



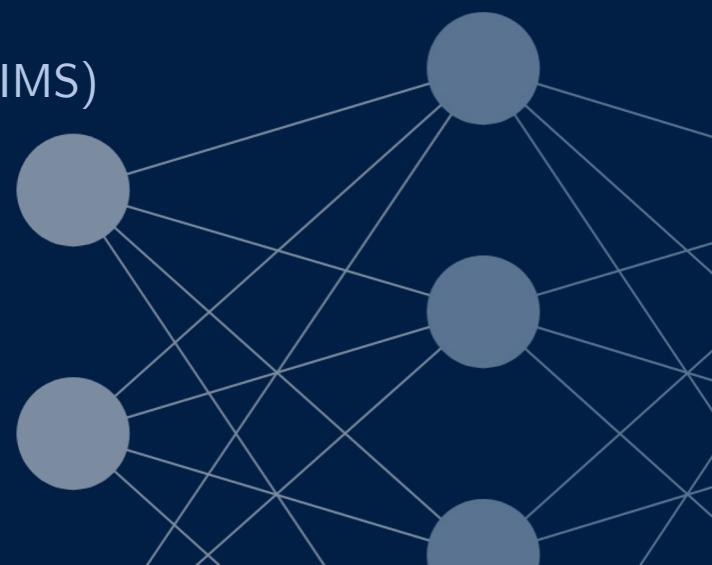
UNIVERSITY OF
OXFORD

Towards Certifiable Machine Learning

Francisco Girbal Eiras

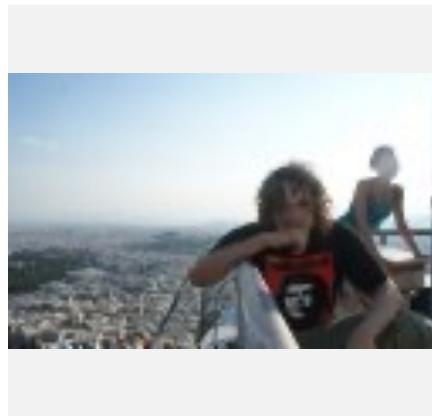
Torr Vision Group · EPSRC CDT Autonomous Intelligent Machines and Systems (AIMS)

July 2023





Francisco Eiras
University of Oxford
eiras@robots.ox.ac.uk



Philip H.S. Torr
University of Oxford
philip.torr@eng.ox.ac.uk



Adel Bibi
University of Oxford
adel.bibi@eng.ox.ac.uk



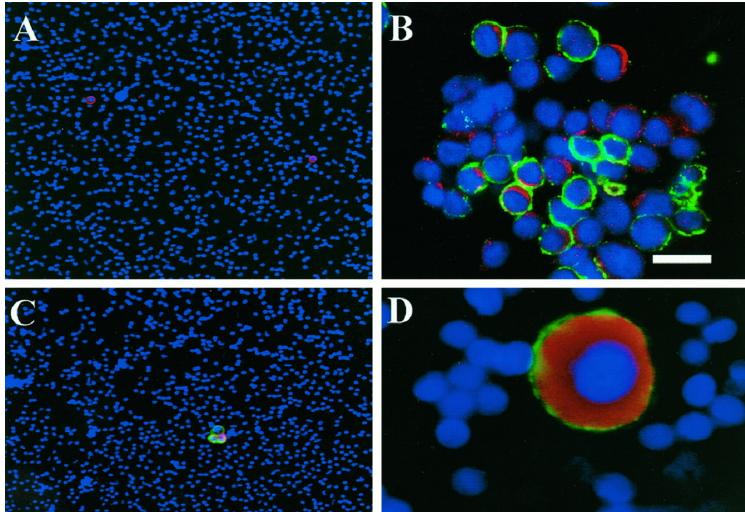
M. Pawan Kumar
Google DeepMind
mpawan@google.com

PhD supervisors

Industry supervisor

Why is Certified Machine Learning *important?*

Safety in Human interaction



Cancer cell detection in the bloodstream



Self-driving cars and autonomous vehicles



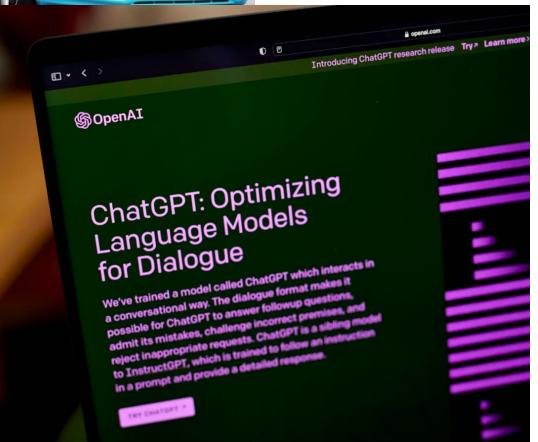
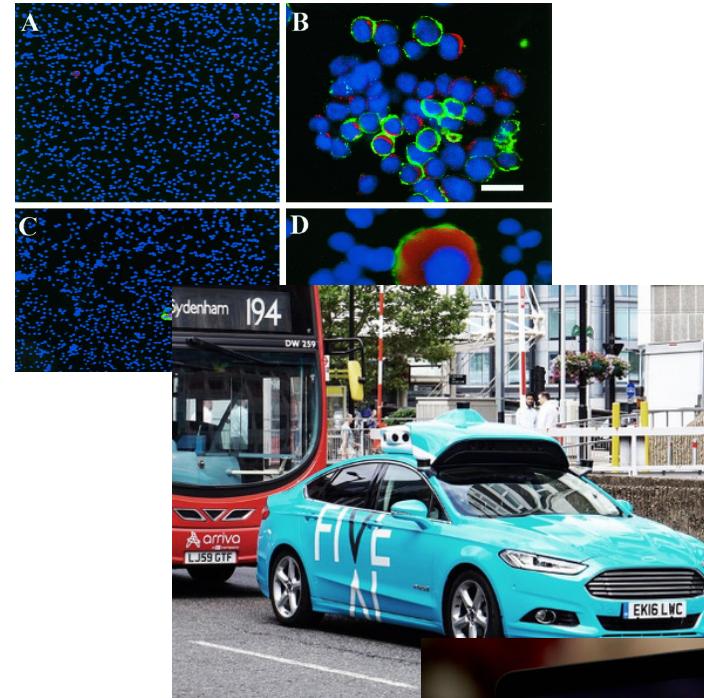
Large generative AI models

Safety in Human interaction

Interaction with humans → makes them **safety-critical** situations

Requires guarantees to be given about their performance!

Minimum guarantee is robustness, *i.e.*, *the ability to tolerate small input perturbations.*



Safety in Human interaction – Adv. Robustness (?)

