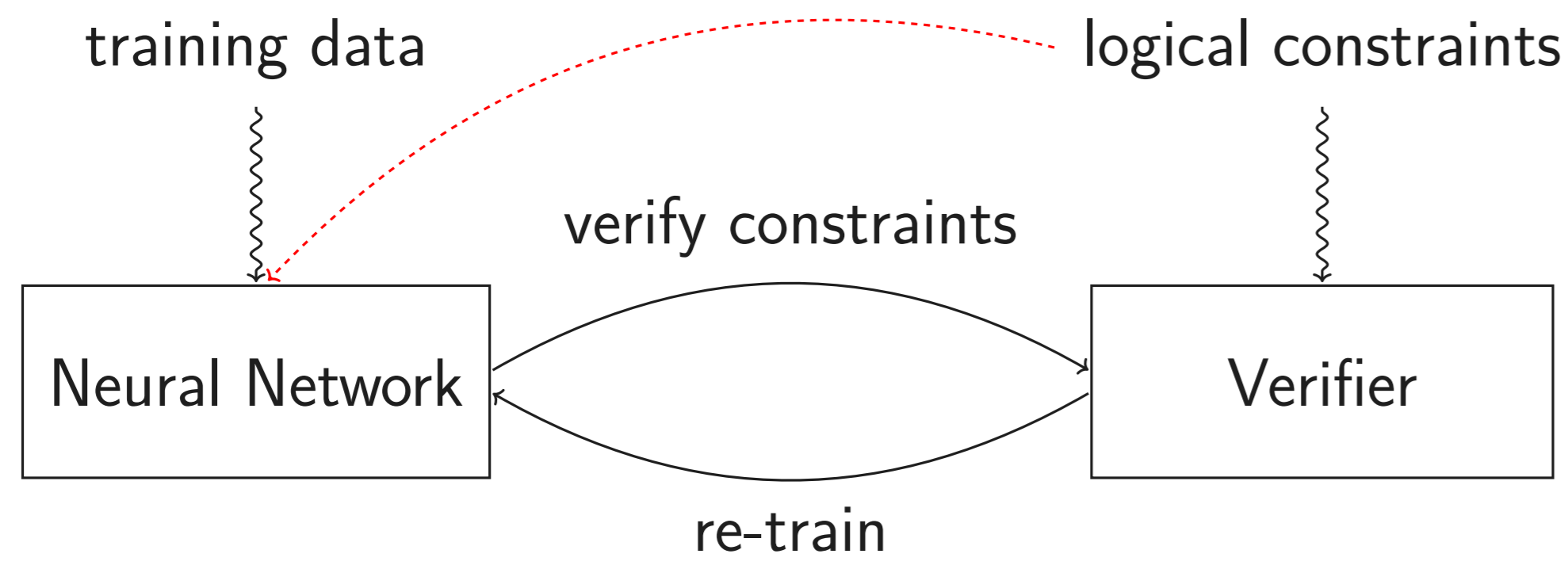


1. Motivation: Correct-by-construction ML models

Task: train a neural network \mathcal{N} to satisfy constraint ϕ .

Train: iteratively update weights using gradient descent and loss.

Verify: with α, β -CROWN, Marabou, NNV, ERAN, ...



