

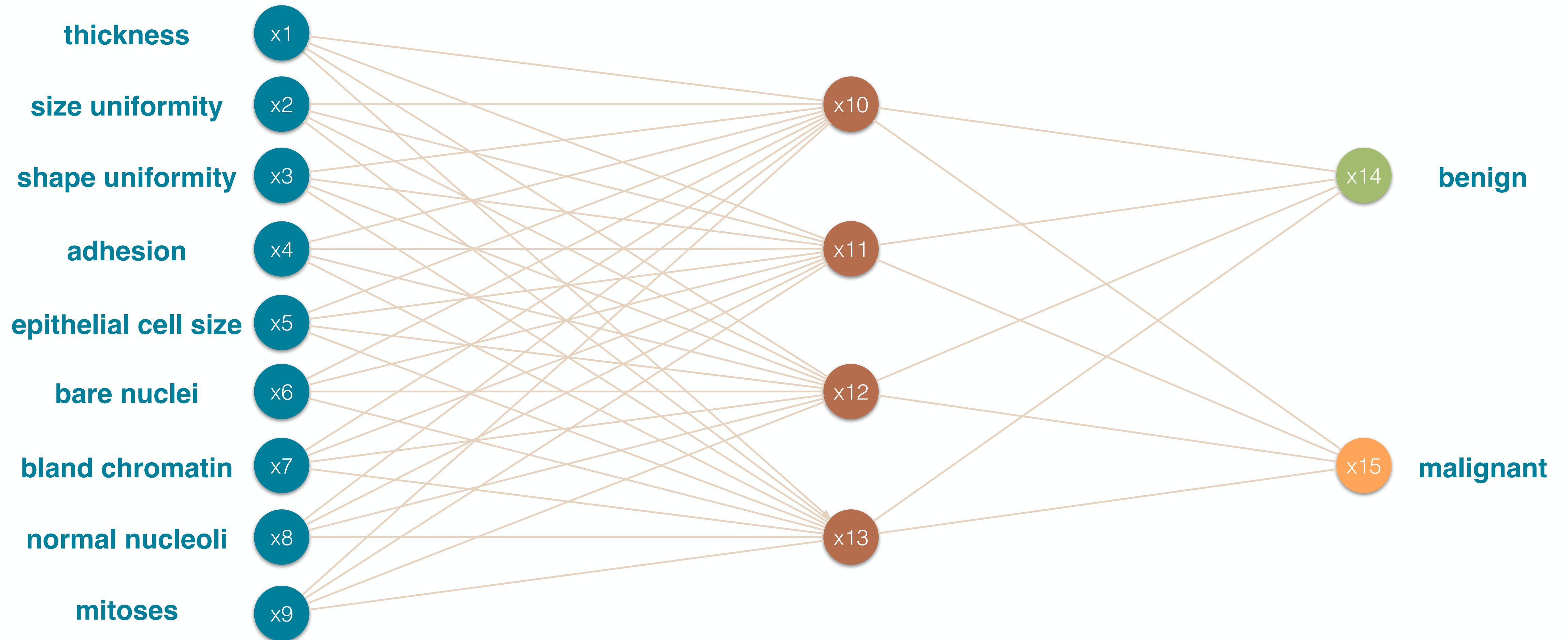
Faster Verified Explanations for Neural Networks

SAIF Days 2026

Caterina Urban (joint work with **Alessandro De Palma** and **Greta Dolcetti**)
Inria & École Normale Supérieure | Université PSL

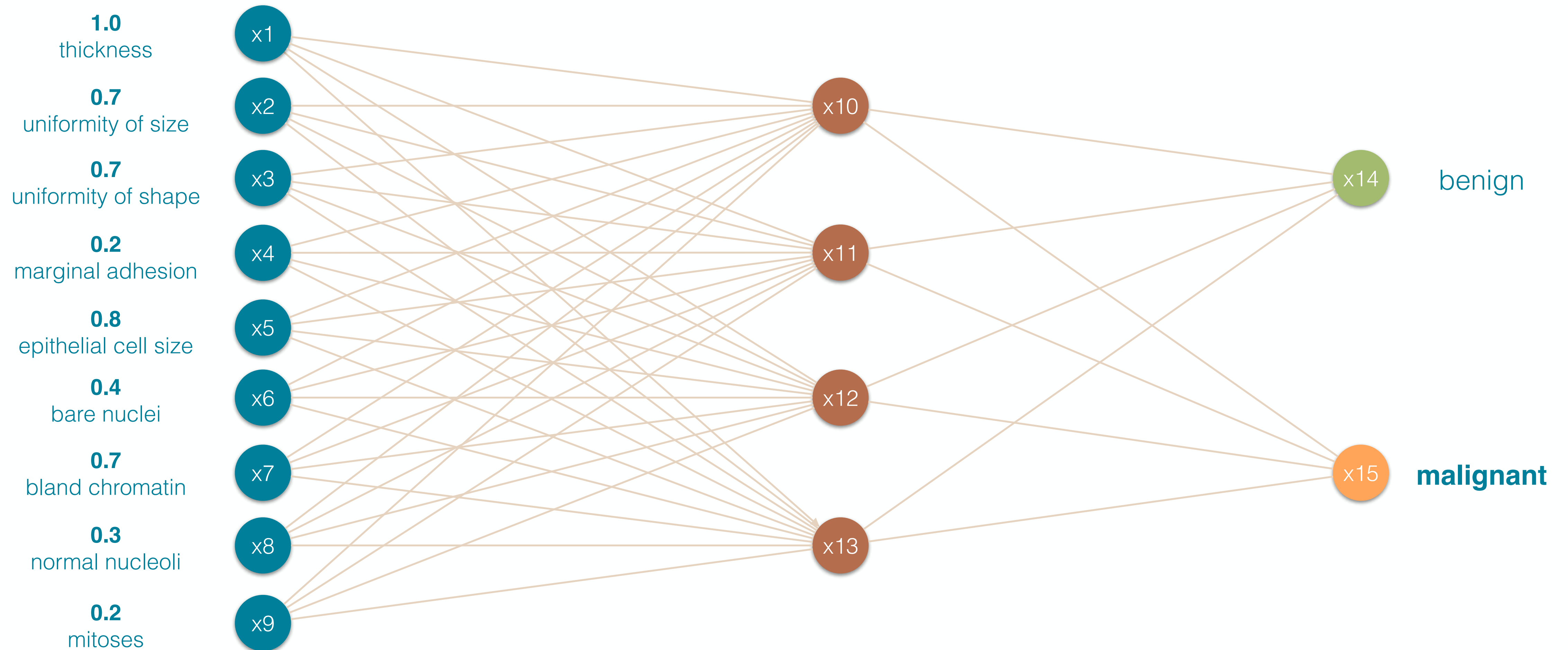
Optimal Robust Explanations

Abductive Explanations (AXps)