- The existence of adversarial examples is well known, particularly in image classification.
- ***Is it a safety problem?***

$$\max_{\delta \in \mathbb{B}_\epsilon} \ell(f_\theta(x + \delta), y) \longrightarrow$$



- They are an interesting phenomenon, potentially attributed to *non-robust/spurious* features.
- **Certification is still important**



A Discussion of **Adversarial Examples Are Not Bugs, They Are Features**

PUBLISHED
Aug. 6, 2019

DOI
10.23915/distill.00019

On May 6th, Andrew Ilyas and colleagues published a paper [\[1\]](#) outlining two sets of experiments. Firstly, they showed that models trained on adversarial examples can transfer to real data, and secondly that models trained on a dataset derived from the representations

01

Safety in Human Interaction/ Understanding Decision-Making Process



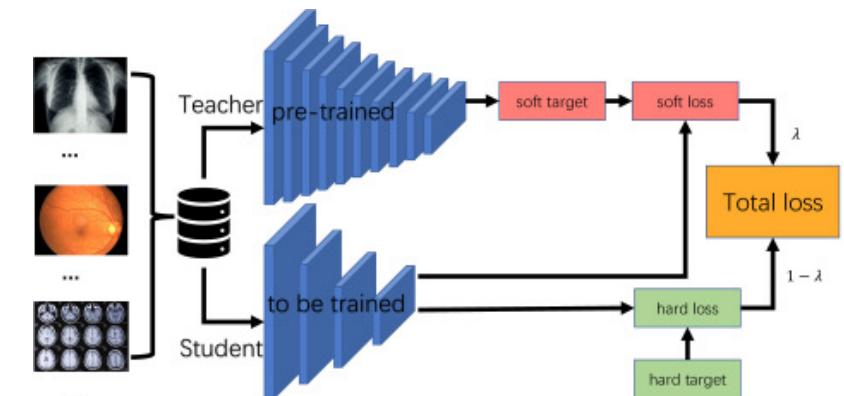
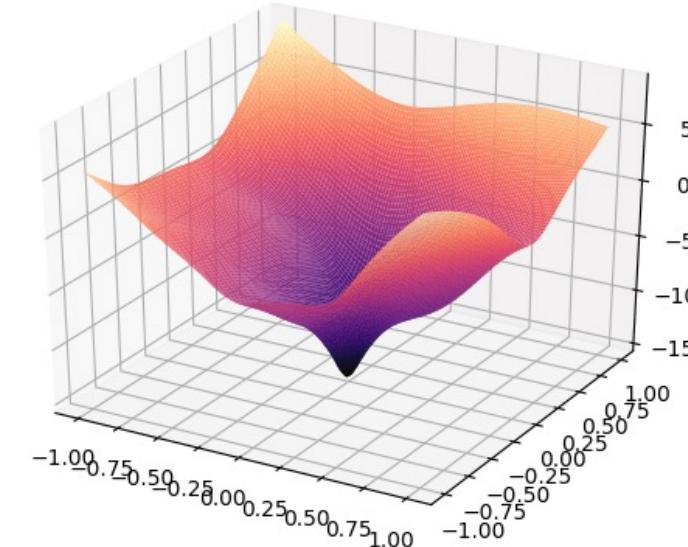
Correctness of Machine Learning based-Systems

ML is used within systems to replace experts/speed-up processes.

E.g., to solve PDEs or knowledge distillation.

Deployed in real-world, but can we trust it to behave **correctly**?

Certified machine learning could help establish reliability/correctness of these components.



01

Safety in Human Interaction/
Understanding Decision-Making Process

02

Correctness of ML-based Systems



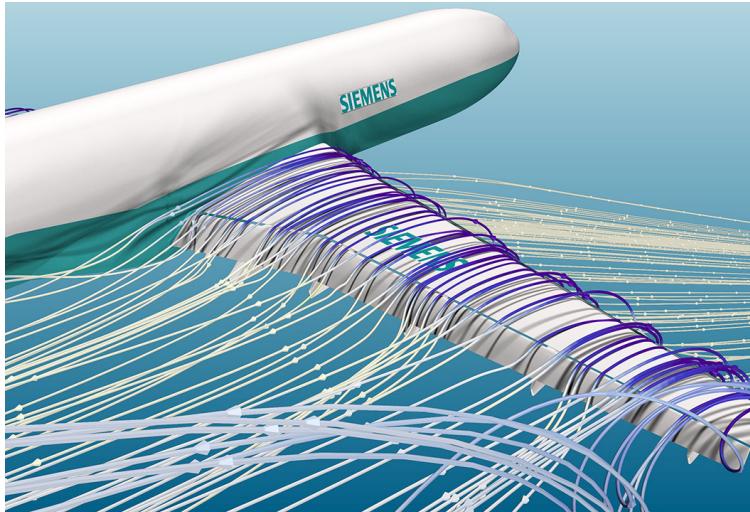
01

Safety in Human Interaction/
Understanding Decision-Making Process

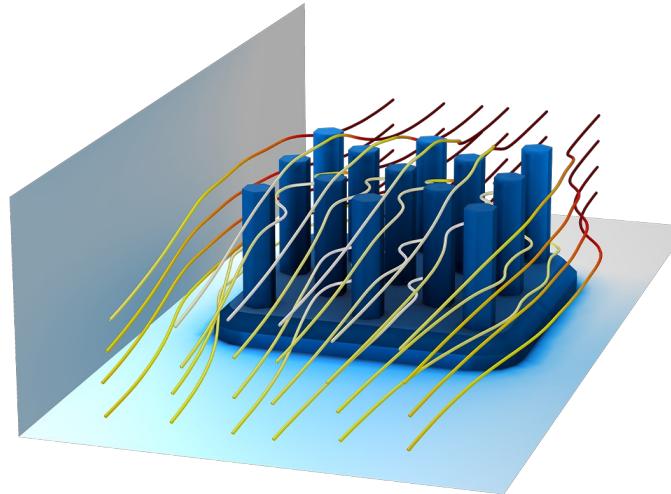
02

Correctness of ML-based Systems

Physical PDEs and Where to Find Them



Aerodynamics (e.g., Euler's equation)



Thermodynamics (e.g., Heat equation)



Earth and Space Sciences
(e.g., turbulence in 2-D Navier-Stokes)

Among many others...

