

Faster Verified Explanations for Neural Networks

Alessandro De Palma, Greta Dolcetti and Caterina Urban

40th European Conference on Object-Oriented Programming (ECOOP 2026)

Neural Networks in High-Stakes Systems



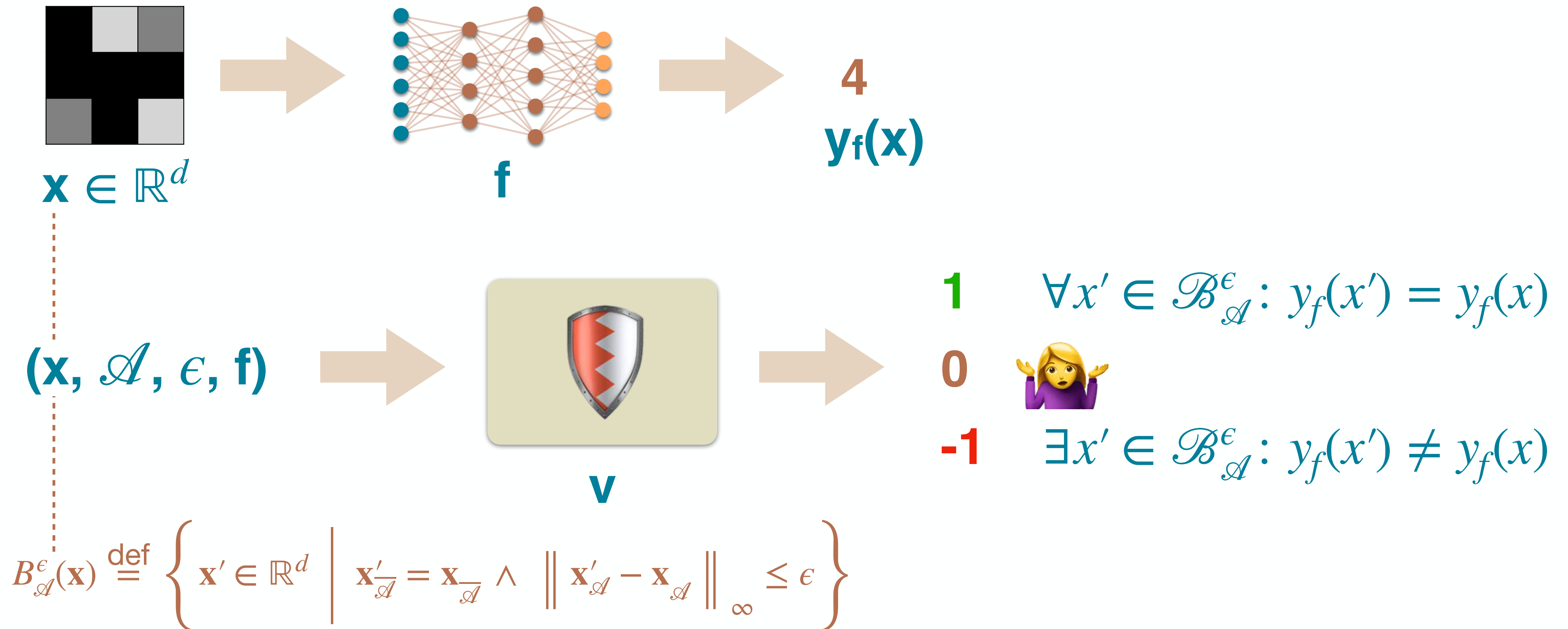
Safety-Critical Applications



Socio-Economic Applications

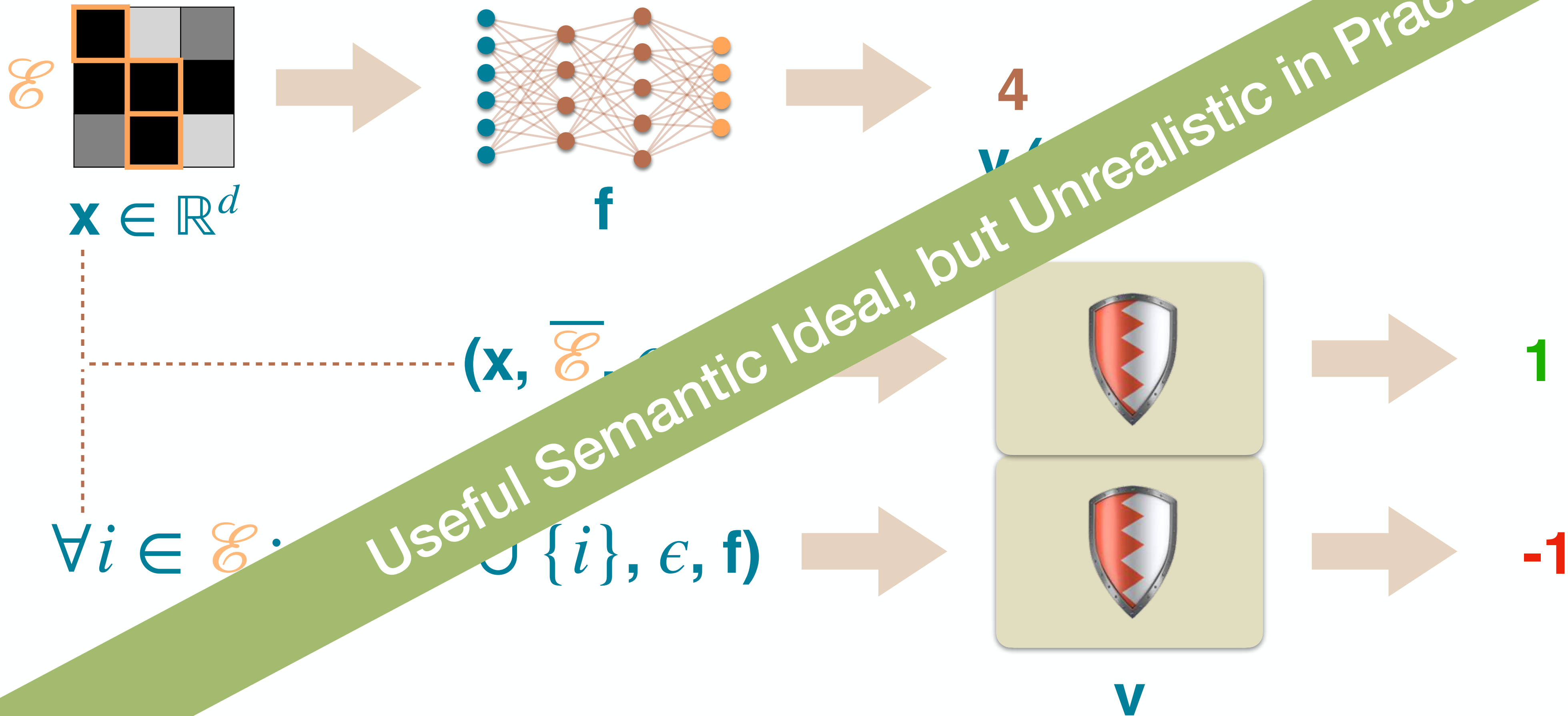
Neural Network Verification

LOCAL ROBUSTNESS



