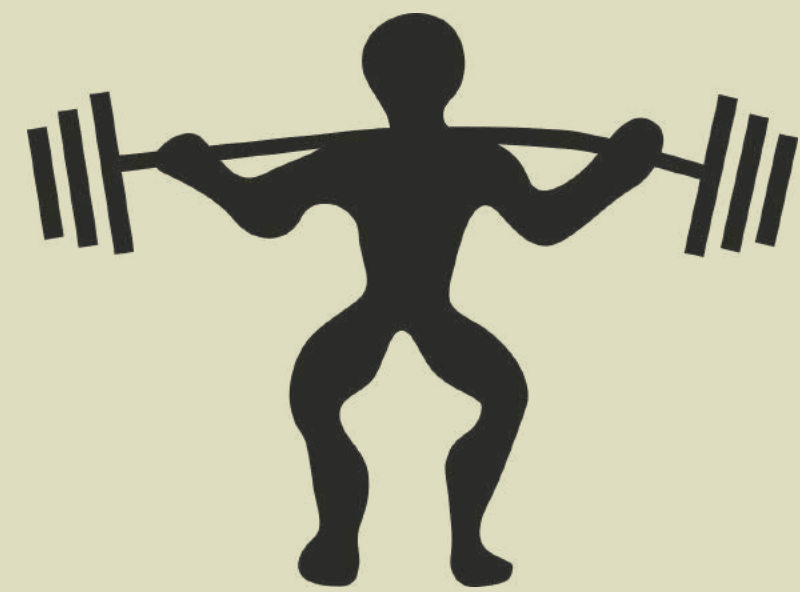
Formal Methods Community



Can we make  
formal methods  
**interesting**  
for the machine  
learning community?

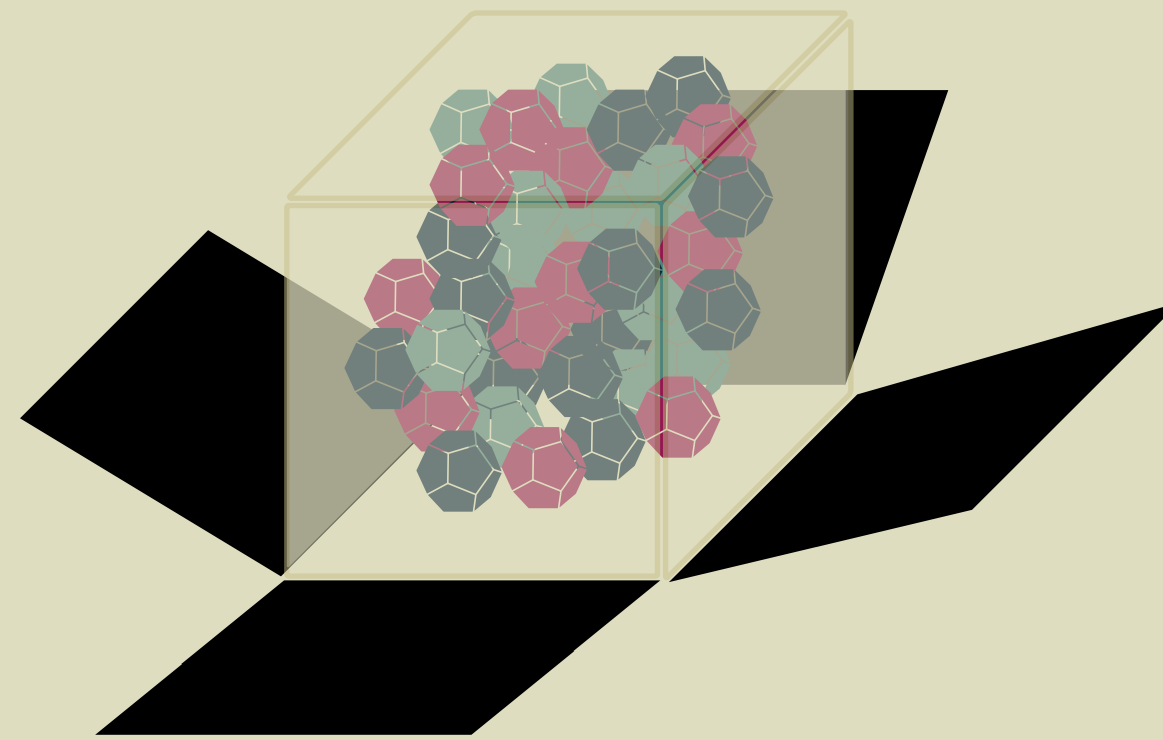






## Training

CIKM 2021



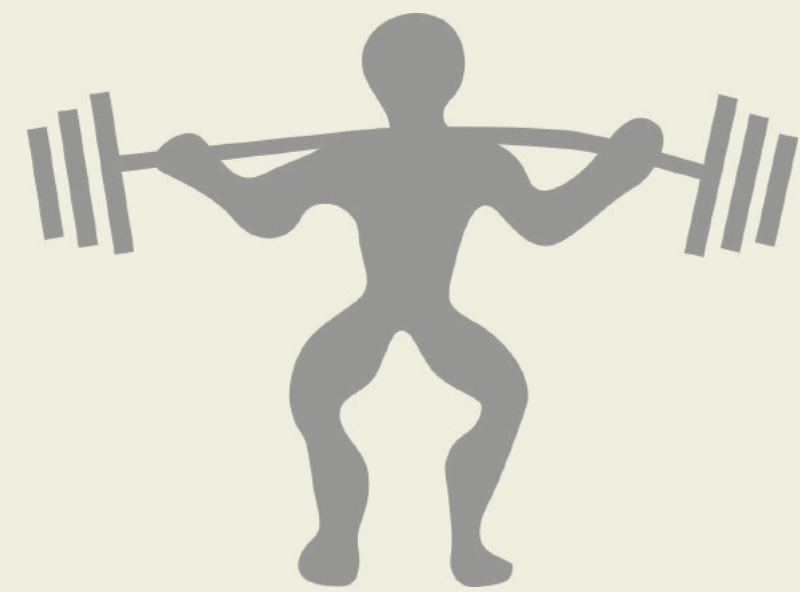
## Interpretability

VMCAI 2024



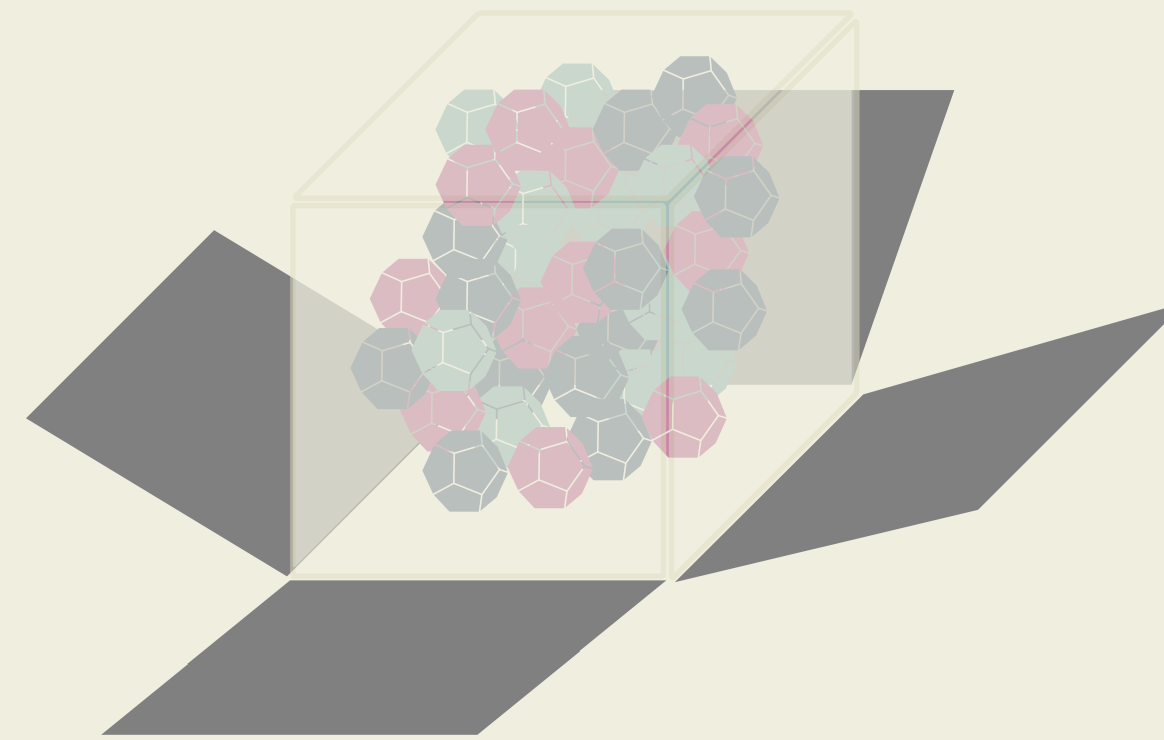
## Verification

NFM 2023



## Training

CIKM 2021



## Interpretability

VMCAI 2024



## Verification

NFM 2023





Using formal methods  
for robustness verification

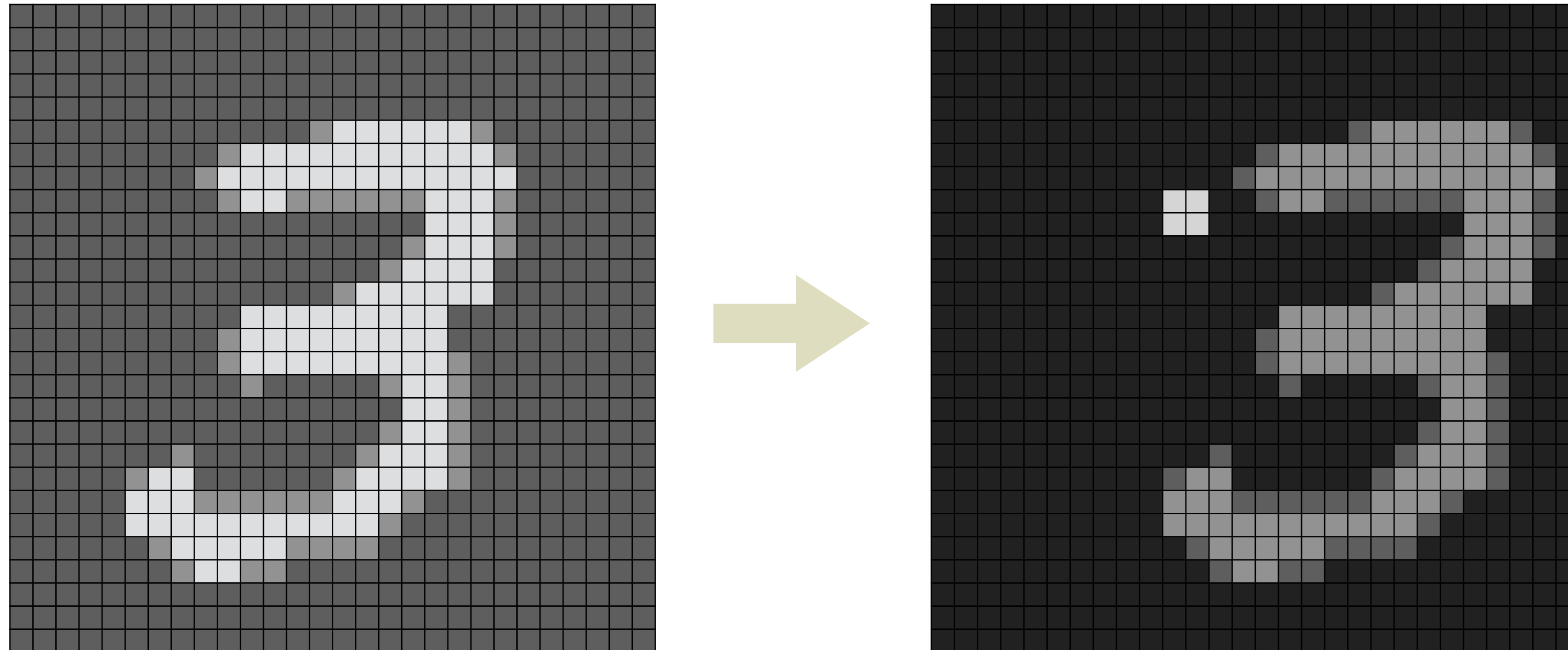


Using formal methods  
**interpretability-aware**  
robustness verification

# Local Robustness Verification

## Combinations of Semantic Perturbations

- Brightness Change
- Patch Placement
- Object Translation

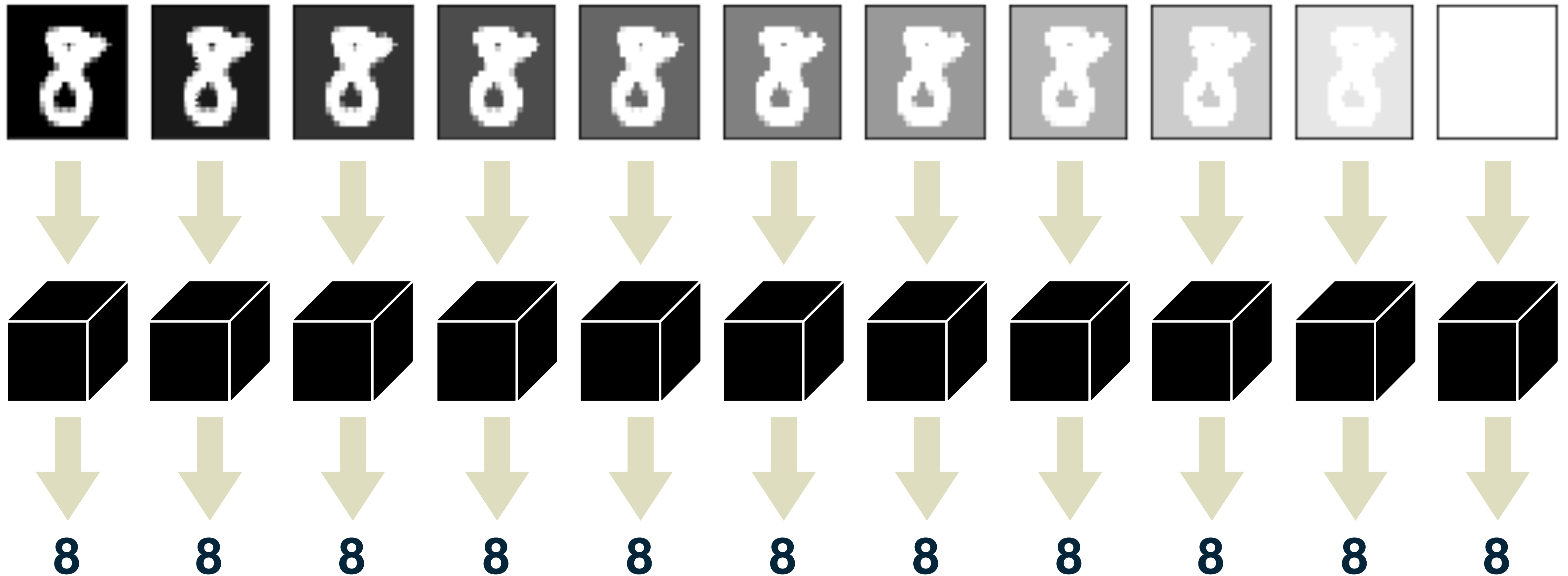


**Goal: Identifying Safe Ranges of Perturbation Parameters**



# Local Robustness Verification

Classification Robustness is **NOT ENOUGH!**

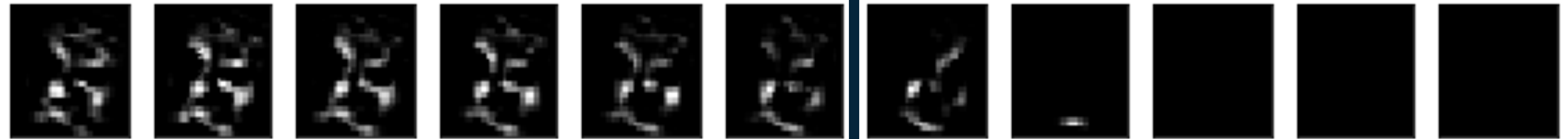


# Local Robustness Verification

Input Image

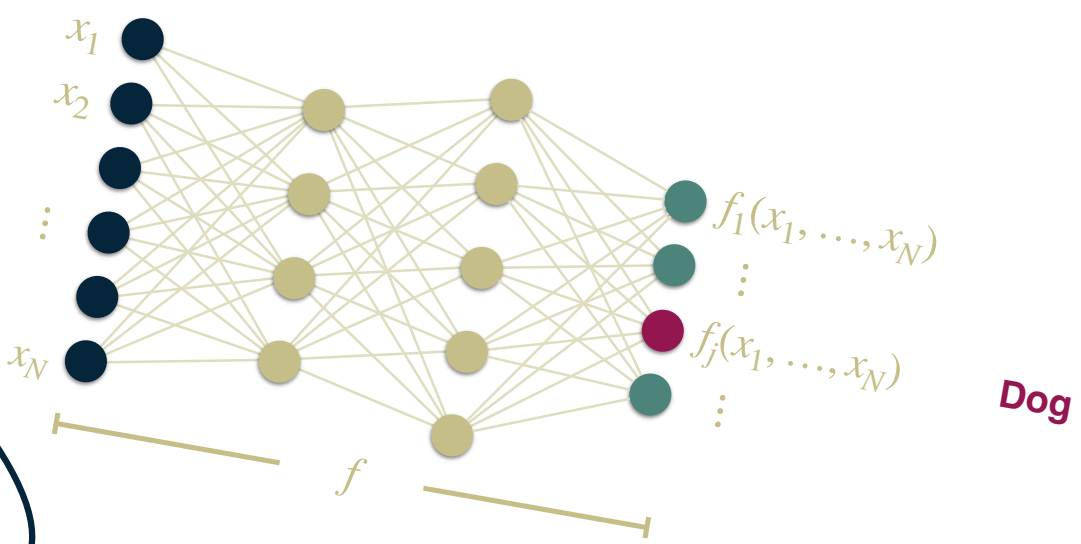
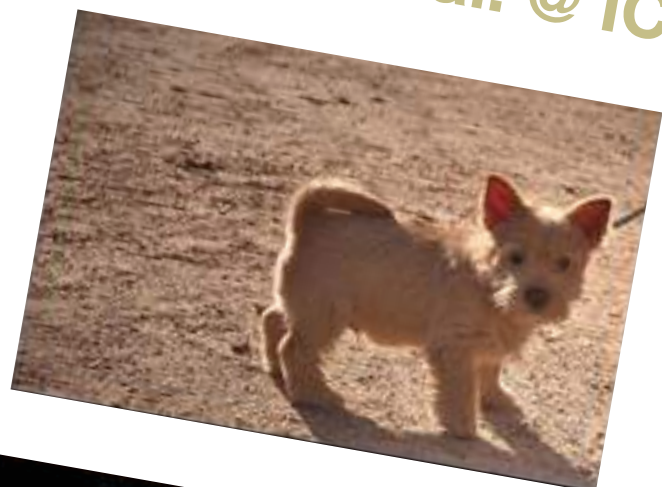


Saliency Map



## Saliency Maps

Simonyan & al. @ ICLR 2014

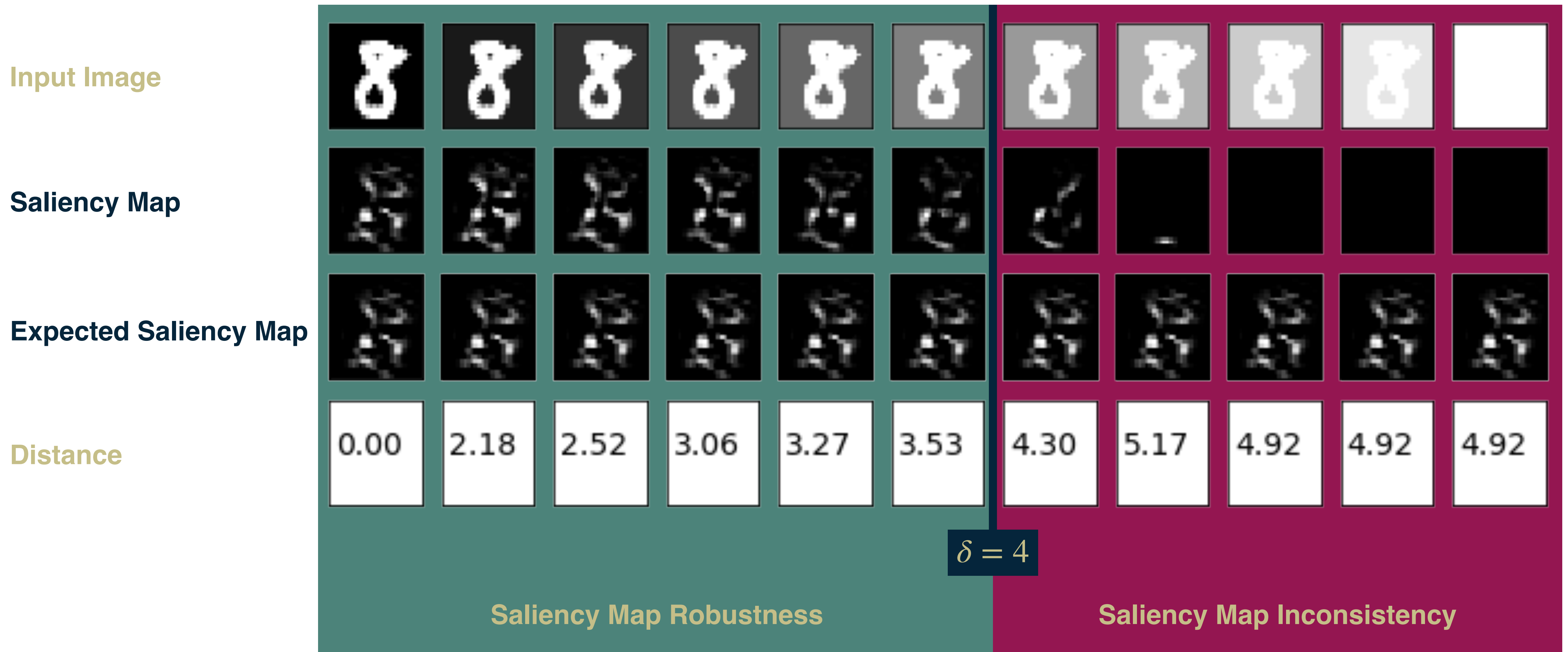


$$\text{map}_j(x) = \left| \frac{\partial f_j(x_1, \dots, x_N)}{\partial x_1} \right| \dots \left| \frac{\partial f_j(x_1, \dots, x_N)}{\partial x_N} \right|$$

