

Property-driven Machine Learning with Differentiable Logics

1. Motivation: Formal Verification of Neural Networks

Many formal verification tools can determine whether a neural network satisfies a given logical property. But how can we train a network to *learn* such a property in the first place?

2. Background: Property-driven ML

Standard ML: Given data x , target y , and loss \mathcal{L} ,

$$\text{minimise}_{(x,y) \sim \mathcal{D}} \mathbb{E} \mathcal{L}(x, y).$$

Approach (based on Adversarial Training): Train network f to satisfy logical constraints of the form $\forall x. P(x) \rightarrow Q(f(x))$ by:

- ① translating P into a *hyper-rectangle* $\llbracket P \rrbracket$
- ② translating Q into *real-valued, differentiable loss* $\llbracket Q \rrbracket$
- ③ finding a counterexample x^* in the input space $\llbracket P \rrbracket$ that does *not satisfy* Q using PGD to maximise $\llbracket Q \rrbracket$:

$$x^* = \arg \max_{x' \in \llbracket P \rrbracket} \llbracket Q \rrbracket(x', y)$$

- ④ using this counterexample in training:

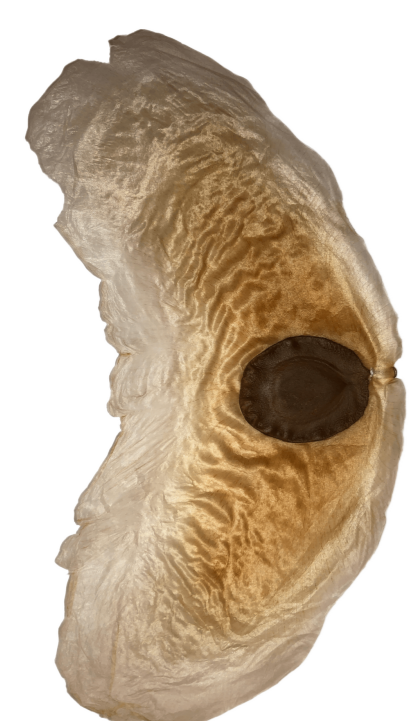
$$\text{minimise}_{(x,y) \sim \mathcal{D}} \mathbb{E} \left[\lambda \underbrace{\mathcal{L}(x, y)}_{\text{prediction loss}} + (1 - \lambda) \underbrace{\llbracket Q \rrbracket(x', y)}_{\text{logical constraint loss}} \right].$$

3. Research Questions

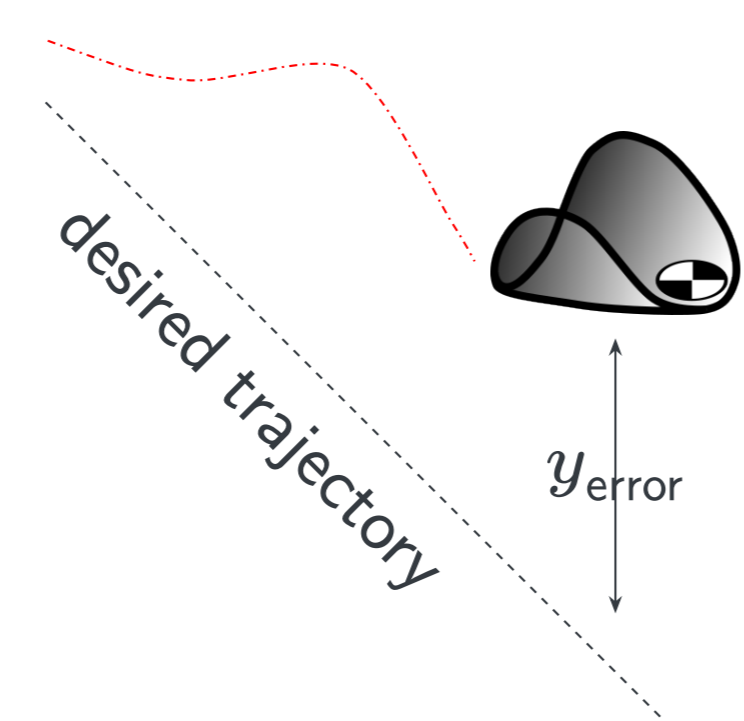
Many differentiable logics (DL2, STL, fuzzy logics, ...) with varying domains and operators, but how do they differ in terms of:

- **Learning behaviour:** gradients and suitability for optimisation.
- **Logical properties:** soundness and internal consistency.
- **Formal guarantees:** empirical performance vs. strong guarantees.