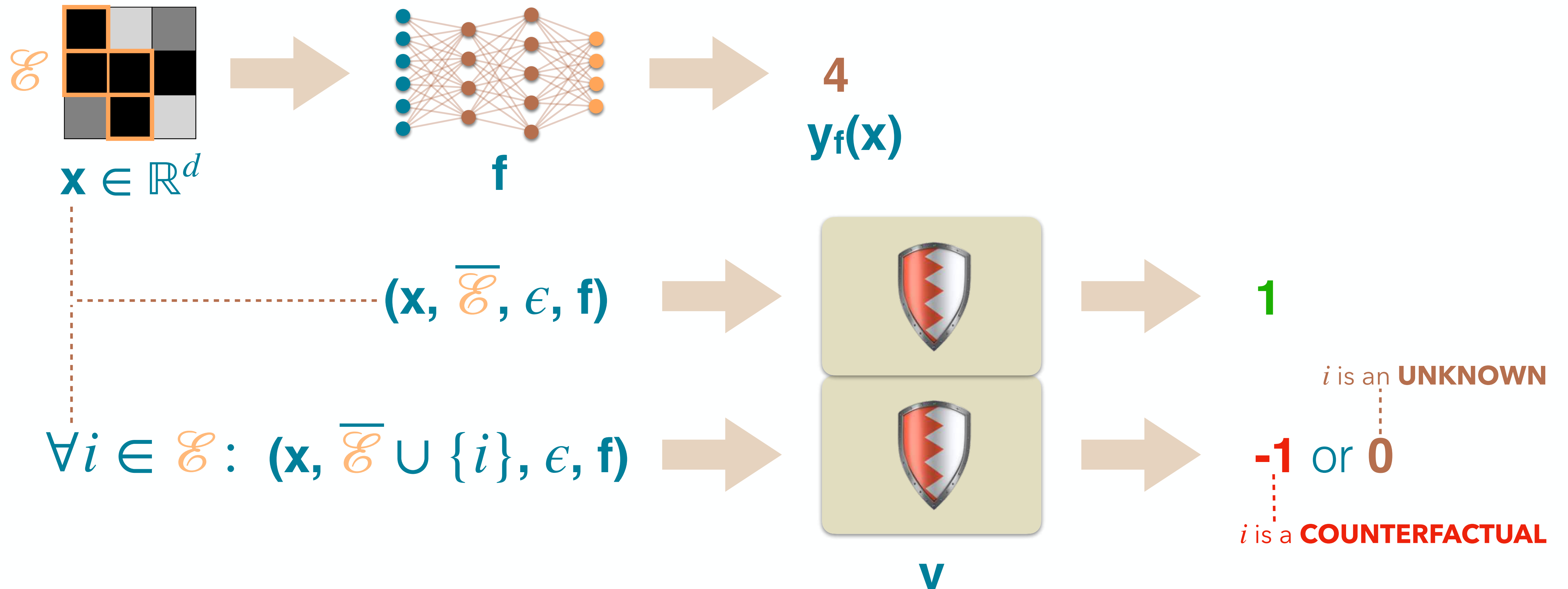
Optimal Robust Explanations

ABDUCTIVE EXPLANATIONS (AX_{ps})



Optimal Robust Explanations

WEAK ABDUCTIVE EXPLANATIONS



Optimal Robust Explanations

WEAK ABDUCTIVE EXPLANATIONS



Computing ~~Optimal~~ Robust Explanations

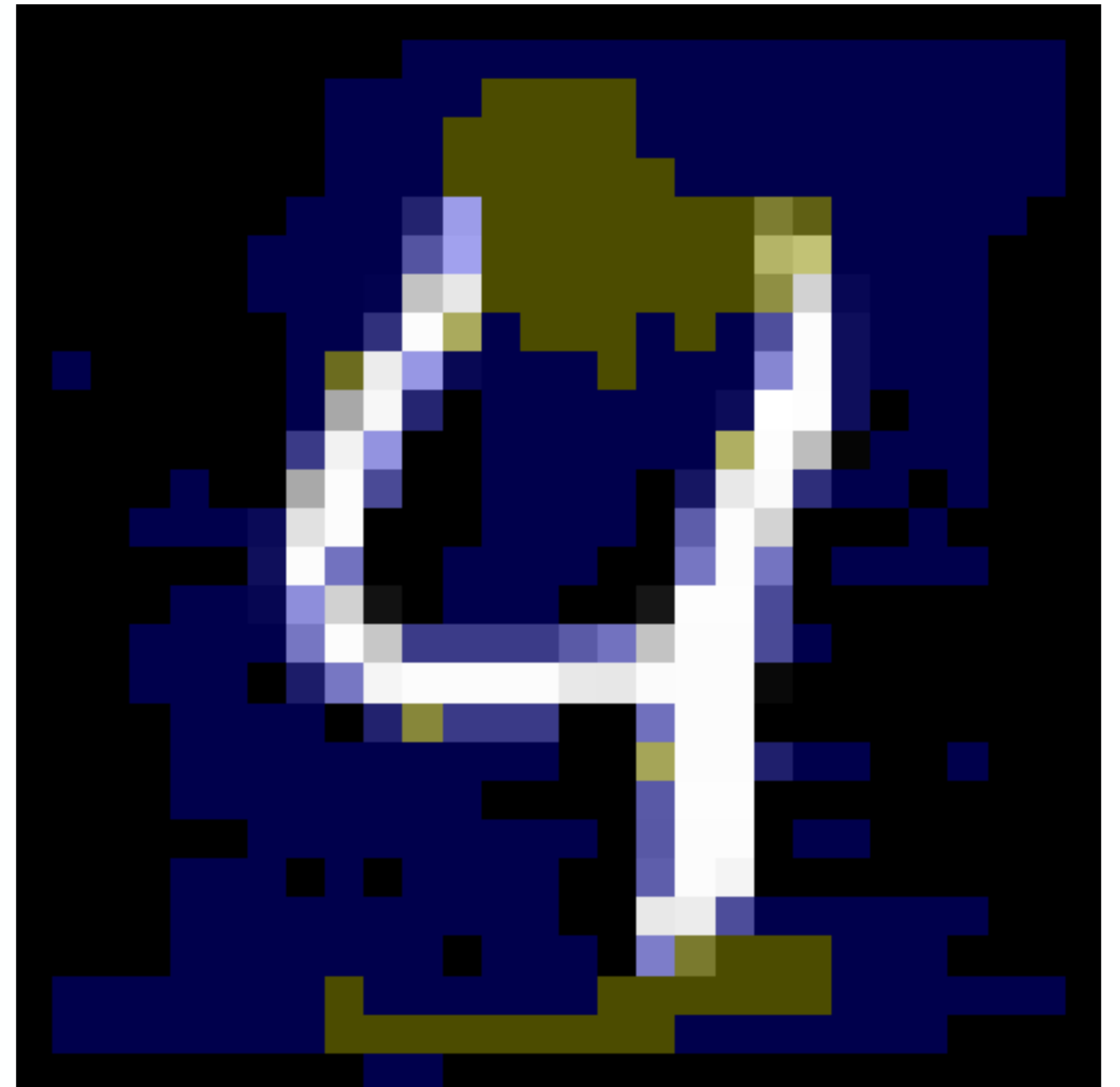
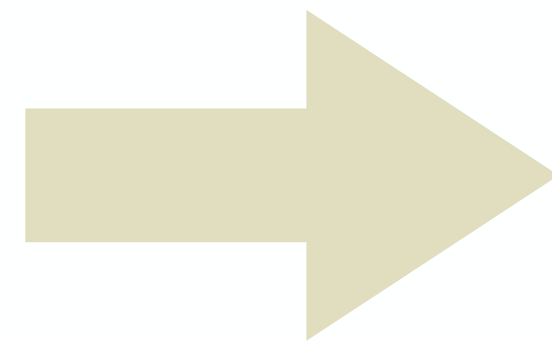
FINDING COUNTERFACTUALS RAPIDLY BECOMES INFEASIBLE