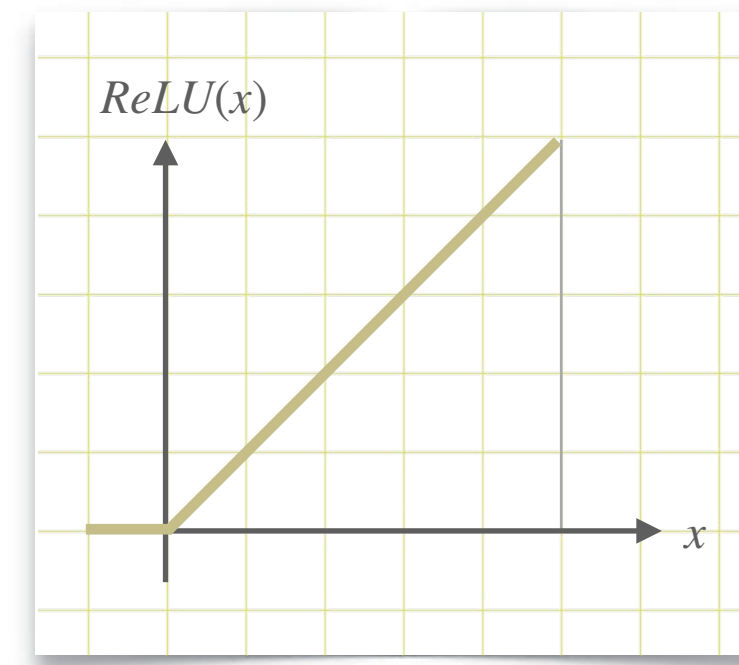


# Local Robustness Verification

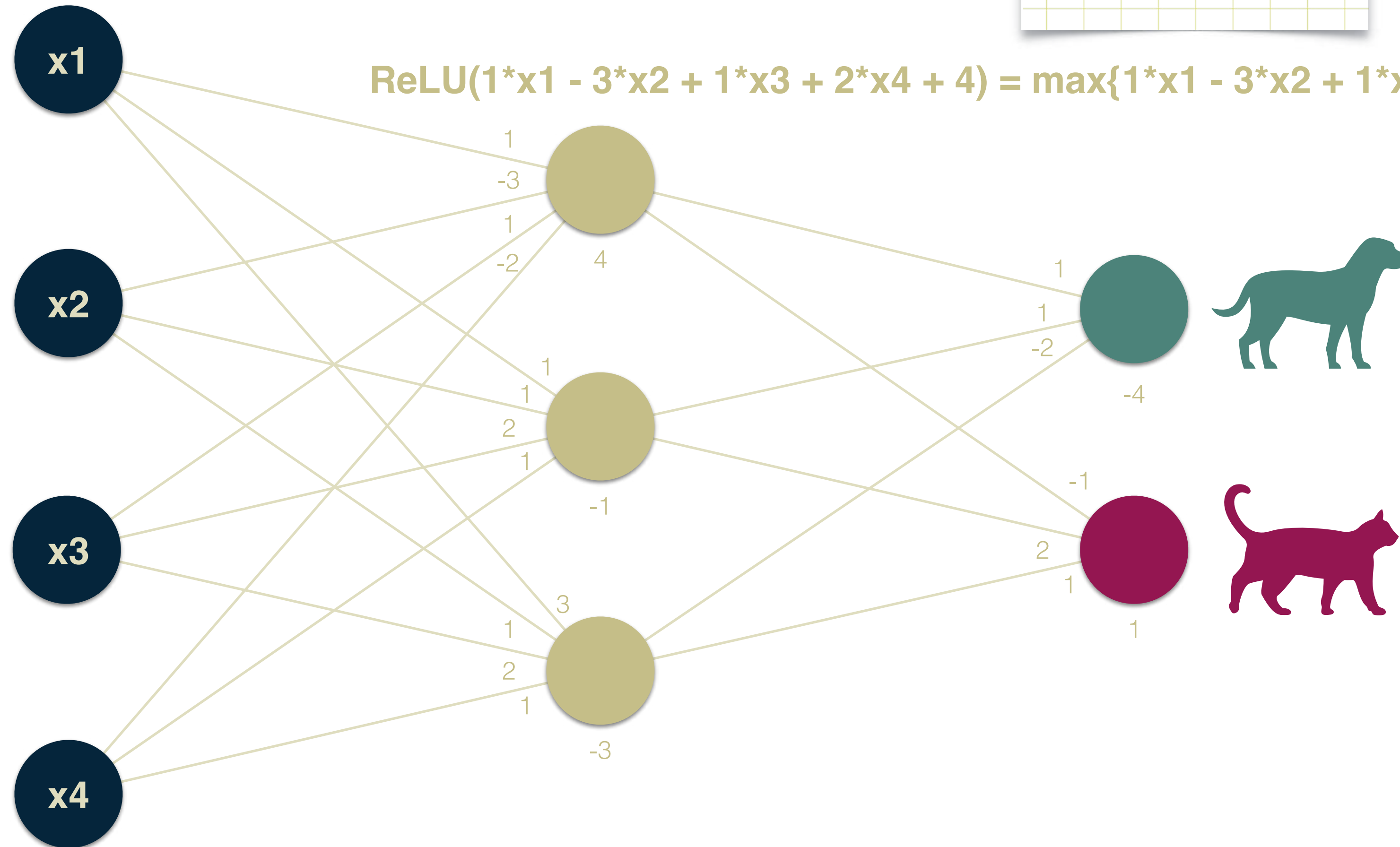
## Saliency Map Robustness



# (A Very Small) Example



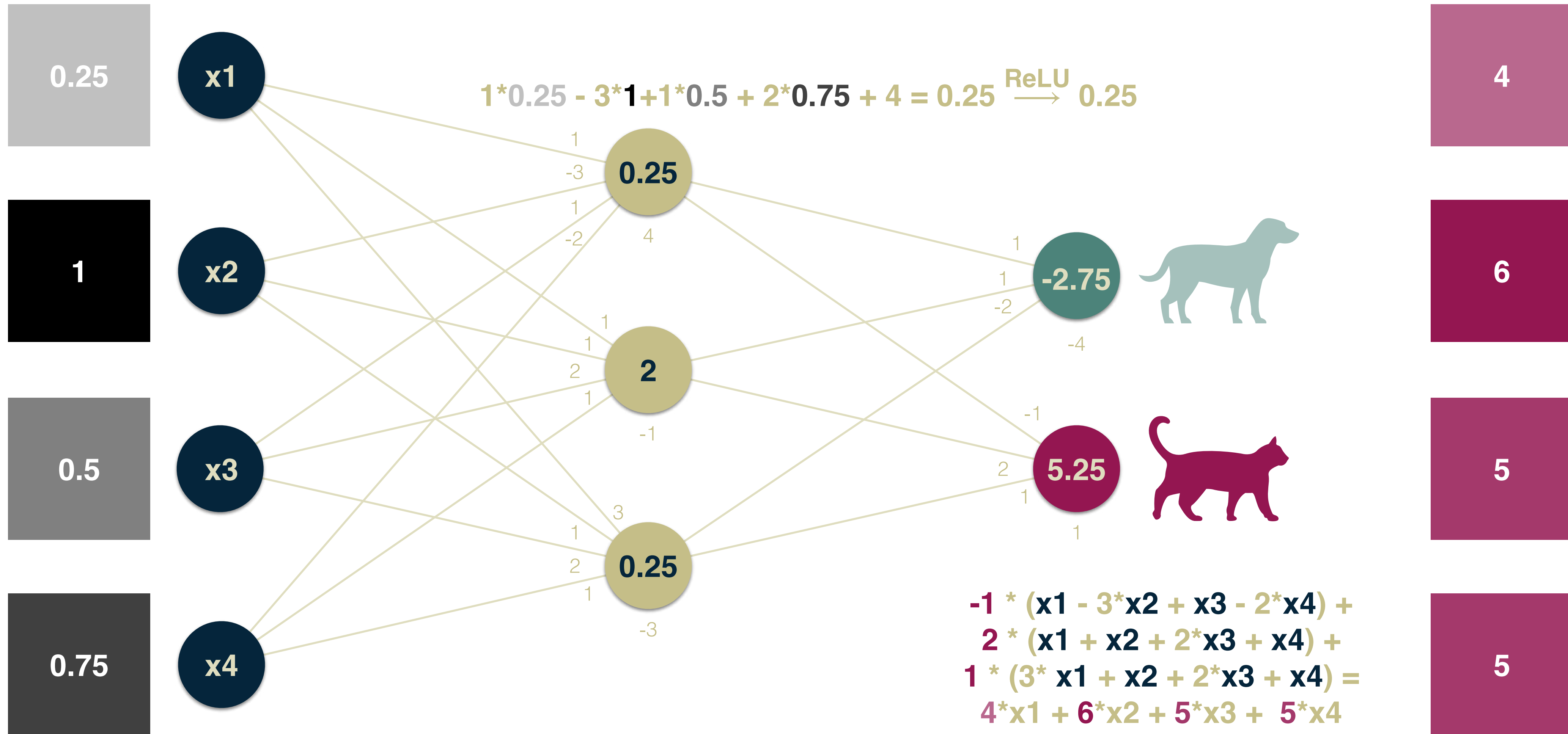
$$ReLU(1 \cdot x_1 - 3 \cdot x_2 + 1 \cdot x_3 + 2 \cdot x_4 + 4) = \max\{1 \cdot x_1 - 3 \cdot x_2 + 1 \cdot x_3 + 2 \cdot x_4 + 4, 0\}$$





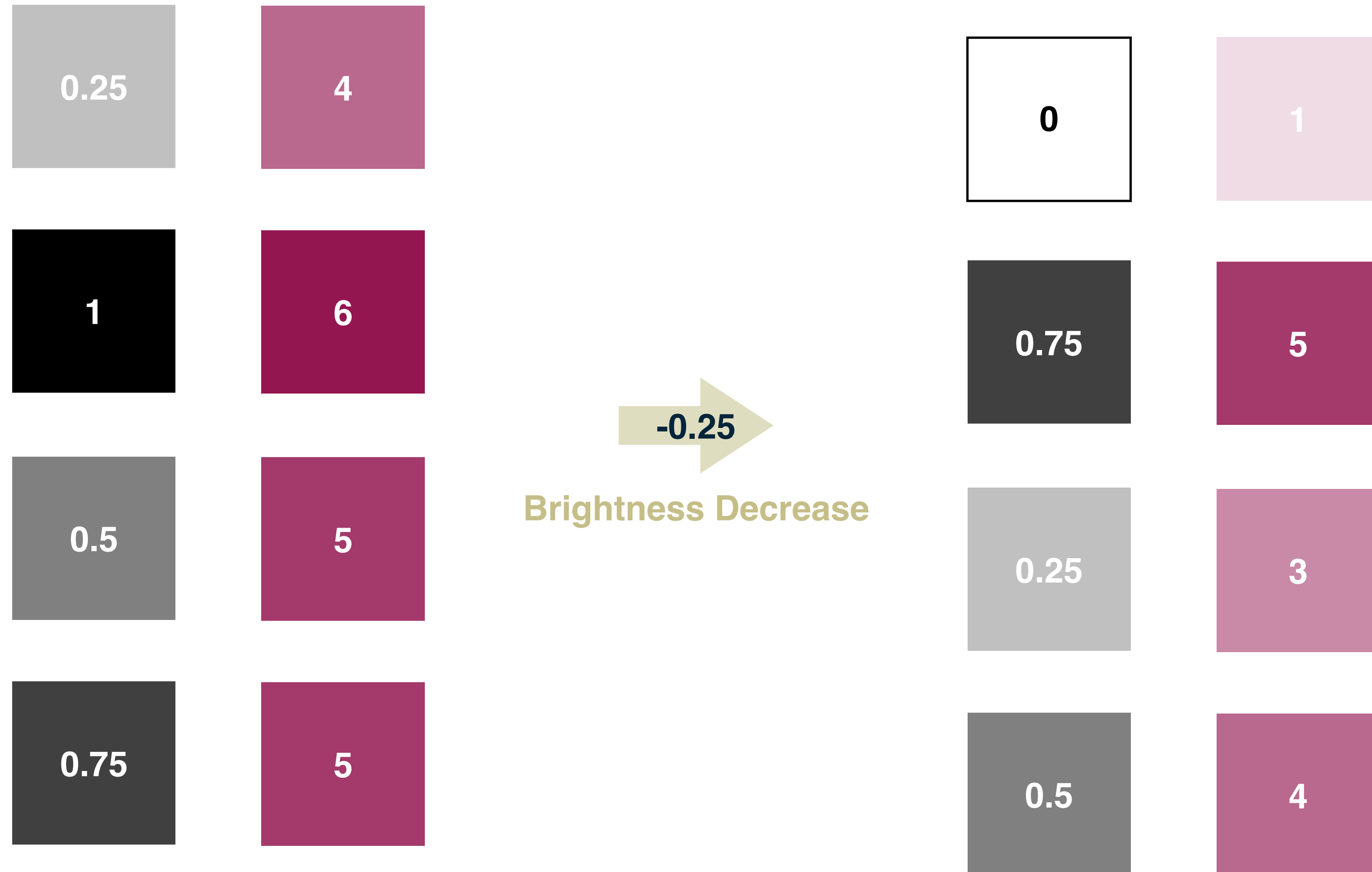
# (A Very Small) Example

## Saliency Maps



# (A Very Small) Example

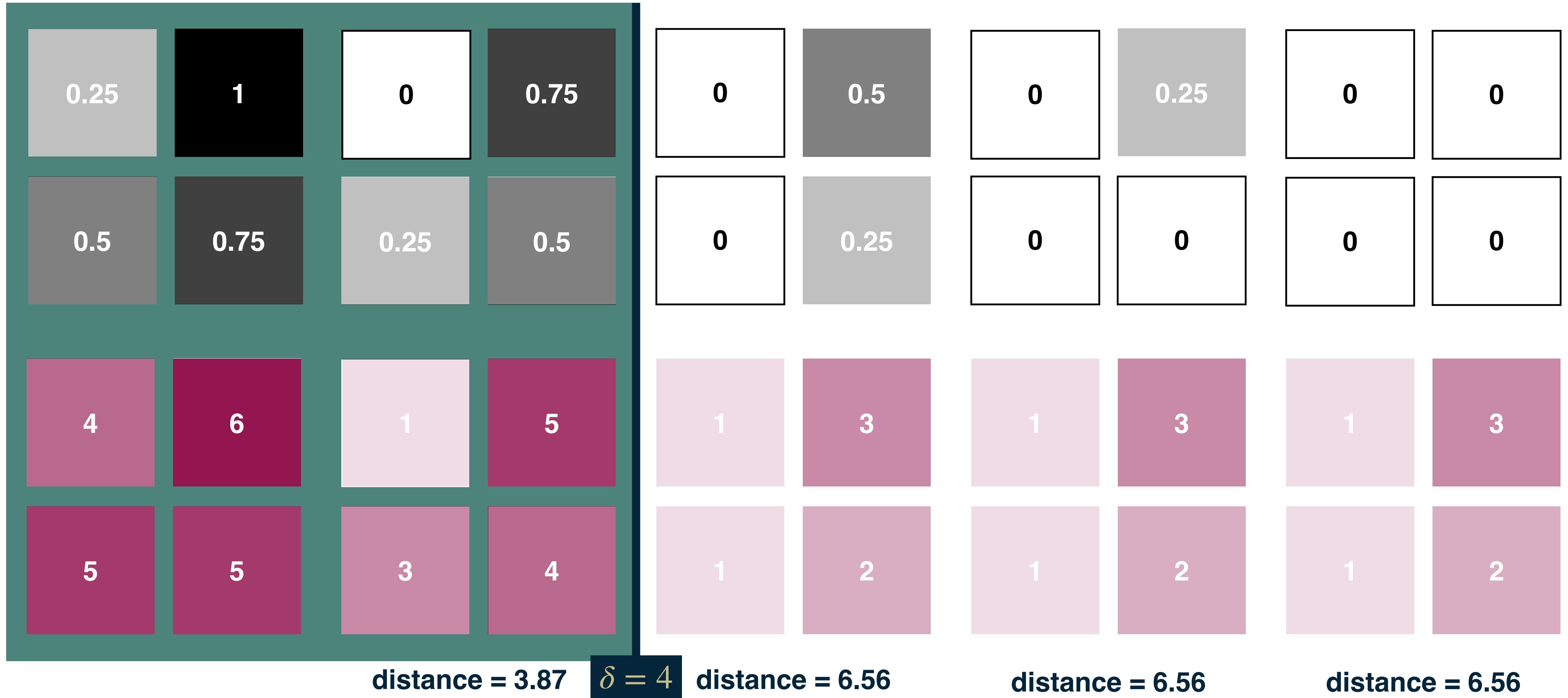
## Semantic Perturbations





# (A Very Small) Example

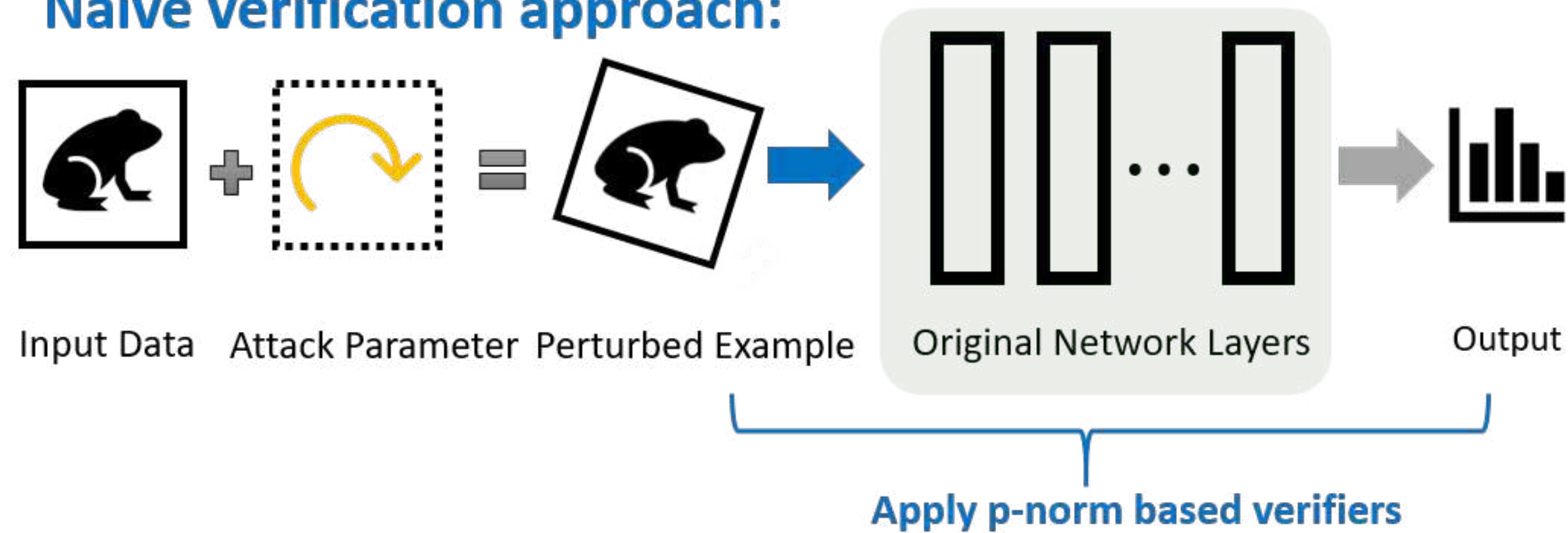
## Saliency Map Robustness



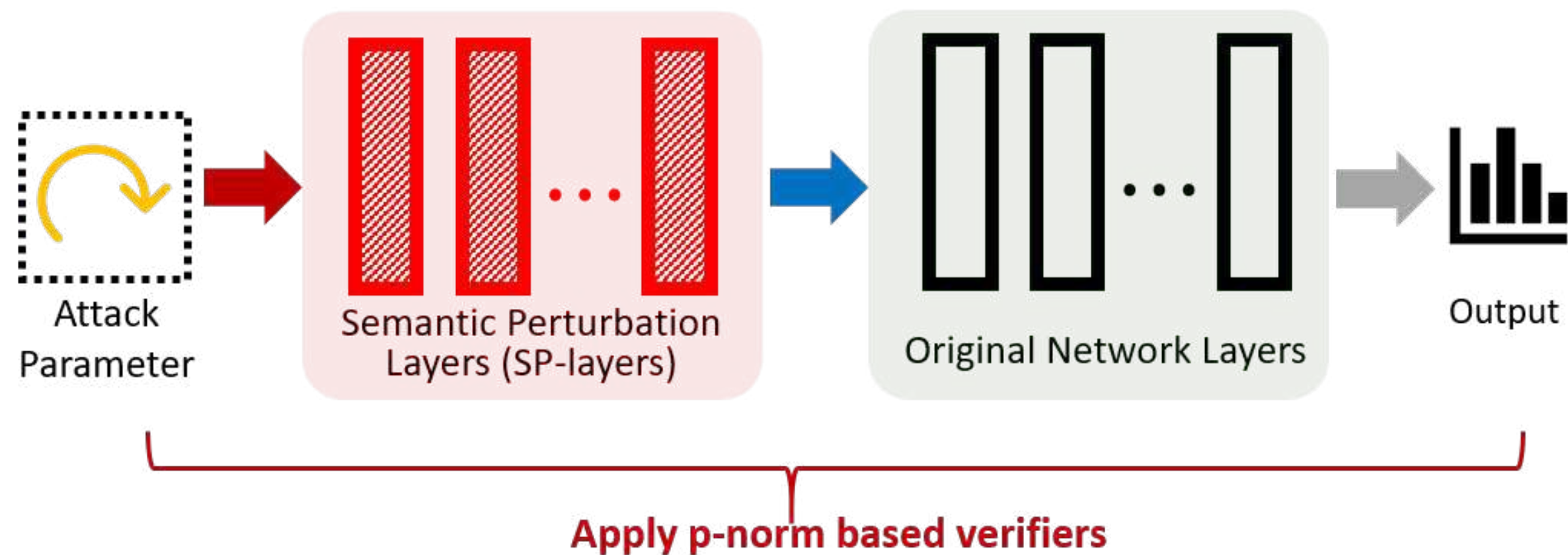
# Encoding Semantic Perturbations

Mohapatra & al. @ CVPR 2020

- Naïve verification approach:



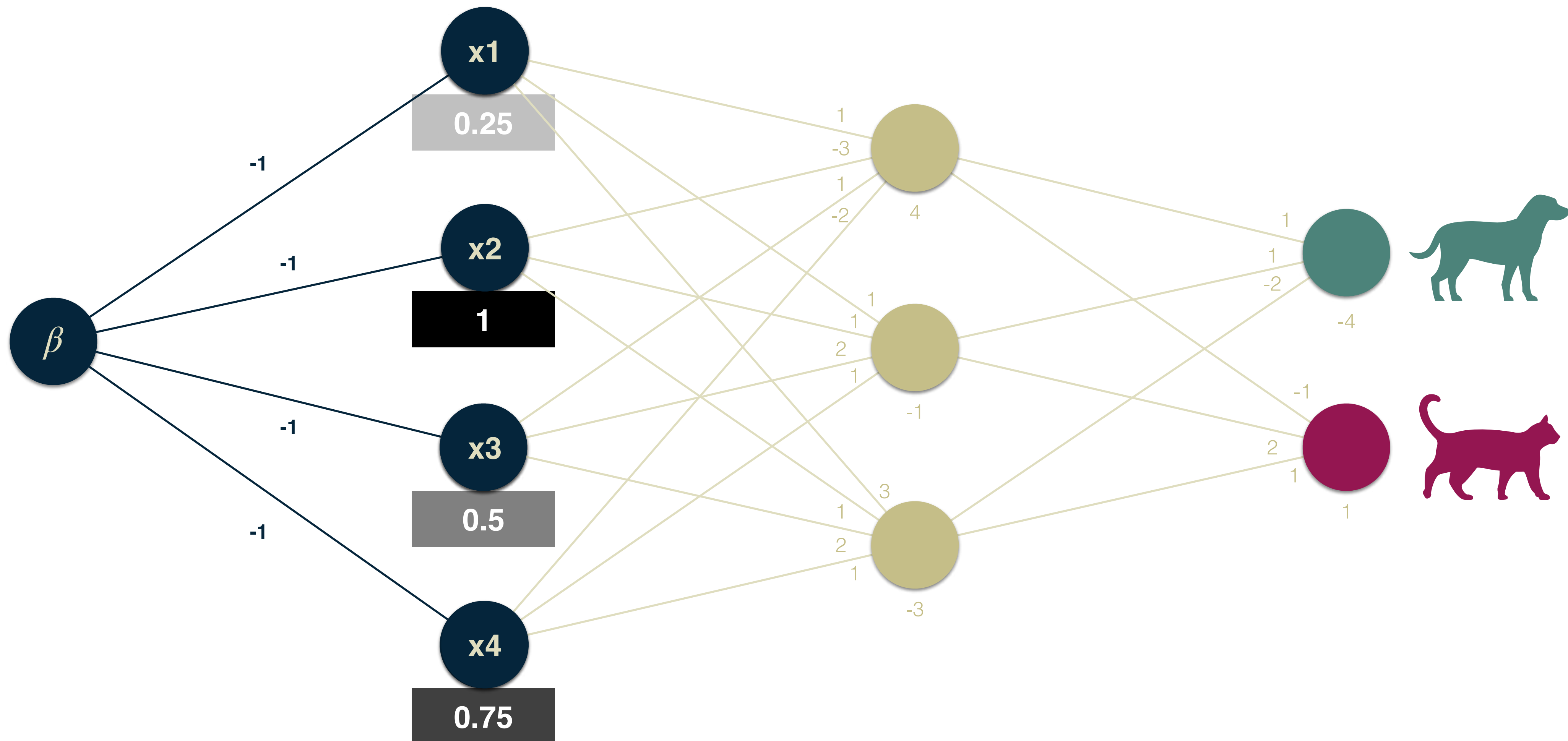
- Our Semantify-NN:





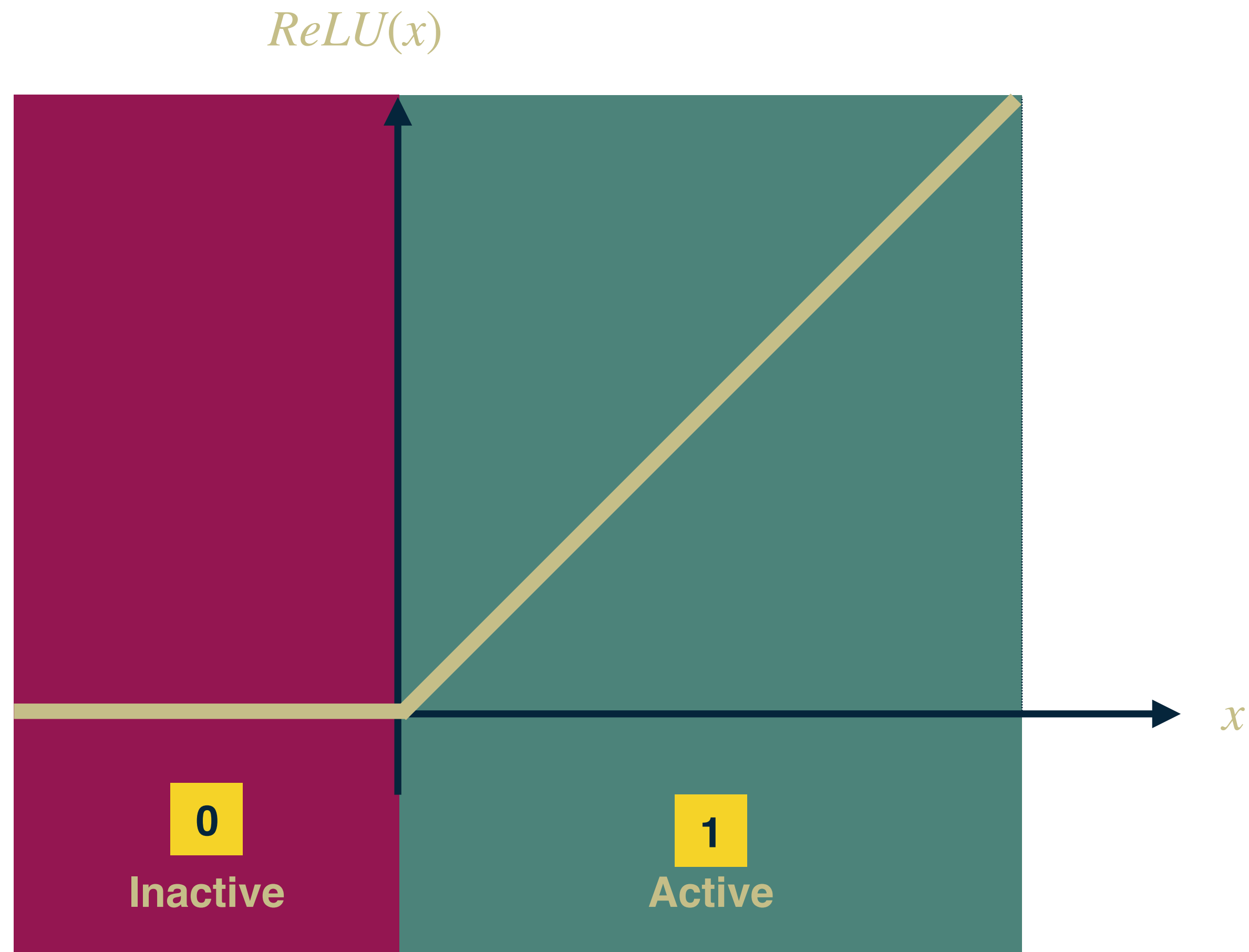
# (A Very Small) Example

## Encoding Semantic Perturbations



# Naïve Breadth-First Search

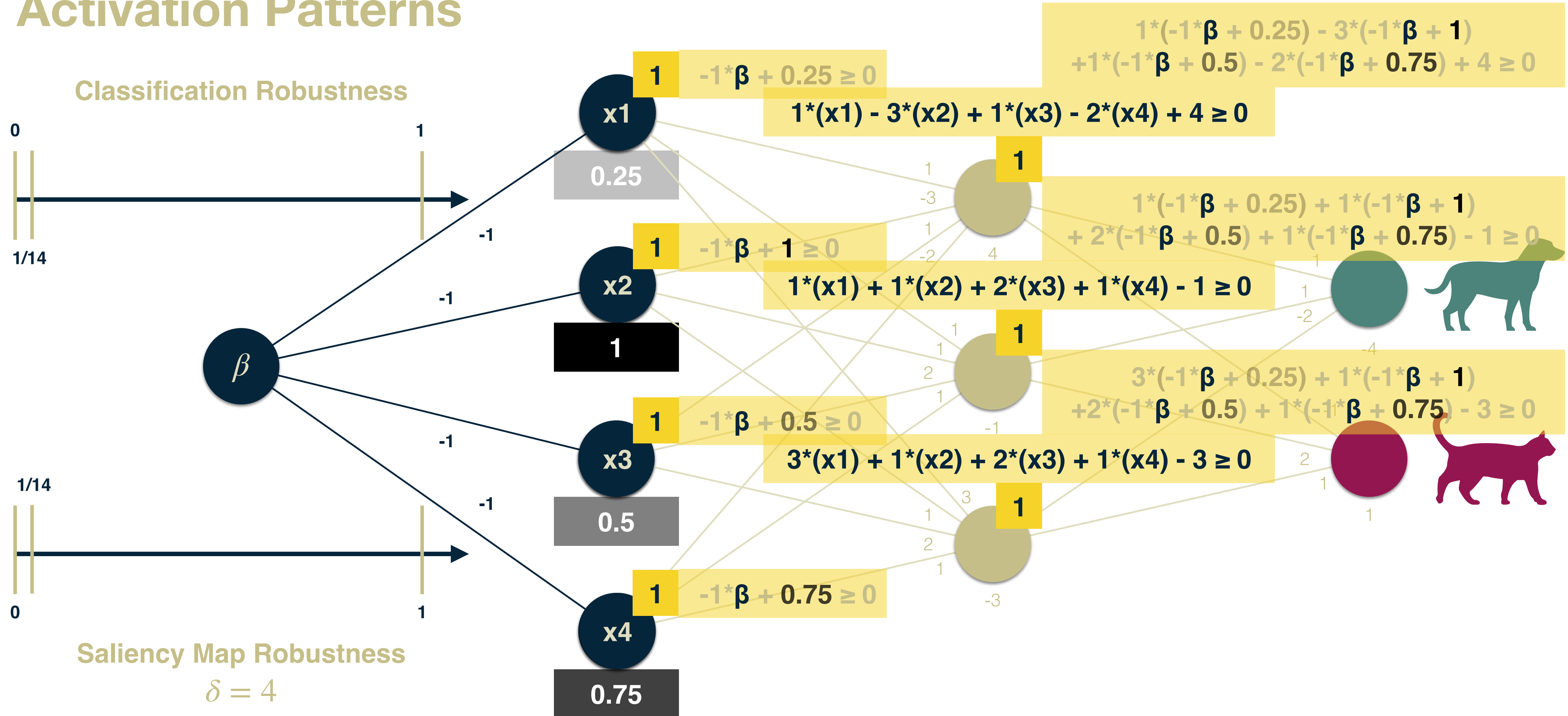
## Activation Patterns





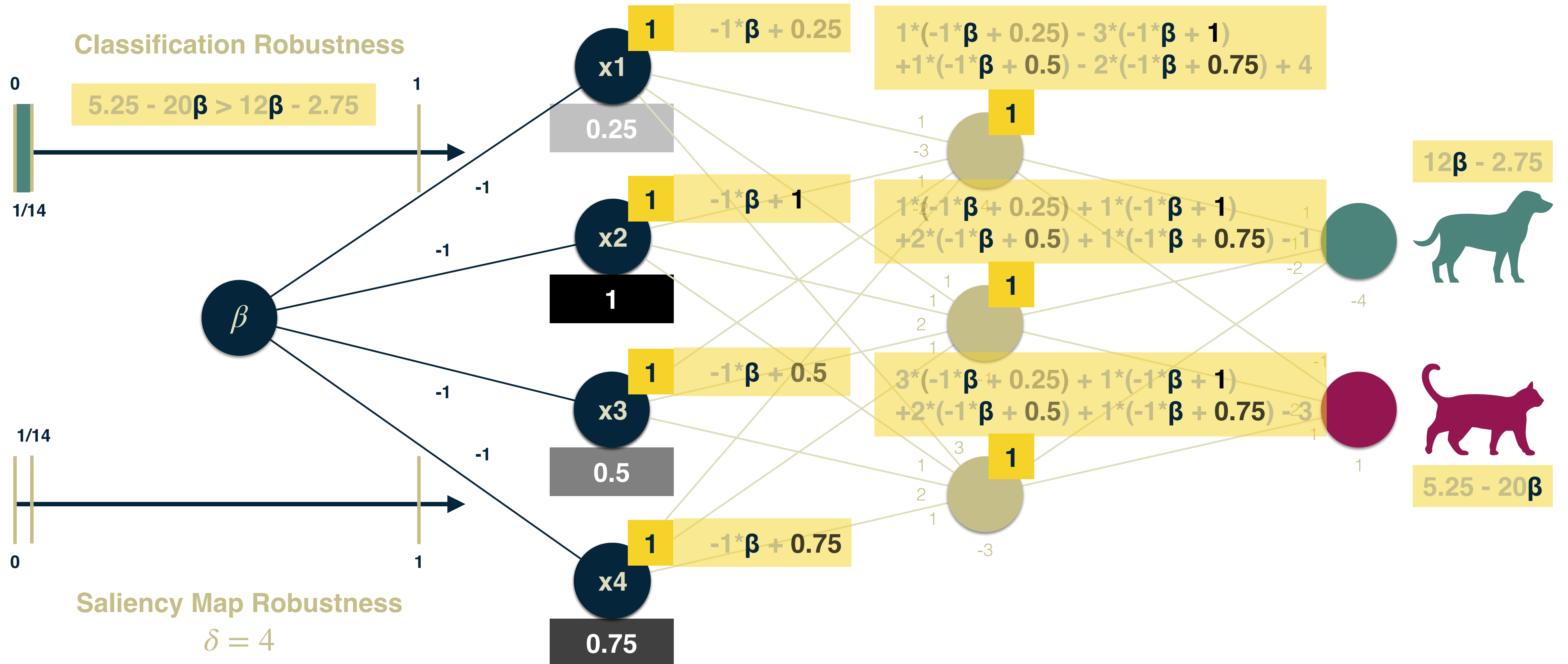
# (A Very Small) Example

## Activation Patterns



# (A Very Small) Example

## Classification Robustness





# (A Very Small) Example

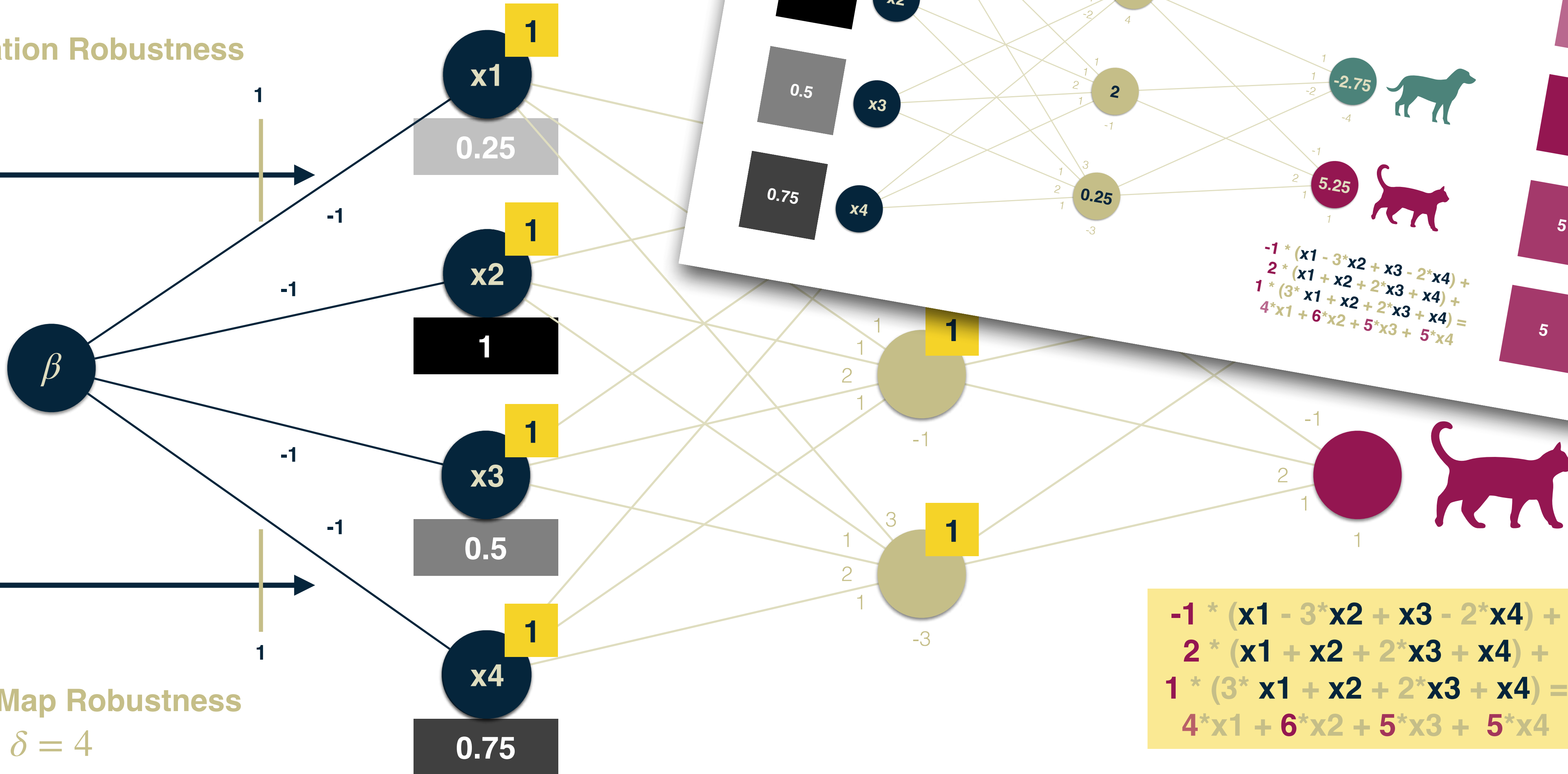
## Saliency Map Robustness

Classification Robustness



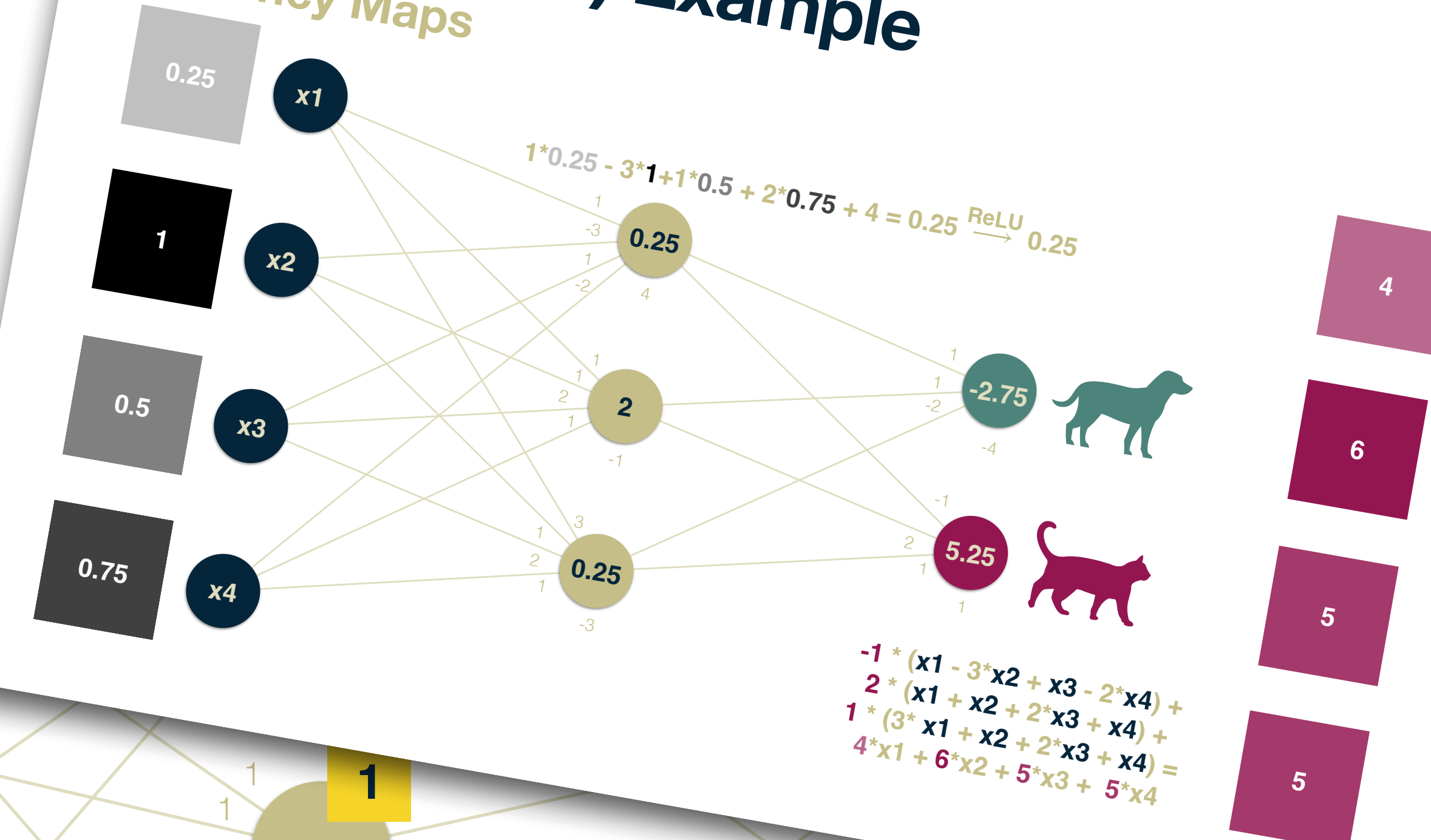
Saliency Map Robustness

$$\delta = 4$$



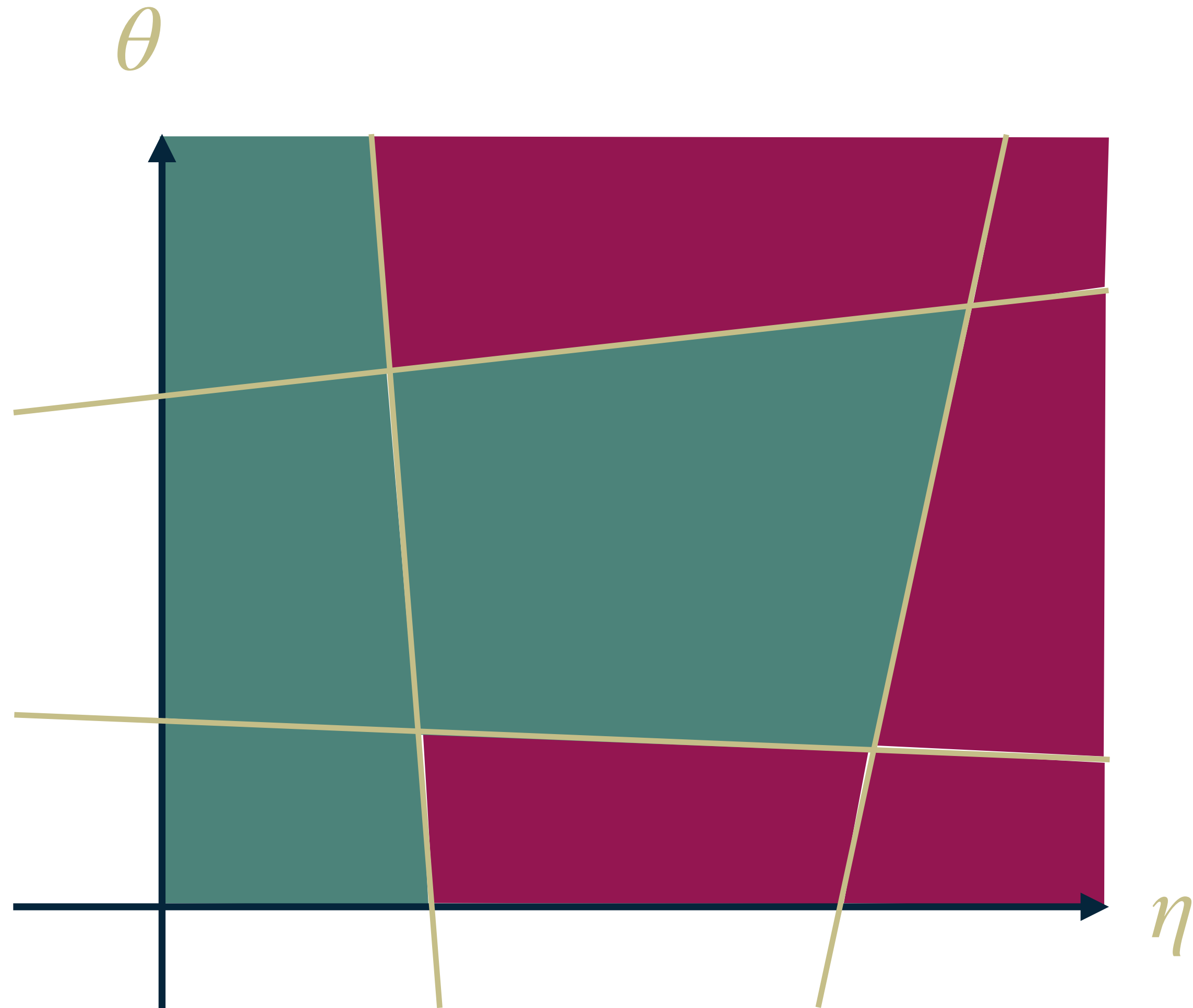
### (A Very Small) Example

#### Saliency Maps



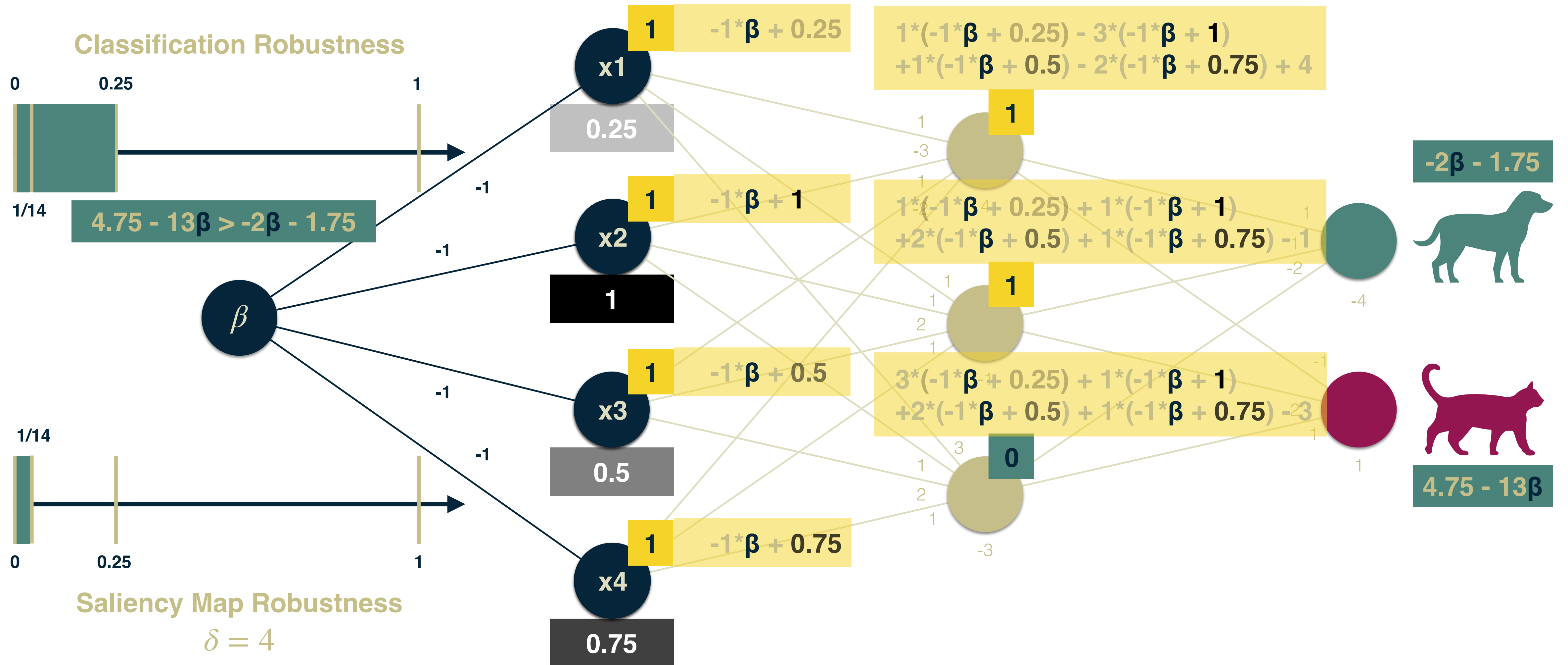
$$-1 \cdot (x_1 - 3 \cdot x_2 + x_3 - 2 \cdot x_4) + 2 \cdot (x_1 + x_2 + 2 \cdot x_3 + x_4) + 1 \cdot (3 \cdot x_1 + x_2 + 2 \cdot x_3 + x_4) = 4 \cdot x_1 + 6 \cdot x_2 + 5 \cdot x_3 + 5 \cdot x_4$$

# Naïve Breadth-First Search



# (A Very Small) Example

## Classification Robustness

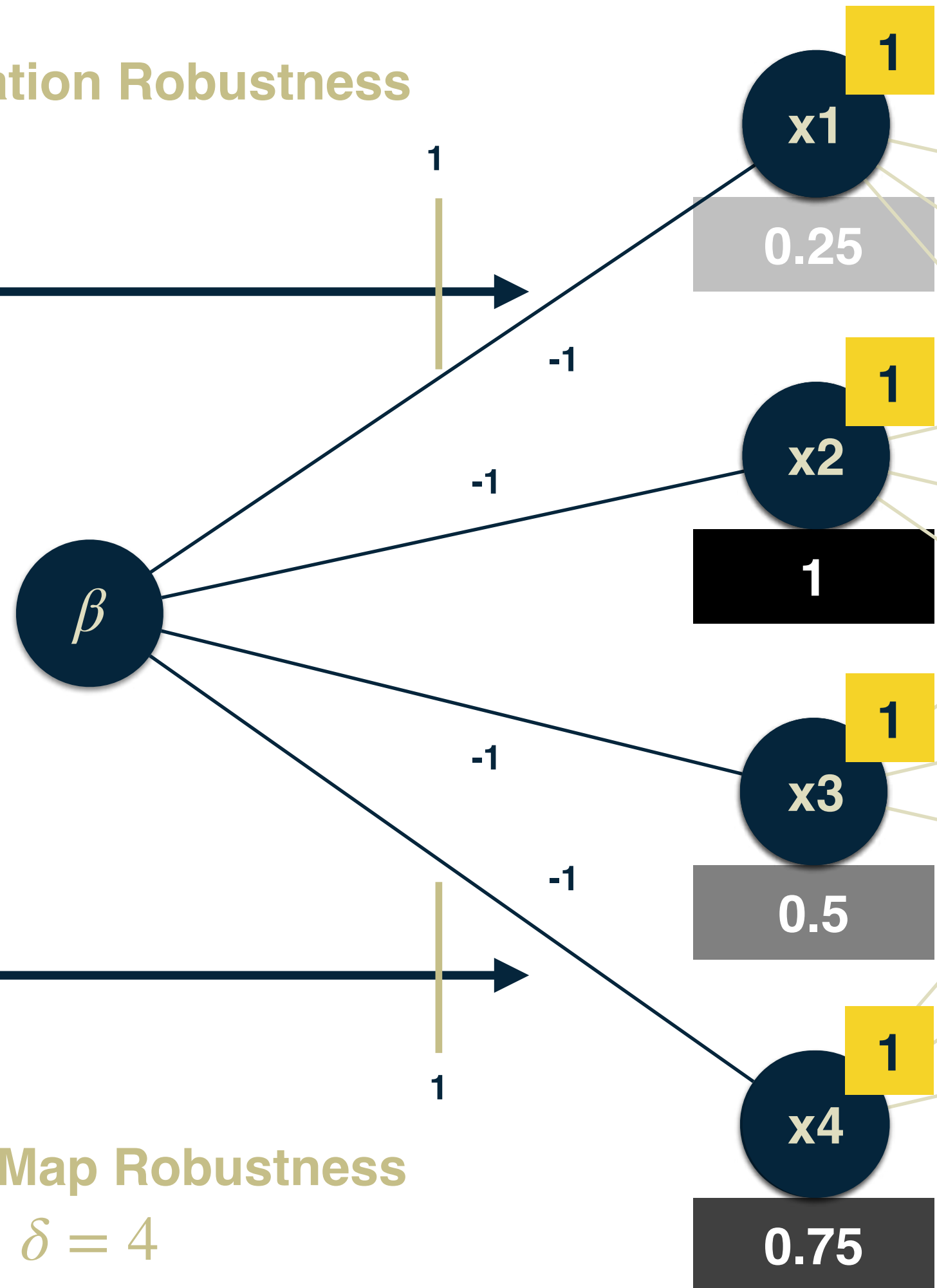
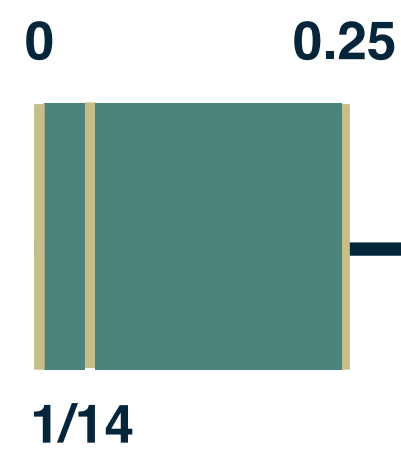




# (A Very Small) Example

## Saliency Map Robustness

Classification Robustness



Saliency Map Robustness

$$\delta = 4$$

### (A Very Small) Example

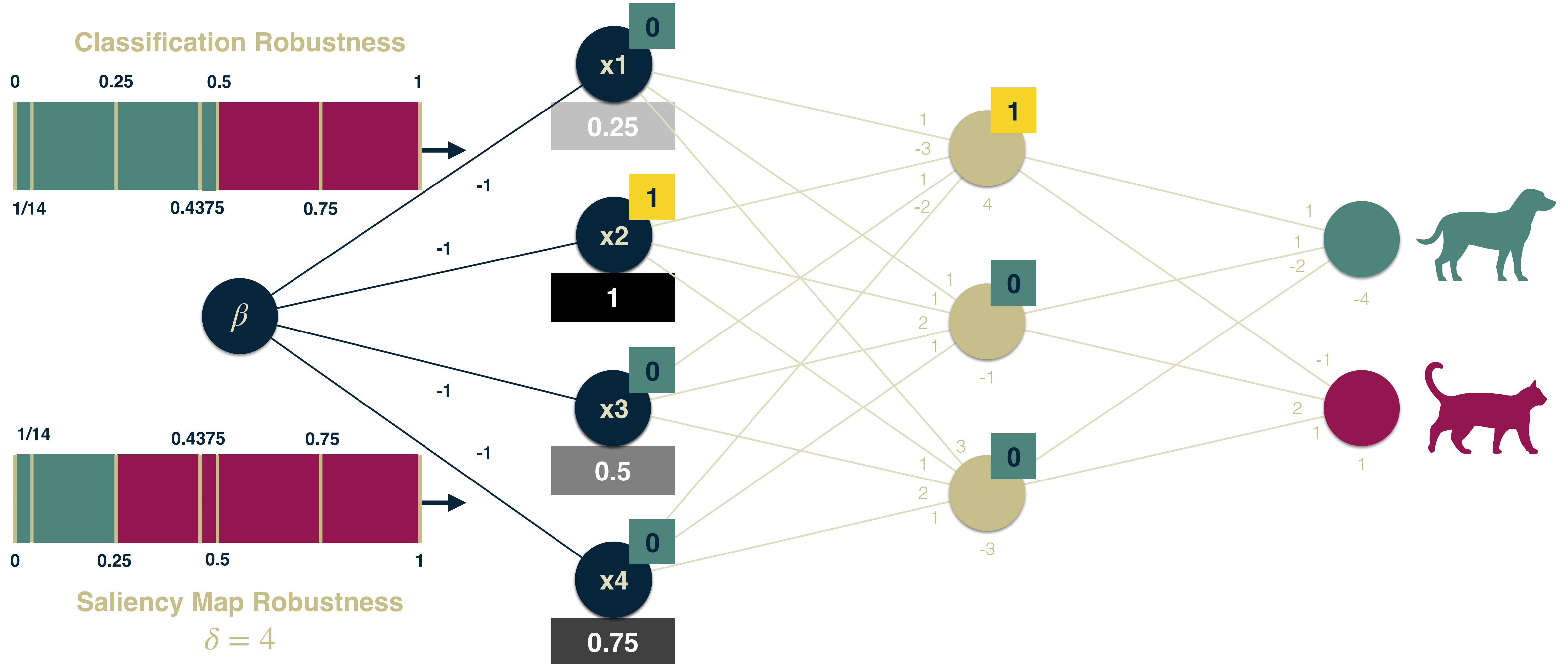
#### Saliency Maps

A neural network diagram showing saliency maps. The input nodes  $x_1, x_2, x_3, x_4$  have saliency maps with values 0.25, 1, 0.5, and 0.75 respectively. The hidden nodes have saliency maps with values 0.25, 2, and 0.25. The output nodes have saliency maps with values -2.75 and 5.25. There are icons of a dog and a cat next to the output nodes. A calculation is shown:  $1 \cdot 0.25 - 3 \cdot 1 + 1 \cdot 0.5 + 2 \cdot 0.75 + 4 = 0.25 \xrightarrow{\text{ReLU}} 0.25$ . Another calculation is shown:  $-1 \cdot (x_1 - 3 \cdot x_2 + x_3 - 2 \cdot x_4) + 2 \cdot (x_1 + x_2 + 2 \cdot x_3 + x_4) + 1 \cdot (3 \cdot x_1 + x_2 + 2 \cdot x_3 + x_4) = 4 \cdot x_1 + 6 \cdot x_2 + 5 \cdot x_3 + 5 \cdot x_4$ . There are also icons of a dog and a cat next to the output nodes.

$$\begin{aligned}
 & -1 \cdot (x_1 - 3 \cdot x_2 + x_3 - 2 \cdot x_4) + \\
 & 2 \cdot (x_1 + x_2 + 2 \cdot x_3 + x_4) + \\
 & 0 \cdot (3 \cdot x_1 + x_2 + 2 \cdot x_3 + x_4) = \\
 & 1 \cdot x_1 + 5 \cdot x_2 + 3 \cdot x_3 + 4 \cdot x_4
 \end{aligned}$$