4. Differentiable Logics in Training



(a) An *Alsomitra macrocarpa* seed, capable of stable flight over long distances.



(b) The desired linear trajectory of the *Alsomitra*-inspired drone.

Constraint: If the drone is pitching down fast, is close to and above the line, the network will always make the drone pitch up.

Logic	RMSE	Constraint Accuracy	Constraint Security
Baseline	3.6111×10^{-4}	0 %	0 %
DL2	1.2287×10^{-3}	100 %	95.3 %
Fuzzy logic	1.1632×10^{-3}	100 %	92.2 %

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

5. Current Work: Smooth Differentiable Logic

Idea: Prefer good gradients over logical soundness—just because PGD cannot find a counterexample, does *not* mean that none exist!

- less-than-or-equal atom as a *smooth approximation of ReLU*:

$$\llbracket x \leq y \rrbracket := \tau \text{softplus}(x - y - \delta/\tau)$$

- disjunction as a *smooth approximation of* $\min\{x_i\}$:

- ① via p -means:

$$\llbracket \bigvee x_i \rrbracket := \left(\frac{1}{n} \sum_{i=1}^n x_i^p \right)^{1/p}, \quad p < 0$$

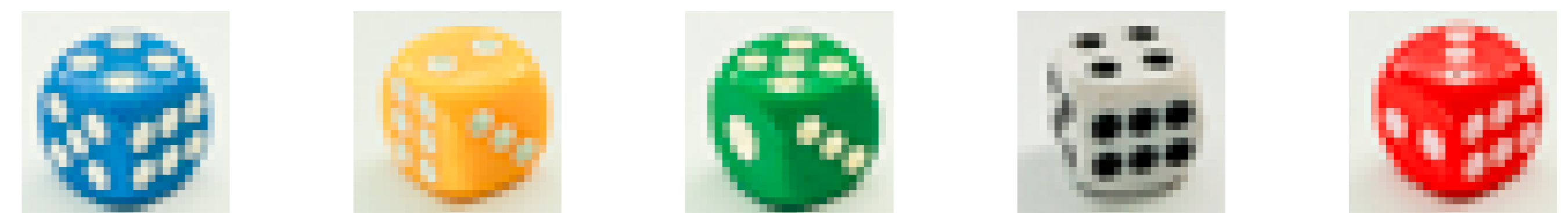
- ② via log-sum-exp:

$$\llbracket \bigvee x_i \rrbracket := \tau \text{LSE}(x_1/\tau, \dots, x_n/\tau), \quad \tau < 0$$

- conjunction as a *smooth approximation of* $\max\{x_i\}$ (dual of disjunction, via p -means with $p > 0$ or log-sum-exp with $\tau > 0$).

6. Property-driven ML and Formal Guarantees?

Example: Multi-label classification of dice faces.