Nonlinear (Physical) Partial Differential Equations

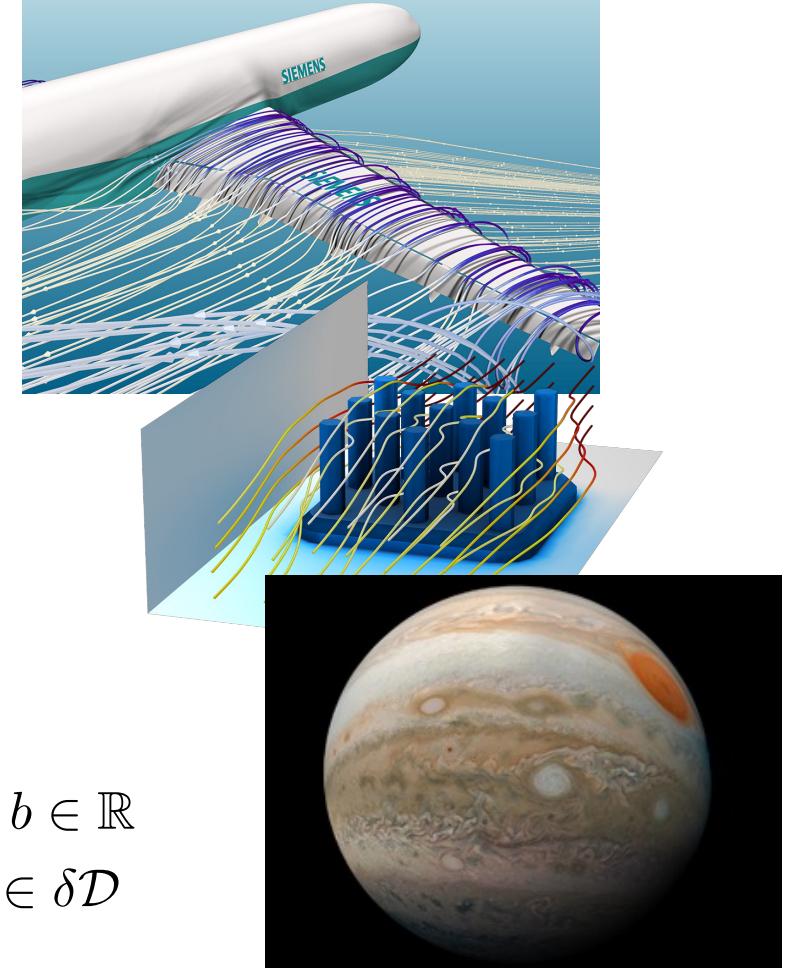
$$\underbrace{\partial_t u(t, x) + \mathcal{N}[u](t, x)}_{\text{residual}} = 0, \quad x \in \mathcal{D}, t \in [0, T]$$

s.t.

1. **Initial condition:** $u(0, x) = u_0(x)$

2. **Robin boundary conditions:** $au(t, x) + b\partial_{\mathbf{n}} u(t, x) = u_b(t, x) \quad a, b \in \mathbb{R}$
 $x \in \delta\mathcal{D}$

domain



PDE example: Diffusion-Sorption Equation

- Applications in groundwater contaminant transportation
- 1D Diffusion-Sorption:

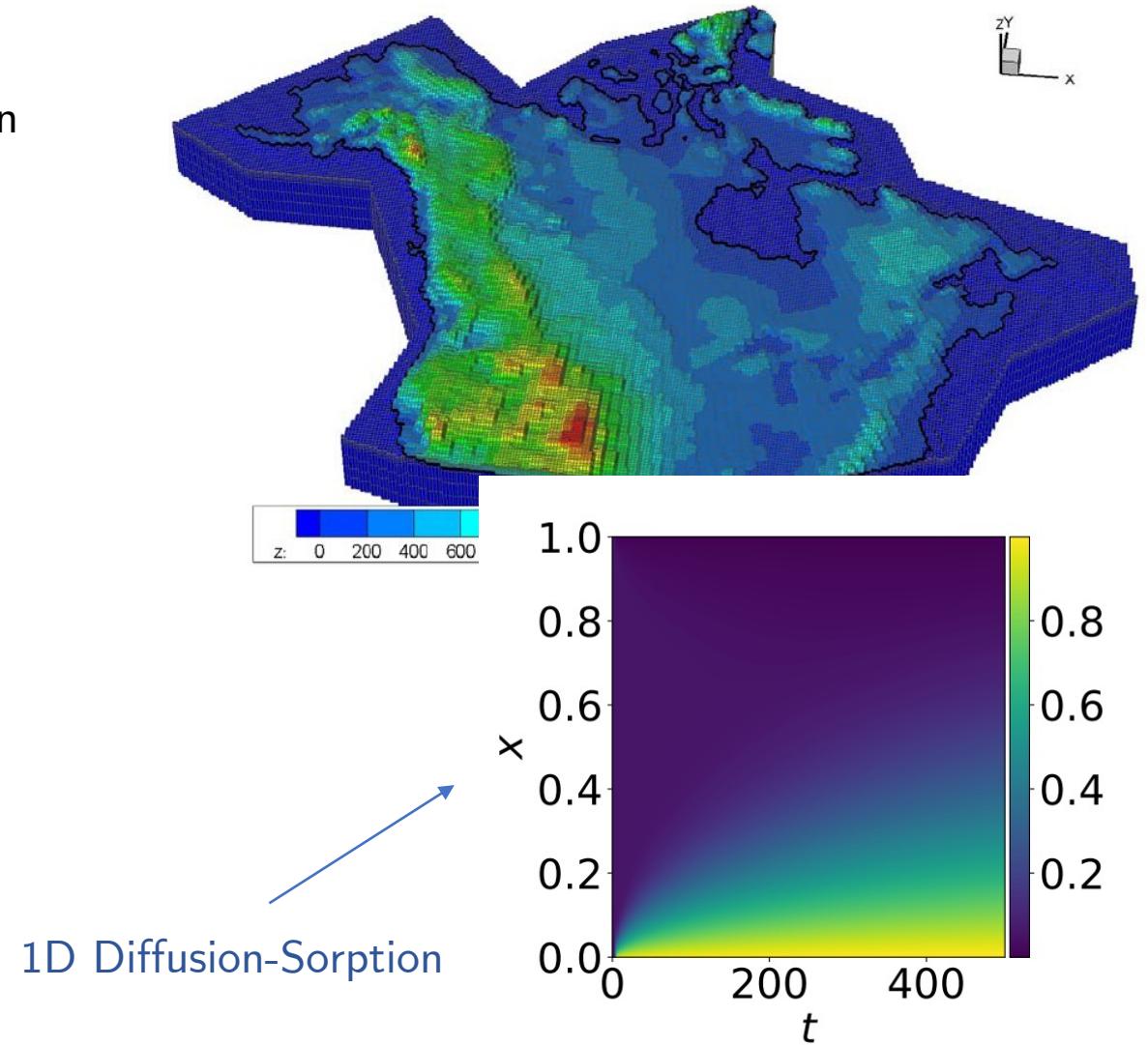
$$\partial_t u(t, x) - D/R(u(t, x)) \partial_{x^2} u(t, x) = 0$$

where

$$R(u(t, x)) = 1 + \frac{(1 - \phi)}{(\phi)} \rho_s k n_f u^{n_f - 1}(t, x)$$

for physical constants D, ϕ, ρ_s, k, n_f , and:

$$u(0, x) = u(t, 0) = 0, \quad u(t, 1) = D \partial_x u(t, 1)$$



Issue: Solving for $u(t, x)$ is computationally expensive

Solution: Use NNs to approximate it

Physics-Informed Neural Networks (PINNs)