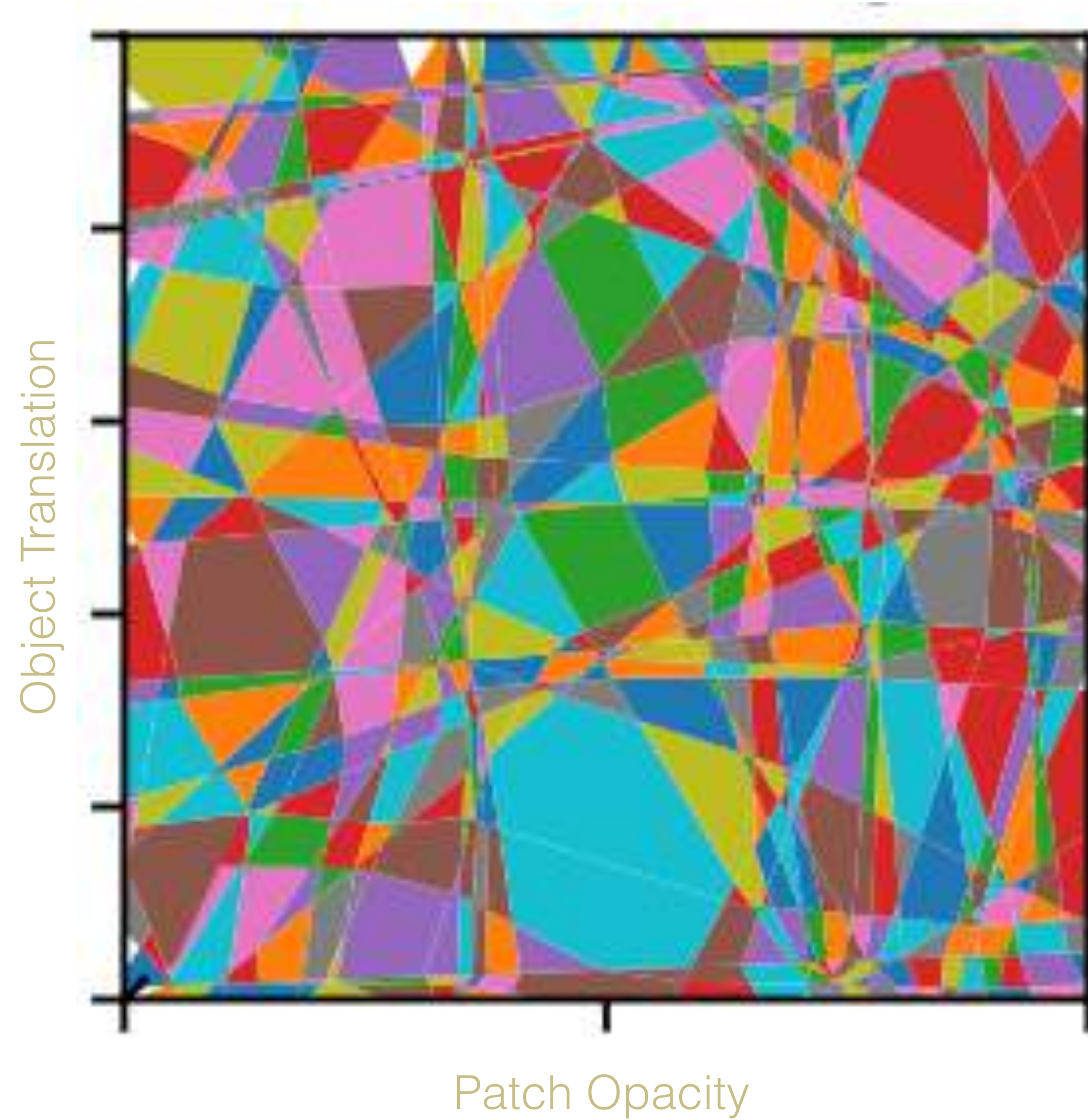
# (A Very Small) Example

## Naïve Breadth-First Search





# Naïve Breadth-First Search

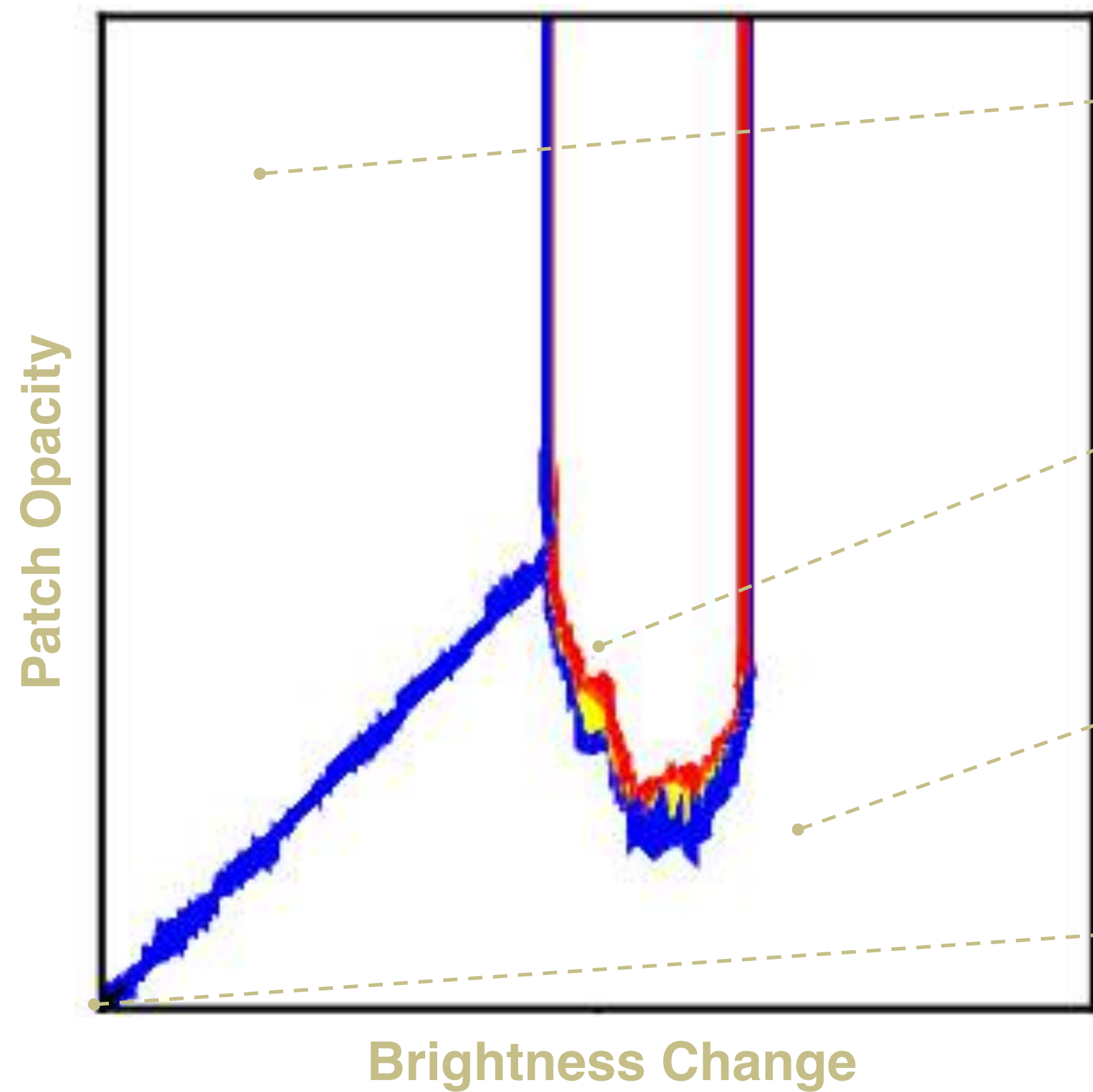


 **TOO MANY ACTIVATION PATTERNS!**

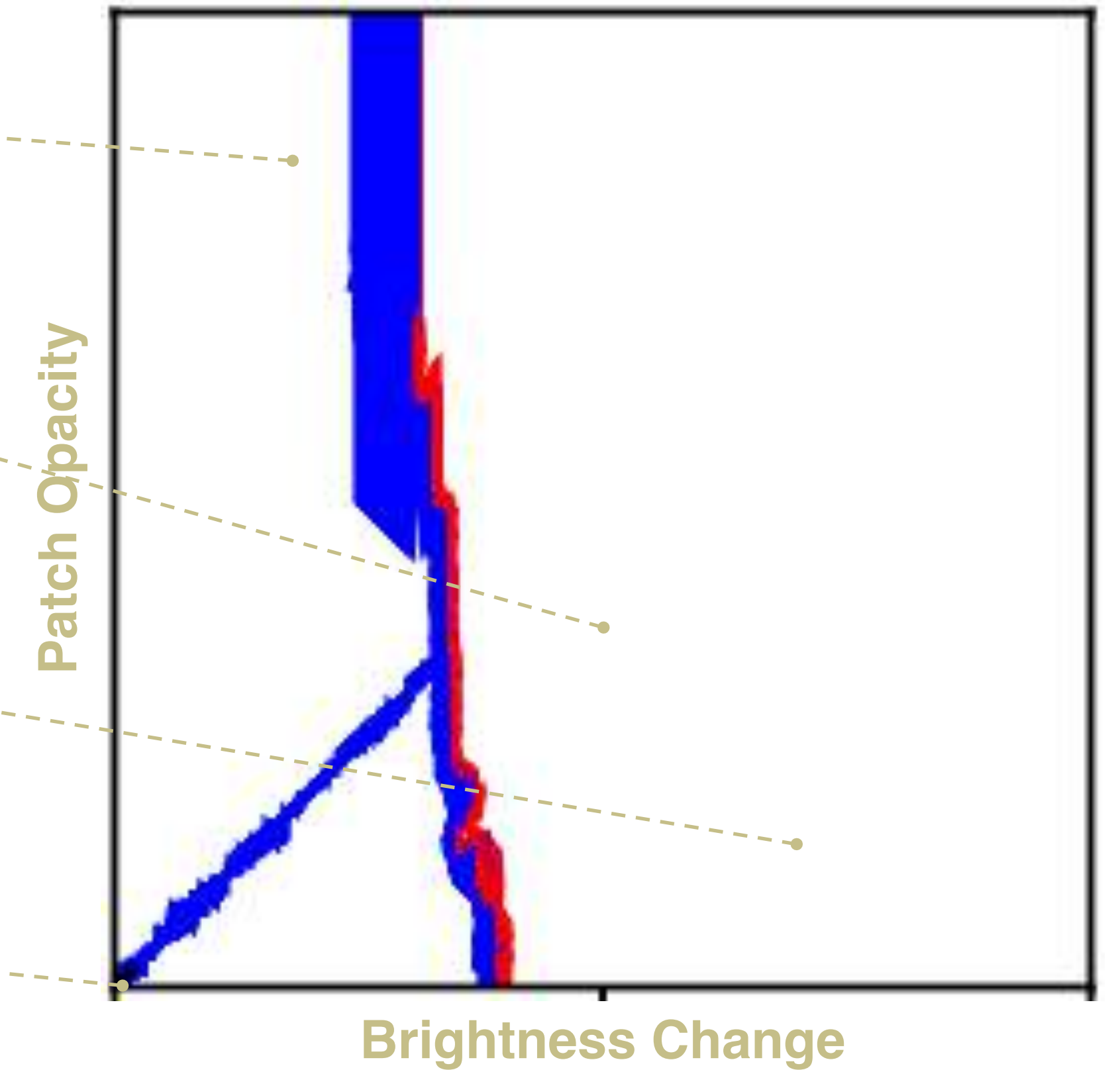


# Geometric Boundary Search

## Classification Robustness

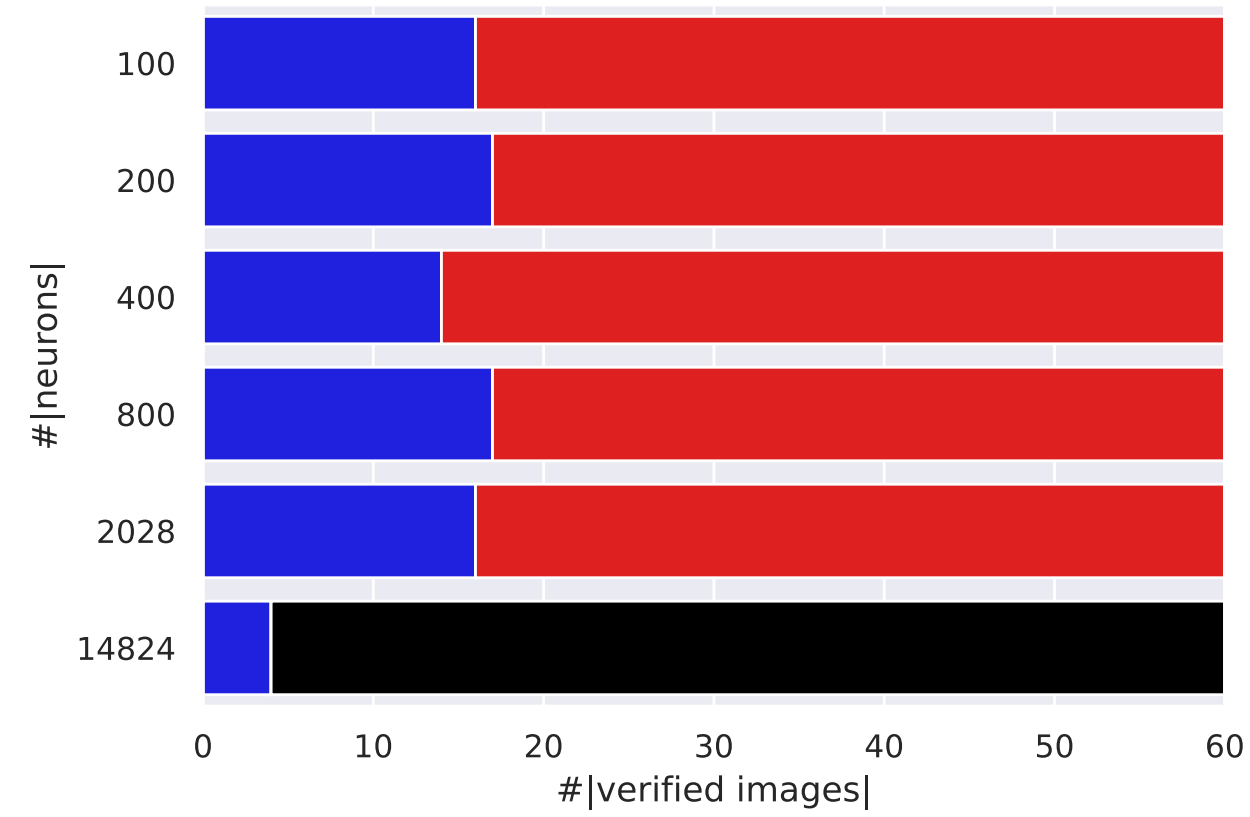


## Saliency Map Robustness

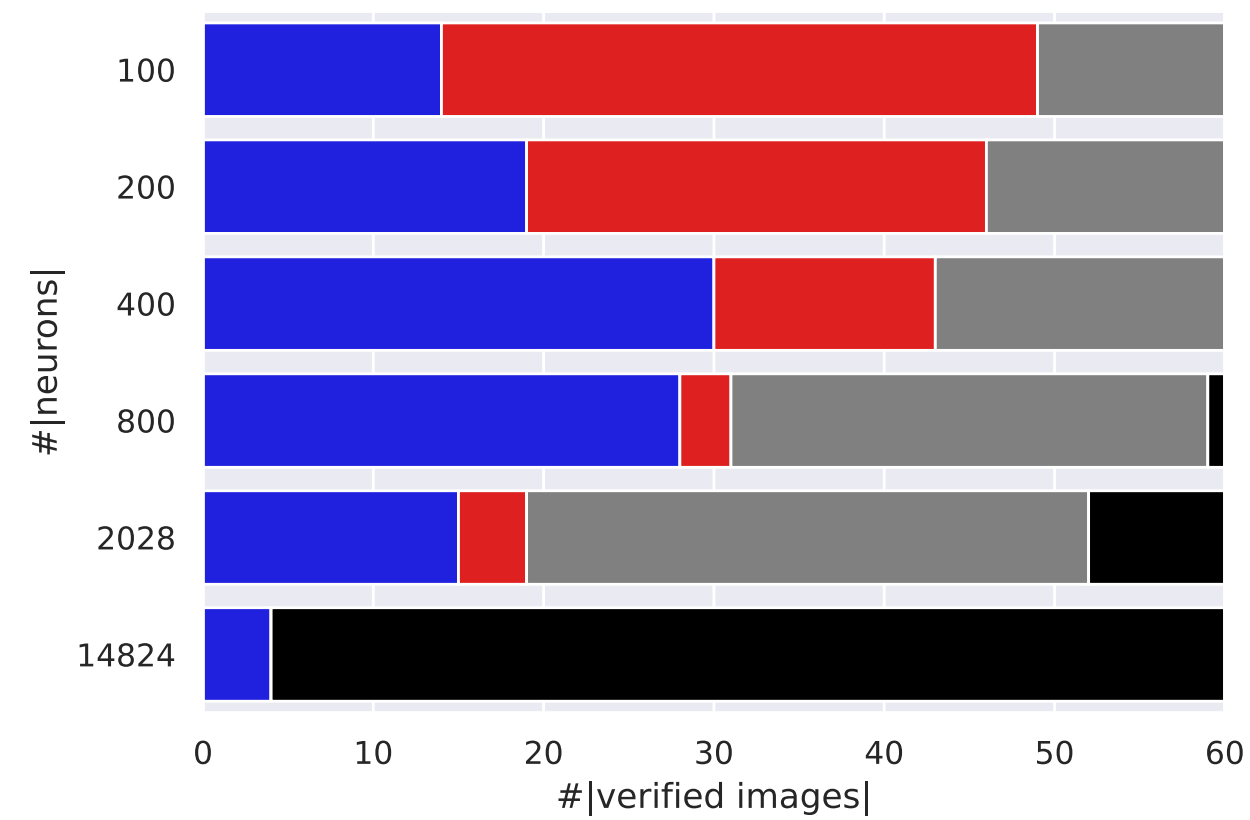


# Geometric Boundary Search

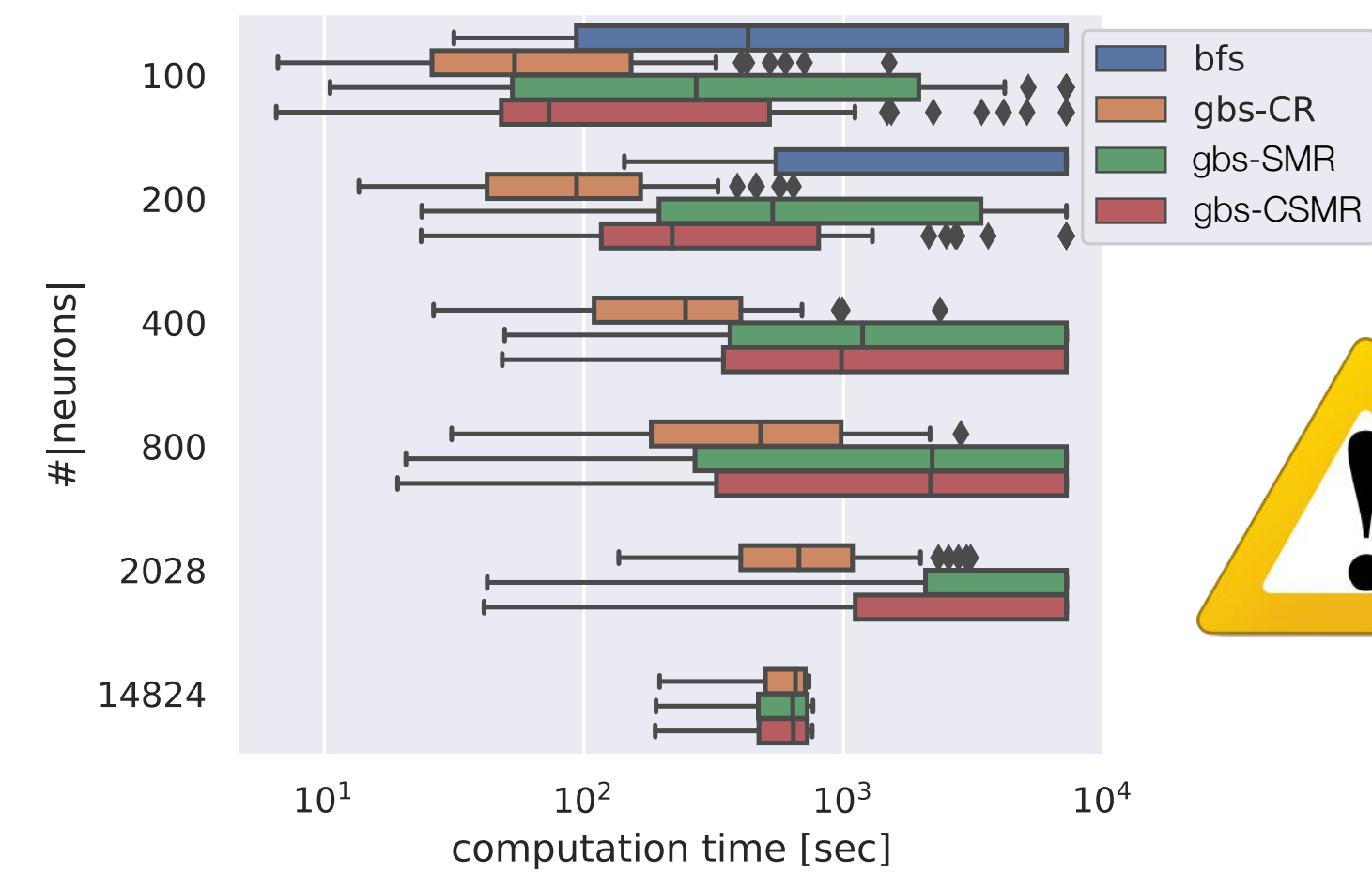
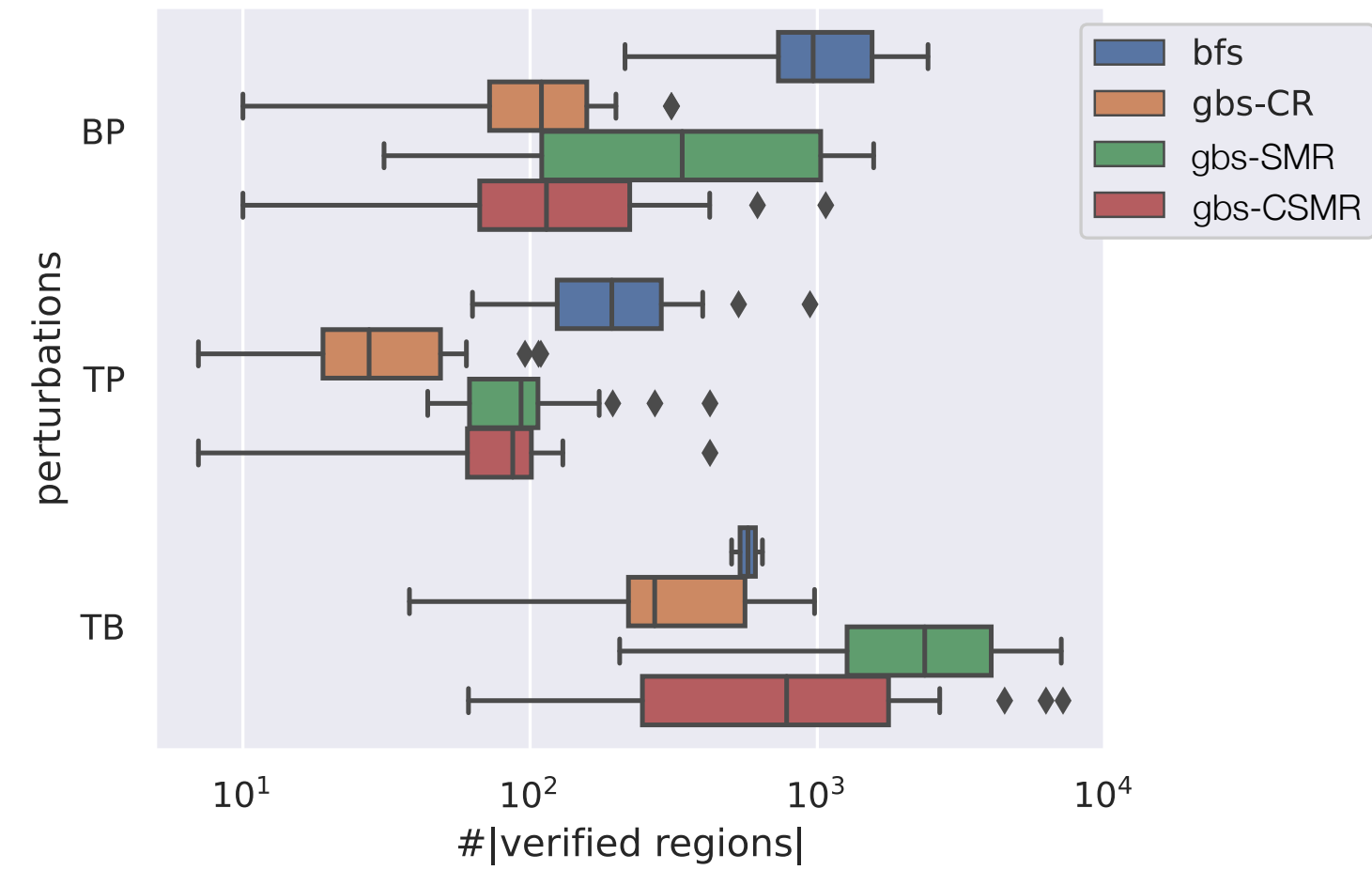
## Experimental Results