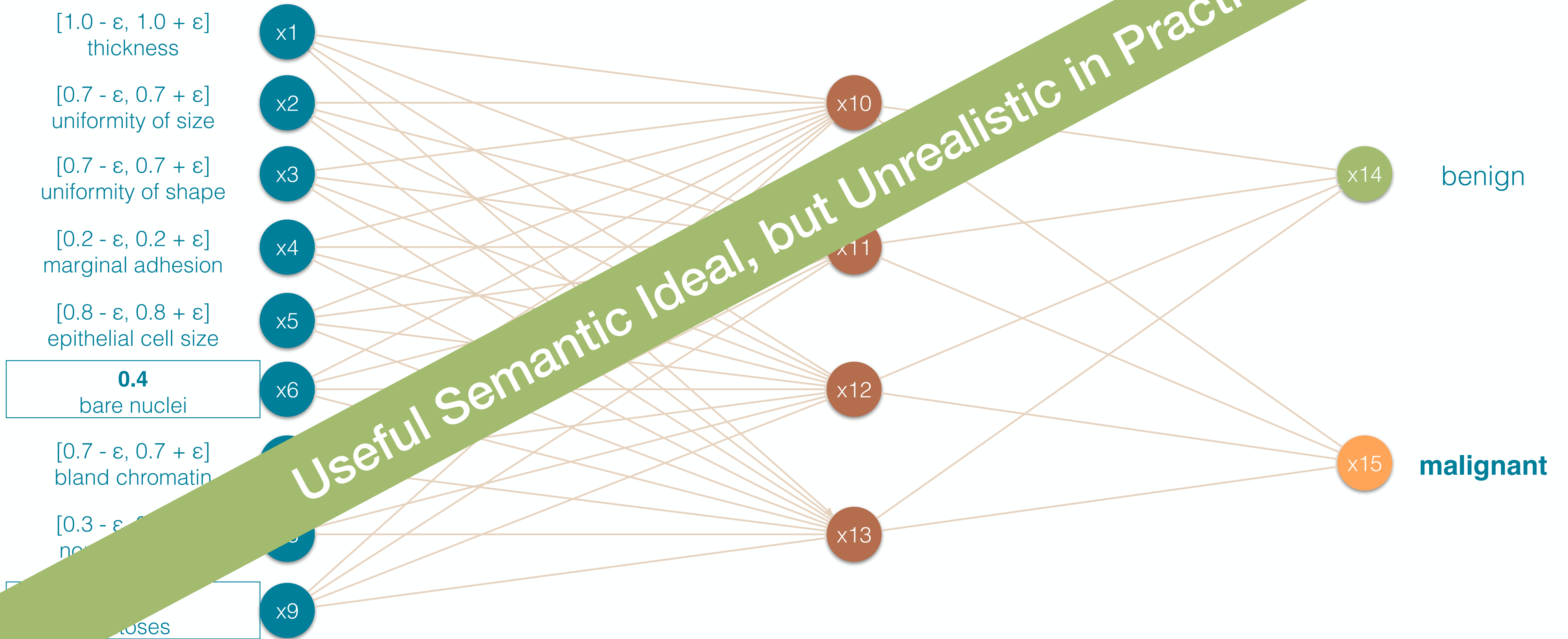
Optimal Robust Explanations

Abductive Explanations (AXps)



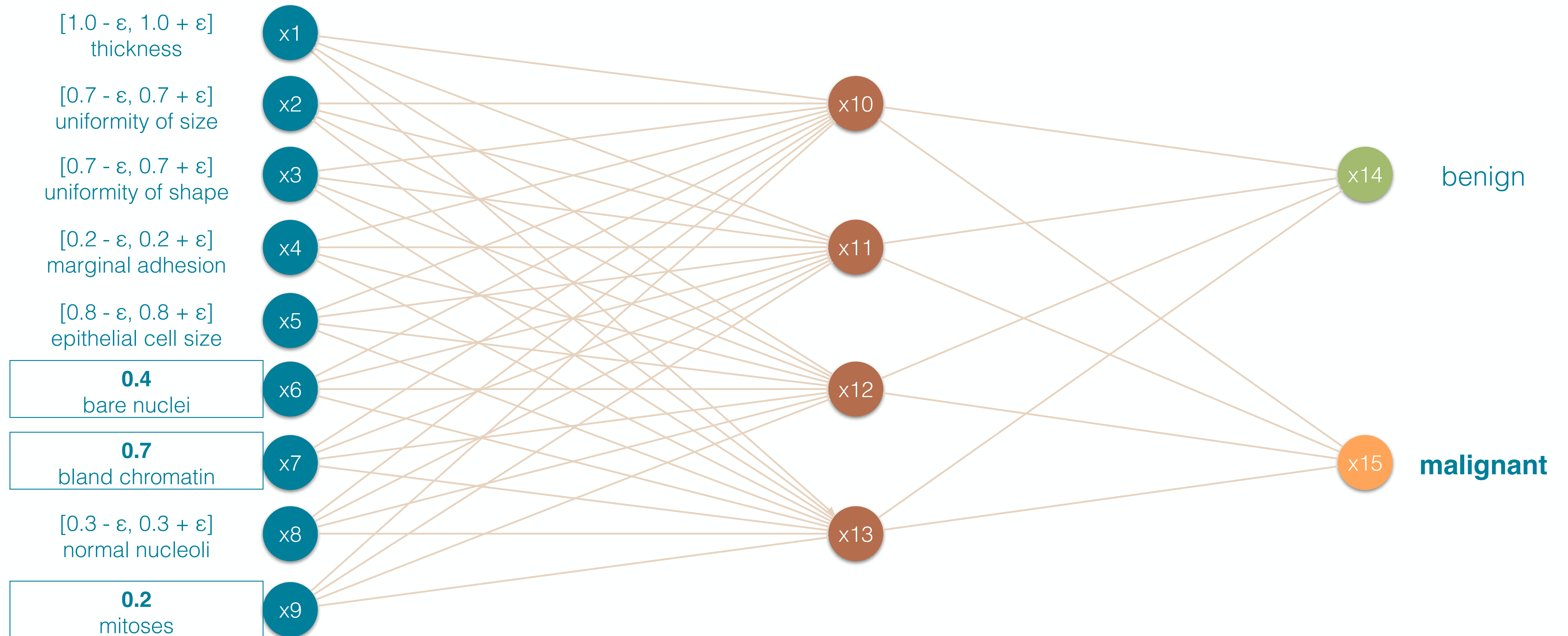
Optimal Robust Explanations

Abductive Explanations (AXps)