- Approximate solution using a neural network, $u_\theta(t, x) \simeq u(t, x)$
- Take the residual evaluated for u_θ as the network $f_\theta(t, x) = \partial_t u_\theta(t, x) + \mathcal{N}[u_\theta](t, x)$
- Train *both networks* jointly using a loss evaluated at collocation points \mathbb{P} (i.e. points in the domain):

$$\mathcal{L} = \underbrace{\sum_{x \in \mathbb{P}_0} |u(0, x) - u_\theta(0, x)|^2}_{\text{initial conditions}} + \underbrace{\sum_{(t, x) \in \mathbb{P}_b} |u(t, x) - u_\theta(t, x)|^2}_{\text{boundary conditions}} + \underbrace{\sum_{(t, x) \in \mathbb{P}_f} |f_\theta(t, x)|^2}_{\text{residual}}$$

PINN

- Evaluate empirically by comparing $u_\theta(t, x)$ to the solution obtained by a numerical solver

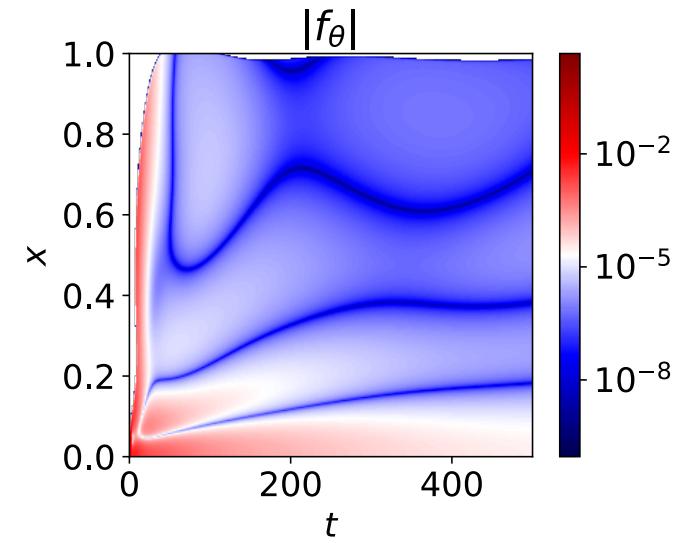
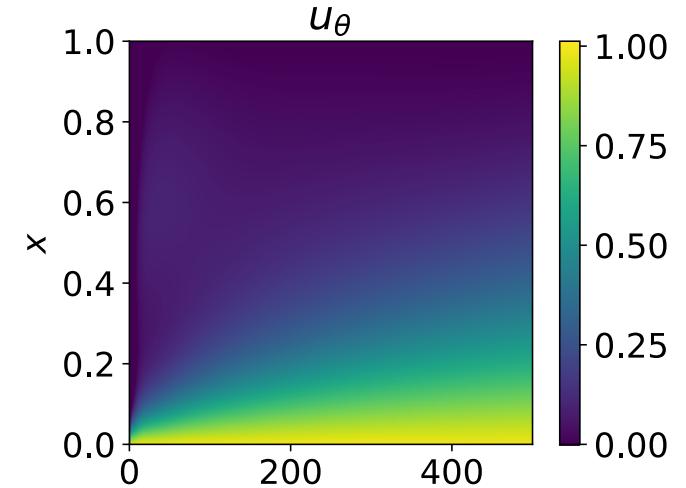
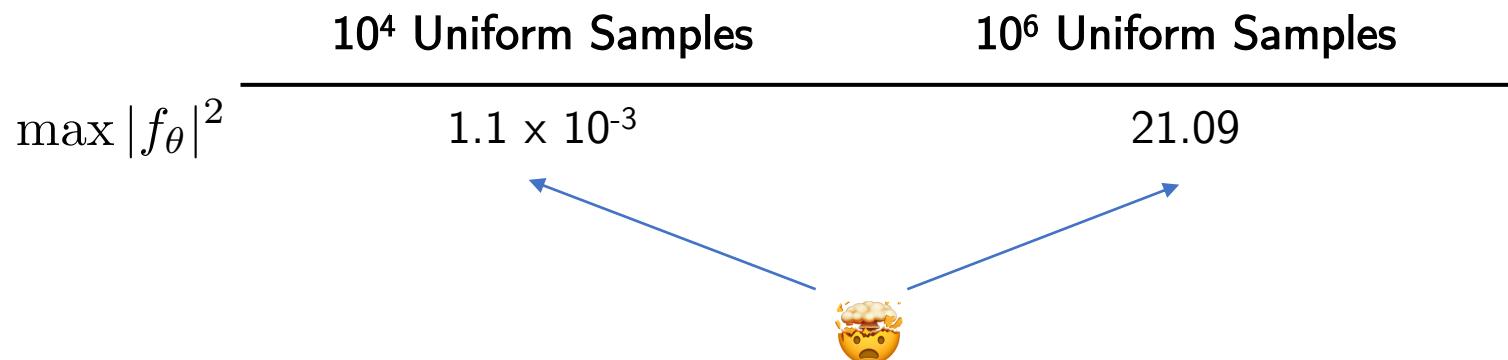
PINN: 1D Diffusion-Sorption Equation

- Inference times on 2 core CPU + 1 GPU (NVIDIA V100):

Numerical Solver	PINN [Takamoto et al. 2022]
59.83 s	2.7×10^{-3} s



- PINN average ℓ_2 solution error is 9.9×10^{-2} compared to numerical solver
- Valid PDE solution **must** satisfy $f_\theta(t, x) = 0$
- Is it satisfied across the domain?*



How can we be confident errors are small enough across the *entire* domain?