### Classification Robustness



### Saliency Map Robustness

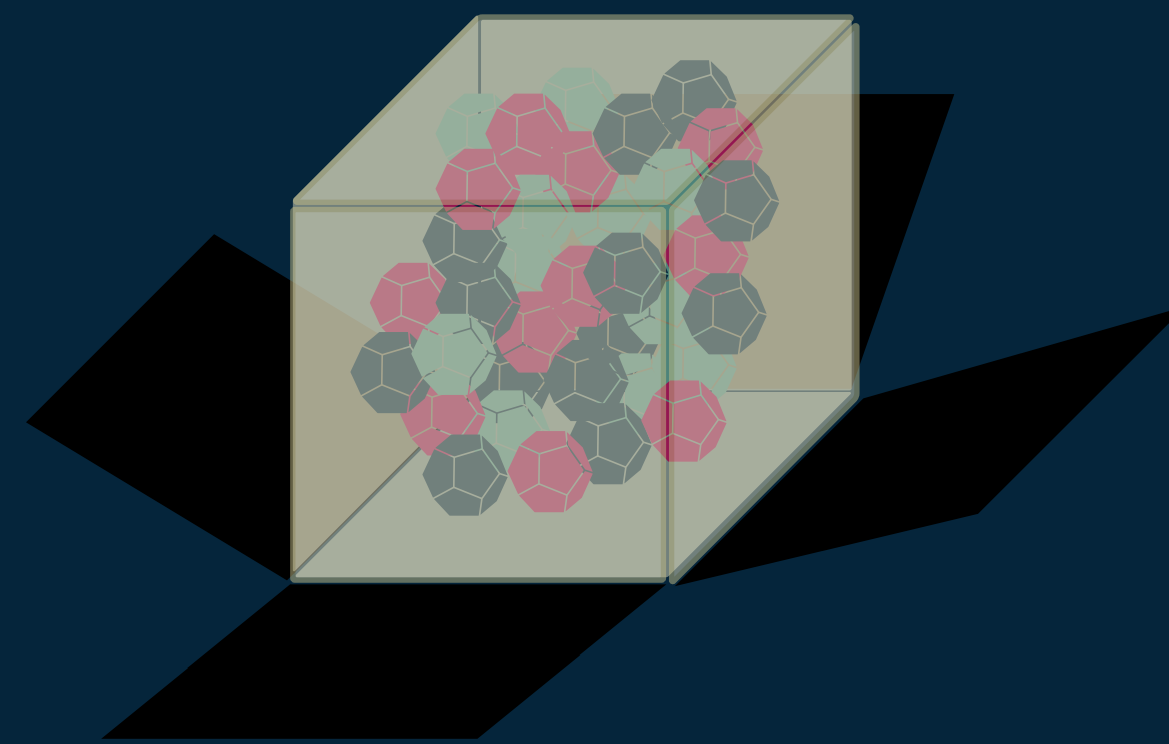


**EXPONENTIAL INCREASE**



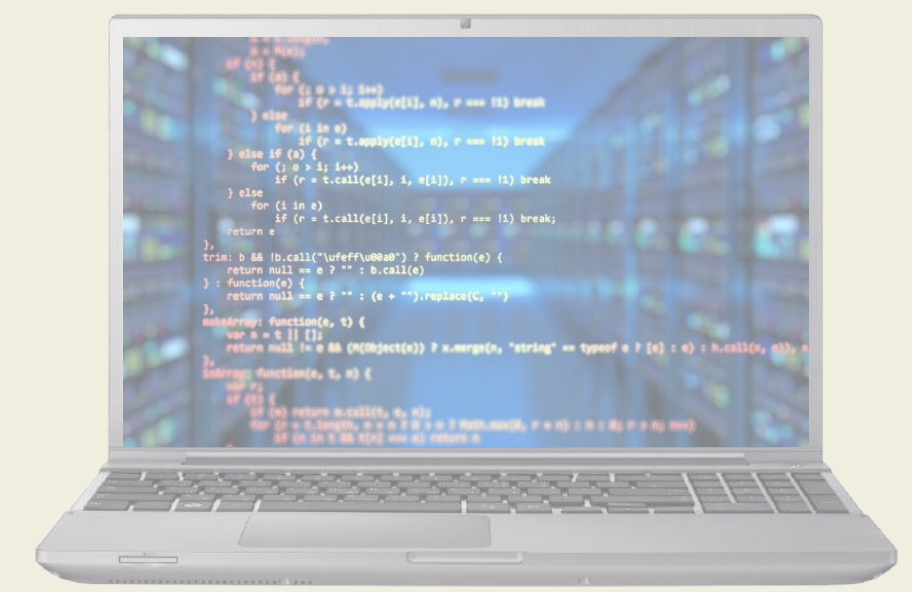
## Training

CIKM 2021



## Interpretability

VMCAI 2024



## Verification

NFM 2023

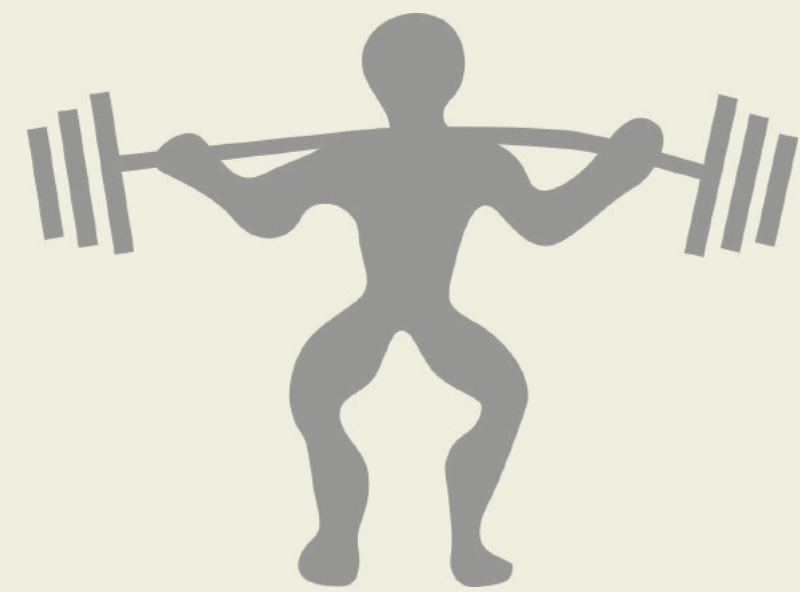




Using formal methods  
for verification

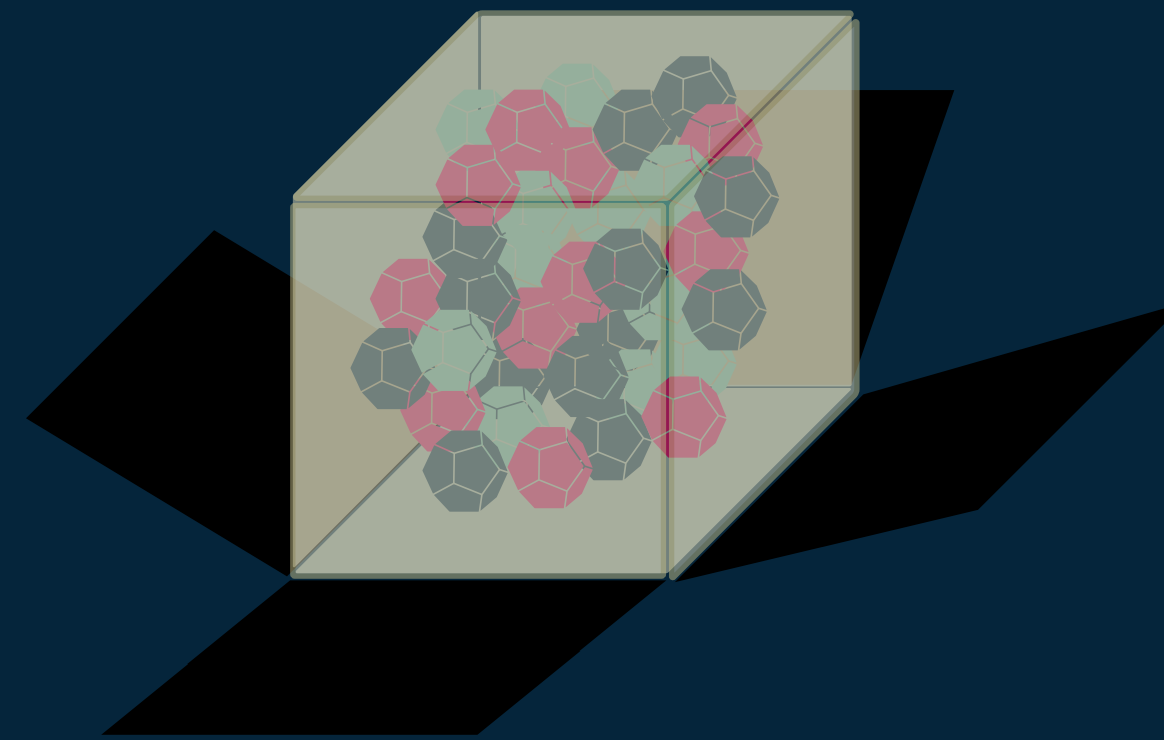


Using formal methods  
**for something else**  
than verification



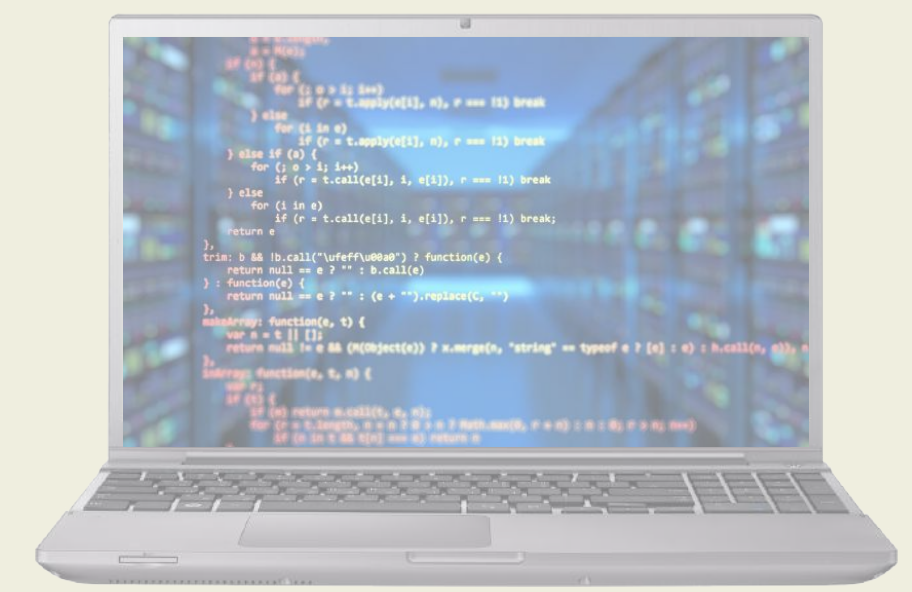
## Training

CIKM 2021



## Interpretability

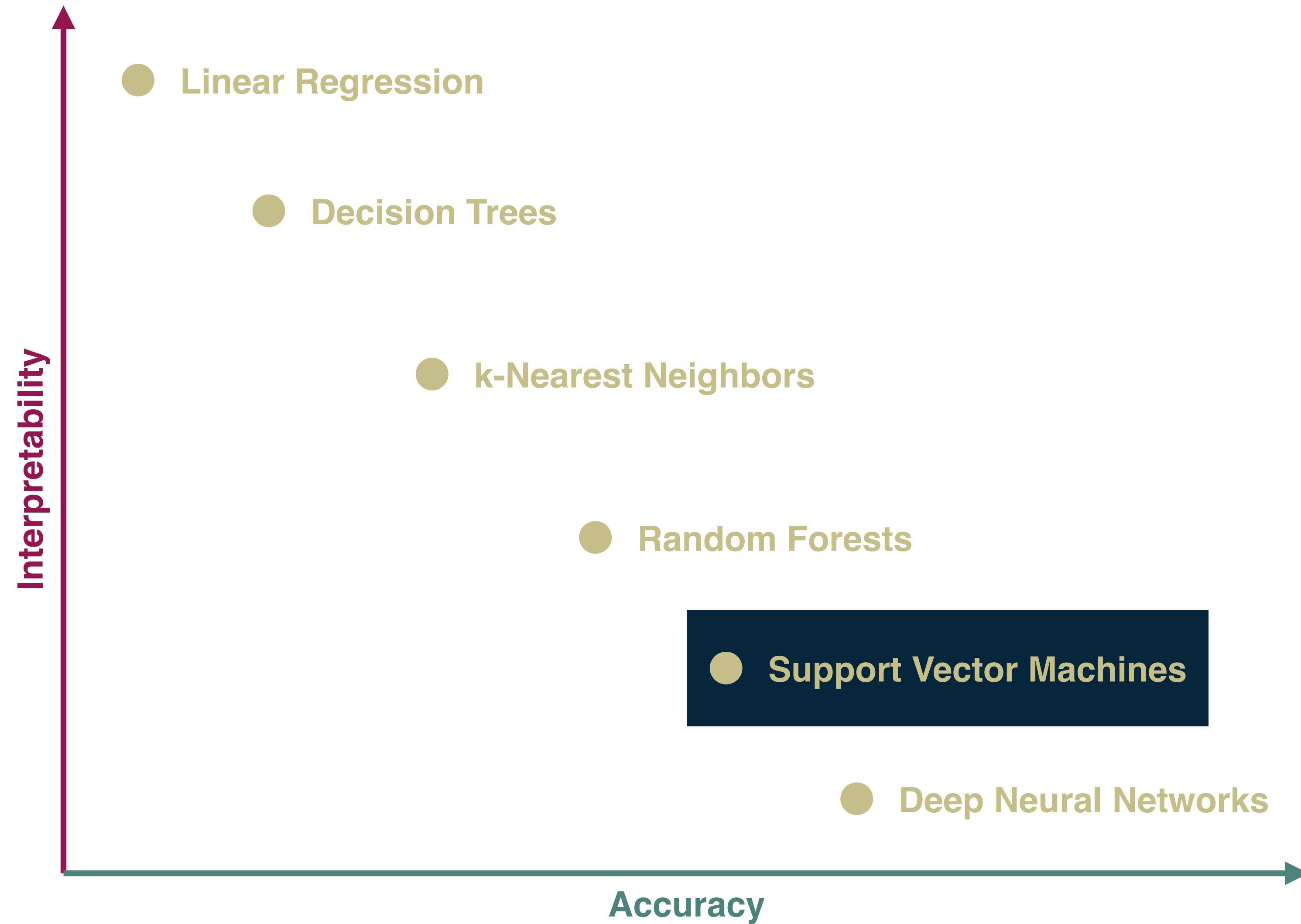
VMCAI 2024



## Verification

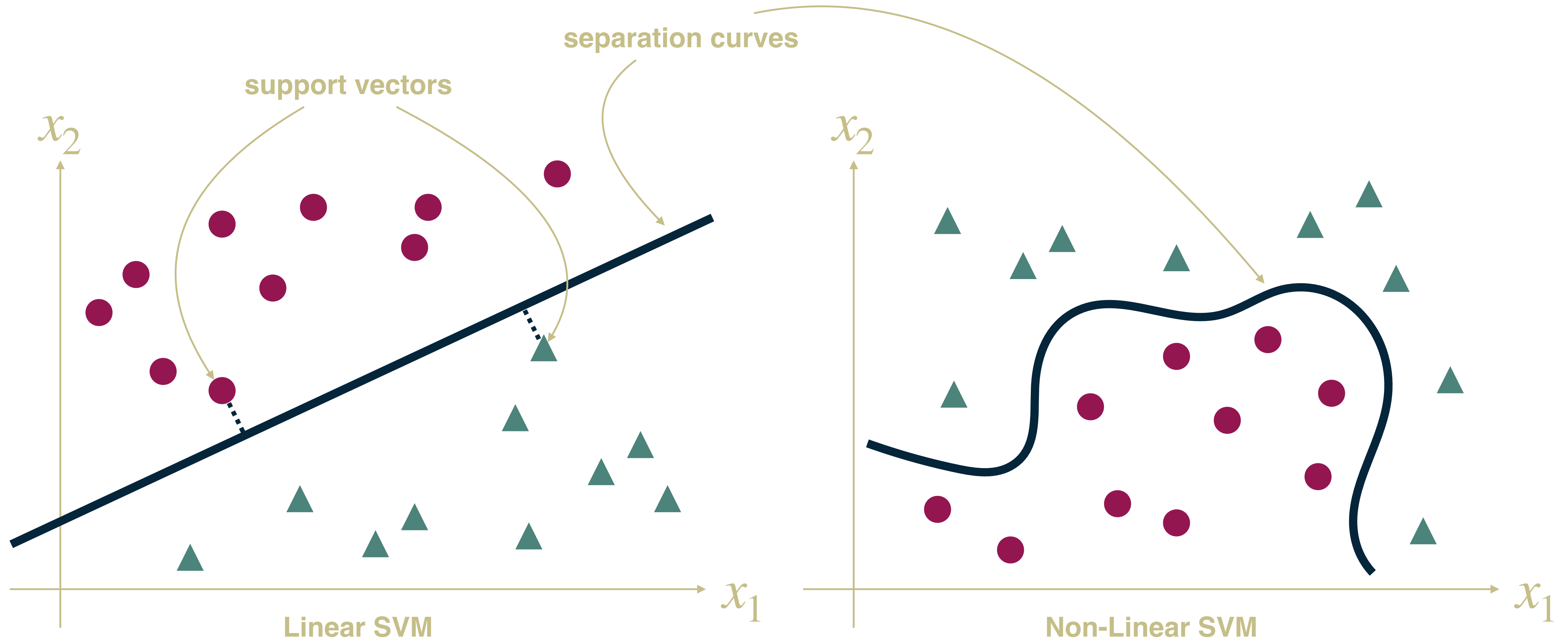
NFM 2023

# Interpretability vs Performance





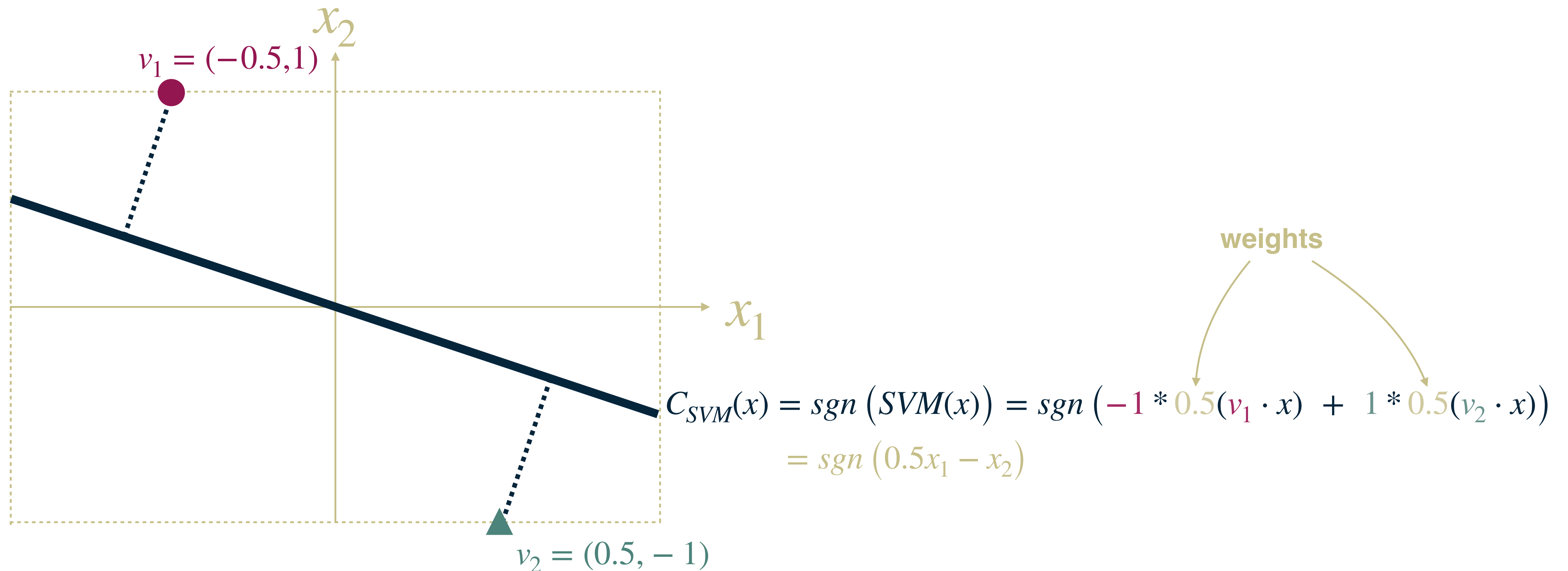
# Support Vector Machines (SVMs)



# Support Vector Machines (SVMs)

## Example

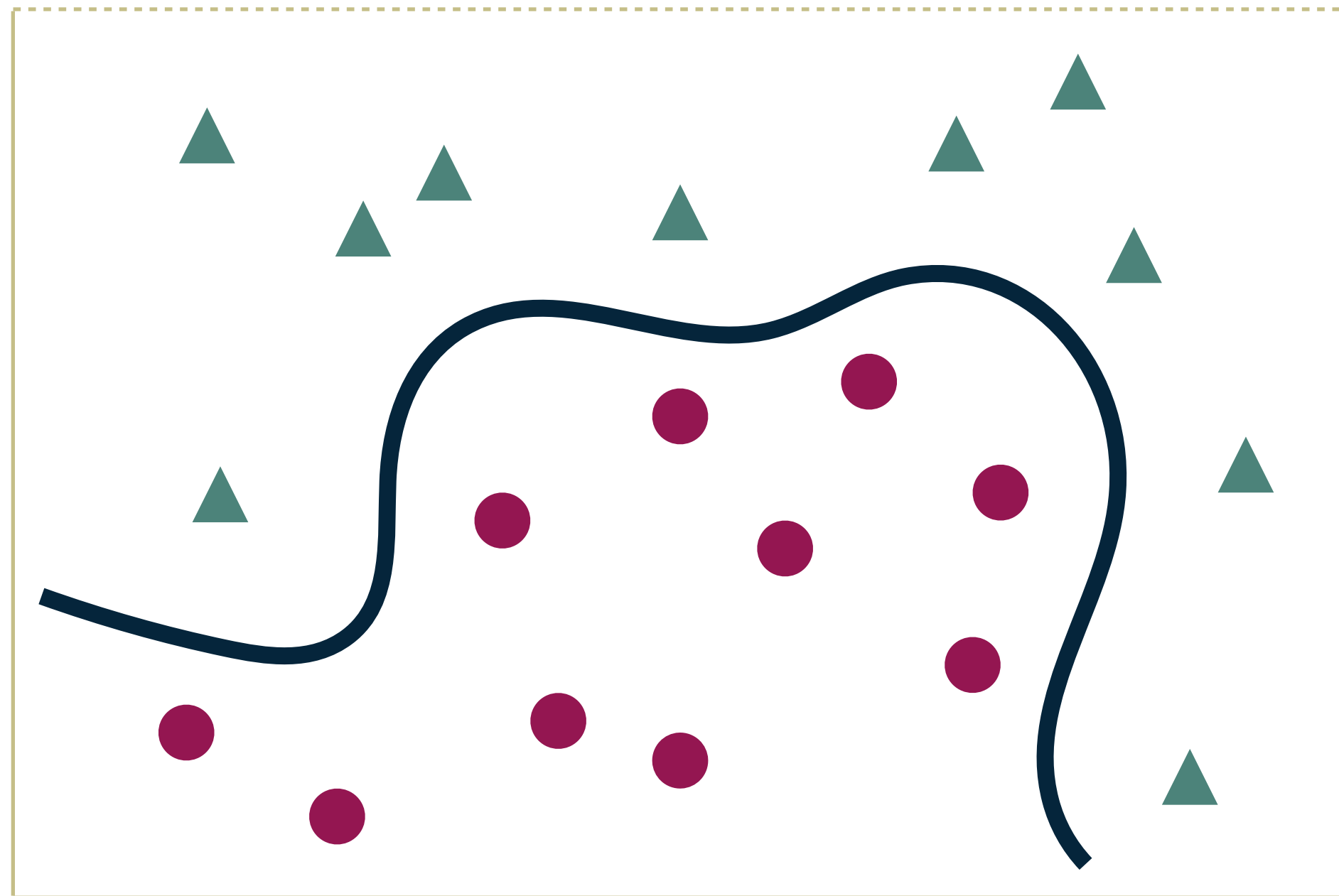
●  $\mapsto -1$   
▲  $\mapsto 1$



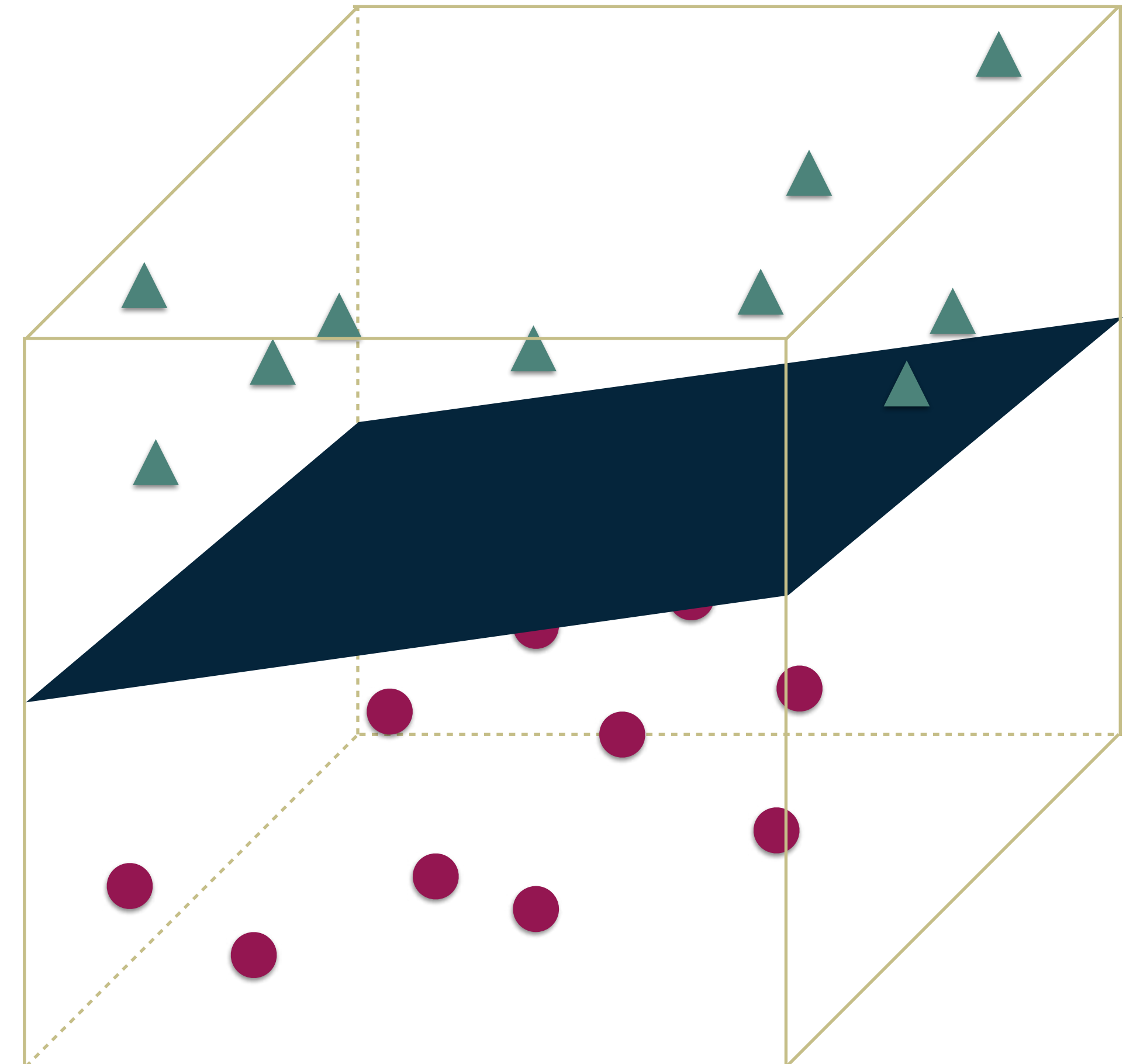
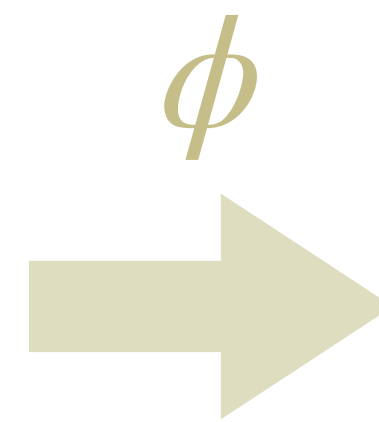
# Non-Linear SVMs

