CWI, Amsterdam

20th May 2025

Provable guarantees for datadriven policy synthesis: a formal methods perspective



Prof. Marta Kwiatkowska Department of Computer Science University of Oxford

Safety and security risks in AI decision making

• Well known that neural networks are unstable to adversarial perturbations









Adversarial example

Physical attack

Patch attack

Real traffic sign

(Red light classified as green after <u>one</u> pixel change)

- For high-stakes applications, need provable guarantees on correctness
- Yet AI/ML community focuses on performance formal verification to the rescue?

Feature-Guided Black-Box Safety Testing of Deep Neural Networks. Wicker et al, In Proc. TACAS, 2018.



20 -

40

60

80

100

Adv

(Red |

٠

aft

F

• Y

Feat

Senior researcher & writer

New research from Which? reveals that more than half of drivers are turning off safety tech in their cars with many finding the tech annoying, distracting or even dangerous. We explain why, and how to get the best from your current car or your next purchase.

Like airbags and crumple zones, various car safety technologies are mandatory on new cars. While airbags are considered 'passive' safety tech – they only activate when you crash – Advanced Driver-Assistance Systems (ADAS) are 'active' and are intended to prevent you from having an accident in the first place.

With driver error a leading cause of road accidents in the UK, the best case scenario with ADAS features is that they prevent avoidable accidents. These include accidents where the driver unintentionally leaves their lane, is driving too fast or hasn't spotted an obstacle ahead of them.

However, Which? has found evidence that these features are being habitually turned off by drivers, with just over half of drivers who have an ADAS feature on their car reporting they turn at least one feature off at least some of the time. And when the tech is off, it isn't protecting anybody.

This highlights that there's a lot of room for improvement in the way these systems are implemented and explained.



Formal verification provides provable guarantees



- Modelling = rigorous, mathematical abstraction
- Verification = proof that the model satisfies specification
- Synthesis = correct-by-construction model/policy from specification
- Automated = algorithmic, implemented in software

Probabilistic Model Checking in Autonomy. Kwiatkowska et al, Ann Rev of Control, Robotics and Aut. Sys. (2022).

Multiple applications and use cases!



Optimal controller synthesis

Formal verification for neural networks (NNs)

- Rigorous formal verification
 - can provide provable guarantees, e.g. that no adversarial examples exist
 - enables robustness certification and correct-by-construction synthesis
 - crucial part of safety assurance
- Neural network models more challenging
 - black box, lack interpretability
 - high-dimensional function
 - interplay between architecture and training (non-linear optimization)

Image classifier is a function f: $R^n \to \{c_1, \ldots c_k\}$ Learnable weights and bias

Approximates human perception from M training examples

• Much progress since 2017: Reluplex, DLV, DeepPoly, ReluVal, CROWN, ...

Safety Verification of Deep Neural Networks. CAV 2017 keynote

This talk: provable guarantees via formal verification

- Overview of recent research
- Focus on (data-driven) NN/RL policies and components
- Brief recap of (local) adversarial robustness certification
 - crucial part of safety assurance, pre-deployment
- A selection of snapshots
 - pre-image approximation
 - quantitative verification
 - exploiting causality
 - handling uncertainty
 - neuro-symbolic models
- Conclusions and future directions

Recap of adversarial robustness

• Focus on local adversarial robustness, for a specific input



- Informally, <u>no</u> perturbation results in a misclassification
- More formally, assume given
 - trained neural network classifier $f:R^m \rightarrow \{c_1, \ldots c_k\}$
 - region η centred at x wrt distance function, e.g. $L^2,\,L^\infty$
- Define local robustness at x wrt η by (SAT friendly)
 - $\exists y \in \eta$ such that $f(x) \neq f(y)$
- Here, focus on computing provable guarantees on correctness, rather than constructing defences



Neural network verification



Typically, exact verification intractable, focus on computing lower/upper bounds

Neural network verification: forward analysis

- Given a neural network $f: \mathbb{R}^n \to \mathbb{R}^m$, the NN verification problem is defined as $(\varphi_{pre}, \varphi_{post})$ requiring that
 - $\ \forall x \ \in \ R^n. x \ \vdash \varphi_{pre} \longrightarrow f(x) \vdash \varphi_{post}$
- Typical approach: forward analysis
 - start from $X = \{x \in \mathbb{R}^n | x \vdash \varphi_{pre}\}$
 - bound the worst case on each layer
 - propagate bounds through layers
 - check whether the predicted labels are preserved



- Computes over-approximation of output set
- Note may result in loose bounds...