Defining Correctness Conditions for PINNs (PINN, 2023)

- Intuitively, a PINN is a correct approximation of the underlying PDE if:
 1. *The solution satisfies the initial conditions* to a reasonable degree
 2. *The solution satisfies the boundary conditions* to a reasonable degree
 3. *The norm of the PINN output is small enough*
- Formally, for a D dimensional spatial input $\hat{\mathbf{x}} \in \mathcal{D}$, and solution/PINN input $\mathbf{x} = (t, \hat{\mathbf{x}}) :$

Definition 1 (Correctness Conditions for PINNs). $u_\theta : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ is a $\delta_0, \delta_b, \varepsilon$ -globally correct approximation of the exact solution $u : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ if:

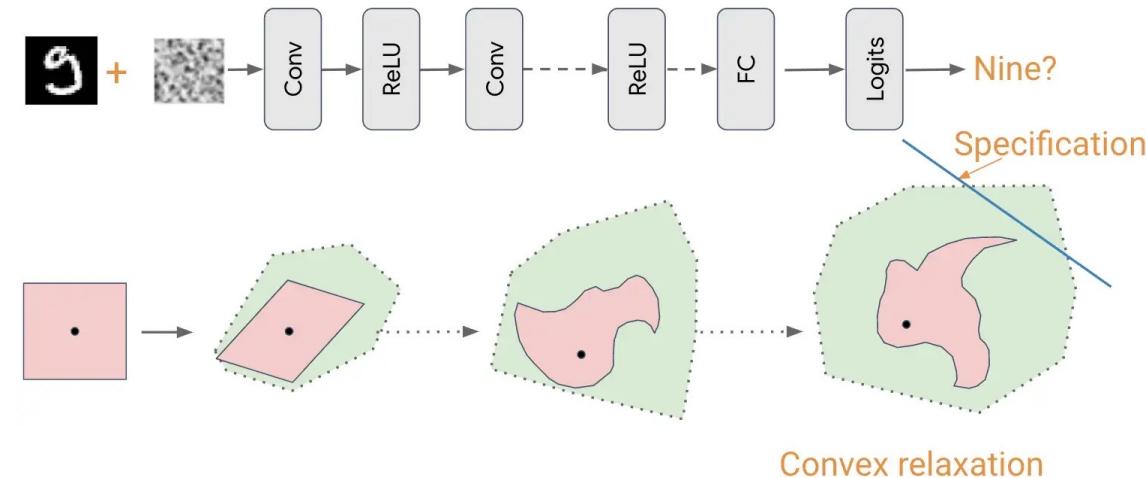
$$\textcircled{1} \quad \max_{\hat{\mathbf{x}} \in \mathcal{D}} |u_\theta(0, \hat{\mathbf{x}}) - u_0(\hat{\mathbf{x}})|^2 \leq \delta_0,$$

$$\textcircled{2} \quad \max_{t \in [0, T], \hat{\mathbf{x}} \in \delta \mathcal{D}} |au_\theta(t, \hat{\mathbf{x}}) + b\partial_{\mathbf{n}} u_\theta(t, \hat{\mathbf{x}}) - u_b(t, \hat{\mathbf{x}})|^2 \leq \delta_b,$$

$$\textcircled{3} \quad \max_{\mathbf{x} \in \mathcal{C}} |f_\theta(\mathbf{x})|^2 \leq \varepsilon.$$

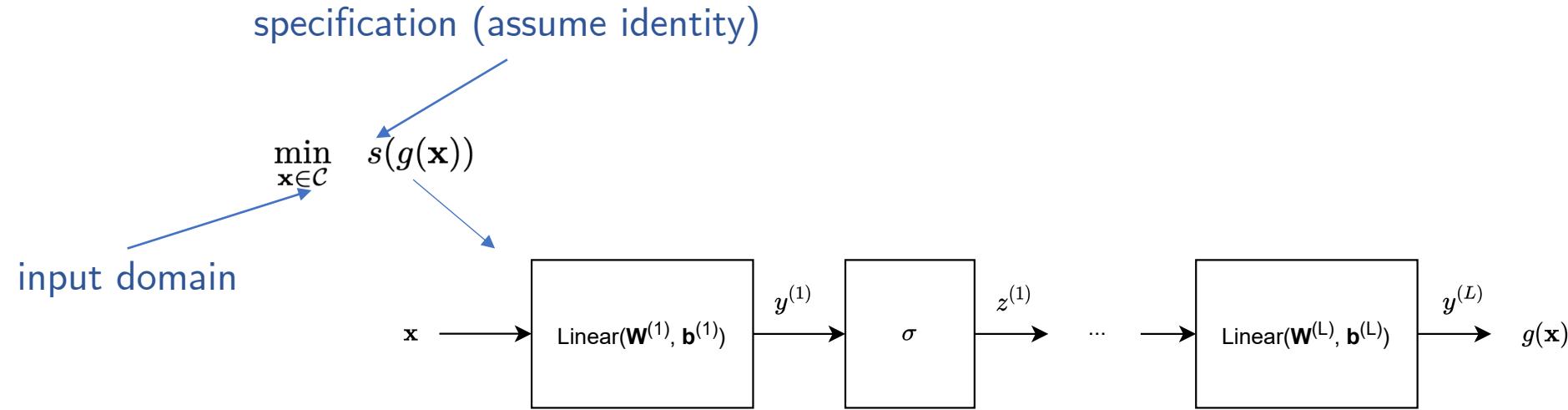
Scalable Adversarial Robustness Verification/Certification

- Complete vs. incomplete methods:
 - *Can it certify a ResNet50? If YES → incomplete.*
- “Sound, but incomplete” methods:
 - Bound propagation (e.g., DiffAI, IBP)
 - **Convex relaxation optimization-based methods** (e.g. CROWN)
 - Smoothing-based methods (e.g. randomized smoothing)



Verification of Neural Networks

- CROWN/α-CROWN [Zhang et. al 2018, Xu et. al 2020]



- Relax network to a linear program → solve it in closed form:
- $$\min_{\mathbf{x} \in \mathcal{C}} \mathbf{A}^L \mathbf{x} + \mathbf{a}^L = l_b \leq \min_{\mathbf{x} \in \mathcal{C}} g(\mathbf{x})$$
- relaxed problem
- Apply it to u_θ and f_θ using the previous specifications

Correctness Certification of PINNs (PINN, 2023)

Definition 1 (Correctness Conditions for PINNs). $u_\theta : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ is a $\delta_0, \delta_b, \varepsilon$ -globally correct approximation of the exact solution $u : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ if:

✓ $\max_{\hat{\mathbf{x}} \in \mathcal{D}} |u_\theta(0, \hat{\mathbf{x}}) - u_0(\hat{\mathbf{x}})|^2 \leq \delta_0,$

✗ $\max_{t \in [0, T], \hat{\mathbf{x}} \in \delta \mathcal{D}} |au_\theta(t, \hat{\mathbf{x}}) + b\partial_{\mathbf{n}} u_\theta(t, \hat{\mathbf{x}}) - u_b(t, \hat{\mathbf{x}})|^2 \leq \delta_b,$