## Kernel Functions

- Polynomial
- Radial Basis Function (RBF)



Input Space



Feature Space



# Feature Importance

## Measuring Contribution of Input Features to Prediction

	Local	Global	Model-		Performance	Effect
			Specific	Agnostic		
Permutation Feature Importance (PFI)		X		X	X	
Partial Dependence (PD) Plots		X		X		X
Individual Conditional Expectation (ICE) Plots		X		X		X
Accumulated Local Effects (ALE) Plots		X		X		X
Local Interpretable Model-Agnostic Explanations (LIME)	X			X		X
SHapley Additive exPlanations (SHAP)	X			X		X
Individual Conditional Importance (ICI) Curves	X			X	X	
Partial Importance (PI) Curves	X			X	X	
Shapley Feature Importance (SFIMP)		X		X	X	
Input Gradients	X			X	X	X
Abstract Feature Importance (AFI)	X	X	X			X

# Abstract Feature Importance

## Why Another Feature Importance Measure?

Permutation Feature Importance (PFI)

- result may greatly vary depending on the dataset
- resource intensive when the number of features is large
- misleading result when features are correlated
- quality of the result heavily depends on the model accuracy

Local Interpretable Model-Agnostic Explanations (LIME)

- requires a neighborhood of similar data points to find a useful optimal neighborhood: requires a large and easily manipulable dataset
- assumes that the decision boundary is linear at the local level, but there is no theoretical guarantee that this is the case

SHapley values (SHAP)

- Shapley values estimations depend on the dataset
- assumes that features are independent
- has a very high computational cost, even for small models

Abstract Feature Importance (AFI)

- yields a formally correct by construction approximation
- does not depend from a dataset nor the accuracy of the model
- extremely fast to compute, whatever the number of features
- supports both linear and non-linear kernel functions

**“Make Sense” but Give No Guarantees**



# 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

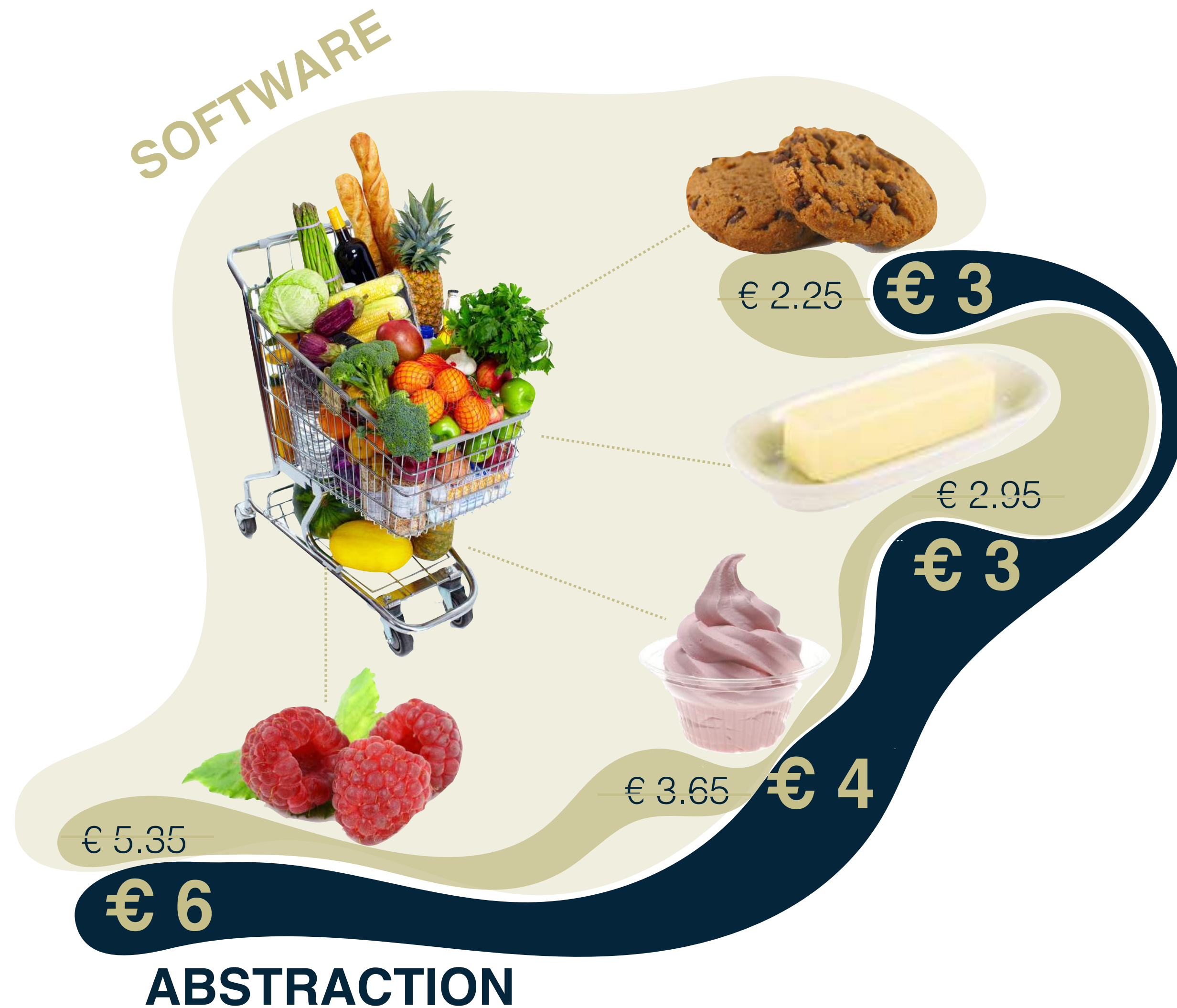


FALSE ALARM

COMPLETENESS



# Abstract Interpretation



Using abstract interpretation for verification

Using abstract interpretation for something else than verification

# Abstract Interpretation of SVMs

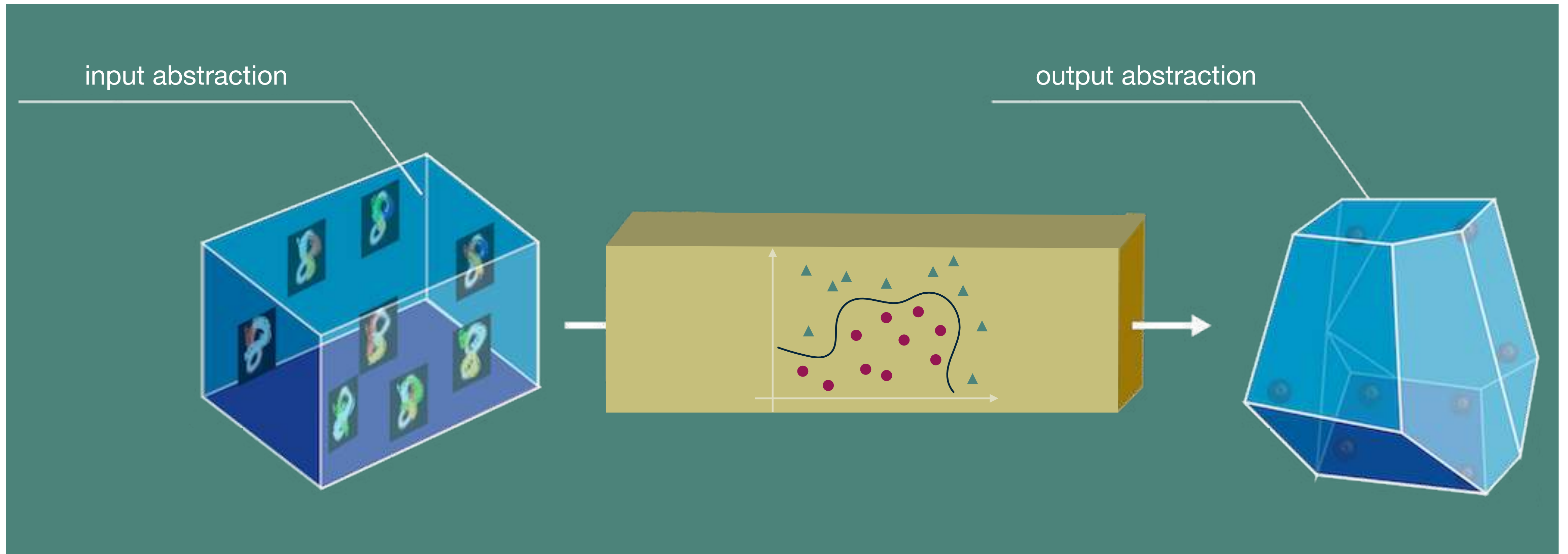
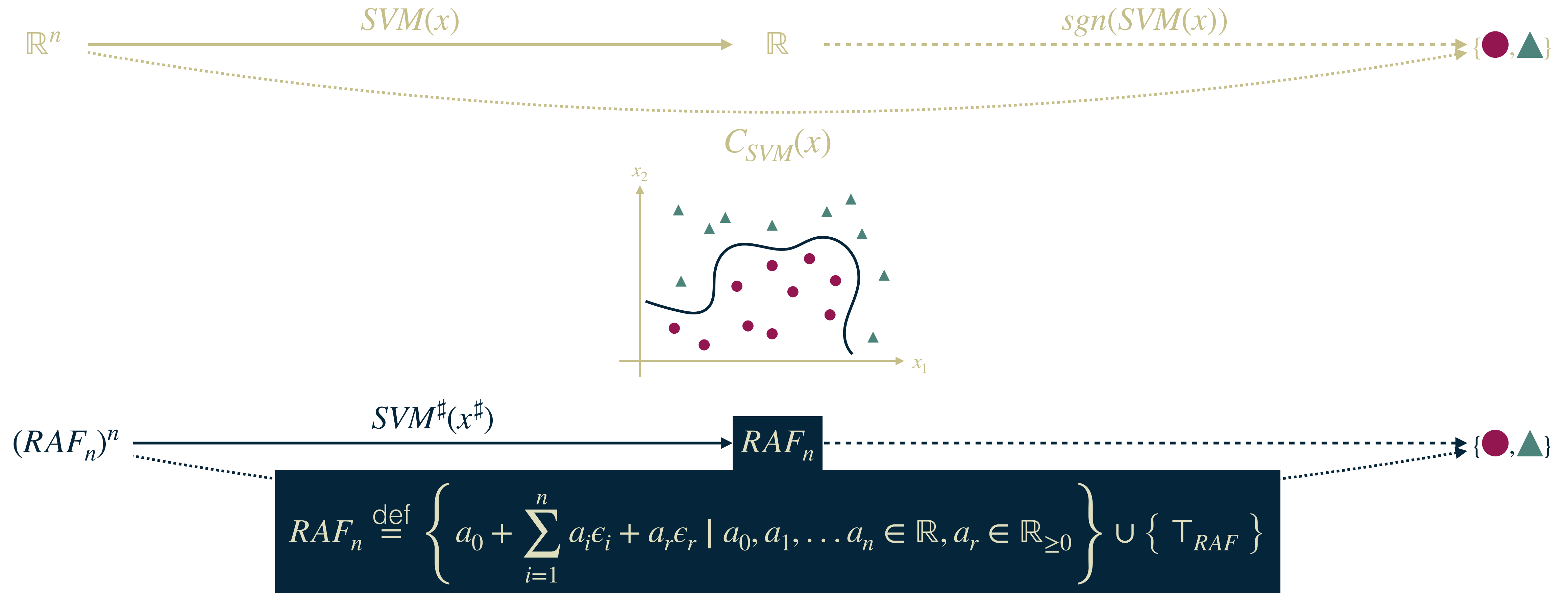


Image taken (and modified) from <http://safeai.ethz.ch>

# Abstract Interpretation of SVMs

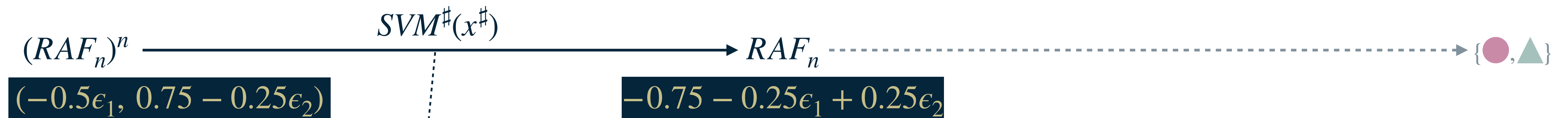
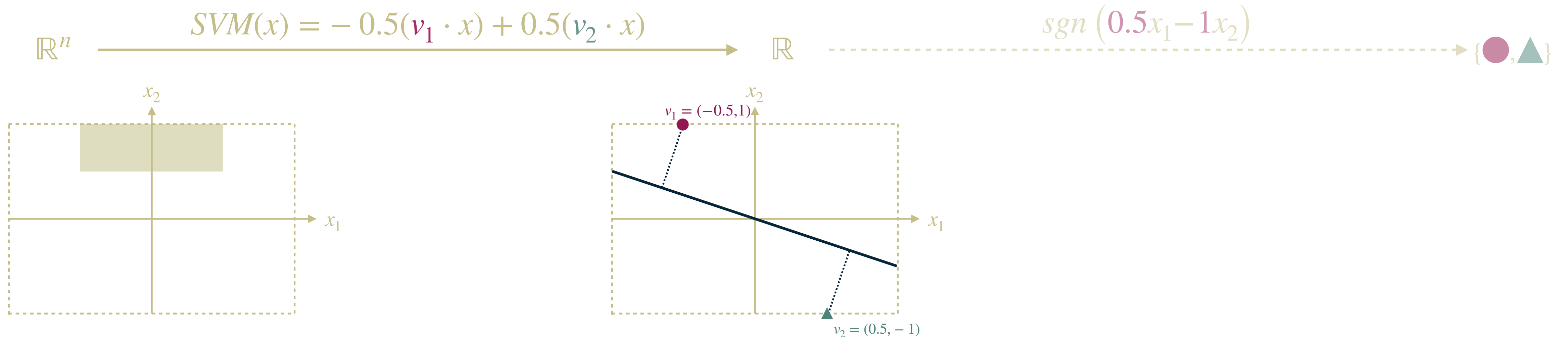
## Reduced Affine Form (RAF) Abstraction





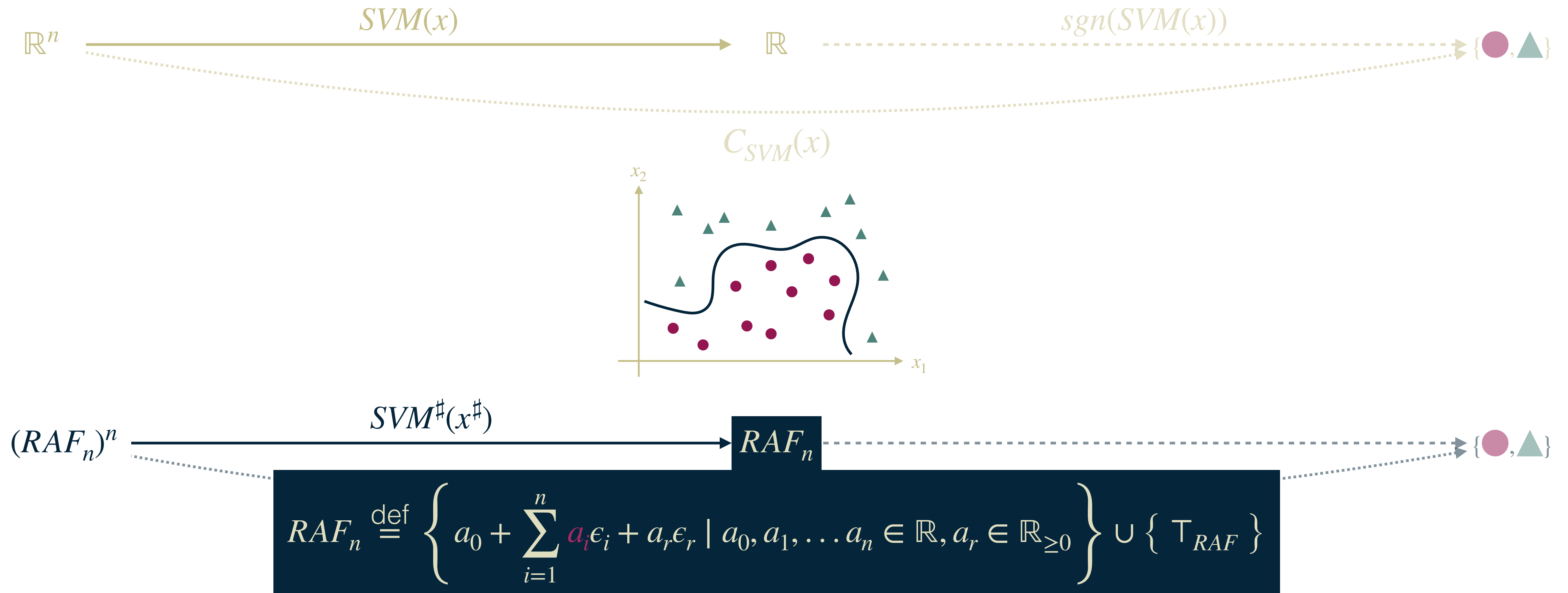
# Abstract Interpretation of SVMs

## Example



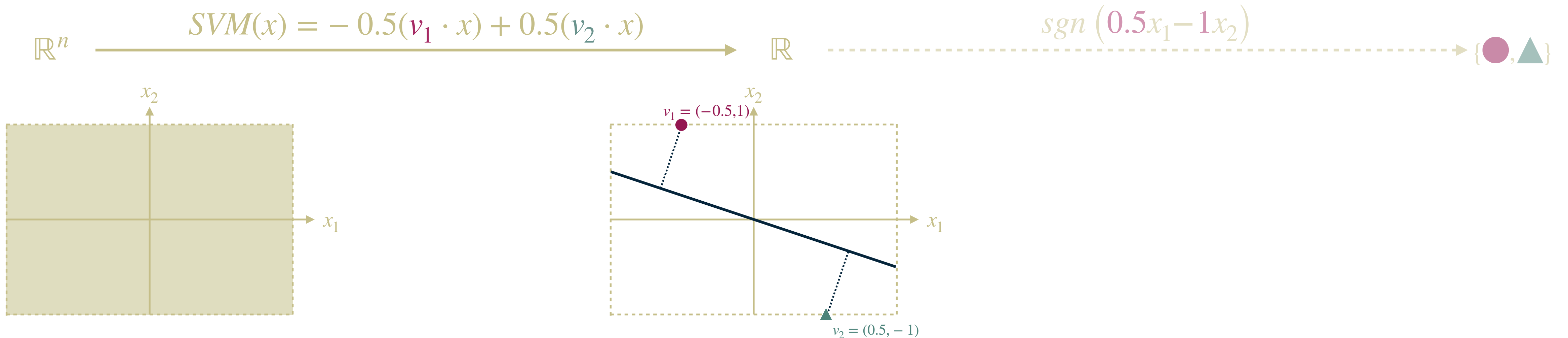
$$\begin{aligned}
 &SVM^\#((-0.5\epsilon_1, 0.75 - 0.25\epsilon_2)) \\
 &= -0.5(-0.5(-0.5\epsilon_1) + 1(0.75 - 0.25\epsilon_2)) + 0.5(0.5(-0.5\epsilon_1) - 1(0.75 - 0.25\epsilon_2)) \\
 &= -0.5(0.75 + 0.25\epsilon_1 - 0.25\epsilon_2) + 0.5(-0.75 - 0.25\epsilon_1 + 0.25\epsilon_2) \\
 &= -0.75 - 0.25\epsilon_1 + 0.25\epsilon_2
 \end{aligned}$$

# Abstract Feature Importance (AFI)



# Abstract Feature Importance (AFI)

## Example



$$\begin{aligned}
 &SVM^\#((\epsilon_1, \epsilon_2)) \\
 &= -0.5(-0.5\epsilon_1 + 1\epsilon_2) + 0.5(0.5\epsilon_1 - 1\epsilon_2) \\
 &= 0.25\epsilon_1 - 0.5\epsilon_2 + 0.25\epsilon_1 - 0.5\epsilon_2 \\
 &= 0.5\epsilon_1 - \epsilon_2
 \end{aligned}$$



# AFI vs PFI

## German Dataset

**Grade for each feature**

<b>Linear</b>	Baseline (13.55s)	<b>5</b>	<b>5</b>	<b>5</b>	<b>6</b>	<b>6</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>8</b>	Distance
	AFI (0.01s)	<b>5</b>	<b>5</b>	<b>5</b>	<b>6</b>	<b>6</b>	<b>7</b>	8	<b>7</b>	<b>7</b>	<b>8</b>	<b>1.0</b>
	PFI (4.07s)	<b>5</b>	<b>5</b>	6	7	7	9	6	6	<b>7</b>	7	3.16
<b>RBF</b>	Baseline (17.98s)	<b>5</b>	<b>5</b>	<b>5</b>	<b>6</b>	<b>6</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>8</b>	<b>8</b>	Distance
	AFI (0.02s)	<b>5</b>	6	<b>5</b>	<b>6</b>	<b>6</b>	8	<b>7</b>	<b>7</b>	<b>8</b>	7	<b>1.73</b>
	PFI (6.23s)	6	7	<b>5</b>	<b>6</b>	7	8	<b>7</b>	6	7	5	4.24
<b>Polynomial</b>	Baseline (15.83s)	<b>5</b>	<b>5</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>7</b>	<b>8</b>	Distance
	AFI (0.01s)	7	6	7	7	5	<b>7</b>	6	6	5	<b>8</b>	<b>4.47</b>
	PFI (4.15s)	6	7	9	7	6	<b>7</b>	5	6	6	6	5.74