Optimal Robust Explanations

Weak Abductive Explanations



Optimal Robust Explanations

Weak Abductive Explanations



Computing ~~Optimal~~ Robust Explanations

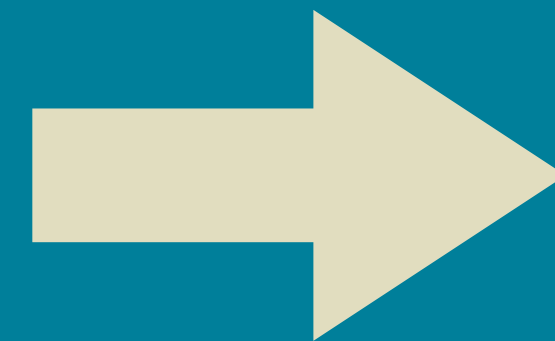
Finding Counterfactuals Rapidly Becomes Infeasible

Model	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	Counterfactuals	Unknowns	Time	Counterfactuals	Unknowns	Time
CNN-3	0.00	247.80	45m	0.00	461.00	2h 30m

Model	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=16/255$		
	Counterfactuals	Unknowns	Time	Counterfactuals	Unknowns	Time
CNN-7	0.00	452.00	3h 59m	0.00	730.67	7h 5m

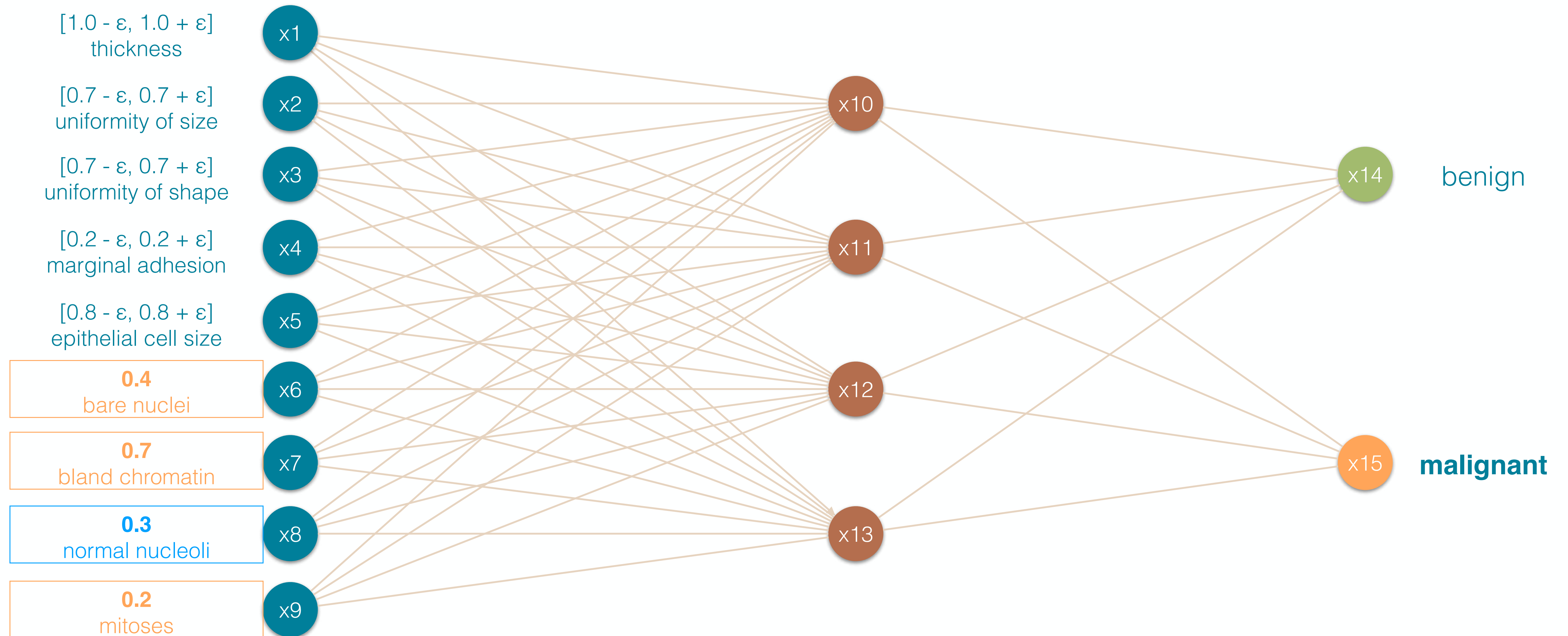
Verifier-Optimal Robust Explanations

Weak Abductive Explanations



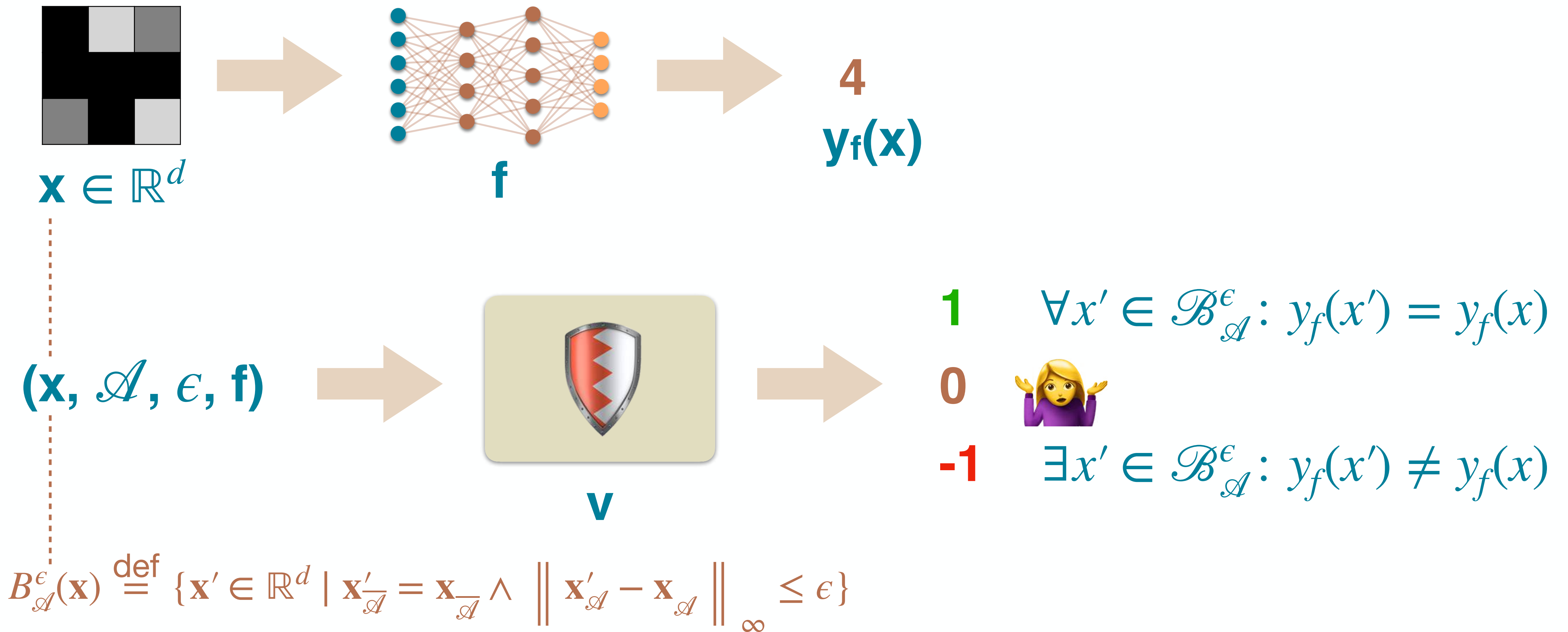
Verifier-Optimal Robust Explanations

Weak Abductive Explanations



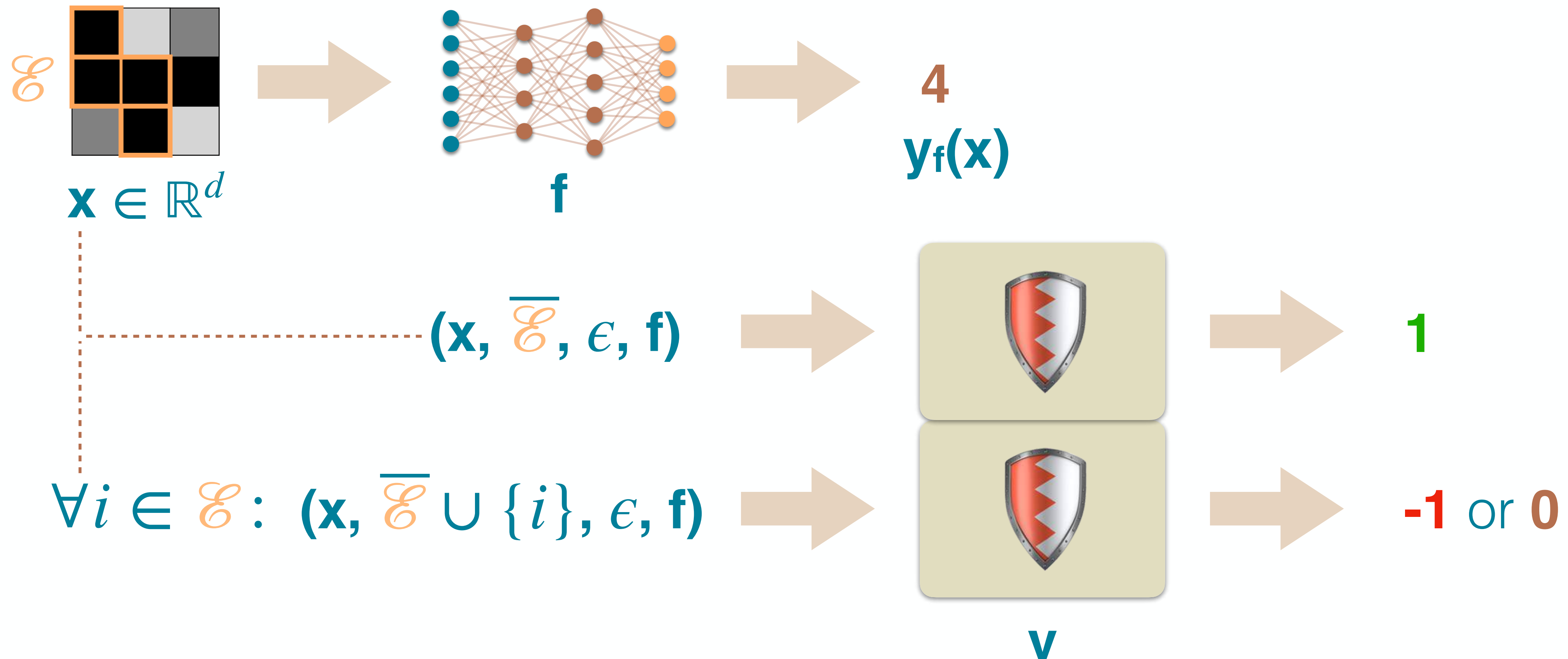
Neural Network Verification

Local Robustness



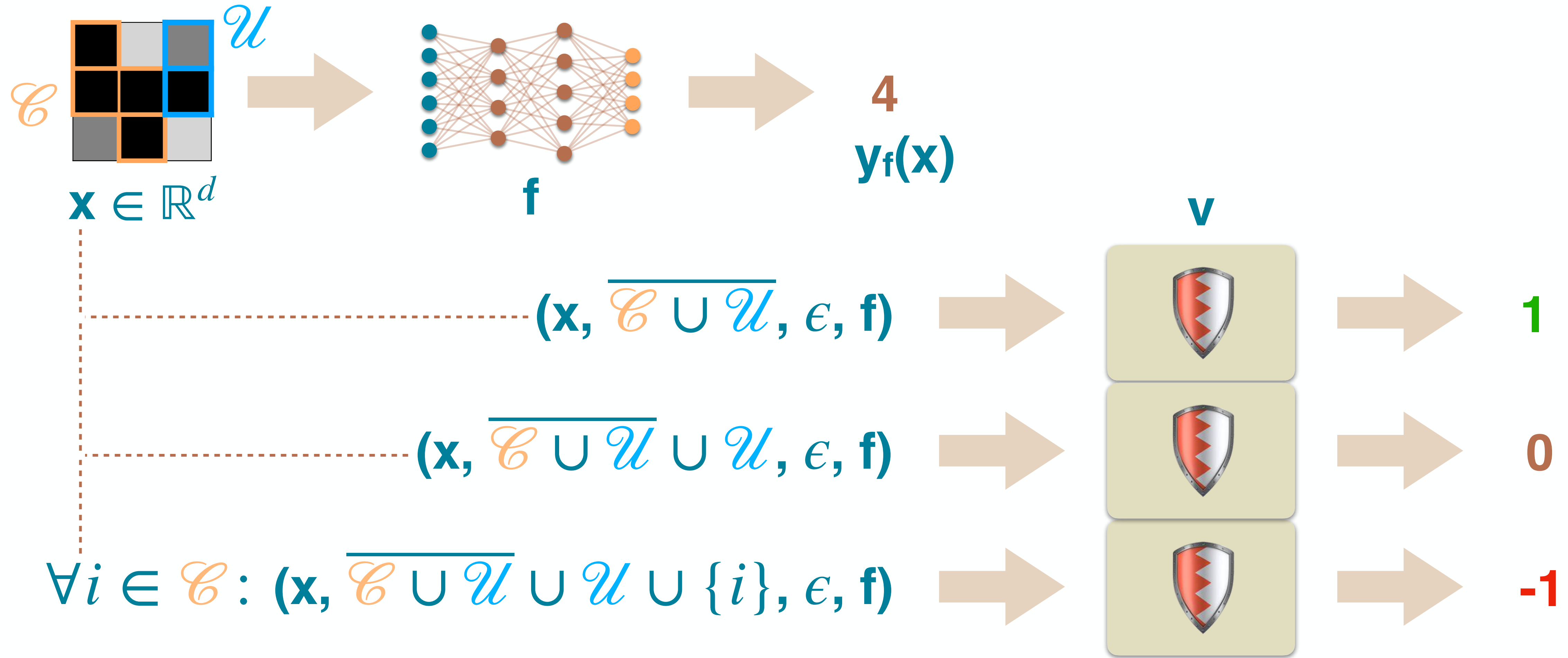