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

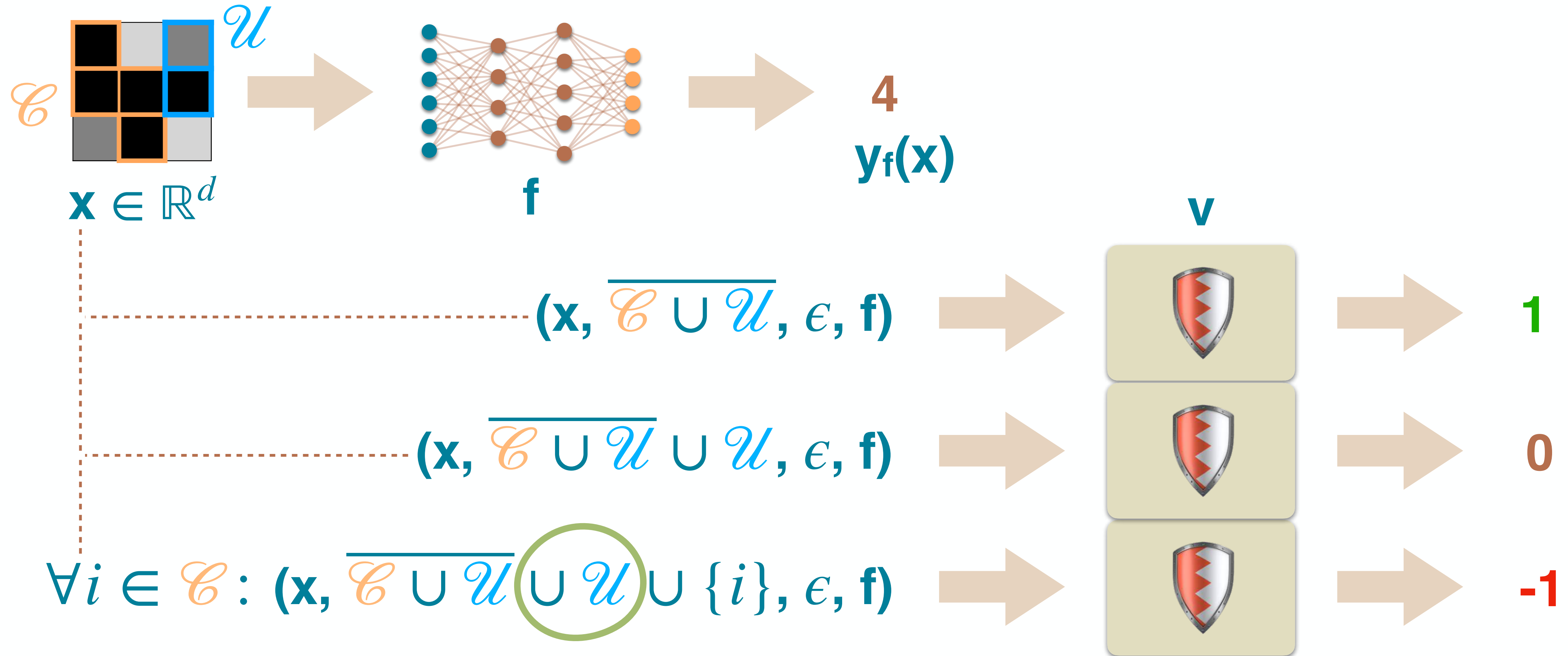
Verifier-Optimal Robust Explanations

OUR 1ST CONTRIBUTION



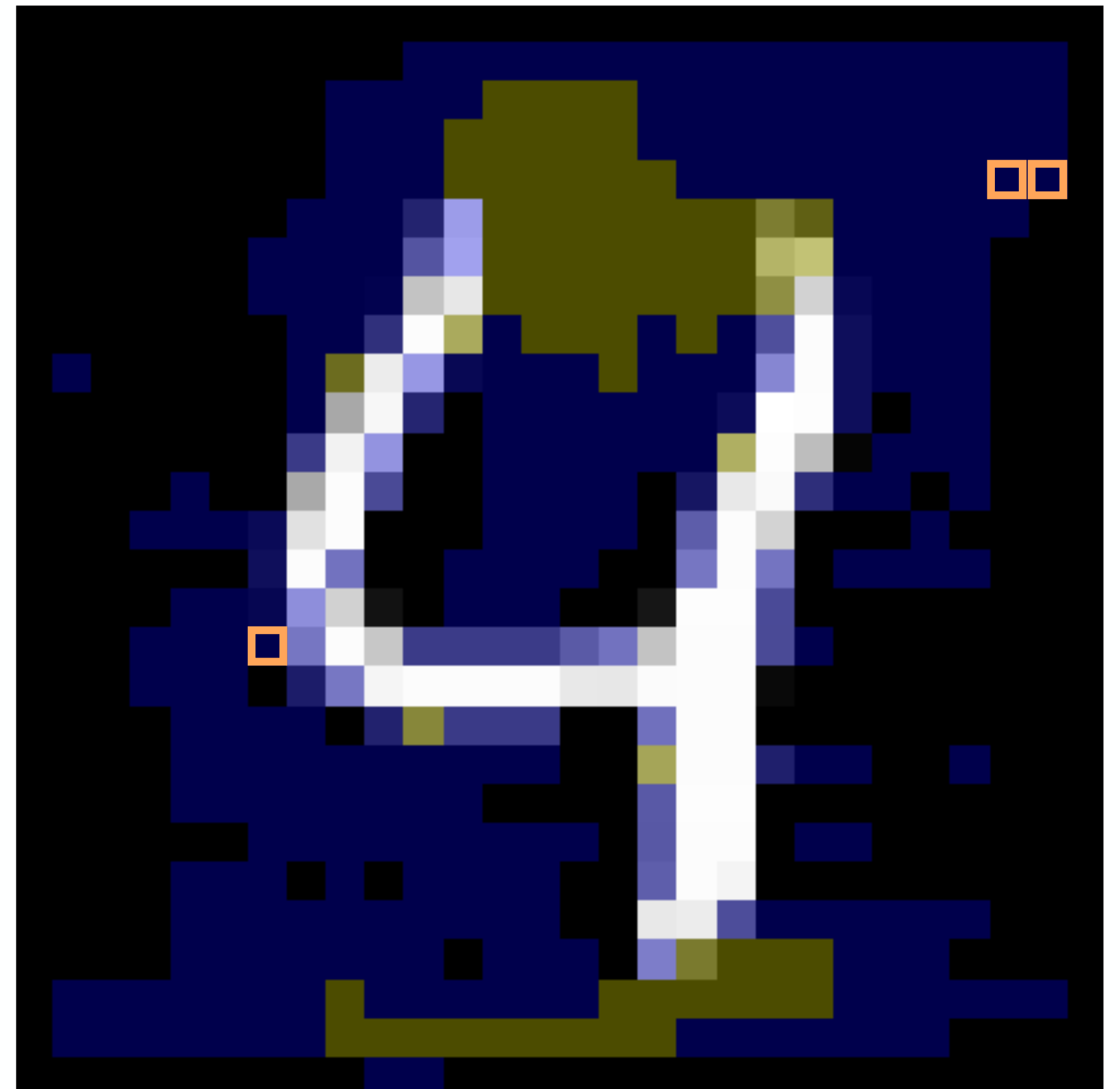
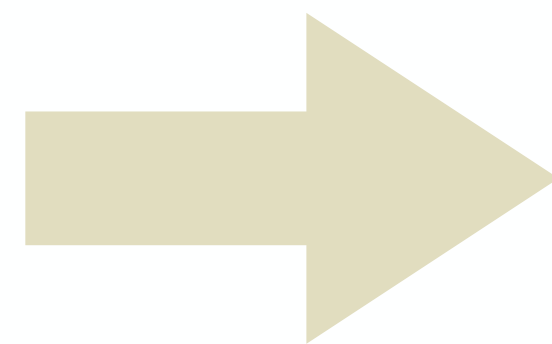
Verifier-Optimal Robust Explanations

WEAK ABDUCTIVE EXPLANATIONS



Verifier-Optimal Robust Explanations

OUR 1ST CONTRIBUTION

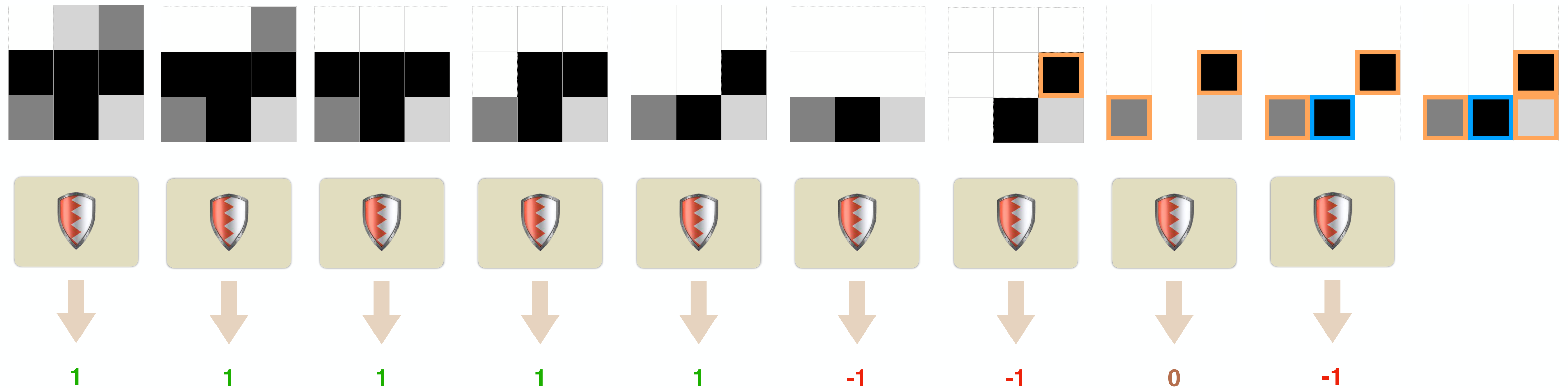


Computing Verifier-Optimal Robust Explanations

DROP (I.E., FREE) INPUT DIMENSIONS WHILE AX_p CONDITION HOLDS

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

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



Computing Robust Explanations

CNN-3

Robust Explanations	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	Counterfactuals $ \mathcal{E} $	Unknowns $ \mathcal{U} $	Time	Counterfactuals $ \mathcal{E} $	Unknowns $ \mathcal{U} $	Time
Optimal	0.00	247.80	45m	0.00	461.00	2h 30m
Verifier-Optimal	129.10	125.50	19m	209.30	253.10	21m

$$129.10 + 125.50 = 254.60$$

$$209.30 + 253.10 = 462.40$$

Computing Robust Explanations

CNN-7

Robust Explanations	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	Counterfactuals $ \mathcal{C} $	Unknowns $ \mathcal{U} $	Time	Counterfactuals $ \mathcal{C} $	Unknowns $ \mathcal{U} $	Time
Optimal	0.00	452.00	3h 59m	0.00	730.67	7h 5m
Verifier-Optimal	142.33	315.00	2h 13m	464.00	269.33	1h 54m