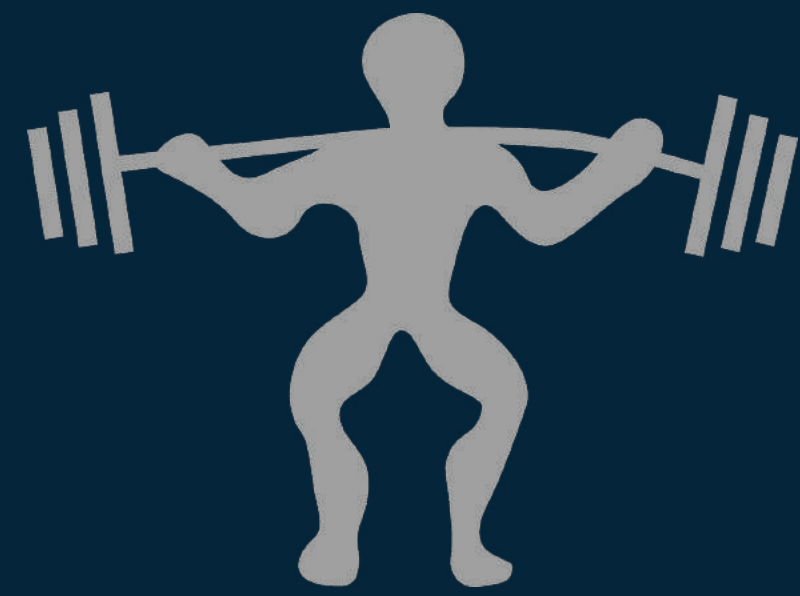
# AFI vs PFI

	<b>Baseline</b>	$N = 2k$ $\epsilon = 0.2$	$N = 10k$ $\epsilon = 0.2$	$N = 2k$ $\epsilon = 0.4$	$N = 10k$ $\epsilon = 0.4$	$N = 2k$ $\epsilon = 0.6$	$N = 5k$ $\epsilon = 0.6$	$N = 10k$ $\epsilon = 0.6$	$N = 2k$ $\epsilon = 0.8$	$N = 5k$ $\epsilon = 0.8$	$N = 10k$ $\epsilon = 0.8$
<b>Adult</b> Linear	AFI (0.27s)	<b>0.0</b>	<b>0.0</b>	<b>1.0</b>	<b>0.0</b>	<b>1.0</b>	<b>1.41</b>	<b>1.0</b>	<b>1.0</b>	<b>1.41</b>	<b>1.0</b>
	PFI (10009s)	2.45	2.45	2.24	2.45	2.24	<b>1.41</b>	2.24	2.24	<b>1.41</b>	2.24
<b>Adult</b> RBF	AFI (0.48s)	<b>1.0</b>	<b>1.41</b>	<b>1.41</b>	<b>1.41</b>	<b>1.73</b>	<b>1.73</b>	<b>1.41</b>	<b>1.41</b>	<b>1.41</b>	<b>1.41</b>
	PFI (25221s)	1.73	2.45	2.45	2.0	2.65	2.65	2.45	2.45	2.45	2.45
<b>Adult</b> Polynomial	AFI (0.44s)	<b>1.0</b>	<b>1.0</b>	<b>0.0</b>	1.41	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
	PFI (9985s)	<b>1.0</b>	<b>1.0</b>	1.41	<b>1.0</b>	1.41	1.41	1.41	1.41	1.41	1.41
<b>Compas</b> Linear	AFI (0.22s)	<b>1.41</b>	<b>1.41</b>	<b>1.73</b>	<b>1.73</b>	<b>1.41</b>	<b>1.73</b>	<b>1.41</b>	<b>1.41</b>	<b>1.41</b>	<b>1.73</b>
	PFI (1953s)	1.73	1.73	2.0	2.0	2.24	2.0	2.24	2.24	2.24	2.83
<b>Compas</b> RBF	AFI (0.27s)	<b>2.0</b>	<b>2.0</b>	<b>2.65</b>	<b>2.65</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>
	PFI (6827s)	<b>2.0</b>	<b>2.0</b>	<b>2.65</b>	<b>2.65</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>	<b>2.83</b>
<b>Compas</b> Polynomial	AFI (0.22s)	4.24	4.24	4.12	4.12	4.24	4.24	4.24	4.24	4.24	4.24
	PFI (2069s)	<b>2.45</b>	<b>2.45</b>	<b>3.0</b>	<b>3.0</b>	<b>3.74</b>	<b>3.74</b>	<b>3.74</b>	<b>3.74</b>	<b>3.74</b>	<b>3.74</b>
<b>German</b> Linear	AFI (0.01s)	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>1.41</b>	<b>1.73</b>	<b>1.41</b>
	PFI (4.07s)	3.16	3.46	3.16	3.16	3.16	3.16	3.16	3.6	3.74	3.0
<b>German</b> RBF	AFI (0.02s)	<b>1.73</b>	<b>1.0</b>	<b>1.73</b>	<b>1.73</b>	<b>2.0</b>	<b>1.41</b>	<b>1.73</b>	<b>1.73</b>	<b>2.0</b>	<b>2.24</b>
	PFI (6.23s)	4.0	3.46	4.24	4.24	4.36	3.61	4.24	4.24	4.36	4.47
<b>German</b> Polynomial	AFI (0.01s)	<b>4.90</b>	<b>4.12</b>	<b>4.47</b>	<b>3.87</b>	<b>3.87</b>	<b>4.24</b>	<b>3.46</b>	<b>3.46</b>	<b>3.46</b>	<b>3.46</b>
	PFI (4.15s)	5.74	5.10	5.74	4.69	4.69	5.0	4.58	4.58	4.58	4.58

# AFI vs LIME

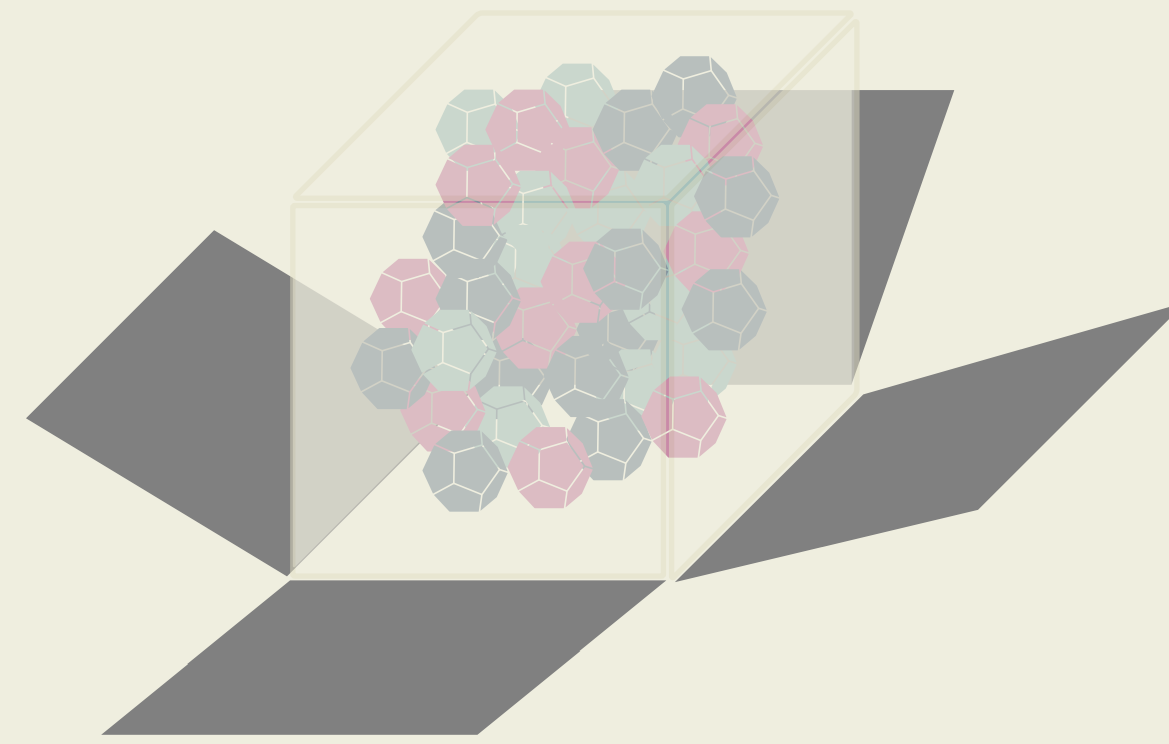
Distance between LIME and ...	Adult			Compas			German		
	Lin.	RBF	Poly	Lin.	RBF	Poly	Lin.	RBF	Poly
AFI ( $\epsilon = 0.1$ )	2.42	2.04	2.98	1.67	1.06	3.05	2.62	2.03	<b>5.31</b>
AFI ( $\epsilon = 0.2$ )	1.68	1.32	2.67	1.63	0.17	2.73	2.21	2.00	5.41
AFI ( $\epsilon = 0.3$ )	1.39	0.51	2.58	<b>1.57</b>	0.14	<b>2.62</b>	1.92	2.05	5.45
AFI (Global)	<b>1.37</b>	<b>0.01</b>	<b>1.01</b>	<b>1.57</b>	<b>0.13</b>	3.16	<b>1.90</b>	<b>1.89</b>	5.53





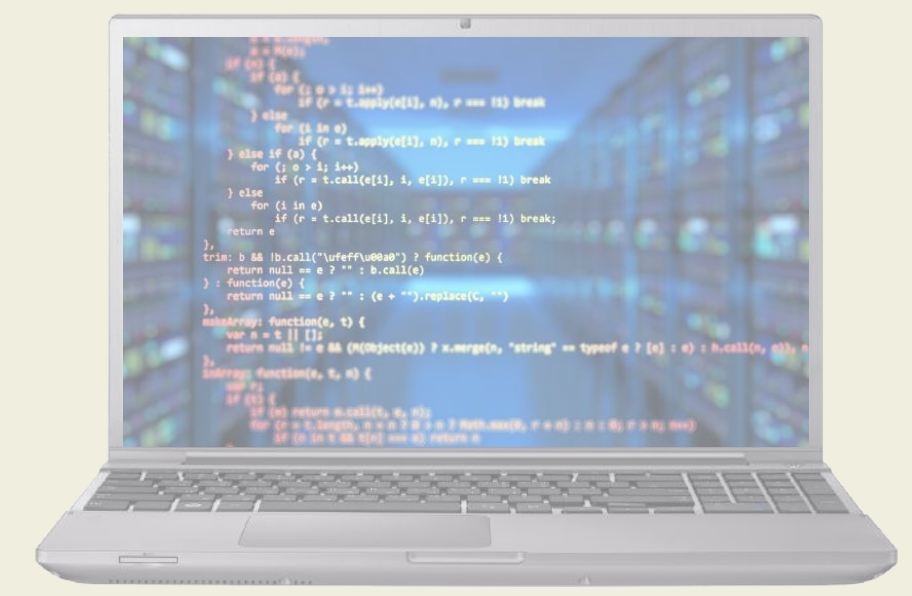
# Training

CIKM 2021



# Interpretability

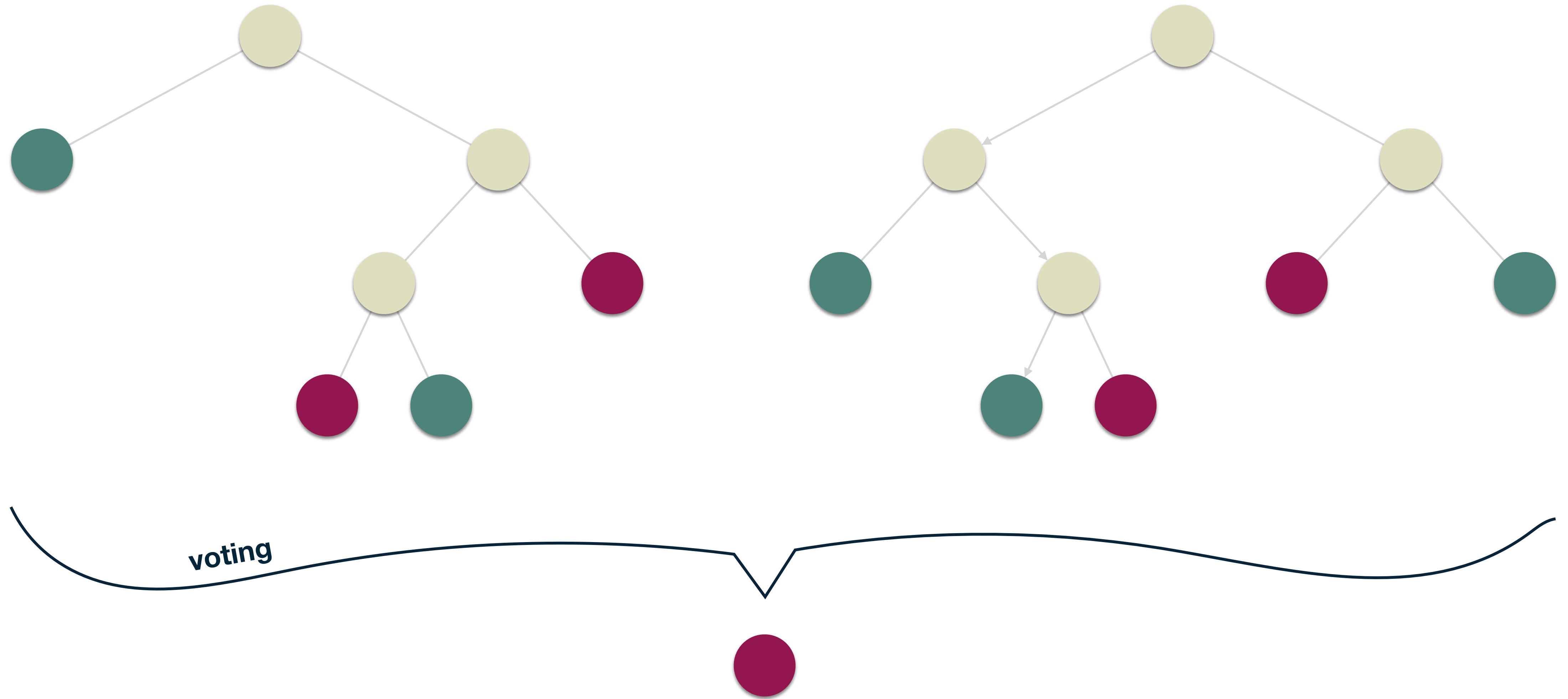
VMCAI 2024



# Verification

NFM 2023

# Random Forests

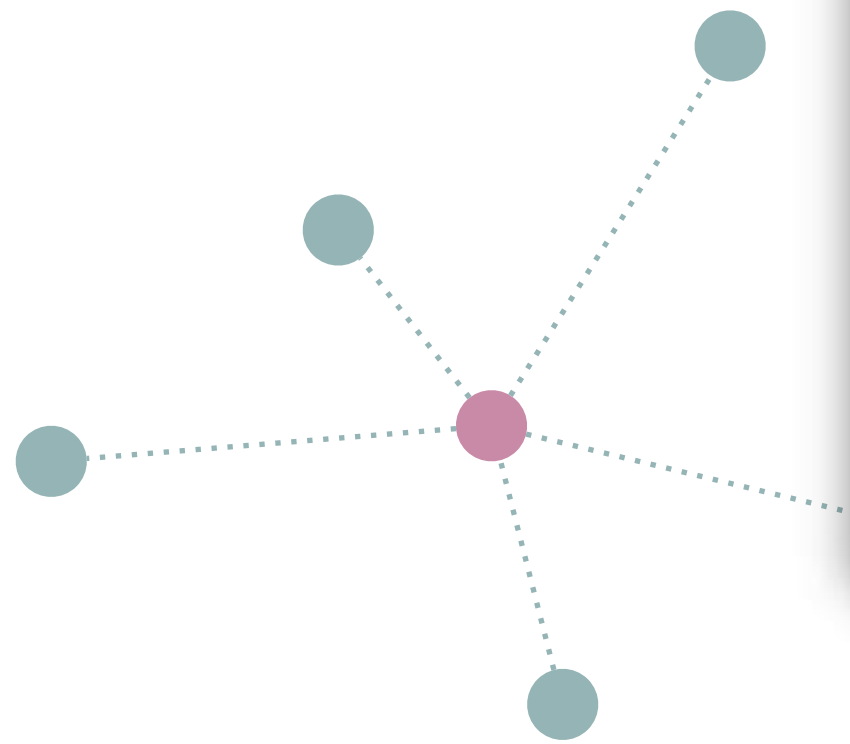


# Robust Training

Minimizing the Worst-Case Loss for Each Input

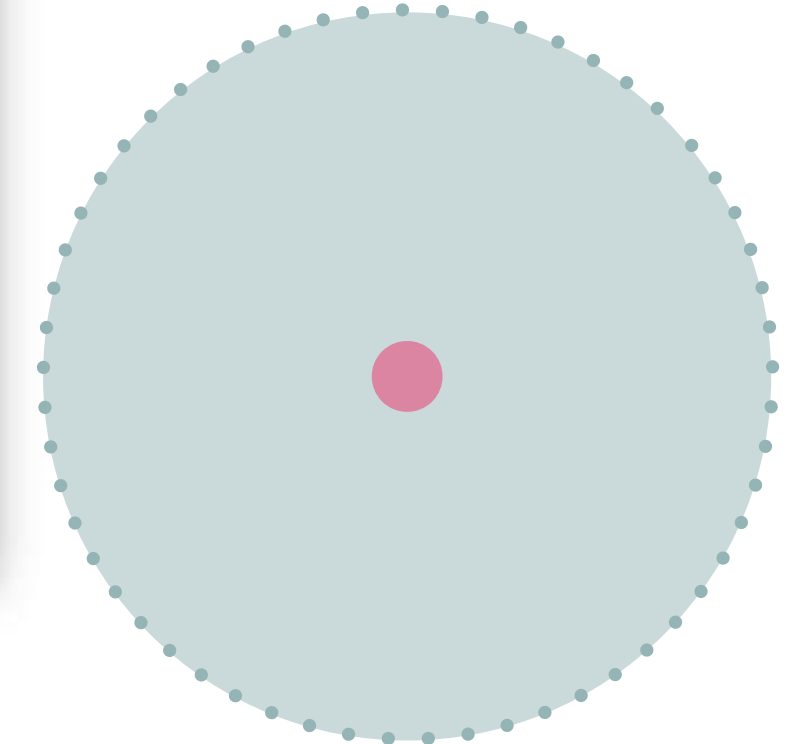
## Adversarial Training

Minimizing a Lower-Bound on the Worst-Case Loss



## Certified Training

Minimizing an Upper-Bound on the Worst-Case Loss



## Hybrid Training

$$(1 - \alpha)\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y) + \alpha \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

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

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

Machine Learning Community

Formal Methods Community

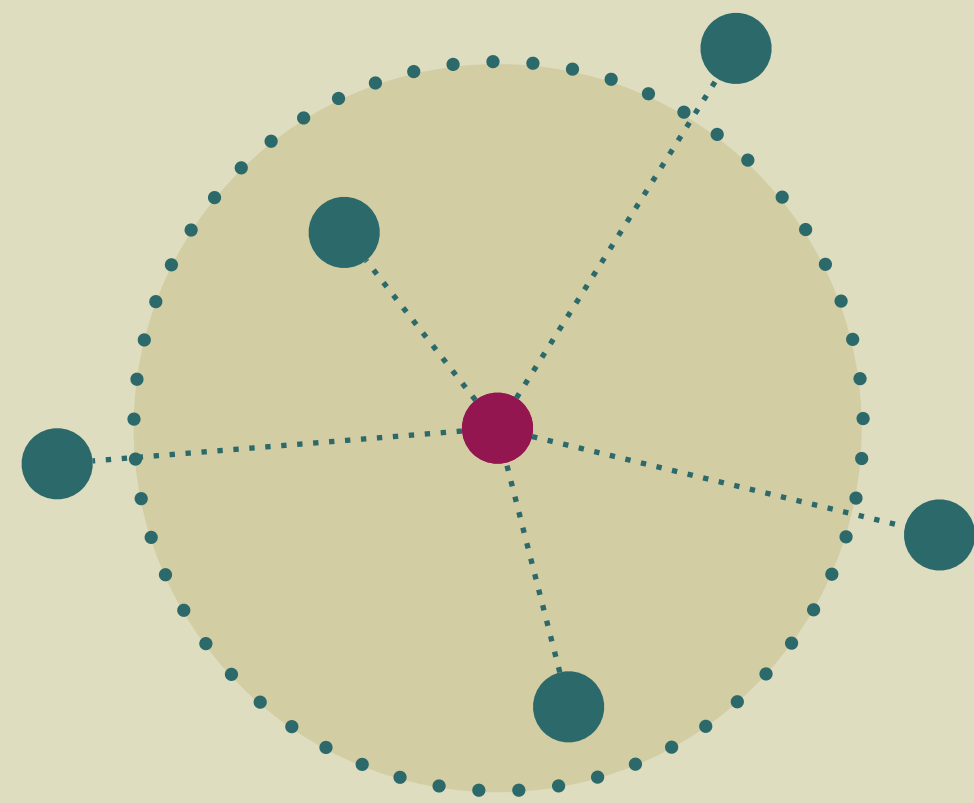


# Hybrid Training

## Random Forests

Dataset	FATT			Natural CART			CART with Hints		
	Accuracy %	Fairness %	Size	Accuracy %	Fairness %	Size	Accuracy %	Fairness %	Size
Adult	80.84	95.21	43	85.32	77.56	270	84.77	87.46	47
Compas	64.11	85.98	75	65.91	22.25	56	65.91	22.25	56
Crime	79.45	75.19	11	77.69	24.31	48	77.44	60.65	8
German	72.00	99.50	2	75.50	57.50	115	73.50	86.00	4
Health	77.87	97.03	84	83.85	79.98	2371	82.25	93.64	100
<b>Average</b>	74.85	<b>90.58</b>	<b>43</b>	<b>77.65</b>	52.32	572	76.77	70.00	<b>43</b>





## Hybrid Training

CIKM 2021

$$a_0 + \sum_{i=1}^n a_i \epsilon_i + a_r \epsilon_r$$

## Verification for Interpretability

VMCAI 2024



## Interpretability for Verification

NFM 2023

THANKS!