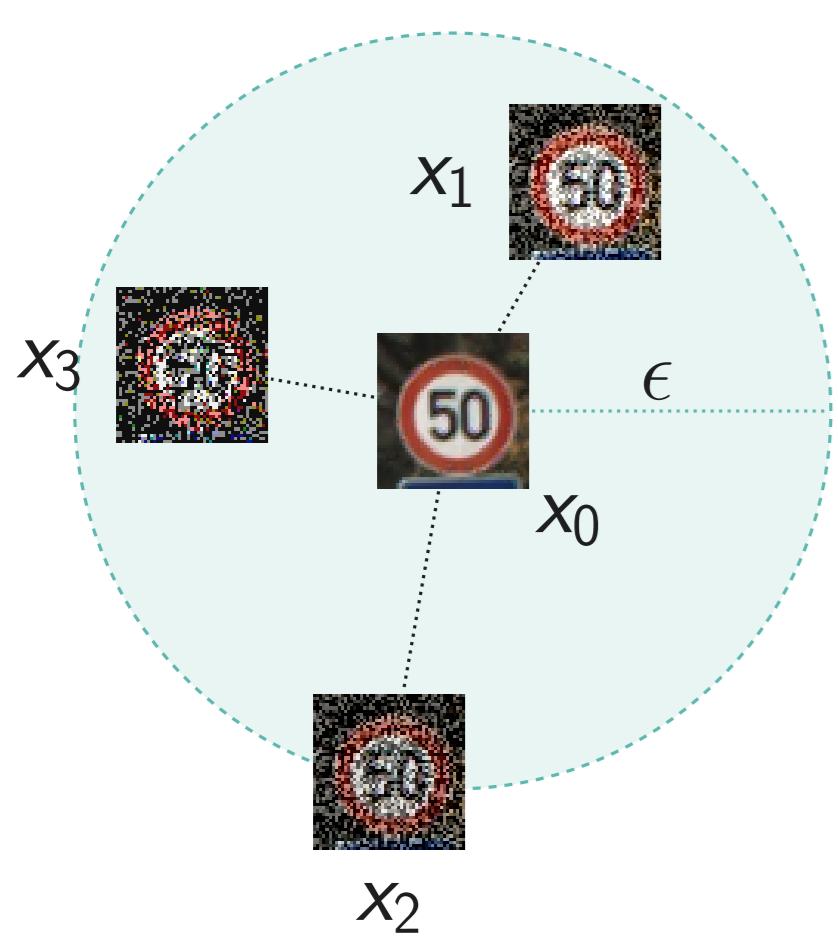
2. Background: Property-driven ML

Standard ML: Given data \mathbf{x} , target \mathbf{y} , and prediction loss \mathcal{L}_P , obtain optimal network weights θ^+ by

$$\theta^+ = \arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \mathcal{L}_P(\mathbf{x}, \mathbf{y}).$$

Insight from DL2 (Deep Learning with Differentiable Logic):

Approximate constraints of the form $\phi: \forall \mathbf{x}. P(\mathbf{x}) \rightarrow Q(\mathbf{x})$ by finding a counterexample \mathbf{x}^* that does *not* satisfy $Q(\mathbf{x})$ in the input space induced by $P(\mathbf{x})$.



- Approximate counterexample *outside* of training set using PGD:

$$\mathbf{x}^* = \arg \max_{\mathbf{x}' \in \mathcal{S}} \mathcal{L}_{\phi}(\mathbf{x}, \mathbf{x}', \mathbf{y})$$

- Use this counterexample in training:

$$\theta^+ = \arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\lambda \mathcal{L}_P(\mathbf{x}, \mathbf{y}) + (1 - \lambda) \mathcal{L}_{\phi}(\mathbf{x}, \mathbf{x}^*, \mathbf{y})].$$

prediction loss logical constraint loss

3. Background: Differentiable Logics

- DL2:** $\llbracket \cdot \rrbracket_{DL2}: \Phi \rightarrow [0, \infty)$, where $\llbracket \top \rrbracket_{DL2} = 0$, and

$$\llbracket x \leq y \rrbracket_{DL2} := \max\{x - y, 0\}$$

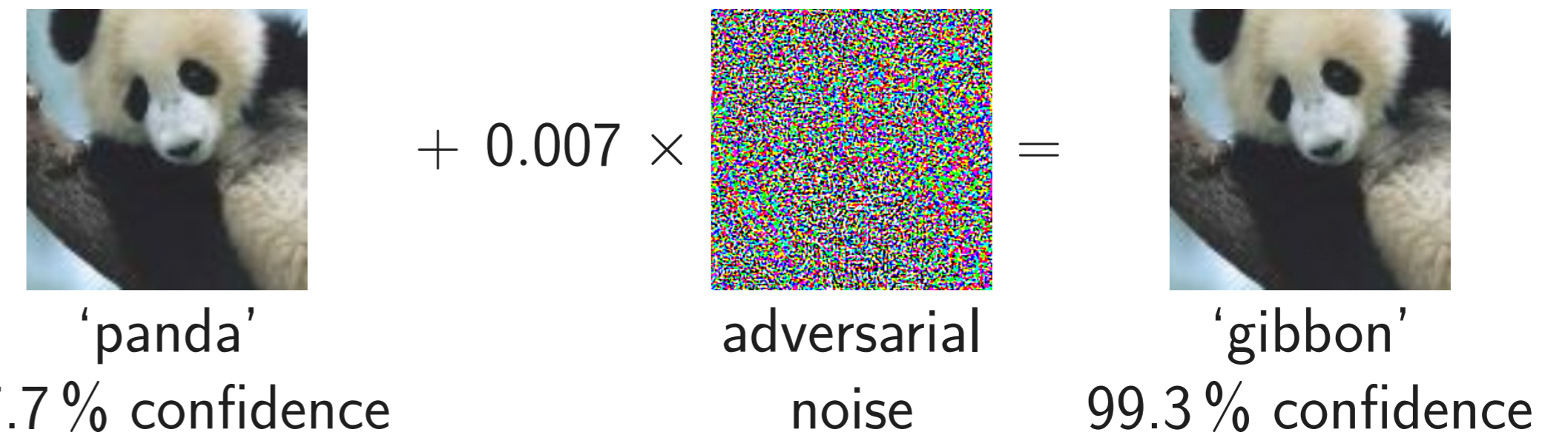
$$\llbracket \phi \wedge \psi \rrbracket_{DL2} := \llbracket \phi \rrbracket_{DL2} + \llbracket \psi \rrbracket_{DL2}$$

$$\llbracket \phi \vee \psi \rrbracket_{DL2} := \llbracket \phi \rrbracket_{DL2} \cdot \llbracket \psi \rrbracket_{DL2}.$$