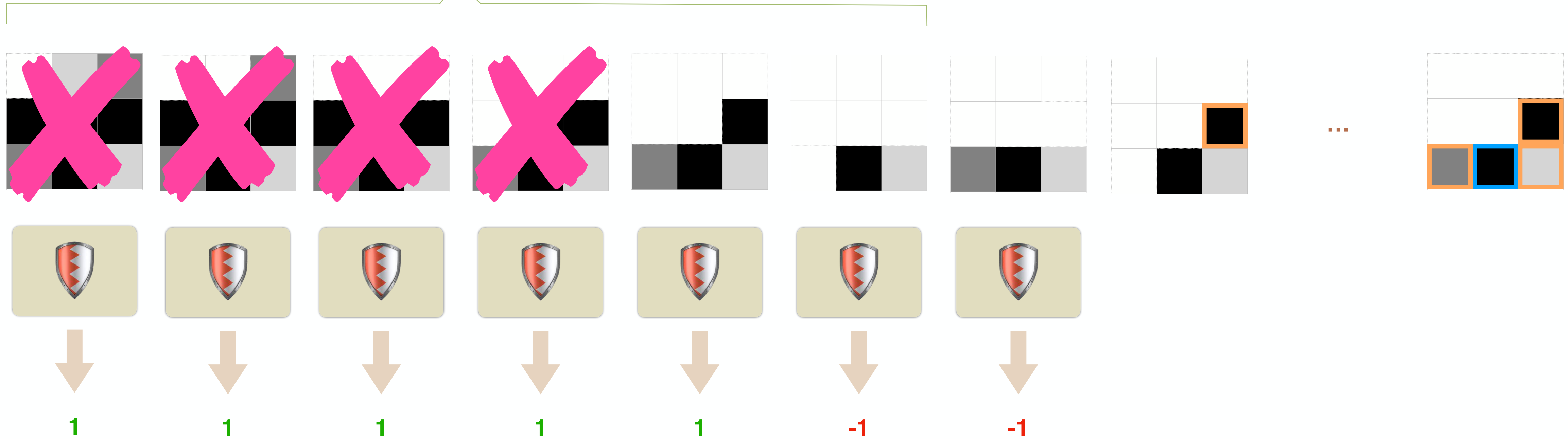
$$142.33 + 315.00 = 457.33$$

$$464.00 + 269.33 = 733.33$$

FaVeX

OUR 2ND CONTRIBUTION

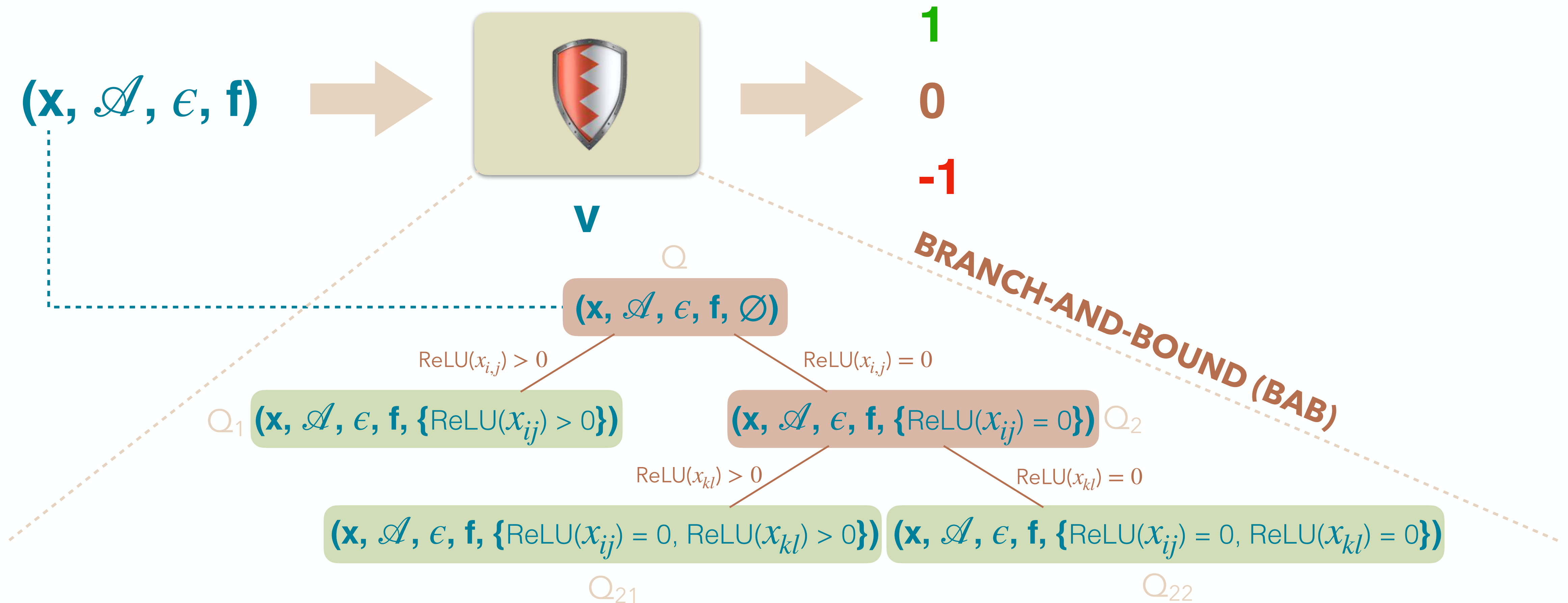