Progress in neural network verification

- Compute provable guarantees by lower/upper bounding the reachable values
- Methods include exact/approximate
 - <u>search-based/Lipschitz</u>, e.g. DLV
 - <u>constraint solving/SMT/MIP</u>, e.g., Reluplex
 - <u>convex relaxation</u>, e.g., interval/linear bound propagation, as in CROWN
 - <u>abstract interpretation</u>, e.g., DeepPoly
 - <u>global optimisation</u>, under assumption of Lipschitz continuity, e.g., DeepGO



- Convex relaxation best performers, see VNN-Comp
- Scaling, loose bounds and complex architectures an issue...



Linear bounding of ReLU activations ReLU(x) := max(0, x)

Neural network verification: backward analysis

- Given the NN verification problem ($\varphi_{pre}, \varphi_{post}$) for a neural network $f: \mathbb{R}^n \to \mathbb{R}^m$, requiring that
 - $\ \forall x \ \in \ R^n. x \ \vdash \varphi_{pre} \longrightarrow f(x) \vdash \varphi_{post}$
- Focus instead on <u>backward analysis</u>
- Characterize the inputs for output constraints Y = {y ∈ R^m | y ⊢φ_{post} }



- Advantages
 - more precise correctness guarantees, particularly under-approximation
- but
 - <u>exact</u> preimage computation is intractable at scale, $O(2^n)$ for n unstable ReLU neurons

<u>Provably bounding neural network preimages</u>. Koha *et al*, In Proc. NeurIPS 2023. <u>Provable Preimage Under-Approximation for Neural Networks</u>. Zhang *et al*, In Proc. TACAS 2024.

Preimage approximation

- Work backwards to generate preimage approximation via convex relaxation in terms of <u>disjoint union</u> of polytopes
- Given output specification $y = f(x) \ge 0$ (any polyhedral property) $a_j^{(i)}$
- Compute symbolic lower/upper bounding functions for activations from output layer to input:

 $-Ax + b \leq f(x) \leq \overline{A}x + \overline{b}$

• Preimage under-approximation as a polytope:

$$- \{ \mathbf{x} \mid \underline{A}\mathbf{x} + \underline{b} \ge \mathbf{0} \} \longrightarrow \{ \mathbf{x} \mid f(\mathbf{x}) \ge \mathbf{0} \}$$

- Also over-approximation
- Method relies on
 - backward propagation
 - preimage refinement through input/ReLU splitting planes
 - heuristics and optimisations, to deal with exponential growth in constraints

PREMAP: A Unifying PREiMage APproximation Framework for Neural Networks, arXiv:2408.09262

 $h_{i}^{(i)}$ $u_j^{(i)}$ Linear bounding of **ReLU** activations

Preimage under/over-approximation

- Anytime algorithm, backward propagation via convex relaxation
- Preimage refinement to handle approximation loss
 - parallel processing of split regions
 - tightening of approximation by optimizing relaxation parameters
 - (novel differential objective)
- Two types of (sound) preimage refinement
 - via input-feature-aligned cutting plane (not shown)
 - via ReLU-aligned cutting plane
 - (unstable ReLU neuron into two stable cases: approximation becomes exact)
- Volume-estimated prioritization of splitting subregions
- Exact volume for final verification





Experimental results: preimage under-approximation

- Method scales to high-dimensional tasks
 - <u>first</u> method to scale to l_{∞} attack (noise in all image pixels) and patch attack





- evaluation on MNIST (GTSRB in progress) with varied size and position of the patch, indicating areas of vulnerability
- provides quantitative coverage results for larger perturbation bounds

L_{∞} attack	$ \#\mathbf{Poly} $	$\operatorname{Cov}(\%)$	$\left \mathbf{Time}(\mathbf{s}) \right $	Patch attack	#Poly	Cov(%)	Time(s)
0.05	2 $ $	100.0	3.107	$3 \times 3(\text{center})$	1	100.0	2.611
0.07	247	75.2	121.661	4×4 (center)	678	38.2	455.988
0.08	522	75.1	305.867	$6 \times 6(\text{corner})$	2	100.0	9.065
0.09	733	16.5	507.116	$7 \times 7(\text{corner})$	7	84.2	10.128

Provable Preimage Under-Approximation for Neural Networks. Zhang et al, In Proc. TACAS 2024.