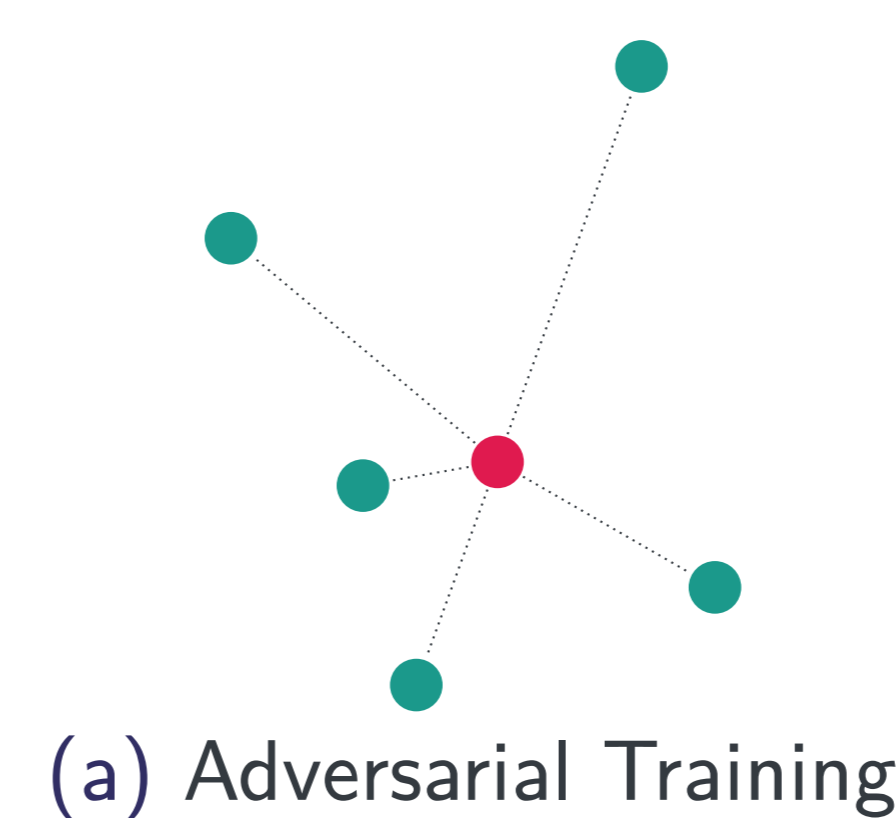
Constraint: Predictions must physically be possible, i.e. do not predict two faces at a time that are on opposite sides of the die. Using Marabou to determine verified constraint satisfaction on the test set consisting of 69 images on networks trained with $\epsilon = 16/255$.

Logic	Prediction Accuracy	Constraint Accuracy	Constraint Security	Verified Satisfaction
Baseline	91.3 %	98.6 %	8.7 %	0 % (0/66)
DL2	91.3 %	100 %	37.7 %	0 % (0/49)
Gödel	91.3 %	100 %	7.3 %	0 % (0/65)
Reichenbach	92.8 %	98.6 %	7.3 %	0 % (0/65)
Our Logic	91.3 %	100 %	100 %	100 % (66/66)

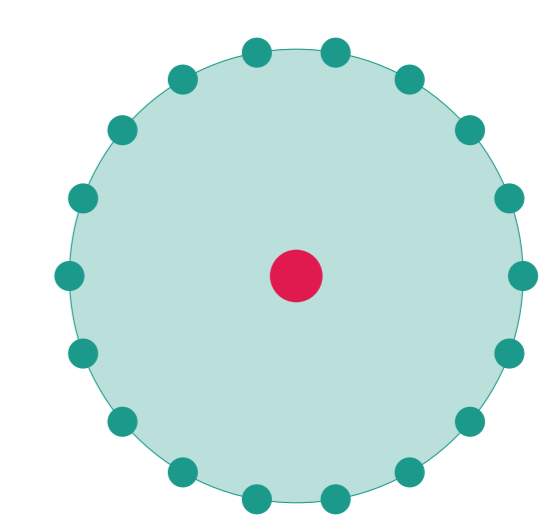
Findings: Property-driven training yields some formal guarantees but often fails to establish strong ones.

7. Future Work: Formal Guarantees & Expressiveness

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



(a) Adversarial Training



(b) Certified Training

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

Immediate next steps: Look into set-based training and train with a set of samples as opposed to using PGD to find counterexamples.



T. Flinkow, B. A. Pearlmutter, R. Monahan: 'Comparing differentiable logics for learning with logical constraints', In *Science of Computer Programming*, 244, 103280 (September 2025). DOI: [10.1016/j.scico.2025.103280](https://doi.org/10.1016/j.scico.2025.103280)

This publication has emanated from research conducted with the financial support of Taighde Éireann – Research Ireland under grant number 20/FFP-P/8853.



T. Flinkow, M. Casadio, C. Kessler, E. Komendantskaya, R. Monahan: 'A General Framework for Property-Driven Machine Learning', (June 2025). PREPRINT: [arXiv:2505.00466](https://arxiv.org/abs/2505.00466)

This publication has emanated from research conducted with the financial support of Taighde Éireann – Research Ireland under grant number 20/FFP-P/8853 and by the ADAPT Research Ireland Centre for AI Driven Digital Content Technology at Maynooth University under grant number 13/RC/2106_P2.