✗ $\max_{\mathbf{x} \in \mathcal{C}} |f_\theta(\mathbf{x})|^2 \leq \varepsilon.$

- Applying CROWN/ α -CROWN to boundary/residual conditions:
 - ✗ Architecture is completely different - f_θ is a nonlinear function of partial derivatives of u_θ
 - ✗ Regression problem – bounds might be too loose to be informative

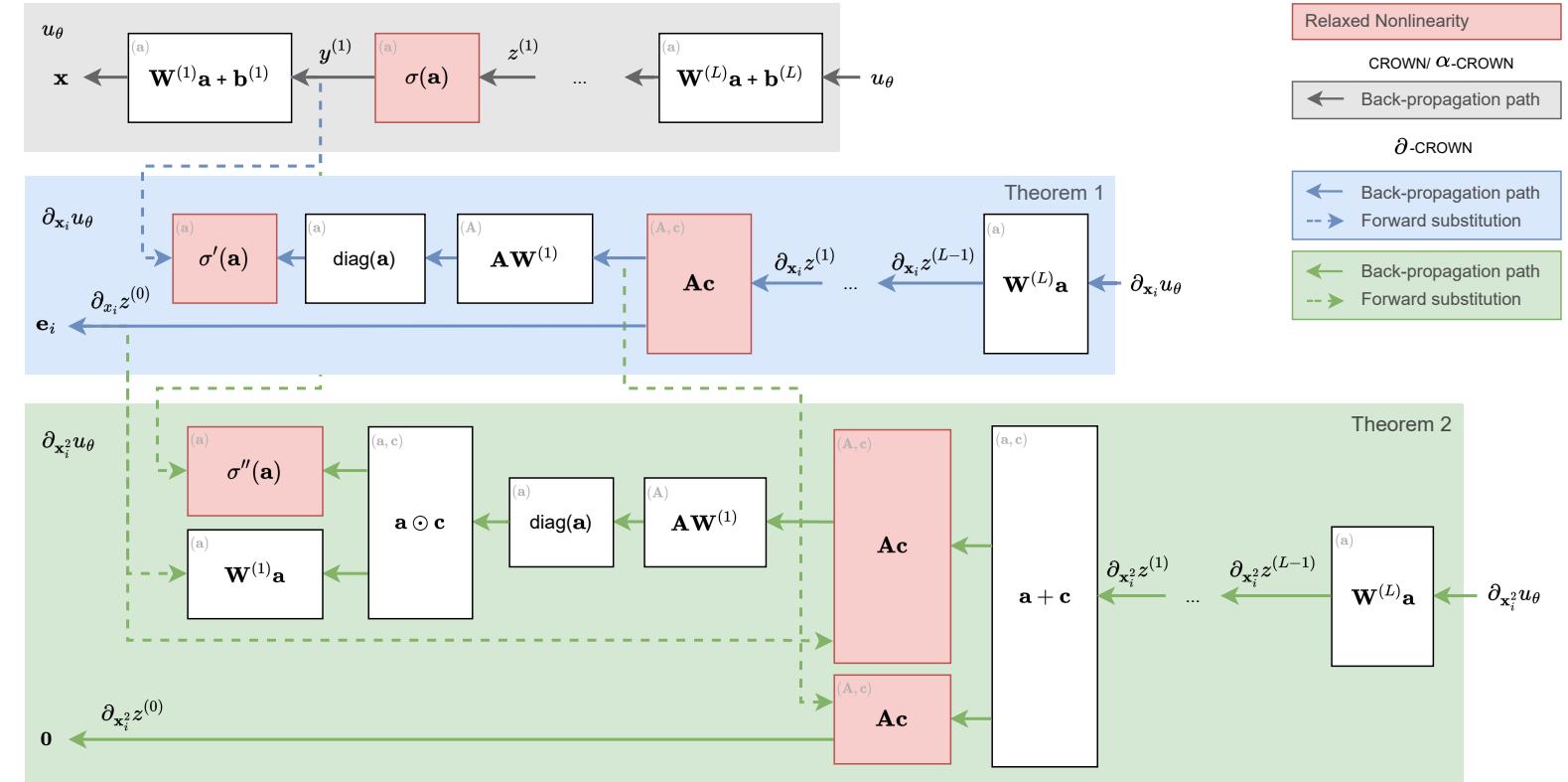
∂ -CROWN: Bounding Derivatives of u_θ and f_θ (PINN, 2023)

- u_θ is bound using CROWN [Zhang et al. 2018]; partial derivatives require purpose-built efficient solution

1st (Theorem 1) and 2nd (Theorem 2) linear bounding of partial derivatives

Hybrid scheme with complexity $\mathcal{O}(L)$ instead of $\mathcal{O}(L^2)$ from [Xu et al. 2020]

Global bounds computed in close-form (similarly to [Zhang et al. 2018])



- f_θ is linearly bounded using McCormick envelopes; global bounds computed in close-form

Correctness Certification of PINNs (PINN, 2023)

Definition 1 (Correctness Conditions for PINNs). $u_\theta : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ is a $\delta_0, \delta_b, \varepsilon$ -globally correct approximation of the exact solution $u : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ if:

✓ $\max_{\hat{\mathbf{x}} \in \mathcal{D}} |u_\theta(0, \hat{\mathbf{x}}) - u_0(\hat{\mathbf{x}})|^2 \leq \delta_0,$

✗ $\max_{t \in [0, T], \hat{\mathbf{x}} \in \delta \mathcal{D}} |au_\theta(t, \hat{\mathbf{x}}) + b\partial_{\mathbf{n}} u_\theta(t, \hat{\mathbf{x}}) - u_b(t, \hat{\mathbf{x}})|^2 \leq \delta_b,$

✗ $\max_{\mathbf{x} \in \mathcal{C}} |f_\theta(\mathbf{x})|^2 \leq \varepsilon.$

- Applying CROWN/ α -CROWN to boundary/residual conditions:
 - ✓ Architecture is completely different - f_θ is a nonlinear function of partial derivatives of u_θ
 - ✗ Regression problem – bounds might be too loose to be informative



Greedy Input Branching

Correctness Certification of PINNs (PINN, 2023)

Definition 1 (Correctness Conditions for PINNs). $u_\theta : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ is a $\delta_0, \delta_b, \varepsilon$ -globally correct approximation of the exact solution $u : [0, T] \times \mathcal{D} \rightarrow \mathbb{R}$ if:

- ✓ $\max_{\hat{\mathbf{x}} \in \mathcal{D}} |u_\theta(0, \hat{\mathbf{x}}) - u_0(\hat{\mathbf{x}})|^2 \leq \delta_0,$
- ✓ $\max_{t \in [0, T], \hat{\mathbf{x}} \in \delta \mathcal{D}} |au_\theta(t, \hat{\mathbf{x}}) + b\partial_{\mathbf{n}} u_\theta(t, \hat{\mathbf{x}}) - u_b(t, \hat{\mathbf{x}})|^2 \leq \delta_b,$
- ✓ $\max_{\mathbf{x} \in \mathcal{C}} |f_\theta(\mathbf{x})|^2 \leq \varepsilon.$