Quantitative neural network verification

- Preimage under-approximation enables quantitative verification
 - i.e. estimating proportion of inputs that satisfy φ_{post}
 - sound and complete
- Useful in cases when verification fails
- Complementary to robustness verifiers, benchmarked against winner of VNN-Comp 2023

Task	$\mid \alpha, \beta$ -CR	OWN	Our			
	Result	Time(s)	$ \mathrm{Cov}(\%) $	#Poly	Time(s)	
Cartpole $(\dot{\theta} \in [-1.642, -1.546])$	yes	3.349	100.0	1	1.137	
Cartpole $(\dot{\theta} \in [-1.642, 0])$	no	6.927	94.9	2	3.632	
MNIST (L_{∞} 0.026)	yes	3.415	100.0	1	2.649	
MNIST (L_{∞} 0.04)	unknown	267.139	100.0	2	3.019	

Provable Preimage Under-Approximation for Neural Networks. Zhang et al, In Proc. TACAS 2024.

Reachability for RL controllers

• Backward reachability analysis, with quantitative guarantees

Task	Property	Config	#Poly		Cov		Time(s)	
			ux	ox	ux	ox	ux	ox
Cartpole (FNN 2×64)	$\{y\in \mathbb{R}^2 \;y_1\geq y_2\}$	$\begin{vmatrix} \dot{\theta} \in [-2, -1] \\ \dot{\theta} \in [-2, -0.5] \\ \dot{\theta} \in [-2, 0] \end{vmatrix}$	$\begin{vmatrix} 25\\42\\66 \end{vmatrix}$	1 8 22	$0.766 \\ 0.750 \\ 0.755$	$1.213 \\ 1.242 \\ 1.246$	$\begin{array}{c} 13.337 \\ 19.732 \\ 30.563 \end{array}$	$2.149 \\ 5.778 \\ 11.476$
$\begin{array}{c} \text{Lunarlander} \\ \text{(FNN } 2 \times 64 \text{)} \end{array}$	$\{y \in \mathbb{R}^4 \wedge_{i \in \{1,3,4\}} y_2 \ge y_i\}$	$\dot{v} \in [-1,0] \ \dot{v} \in [-2,0] \ \dot{v} \in [-4,0]$	$ 18 \\ 67 \\ 97 $	1 23 90	$\begin{array}{c c} 0.754 \\ 0.751 \\ 0.751 \end{array}$	$1.068 \\ 1.246 \\ 1.249$	$\begin{array}{c c} 14.453 \\ 48.455 \\ 76.234 \end{array}$	$2.381 \\ 19.210 \\ 72.285$
Dubinsrejoin (FNN 2×256)	$egin{array}{lll} \{y\in \mathbb{R}^8 \wedge_{i\in [2,4]} & y_1\geq y_i \ & igwedge \wedge_{i\in [6,8]} & y_5\geq y_i \} \end{array}$	$\begin{vmatrix} x_v \in [-0.1, 0.1] \\ x_v \in [-0.2, 0.2] \\ x_v \in [-0.3, 0.3] \end{vmatrix}$	$ \begin{array}{ c c c } 211 \\ 409 \\ 677 \\ \end{array} $	$20 \\ 23 \\ 43$	$\begin{array}{c c} 0.751 \\ 0.750 \\ 0.750 \end{array}$	$1.242 \\ 1.241 \\ 1.244$	$\begin{array}{c} 182.821 \\ 323.839 \\ 589.939 \end{array}$	$\begin{array}{c} 18.666 \\ 24.788 \\ 41.502 \end{array}$

- Efficient, often bounding with few polytopes
- Over- and under-approximation

PREMAP: A Unifying PREiMage APproximation Framework for Neural Networks, arXiv:2408.09262