- Fuzzy Logics:** $\llbracket \cdot \rrbracket_L: \Phi \rightarrow [0, 1]$, where $\llbracket \top \rrbracket_L = 1$ and $\llbracket \perp \rrbracket_L = 0$

Logic	Conjunction	Disjunction	Implication
Gödel	$\min\{x, y\}$	$\max\{x, y\}$	$\begin{cases} 1, & \text{if } x < y, \\ y, & \text{else.} \end{cases}$
Kleene-Dienes			$S(N(x), y)$
Łukasiewicz	$\max\{0, x + y - 1\}$	$\min\{1, x + y\}$	$S(N(x), y)$
Reichenbach	xy	$x + y - xy$	$\begin{cases} 1, & \text{if } x < y, \\ y^x, & \text{else.} \end{cases}$
Goguen			y^x

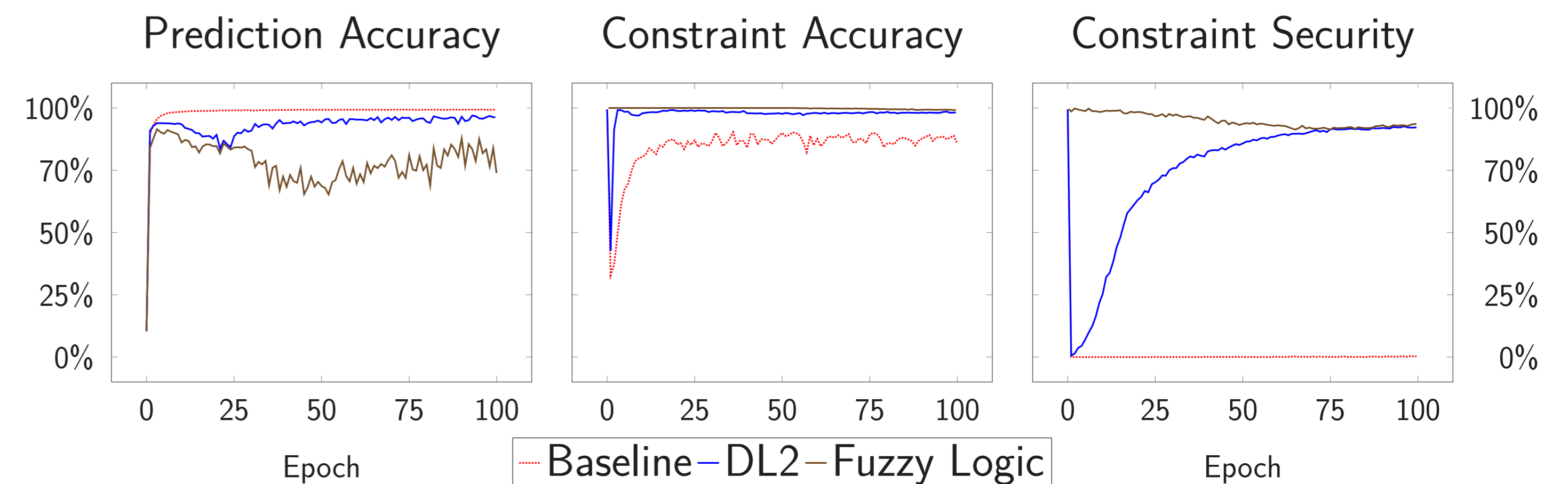
4. Training Experiment



Constraint: A neural network is *locally robust* in input \mathbf{x} , if

$$\forall \mathbf{x}'. \|\mathbf{x}' - \mathbf{x}\|_{\infty} \leq \epsilon \quad \text{implies} \quad \|\mathcal{N}(\mathbf{x}') - \mathcal{N}(\mathbf{x})\|_{\infty} \leq \delta$$

all elements in the input space close to \mathbf{x} the classification is roughly the same



Main finding: Property-driven training with *any* differentiable logic leads to significantly improved constraint satisfaction.