Optimal Robust Explanations

Weak Abductive Explanations



Verifier-Optimal Robust Explanations

Weak Abductive Explanations



Computing Verifier-Optimal Robust Explanations

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

ADD TO $\overline{\mathcal{C} \cup \mathcal{U}}$

LOCAL ROBUSTNESS IN $B_{\overline{\mathcal{C} \cup \mathcal{U}}}^\epsilon(\mathbf{x})$

x_1	●	✓	✓	✓	✓	✓	✓	✓	✓	✓
x_2		●	✓	✓	✓	✓	✓	✓	✓	✓
x_3			●	✓	✓	✓	✓	✓	✓	✓
x_4				●	✓	✓	✓	✓	✓	✓
x_5					●	✓	✓	✓	✓	✓
x_6						●	×	×	×	×
x_7							●	U	U	U
x_8								●	×	×
x_9									●	×

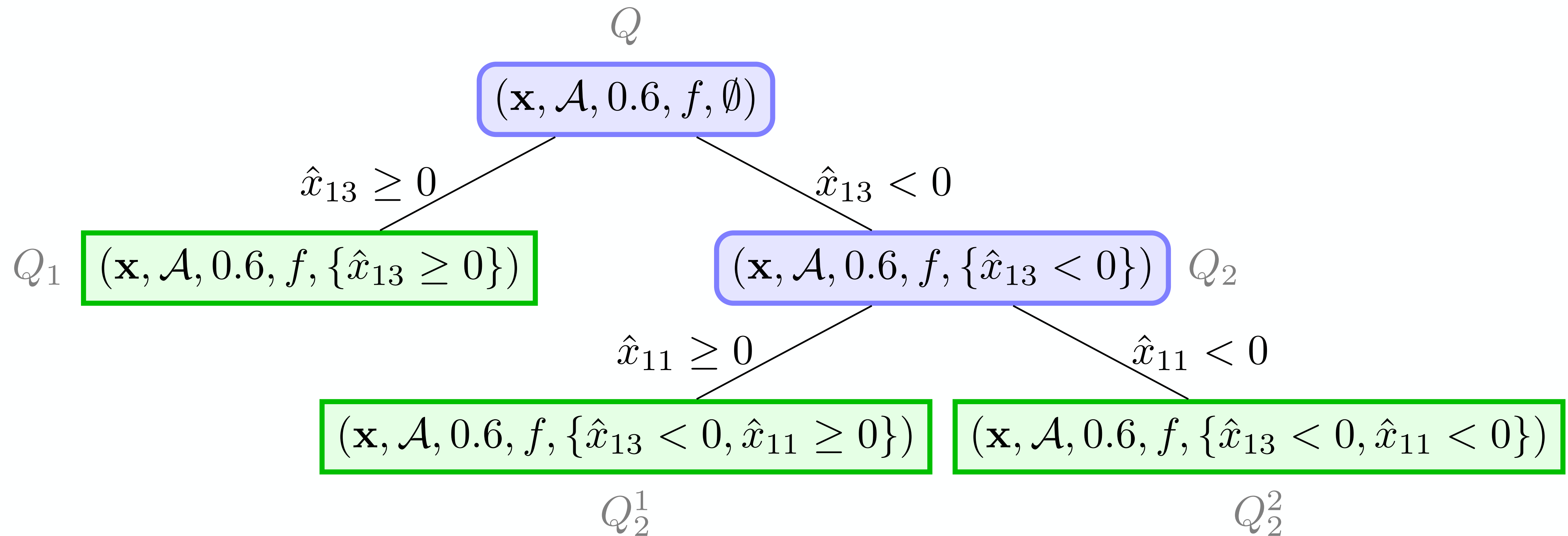
Computing Verifier-Optimal Robust Explanations

FaVeX (Simplified)

```
1: function FAVEX( $f, \mathbf{x}, \epsilon, v, \vec{\mathcal{A}}$ )
2:    $\mathcal{R}_x, \mathcal{U}_x, \mathcal{C}_x \leftarrow \emptyset, \emptyset, \emptyset$ 
3:   batches  $\leftarrow \{\vec{\mathcal{A}}\}$ 
4:   for  $B \in$  batches do