- Applying CROWN/ α -CROWN to boundary/residual conditions:
 - ✓ Architecture is completely different - f_θ is a nonlinear function of partial derivatives of u_θ
 - ✓ Regression problem – bounds might be too loose to be informative

Experiments: Certifying with ∂ -CROWN (PINN, 2023)

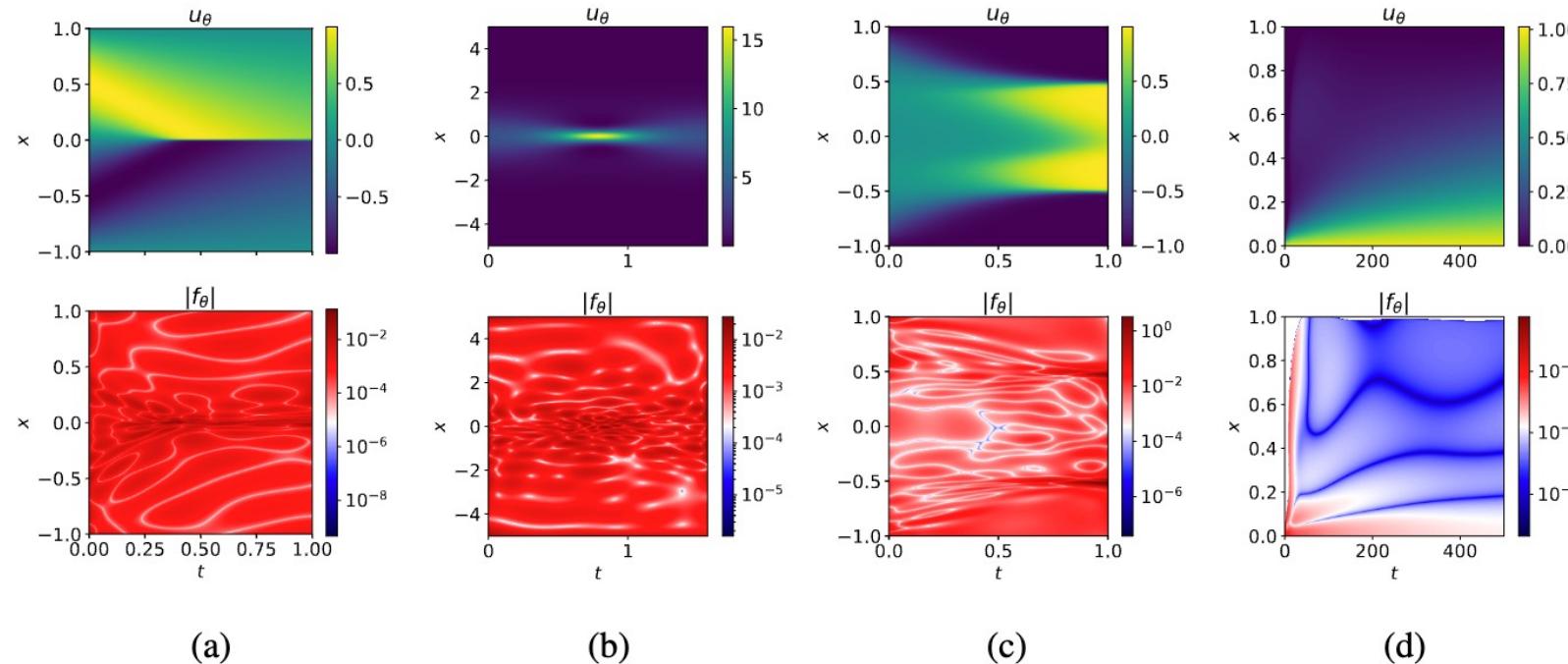


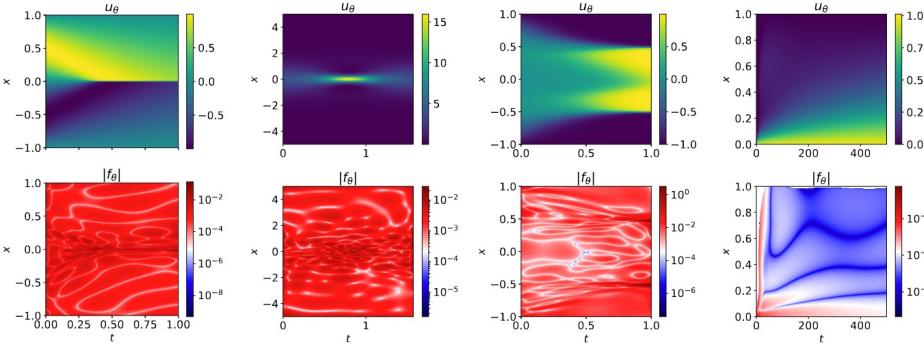
Figure 2: Certifying with ∂ -CROWN: visualization of the time evolution of u_θ , and the residual errors as a function of the spatial temporal domain (log-scale), $|f_\theta|$, for (a) Burgers' equation [Raissi et al., 2019b], (b) Schrödinger's equation [Raissi et al., 2019b], (c) Allan-Cahn's equation [Monaco and Apiletti, 2023], and (d) the Diffusion-Sorption equation [Takamoto et al., 2022].

Raissi, Maziar, Paris Perdikaris, and George E. Karniadakis. "Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations." *Journal of Computational physics* 378 (2019): 686-707.

Monaco, Simone, and Daniele Apiletti. "Training physics-informed neural networks: One learning to rule them all?." *Results in Engineering* 18 (2023): 101023.

Takamoto, Makoto, et al. "PDEBench: An extensive benchmark for scientific machine learning." *Advances in Neural Information Processing Systems* 35 (2022): 1596-1611.