Decision policies with optimality guarantees

- The data generating model can be expressed as follows: • $R \coloneqq f_r(C, A) + \epsilon, \quad \mathbb{E}[\epsilon] = 0, \quad \mathbb{E}[\epsilon|A, C] \neq 0,$ It has been shown that we cannot learn the causal effect ٠ of actions in the presence of hidden confounders without structural assumptions. A: action Thus, we assume access to instrumental variables (IVs) • R: outcome Assumption 2.1. (a) ϵ is additive to R and $\mathbb{E}[\epsilon] = 0$; (b) C: context $\mathbb{E}[\epsilon|C, Z] = 0$; and (c) $\mathbb{P}(A|C, Z)$ is not constant in Z. **Popular conferences** Z: Instrument Supply side costs Airplane ticket price Airplane ticket sales e.q., jet fuel price
- Geographical distance from a medical facility can be an instrument for medical treatment

Learning Decision Policies With Instrumental Variables Through Double Machine Learning. Shao *et al*, In Proc. ICML 2024.

DML-IV

- Double Machine Learning (DML) is a statistical technique that debiases twostage estimators and provides fast convergence rate guarantees of for general two-stage regressions, where N is the sample size
- We propose DML-IV, a novel IV regression algorithm utilising the DML framework
 - Derive a Neyman Orthogonal score function that describes the IV regression problem: $\psi(\mathcal{D}; h, (s, g)) = (s(c, z) g(h, c, z))^2$
 - Design a k-fold cross-fitting learning algorithm.
- We are able to prove that the DML-IV estimator converges to the true causal effect function at rate $O(N^{-1/2})$ under mild regularity assumptions.
- Furthermore, the suboptimality of the induced decision policy is also $O(N^{-1/2})$

Learning Decision Policies With Instrumental Variables Through Double Machine Learning. Shao *et al*, In Proc. ICML 2024.

Experiments

- Evaluate DML-IV on benchmarks and semi-synthetic real-world datasets (Infant development and cardiovascular mortality rate datasets).
- Compare the error of the learned causal effect of actions (lower is better):



Experiments

• Compare the expected reward of the induced decision policy from the learned causal effect function (higher is better).



Extensions

- Extension to imitation learning incorporating causal inference
 - <u>A Unifying Framework for Causal Imitation Learning with Hidden Confounders</u>. Shao et al, In Workshop on Spurious Correlation and Shortcut Learning: Foundations and Solutions, Workshop at The International Conference on Learning Representations (ICLR) 2025
- Learning of optimal policies from temporal logic specifications
 - <u>Sample Efficient Model-free Reinforcement Learning from LTL Specifications with</u> <u>Optimality Guarantees.</u> Shao *et al*, Proc. IJCAI 2023
 - Converts LTL to limit-deterministic Buchi automata
- Learning temporal logic specifications to debug/explain RL policies
 - <u>Learning Probabilistic Temporal Logic Specifications for Stochastic Systems</u>. Roy *et al*, Proc. IJCAI 2025
 - Learns concise probabilistic LTL from positive and negative examples

But what about uncertainty?

- Autonomy, yet highly uncertain scenarios!
- Main focus so far on deterministic neural networks
 - deterministic outcomes, fixed (trained) weights
 - potentially overconfident predictions
- Probabilistic verification
 - enables reasoning, planning, etc, in the presence of uncertainty
 - many tools (PRISM, Storm, etc)
 - limited to state-based models





· Can we extend probabilistic verification to neural network settings?

Probabilistic Model Checking for Strategic Equilibria-based Decision Making: Advances and Challenges, Kwiatkowska 36 et al., In *Proc* MFCS 2022

Uncertainty in decision making

- Requiring that no adversarial examples exist too strict!
- Uncertainty as a <u>first-class citizen</u>: Bayesian neural networks (BNNs)
 - allow for 'don't know' answers, can increase trust in decisions
 - pass distributions through softmax
- Define safety with prob $1-\varepsilon$ $Prob(\exists y \in \eta \text{ s.t. } f(x) \neq f(y) \mid D) \leq \varepsilon$
- i.e., conditioned on training data D



- Aim to provide provable probabilistic guarantees for BNNs
 - certified bounds on decision probability
 - also under adverse conditions
 - (certifiable adversarial robustness training)