5:     batches  $\leftarrow$  batches  $\setminus \{B\}$ 
6:     if  $|B| > 1$  then  $\triangleright$  batch robustness query
7:       result  $\leftarrow v(\mathbf{x}, \mathcal{R}_x \cup \mathcal{U}_x \cup B, \epsilon, f)$ 
8:       if result == 1 then  $\mathcal{R}_x \leftarrow \mathcal{R}_x \cup B$ 
9:       else
10:         $\vec{\mathcal{A}}_1, \vec{\mathcal{A}}_2 \leftarrow$  HALVE( $\vec{\mathcal{A}}$ )
11:        batches  $\leftarrow$  batches  $\cup \{\vec{\mathcal{A}}_1, \vec{\mathcal{A}}_2\}$ 
12:     else  $\triangleright$  single robustness query
13:       result  $\leftarrow v(\mathbf{x}, \mathcal{R}_x \cup \mathcal{U}_x \cup B, \epsilon, f)$ 
14:       if result == 1 then  $\mathcal{R}_x \leftarrow \mathcal{R}_x \cup B$ 
15:       else if result = -1 then  $\mathcal{C}_x \leftarrow \mathcal{C}_x \cup B$  else  $\mathcal{U}_x \leftarrow \mathcal{U}_x \cup B$ 
16:   return  $\mathcal{U}_x, \mathcal{C}_x$ 
```