BATCH QUERIES



SINGLE QUERIES

FaVeX

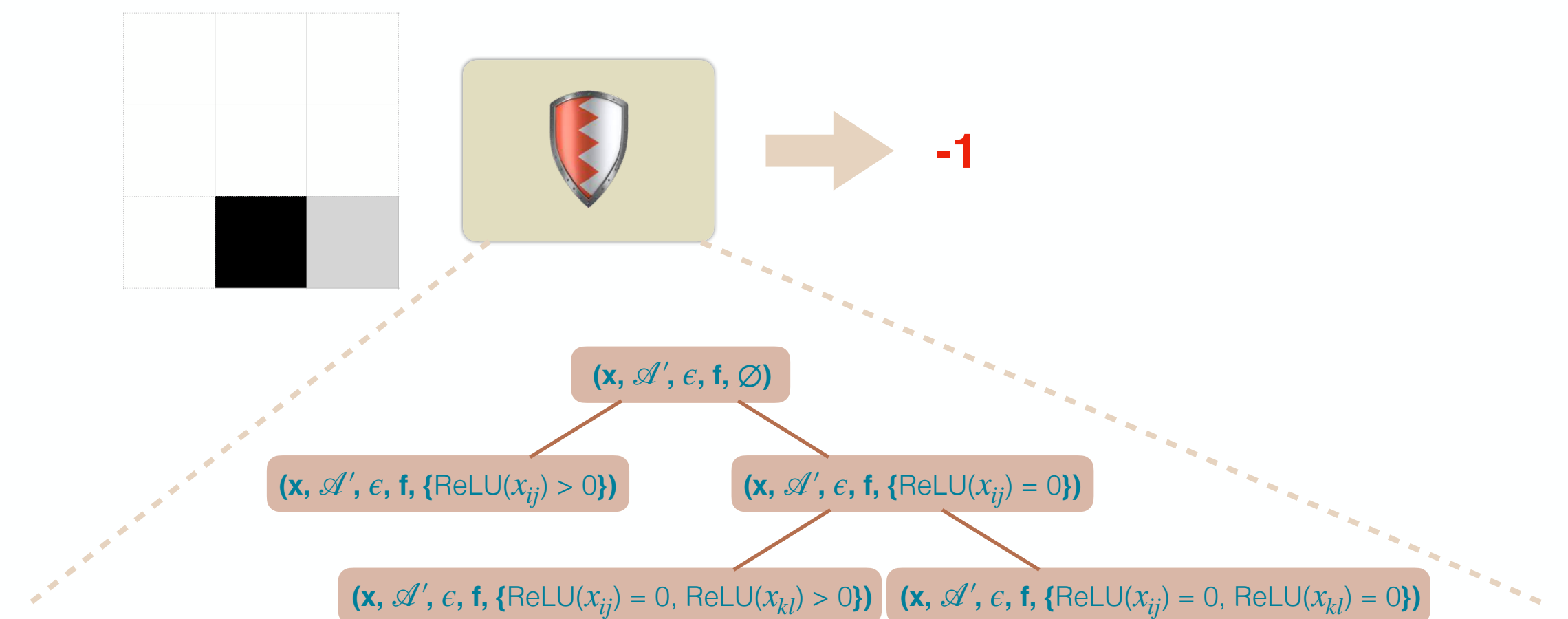
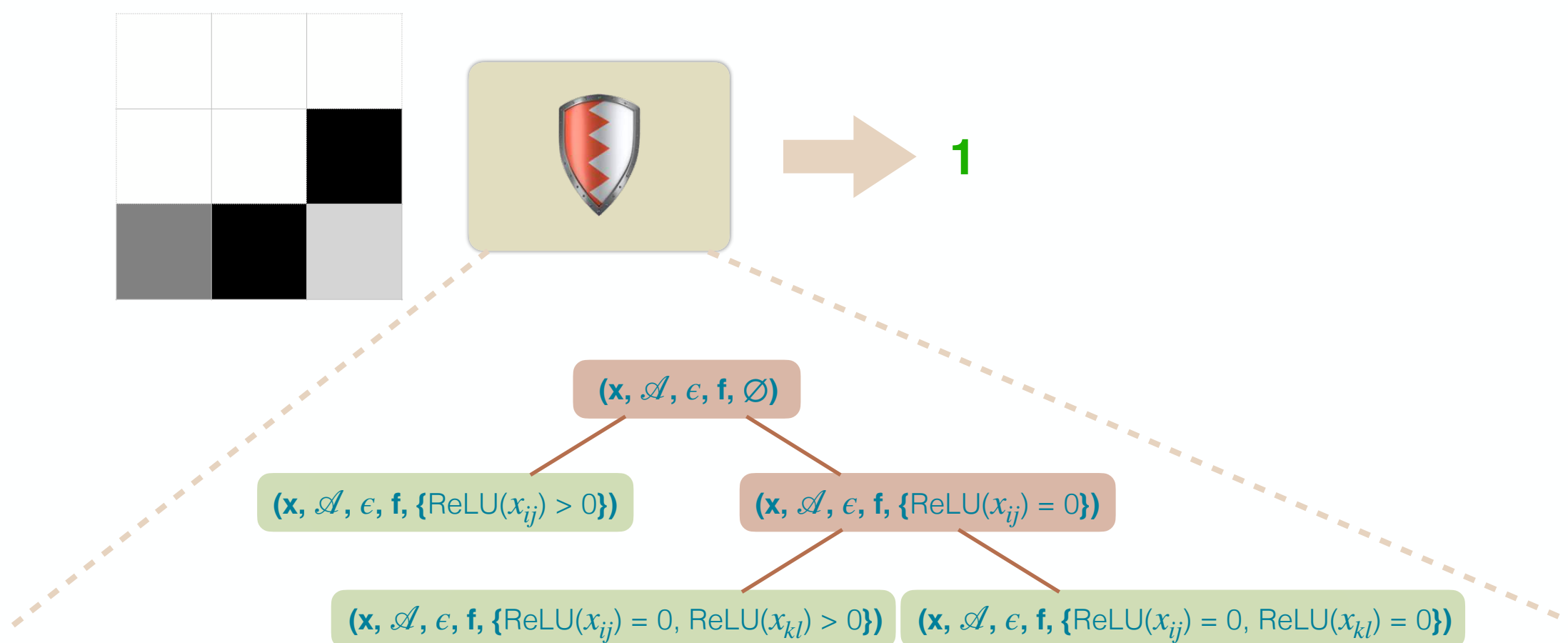
OUR 2ND CONTRIBUTION