<u>Statistical Guarantees for the Robustness of Bayesian Neural Networks</u>. Cardelli *et al*, In Proc. IJCAI 2019. <u>Probabilistic Safety for Bayesian Neural Networks</u>, Wicker et al., In *Proc* UAI 2020

Certification of Bayesian CNN on medical images



- · Limited scalability, more progress needed!
- Also need stronger methodologies beyond certification: correctness by design

Adversarial Robustness Certification for Bayesian Neural Networks. FM 2024

Risk-averse certification for BNNs

- Averaging over full input distributions does not suffice
- Risk-averse certification: robustness under adverse conditions, say 5% of most adversarially unstable cases
- Principled approach incorporating CVar (Conditional Value at Risk)
- Compute certified CVar bounds with probabilistic guarantees



Output support set (empirical)

Tasks	CVaR Level	L_∞ noise	Rotation	Contrast
MNIST	$lpha=1 \ lpha=0.5 \ lpha=0.25$	$\begin{array}{c} \textbf{-0.999} \pm 0.1 \\ \textbf{-0.999} \pm 0.1 \\ \textbf{-0.999} \pm 0.1 \end{array}$	$\begin{array}{c} \text{-0.998} \pm 0.1 \\ \text{-0.998} \pm 0.1 \\ \text{-0.997} \pm 0.1 \end{array}$	$\begin{array}{c} \text{-0.997} \pm 0.1 \\ \text{-0.993} \pm 0.1 \\ \text{-0.986} \pm 0.1 \end{array}$
FASHION	$lpha = 1 \ lpha = 0.5 \ lpha = 0.25$	$\begin{array}{c} \text{-0.191} \pm 0.1 \\ 0.443 \pm 0.1 \\ 0.885 \pm 0.1 \end{array}$	$\begin{array}{c} \text{-0.119} \pm 0.1 \\ 0.525 \pm 0.1 \\ 0.964 \pm 0.1 \end{array}$	$\begin{array}{c} \text{-0.152} \pm 0.1 \\ 0.562 \pm 0.1 \\ 0.849 \pm 0.1 \end{array}$

Table 2: Certified CVaR bounds for different attacks on classification tasks.

Risk-Averse Certification of Bayesian Neural Networks. arXiv:2411.19729

Probabilistic guarantees for autonomous controllers

- So far only one-step predictions; now consider iterated prediction steps
- Here, environment specified as a Bayesian neural network



- Synthesised controllers and reach-avoid probabilities
 - controller benchmarks
 - obstacles





Probabilistic reach-avoid for Bayesian Neural Networks. Wicker et al, Artificial Intelligence, 2024.

Neuro-symbolic games

- Agents endowed with neural perception and symbolic decision making
 - here: NN classifiers (or other machine learning) for perception tasks
 - constrained interface: convert inputs such as images to symbolic percepts
 - plus: local strategies for control decisions
- Neuro-symbolic games (two players/coalitions)
 - finite-state agents + continuous-state environment E

 $\cdot S = (Loc_1 \times Per_1) \times (Loc_2 \times Per_2) \times S_E$

- agents only use a (learnt) perception function to observe E
 - $obs_i : (Loc_1 \times Loc_2) \times S_E \rightarrow Per_i$
- joint actions update state probabilistically
- Example: dynamic vehicle parking
 - NN maps exact vehicle position to perceived grid cell
 - stochasticity from, e.g., motion imprecision







Strategy synthesis for neuro-symbolic games

- Consider zero-sum (discounted) expected reward over infinite horizon
 - for now, we assume full observability
 - value exists under Borel assumptions, fixed point of minimax
 - but optimal value may not be finitely representable
- Value iteration (VI) approach, exploit structure
 - continuous state-space decomposed into regions with the same percept (and reward)
 - further subdivision at each iteration
 via Borel decomposition, under assumptions
 - abstraction based on piecewise-continuous value functions, preserved by NNs and VI
- Implementation
 - pre-image computations of NNs
 - polytope representations of regions (ReLU)
 - LPs to solve zero-sum games at each step

Dynamic vehicle parking with larger (8x8) grid and simpler (regression) perception



Strategy synthesis for zero-sum neuro-symbolic concurrent stochastic games, Inf & Comp, 2024

Neuro-symbolic POMDPs

Need partial observability for neural perception, not just continuous environment!