Computing Verifier-Optimal Robust Explanations

Branch-and-Bound Verification



Computing Verifier-Optimal Robust Explanations

Branch-and-Bound Verification with Branching Save/Reuse

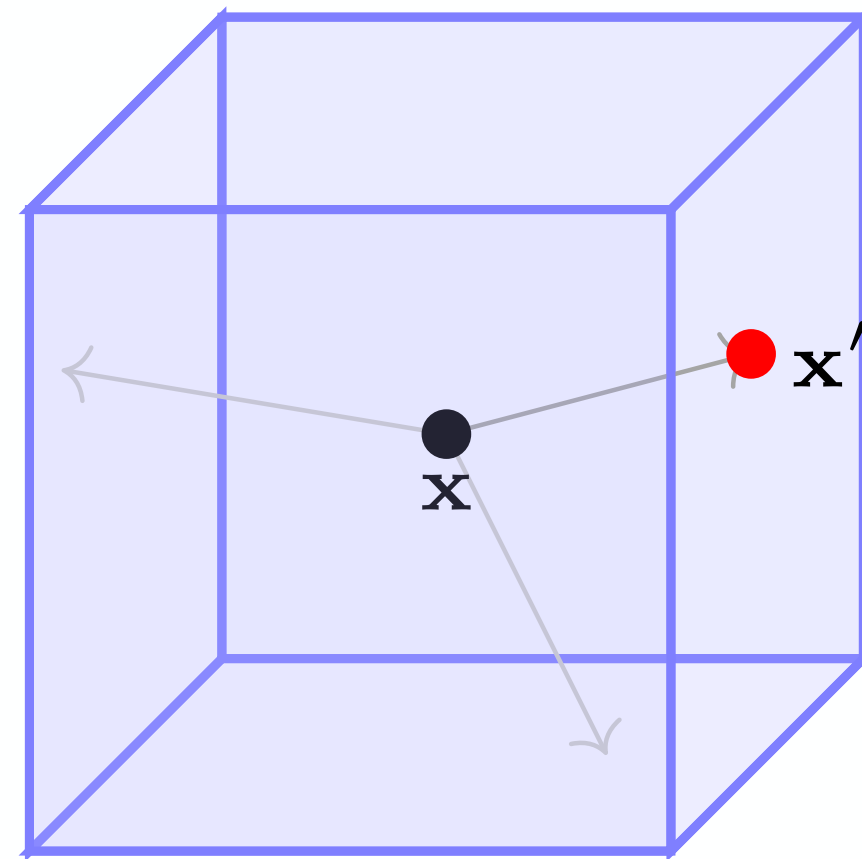
```
1: function BAB( $a, (\mathbf{x}, \mathcal{A}, \epsilon, f), \mathbb{C}$ )  $\triangleright a: \hat{Q} \rightarrow \mathbb{R}, \mathbb{C} \in \mathcal{P}(\mathcal{P}(\mathcal{C}))$ 
2:   for  $C \in \mathbb{C}$  do unresolved  $\leftarrow \{(\mathbf{x}, \mathcal{A}, \epsilon, f, C)\}$   $\triangleright Reuse$ 
3:    $\mathbb{C} \leftarrow \emptyset$ 
4:   for  $Q \in$  unresolved do
5:     if CEX( $\mathbf{x}, \mathcal{A}, \epsilon, f$ ) then
6:       for  $Q' \in$  unresolved do  $\mathbb{C} \leftarrow \mathbb{C} \cup \{Q'_C\}$   $\triangleright Save$  unresolved
7:       return  $-1, \mathbb{C}$ 
8:   unresolved  $\leftarrow$  unresolved  $\setminus \{Q\}$ 
9:   if  $a(Q) > 0$  then  $\mathbb{C} \leftarrow \mathbb{C} \cup \{Q_C\}$   $\triangleright Save$   $Q$ 
10:  else if  $\neg$ TIMEOUT then
11:     $Q_1, Q_2 \leftarrow$  SPLIT( $Q$ )
12:    unresolved  $\leftarrow$  unresolved  $\cup \{Q_1, Q_2\}$ 
13:  else
14:    for  $Q' \in$  unresolved do  $\mathbb{C} \leftarrow \mathbb{C} \cup \{Q'_C\}$   $\triangleright Save$  unresolved
15:    return  $0, \mathbb{C}$ 
16:  return  $1, \mathbb{C}$ 
```

Computing Verifier-Optimal Robust Explanations

Restricted-Space Counterfactual Search

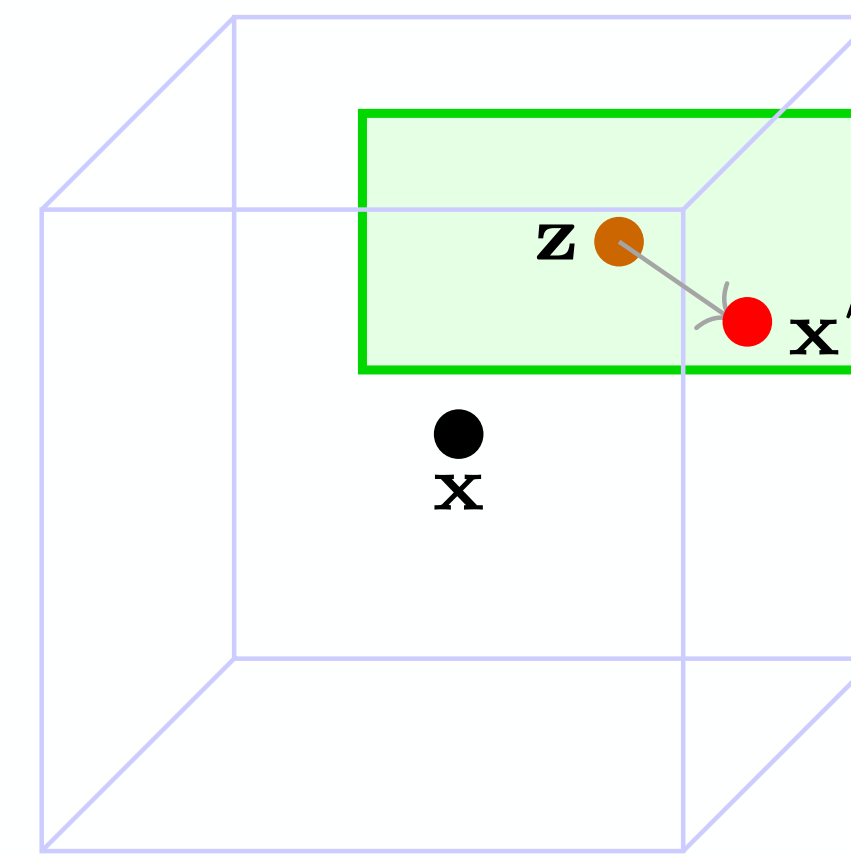
Full-Space Search

$$B_{\mathcal{A} \cup \mathcal{B}}^\epsilon(\mathbf{x})$$



Restricted-Space Search

$$\text{subset of } B_{\mathcal{A} \cup \mathcal{B}}^\epsilon(\mathbf{x})$$



■ **Figure 7** Illustration of full-space vs. restricted-space' counterfactual search. Full-space search (left) explores the entire perturbation region $B_{\mathcal{A} \cup \mathcal{B}}^\epsilon(\mathbf{x})$. Restricted-space search (right) begins from the output \mathbf{z} of the previous search (which may not be a counterfactual) and varies only the newly added feature(s), thus exploring a much smaller subset of $B_{\mathcal{A} \cup \mathcal{B}}^\epsilon(\mathbf{x})$.