FaVeX

OUR 2ND CONTRIBUTION

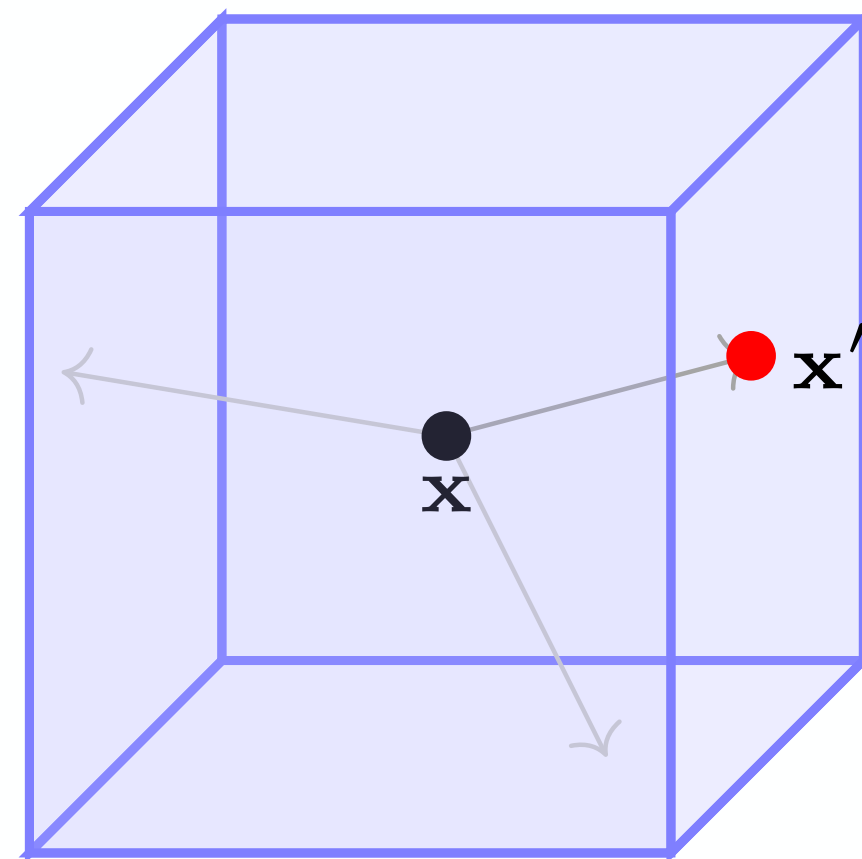
INCREMENTAL BAB VERIFICATION



FaVeX

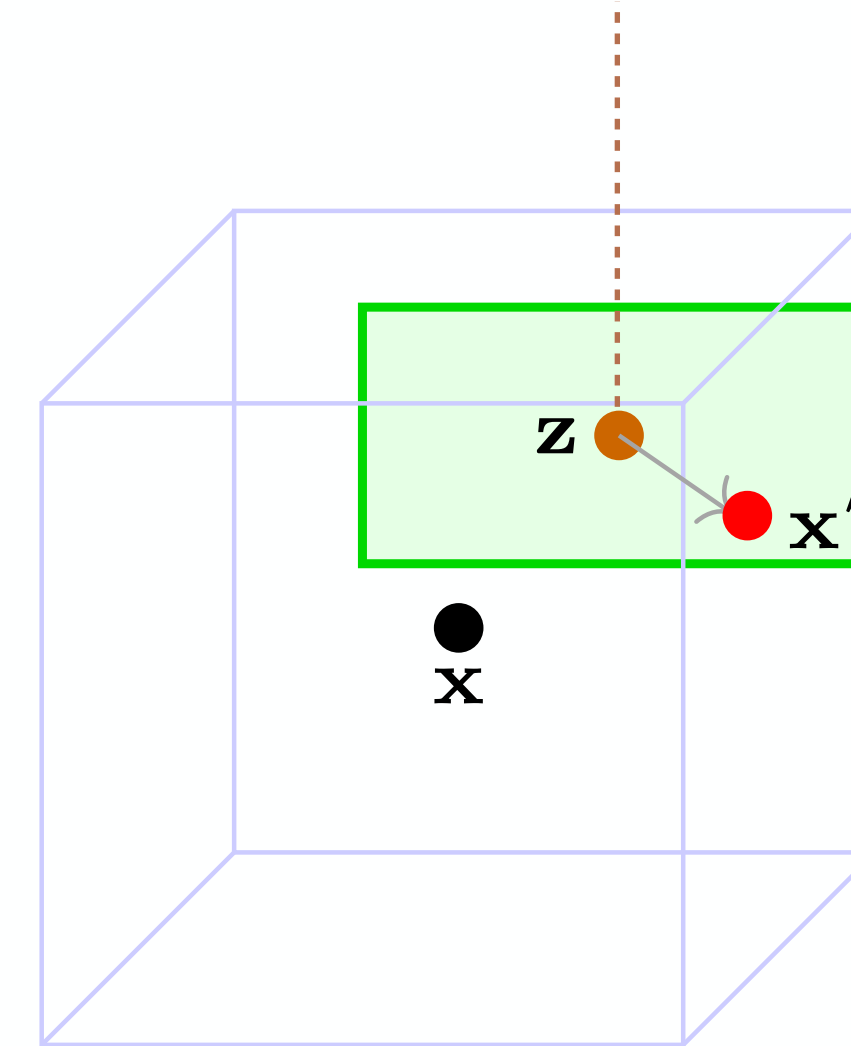
OUR 2ND CONTRIBUTION

output of the previous counterfactual search



$$B_{\mathcal{C} \cup \mathcal{U}}^{\epsilon}(\mathbf{x})$$

Full-Space Search



subset of $B_{\mathcal{C} \cup \mathcal{U}}^{\epsilon}(\mathbf{x})$

Reduced Space Search

Computing Verifier-Optimal Robust Explanations

CNN-3

Verifier-Optimal Explanations	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	$ \mathcal{E} $	$ \mathcal{U} $	Time	$ \mathcal{E} $	$ \mathcal{U} $	Time
Single Queries	129.10	125.50	19m	209.30	253.10	21m
FaVeX	160.30	94.40	10m	210.40	251.70	19m

Computing Verifier-Optimal Robust Explanations

CNN-7

Verifier-Optimal Explanations	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=16/255$		
	$ \mathcal{E} $	$ \mathcal{U} $	Time	$ \mathcal{E} $	$ \mathcal{U} $	Time
Single Queries	142.33	315.00	2h 13m	464.00	269.33	1h 54m
FaVeX	207.33	249.33	1h 14m	467.00	266.33	1h 49m

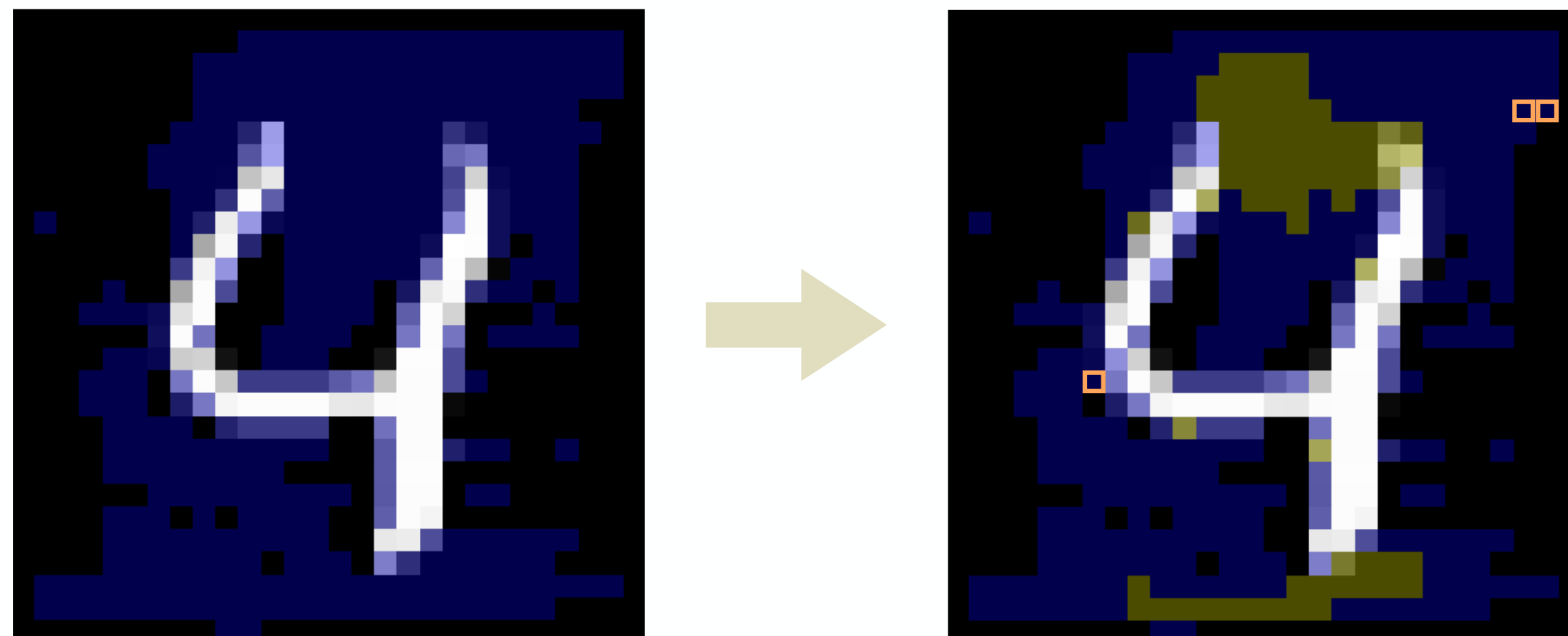
Computing Verifier-Optimal Robust Explanations