5. Verification Experiment

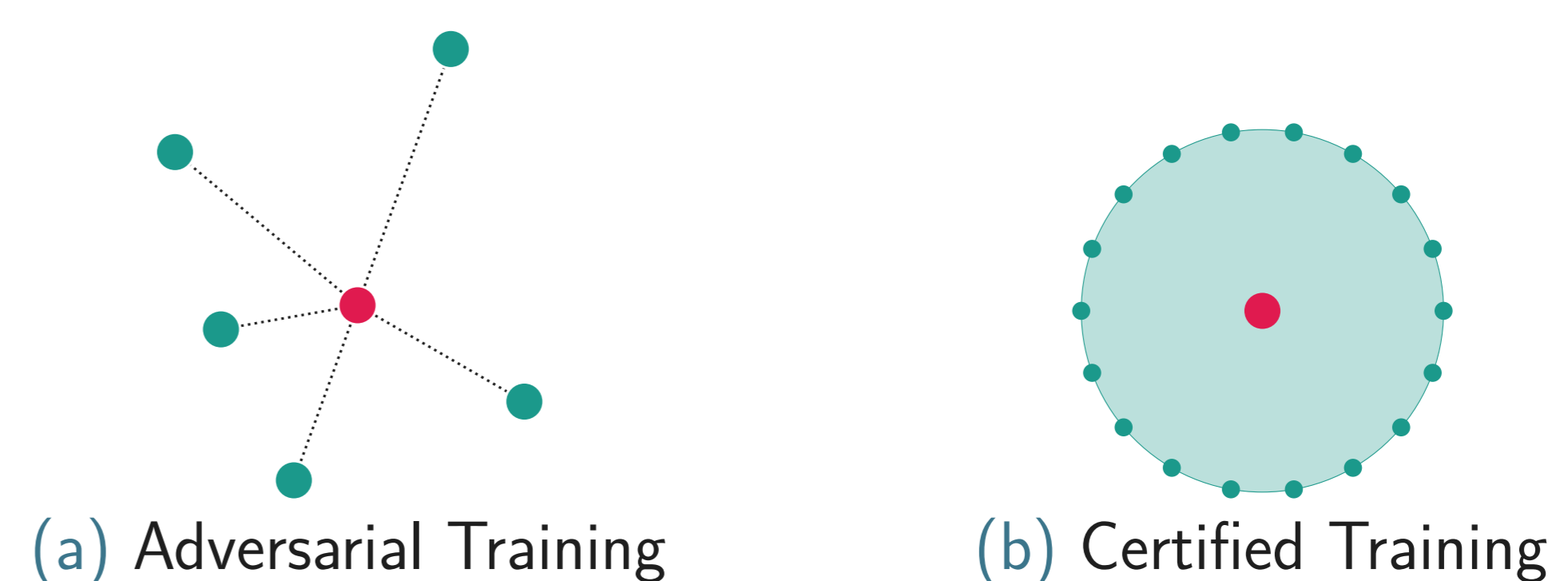
Verified constraint satisfaction on 500 randomly chosen images of the MNIST test set on networks trained with $\epsilon = 0.4$.

Logic	Prediction Accuracy	Constraint Security	Verified Satisfaction		
			$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
Baseline	96.50 %	79.68 %	0.68 % (3/444)	0 % (0/500)	0 % (0/500)
DL2	93.07 %	100 %	92.98 % (384/413)	55.29 % (183/331)	20.51 % (73/356)
Fuzzy Logic	94.87 %	100 %	92.70 % (368/397)	52.16 % (157/301)	9.22 % (27/293)

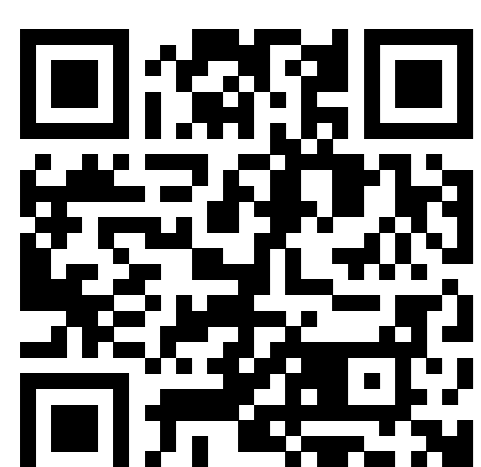
Main finding: Property-driven training leads to some formal guarantees, but significantly lower than what is reported in training.

6. Future Work: Formal Guarantees

Instead of *minimising* an upper bound on the loss, ...



... *maximise* a lower bound on robustness!



Thomas Flinkow, Marco Casadio, Colin Kessler, Rosemary Monahan, Ekaterina Komendantskaya: 'A Generalised Framework for Property-Driven Machine Learning', Submitted to *8th International Symposium on AI Verification (SAIV)* (May 2025). preprint: [arXiv:2505.00466](https://arxiv.org/abs/2505.00466)



Thomas Flinkow, Barak A. Pearlmutter, Rosemary Monahan: 'Comparing differentiable logics for learning with logical constraints', In *Science of Computer Programming*, **244**, 103280 (Sep 2025). DOI: [10.1016/j.scico.2025.103280](https://doi.org/10.1016/j.scico.2025.103280)