Computing Verifier-Optimal Robust Explanations

CNN-3

	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	Counterfactuals	Unknowns	Time	Counterfactuals	Unknowns	Time
Standard	0.00	247.80	45m	0.00	461.00	2h 30m
Verifier-Optimal	160.30	94.40	10m	210.40	251.70	19m

$$\begin{aligned} &160.30 + \\ &94.40 = \\ &254.70 \end{aligned}$$

$$\begin{aligned} &210.40 + \\ &251.70 = \\ &462.10 \end{aligned}$$

Computing Verifier-Optimal Robust Explanations

CNN-3

Ablations

■ **Table 6** On CNN-3, FAVEX computes OVAL-optimal robust explanations faster than previous computation strategies [65, 64] while at the same time finding more counter-factuals. Results are averaged over the first 10 images of the test set. Bold entries correspond to the smallest runtime and the largest number of provided counterfactuals.

Method	MNIST, $\epsilon = 0.25$			CIFAR-10, $\epsilon = \frac{8}{255}$		
	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]
SEQUENTIAL	129.10	125.50	1164.89	209.30	253.10	1287.96
BINARY SEARCH	134.20	120.40	1026.67	209.60	252.70	1367.25
BINS + INCR	135.20	119.50	1040.82	209.60	252.50	1374.40
BINS + INCR + RSA	157.70	96.60	727.94	211.90	250.20	1334.00
FAVEX	160.30	94.40	612.75	210.40	251.70	1150.82

Computing Verifier-Optimal Robust Explanations

CNN-7

	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=16/255$		
	Counterfactuals	Unknowns	Time	Counterfactuals	Unknowns	Time
Standard	0.00	452.00	239m	0.00	730.67	7h 5m
Verifier-Optimal	207.33	249.33	74m	467.00	266.33	1h 49m

207.33 +
249.33 =
456.66

467.00 +
266.33 =
733.33

Computing Verifier-Optimal Robust Explanations

CNN-7

Ablations

■ **Table 7** On CNN-7, FAVEX computes OVAL-optimal robust explanations faster than previous computation strategies [65, 64] while at the same time finding more counter-factuals. Results are averaged over the first 3 images of the test set. Bold entries correspond to the smallest runtime and the largest number of provided counterfactuals.

Method	MNIST, $\epsilon = 0.25$			CIFAR-10, $\epsilon = \frac{16}{255}$		
	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]
SEQUENTIAL	142.33	315.00	7979.61	464.00	269.33	6868.53
BINARY SEARCH	148.33	308.67	9165.06	464.33	269.00	8284.12
BINS + INCR	151.33	305.00	9116.13	467.00	266.33	8231.01
BINS + INCR + RSA	196.33	260.33	5406.70	468.00	265.33	8304.52
FAVEX	207.33	249.33	4446.53	467.00	266.33	6532.98

Computing Verifier-Optimal Robust Explanations

Traversal Strategies

■ **Table 10** Comparison of different traversal strategies for OVAL-optimal robust explanations on CNN-3. Results are averaged over the first 10 images of the test set. Bold entries correspond to the smallest $\mathcal{C}_x \cup \mathcal{U}_x$.

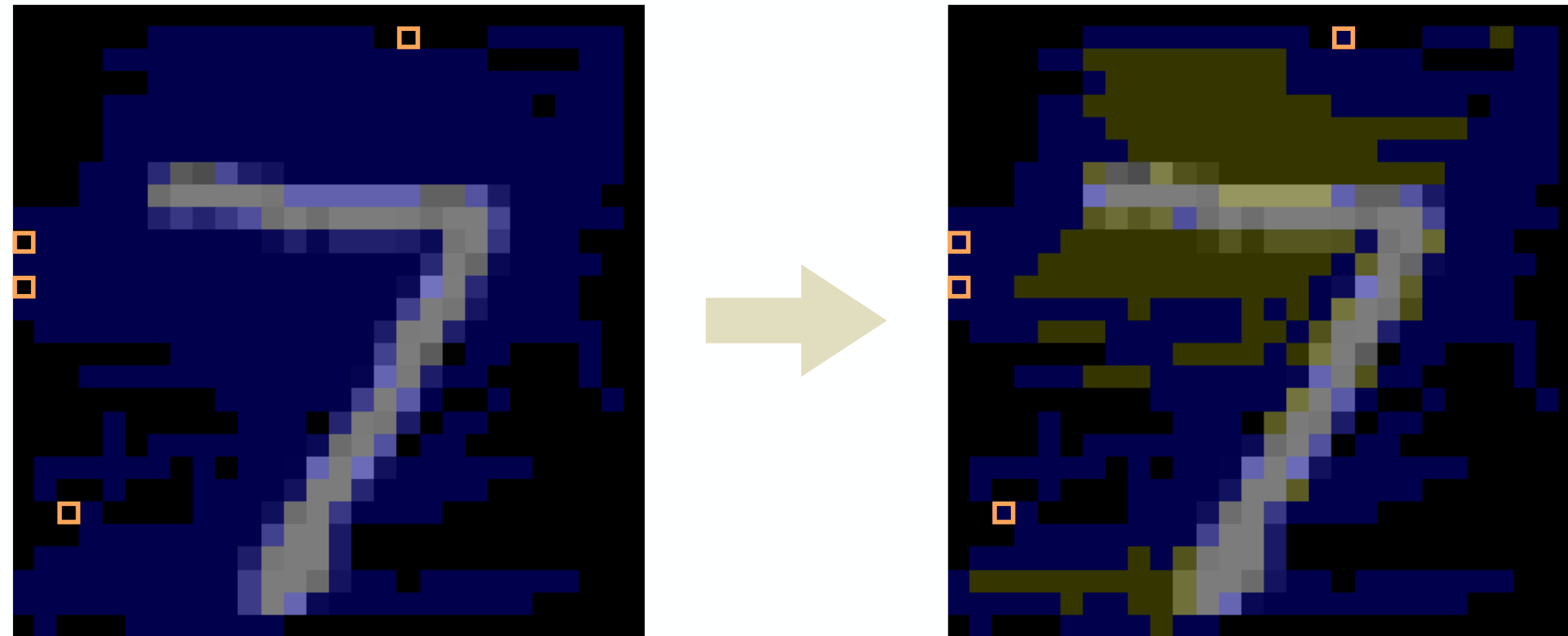
Method	MNIST, $\epsilon = 0.25$				CIFAR-10, $\epsilon = \frac{8}{255}$			
	$ \mathcal{C}_x \cup \mathcal{U}_x $	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]	$ \mathcal{C}_x \cup \mathcal{U}_x $	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]
VERIX	282.80	160.50	122.30	974.35	765.50	437.20	328.30	1925.09
VERIX+	247.40	155.20	92.20	576.67	468.90	262.00	206.90	915.45
α -FAVEX	289.90	181.20	108.70	696.31	462.10	210.40	251.70	1150.82
FAVEX-IBP	254.70	160.30	94.40	612.75	466.40	250.90	215.50	941.95