TRAVERSAL STRATEGIES

Model	Traversal	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
		$ \mathcal{E} $	$ \mathcal{U} $	Time	$ \mathcal{E} $	$ \mathcal{U} $	Time
CNN-3	VeriX	160.50	122.30	16m	437.20	328.30	32m
	VeriX+	155.20	92.20	10m	262.00	206.90	15m
	α -FAVEX	181.20	108.70	12m	210.40	251.70	19m
	FaVeX-IBP	160.30	94.40	10m	250.90	215.50	16m
Model	Traversal	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=16/255$		
		$ \mathcal{E} $	$ \mathcal{U} $	Time	$ \mathcal{E} $	$ \mathcal{U} $	Time
CNN-7	VeriX	123.33	423.67	2h 9m	728.67	216.33	1h 23m
	VeriX+	196.67	232.67	1h 13m	522.00	213.67	1h 25m
	α -FAVEX	234.67	317.67	1h 29m	467.00	266.33	1h 49m
	FaVeX-IBP	207.33	249.33	1h 14m	512.67	216.67	1h 25m

Verifier-Optimal Robust Explanations

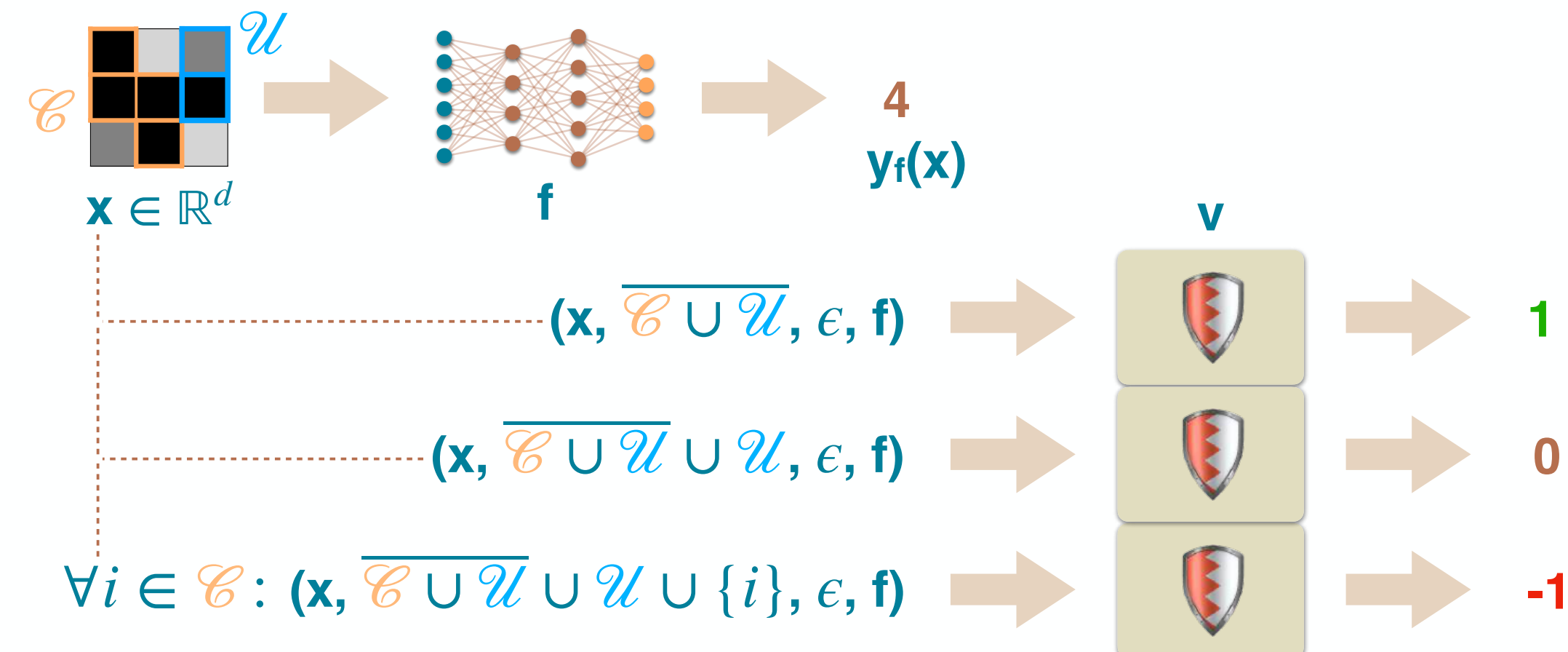
OUR 1ST CONTRIBUTION



10

Verifier-Optimal Robust Explanations

WEAK ABDUCTIVE EXPLANATIONS



9

FaVeX

OUR 2ND CONTRIBUTION

Robust Explanations	MNIST, $\epsilon=0.25$			CIFAR-10, $\epsilon=8/255$		
	$ \mathcal{C} $	$ \mathcal{U} $	Time	$ \mathcal{C} $	$ \mathcal{U} $	Time
Optimal	0.00	452.00	3h 59m	0.00	730.67	7h 5m
Verifier-Optimal (Single Queries)	142.33	315.00	2h 13m	464.00	269.33	1h 54m
Verifier-Optimal (FaVeX)	207.33	249.33	1h 14m	467.00	266.33	1h 49m



OUR PAPER



OUR TOOL

THANKS!