Neuro-symbolic POSGs

Restrict to one-sided variant, new subclass of hybrid-state POSGs



NS-POSGs (same syntax as NS-CSGs)

- finite-state agents + continuous-state environment E
- Agent 1 uses a (learnt) perception function to observe E
 - $obs_1 : (Loc_1 \times S_E) \rightarrow Per_1$
- and transitions on percepts and local states δ : (Loc₁×Per₁×A) \rightarrow Dist(Loc)
- Agent 2 fully informed

Partially observable stochastic games with neural perception mechanisms, In Proc FM 2024

discounted reward:

 $Y(\pi) = \sum\nolimits_{k=0}^{\infty} \beta^k r(\pi(k), \pi[k])$

zero-sum

71

Example: pedestrian-vehicle interaction

Autonomous vehicle

- partially informed
- aims to predict pedestrian's intention
- using NN trained from video data



• Pedestrian

- fully informed for worst-case analysis
- decides whether to cross or return to sidewalk
- Goal: synthesise strategy for vehicle to minimize likelihood of crash (opposite for pedestrian)

https://data.nvision2.eecs.yorku.ca/PIE_dataset/



Pedestrian-vehicle interaction as NS-POSG



- Agent 1: vehicle
 - *loc*₁: speed
 *per*₁: pedestrian intention
 - $-a_1$: acceleration (e.g. +3,-3)
- Agent 2: pedestrian
 - $-a_2$: cross, back
- Environment E
 - two successive pedestrian positions (x_1, y_1, x_2, y_2)



Strategy synthesis for neuro-symbolic POSGs

- Consider zero-sum (discounted) expected reward over infinite horizon
 - one sided, so Agent 2 can recover beliefs of Agent 1
 - assume determined, as value may not exist
- HSVI approach (extend Horak *et al* 2023)
 - continuous state-space decomposed into regions
 - further subdivision at each iteration
 - work with a class of piecewise-continuous α -functions, + closure properties
 - anytime
- Implementation
 - polyhedral pre-image computations of NNs
 - LPs to compute lower/upper bound and minimax values

Partially observable stochastic games with neural perception mechanisms, In Proc FM 2024



PWC α -function polyhedra + value vector



Efficient online minimax strategies

• How to synthesize strategies based on the lower and upper bound functions



NS-HSVI continual re-solving for Ag_1



Online continual resolving

- · keeps track of belief and counterfactual values
- builds and solves a game without storing complete strategy

Our variant

- precomputes HSVI lower bound
- · keeps track of belief and PWC function α_1
- solves a single LP at each stage

HSVI-based Online Minimax Strategies for Partially Observable Stochastic Games with Neural Perception Mechanisms, 86 In *Proc* L4DC 2024

Safe planning in a crowd

- Robotic agent modelled as a partially observable MDP (POMDP)
 - partial observability (e.g., perception inaccuracy)
- Environment is populated with pedestrians
- Pedestrian trajectory prediction
 - data-driven trajectory predictor
 - uncertainty quantification via
 - (statistical) adaptive conformal prediction (ACP)
- Safe online planning via shielding
 - on-the-fly safety shield construction
- Safety guarantee, given any probability threshold
 - outperforms state of the art on real data





<u>Safe POMDP Online Planning Among Dynamic Agents via Adaptive Comformal Prediction</u>. Sheng *et al*, RAL 2024.