■ **Table 11** Comparison of different traversal strategies for OVAL-optimal robust explanations on CNN-7. Results are averaged over the first 3 images of the test set. Bold entries correspond to the smallest $\mathcal{C}_x \cup \mathcal{U}_x$.

Method	MNIST, $\epsilon = 0.25$				CIFAR-10, $\epsilon = \frac{16}{255}$			
	$ \mathcal{C}_x \cup \mathcal{U}_x $	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]	$ \mathcal{C}_x \cup \mathcal{U}_x $	$ \mathcal{C}_x $	$ \mathcal{U}_x $	time [s]
VERIX	547.00	123.33	423.67	7763.91	945.00	728.67	216.33	4974.69
VERIX+	429.33	196.67	232.67	4394.67	735.67	522.00	213.67	5105.86
α -FAVEX	552.33	234.67	317.67	5344.59	733.33	467.00	266.33	6532.98
FAVEX-IBP	456.67	207.33	249.33	4446.53	729.33	512.67	216.67	5099.98

Verifier-Optimal Robust Explanations

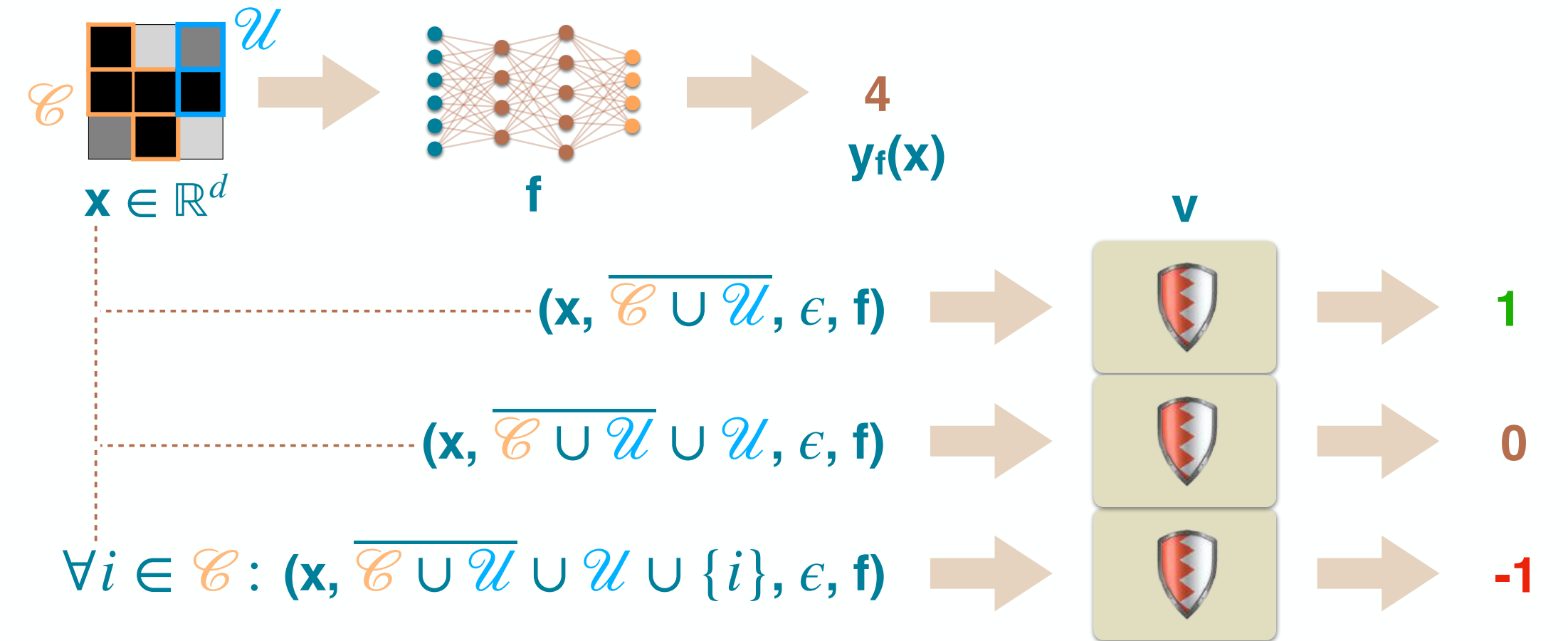
Weak Abductive Explanations



x

Verifier-Optimal Robust Explanations

Weak Abductive Explanations



11

Computing Verifier-Optimal Robust Explanations

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

ADD TO $\overline{\mathcal{C} \cup \mathcal{U}}$

LOCAL ROBUSTNESS IN $B_{\overline{\mathcal{C} \cup \mathcal{U}}}^\epsilon(\mathbf{x})$

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
x_1	•	✓	✓	✓	✓	✓	✓	✓	✓
x_2		•	✓	✓	✓	✓	✓	✓	✓
x_3			•	✓	✓	✓	✓	✓	✓
x_4				•	✓	✓	✓	✓	✓
x_5					•	✓	✓	✓	✓
x_6						•	✗	✗	✗
x_7							•	U	U
x_8								•	✗
x_9									•

12

Computing Verifier-Optimal Robust Explanations

CNN-3

	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	Counterfactuals	Unknowns	Time	Counterfactuals	Unknowns	Time
Standard	0.00	247.80	45m	0.00	461.00	2h 30m
Verifier-Optimal	160.30	94.40	10m	210.40	251.70	19m

160.30 +
94.40 =
254.70

210.40 +
251.70 =
462.10

THANKS!