Experimental evaluation using real-world data







(a) ETH

(b) Hotel TABLE I: Experiment Results

			ETH				Hotel				GC	
Method	N	Safety Rate	Time (s)	Min Distance	N	Safety Rate	Time (s)	Min Distance	$\mid N$	Safety Rate	Time (s)	Min Distance
No Shield Shielding without ACP Shielding with ACP	45	0.893 0.943 0.974	21.1 21.5 22.1	0.28 ± 0.19 0.39 ± 0.29 0.51 ± 0.29	35	0.944 0.969 0.988	20.1 20.3 20.6	$0.42{\pm}0.21$ $0.54{\pm}0.3$ $0.8{\pm}0.49$	160	0.91 0.943 0.963	39.3 67.2 71.4	0.22 ± 0.13 0.23 ± 0.13 0.28 ± 0.16
No Shield Shielding without ACP Shielding with ACP	55	0.891 0.951 0.975	21.0 21.8 22.4	0.26 ± 0.17 0.41 ± 0.25 0.53 ± 0.37	45	0.931 0.959 0.982	20.1 20.1 20.6	0.38 ± 0.24 0.48 ± 0.24 0.62 ± 0.27	180	0.904 0.938 0.953	39.9 66.8 71.1	0.2 ± 0.1 0.23 ± 0.12 0.24 ± 0.16
No Shield Shielding without ACP Shielding with ACP	65	0.872 0.943 0.967	21.2 21.9 22.6	0.24 ± 0.13 0.36 ± 0.2 0.42 ± 0.26	55	0.921 0.957 0.982	20.2 20.3 20.3	$0.36 {\pm} 0.18$ $0.48 {\pm} 0.29$ $0.6 {\pm} 0.24$	200	0.895 0.931 0.951	39.6 65.7 74.3	0.22 ± 0.11 0.2 ± 0.13 0.25 ± 0.15

<u>Safe POMDP Online Planning Among Dynamic Agents via Adaptive Comformal Prediction</u>. Sheng *et al*, RAL 2024.

Multiple applications and NN verification use cases!



http://fun2model.org/

Beyond certification: robust learning

- So far, consider <u>trained</u> neural networks, need to retrain if verification fails
- Can we instead efficiently robustly learn? (correct-by-construction synthesis)
- Similar to PAC framework (polynomial sample complexity), except
 - for concept c and hypothesis h, use robust risk $ProbR(\exists z \in \eta(x, \varrho) \text{ s.t. } c(z) \neq h(z))$
 - instead of standard risk $Prob(c(x) \neq h(x))$
 - NB learning is exact in region η , different from previously required invariance/stability over η (constant)
- Show that no non-trivial concept can be learnt in the distribution-free setting
- For simple concepts, can efficiently *e*-robustly learn under classes of distributions (e.g., uniform)

On the Hardness of Robust Classification. Gourdeau *et al*, In Proc. *NeurIPS 2019*, extended *JMLR, 22(273)* 2021 When are Local Queries Useful for Robust Learning? Gourdeau *et al*, In Proc. *NeurIPS 2022*

Beyond adversaries: strategyproof robustness

- So far, consider only adversarial robustness to <u>individual</u> perturbations, but AI agents can behave strategically
- Can we instead strategyproof policy learning? (correctness by design)
- Consider RLHF (reinforcement learning from human feedback)
 - multiple agents, diverse preferences, leading to potential bias in learnt policy decisions
 - but agents can also strategically manipulate the decisions in their favour by misreporting their preferences
 - existing RLHF methods not strategyproof...
- Aim to devise strategyproof RLHF through mechanism design
 - how? incentivise truthful reporting
 - can provide an algorithm that is <u>approximately</u> strategyproof and <u>converges</u> to the optimal policy as the number of individuals and samples increases

Strategyproof Reinforcement Learning from Human Feedback. Kleine Buening et al, arXiv:2503.09561v1

Concluding remarks

- Range of techniques developed in the AI/ML and formal methods communities
 - robustness guarantees needed for high-stakes decisions
 - optimality, explainability of policies desirable
 - but likely to need human involvement in decisions and act as assistants
 - ML models increasing in complexity, take up of certification lagging behind
- Despite progress, major challenges remain
 - scalability to complex architectures and properties
 - foundational understanding needed
 - ideally, semantic methods, not pixel-based perturbations
 - need support for interactions with human decision makers
 - robust learning for correct-by-construction models and policies
- Need integrated processes for validation and safety assurance, not just (probabilistic) verification

Acknowledgements

- My group and collaborators in this work
- Project funding
 - ERC Advanced Grant fun2model
 - EPSRC project FAIR: Framework for responsible adoption of artificial intelligence in the financial services industry, <u>https://www.turing.ac.uk/research/research-</u> <u>projects/project-fair-framework-responsible-adoption-artificial-intelligence</u>
 - ELSA European Lighthouse on Secure and Safe AI, https://www.elsa-ai.eu/
 - UKRI AI Hub on Mathematical foundations of intelligence: an 'Erlangen Programme' for AI
- See also
 - PRISM <u>www.prismmodelchecker.org</u>