Experiments: Certifying with ∂ -CROWN (PINN, 2023)



	MC max (10^4)	MC max (10^6)	∂ -CROWN u_b (time [s])
(a) Burgers [Raissi et al., 2019b]			
① $ u_\theta(0, x) - u_0(x) ^2$	1.59×10^{-6}	1.59×10^{-6}	2.63×10^{-6} (116.5)
② $ u_\theta(t, -1) ^2$	8.08×10^{-8}	8.08×10^{-8}	6.63×10^{-7} (86.7)
③ $ u_\theta(t, 1) ^2$	6.54×10^{-8}	6.54×10^{-8}	9.39×10^{-7} (89.8)
④ $ f_\theta(x, t) ^2$	1.23×10^{-2}	1.80×10^{-2}	1.03×10^{-1} (2.8×10^5)
(b) Schrödinger [Raissi et al., 2019b]			
① $ u_\theta(0, x) - u_0(x) ^2$	7.06×10^{-5}	7.06×10^{-5}	8.35×10^{-5} (305.2)
② $ u_\theta(t, 5) - u_\theta(t, -5) ^2$	7.38×10^{-7}	7.38×10^{-7}	5.73×10^{-6} (545.4)
③ $ \partial_x u_\theta(t, 5) - \partial_x u_\theta(t, -5) ^2$	1.14×10^{-5}	1.14×10^{-5}	5.31×10^{-5} (2.4×10^3)
④ $ f_\theta(x, t) ^2$	7.28×10^{-4}	7.67×10^{-4}	5.55×10^{-3} (1.2×10^6)
(c) Allen-Cahn [Monaco and Apiletti, 2023]			
① $ u_\theta(0, x) - u_0(x) ^2$	1.60×10^{-3}	1.60×10^{-3}	1.61×10^{-3} (52.7)
② $ u_\theta(t, -1) - u_\theta(t, 1) ^2$	5.66×10^{-6}	5.66×10^{-6}	5.66×10^{-6} (95.4)
③ $ f_\theta(x, t) ^2$	10.74	10.76	10.84 (6.7×10^5)
(d) Diffusion-Sorption [Takamoto et al., 2022]			
① $ u_\theta(0, x) ^2$	0.0	0.0	0.0 (0.2)
② $ u_\theta(t, 0) - 1 ^2$	4.22×10^{-4}	4.39×10^{-4}	1.09×10^{-3} (72.5)
③ $ u_\theta(t, 1) - D\partial_x u_\theta(t, 1) ^2$	2.30×10^{-5}	2.34×10^{-5}	2.37×10^{-5} (226.4)
④ $ f_\theta(x, t) ^2$	1.10×10^{-3}	21.09	21.34 (2.4×10^6)

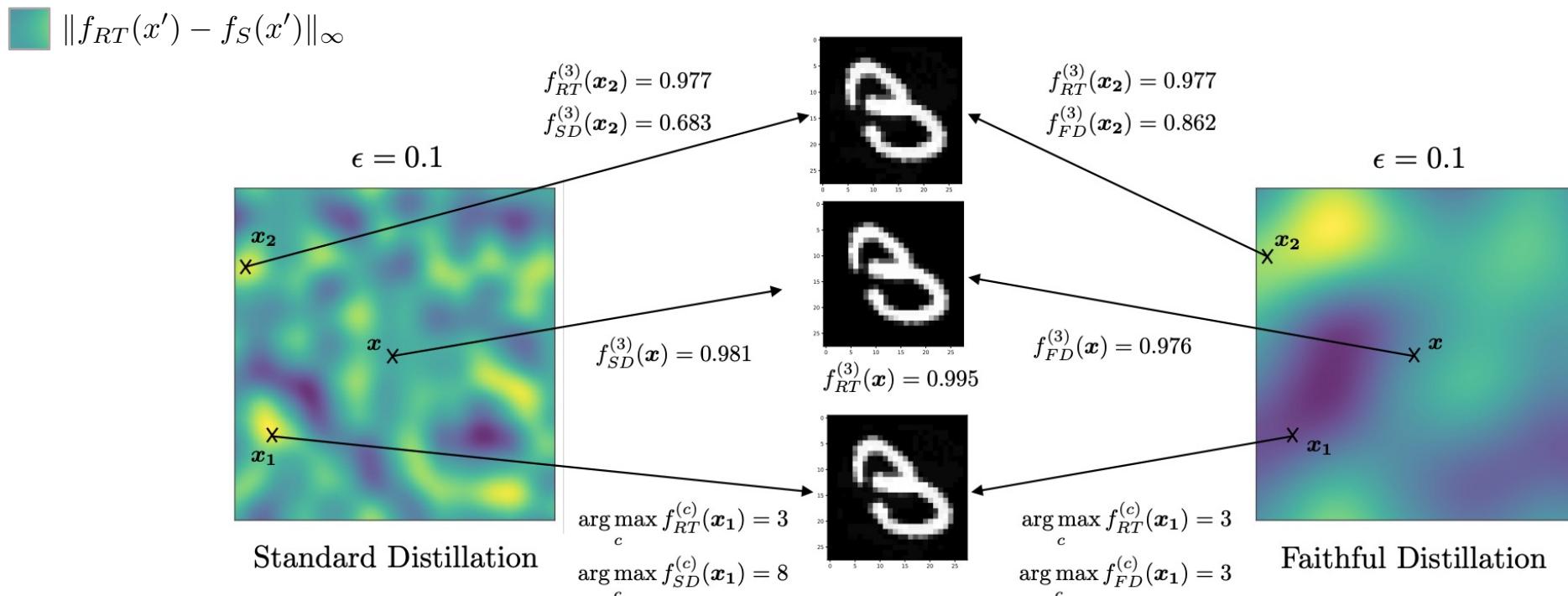
Raissi, Maziar, Paris Perdikaris, and George E. Karniadakis. "Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations." Journal of Computational physics 378 (2019): 686-707.

Monaco, Simone, and Daniele Apiletti. "Training physics-informed neural networks: One learning to rule them all?." Results in Engineering 18 (2023): 101023.

Takamoto, Makoto, et al. "PDEBench: An extensive benchmark for scientific machine learning." Advances in Neural Information Processing Systems 35 (2022): 1596-1611.

Faithful Knowledge Distillation (FKD, 2023)

- *Can we be sure when we distil a teacher network that the student will match the **teacher's output and confidence**?*
- **Intuition:**



Faithful Knowledge Distillation (FKD, 2023)

Definition 2 (Faithful Imitator) We say that $\hat{f} : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is an (ϵ, δ) -faithful imitation of $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ around $\mathbf{x}_0 \in \mathbb{R}^n$ if:

$$d_f(f(\mathbf{x}), \hat{f}(\mathbf{x})) \leq \delta, \quad \forall \mathbf{x} \in B_\epsilon(\mathbf{x}_0) \quad (7)$$

where $B_\epsilon(\mathbf{x}_0) = \{\mathbf{x}' \in \mathbb{R}^n \mid d_{\mathbf{x}}(\mathbf{x}_0, \mathbf{x}') \leq \epsilon\}$, $d_f : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ is a chosen metric function in the output space, and $d_{\mathbf{x}} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ is a metric function in the input space. We refer to any δ that bounds $d_f(f(\mathbf{x}'), \hat{f}(\mathbf{x}'))$ as a faithfulness bound for a given ϵ .

- **Faithful Imitator for Knowledge Distillation:** underestimate δ empirically using a PGD-inspired method, overestimate it using an extension to CROWN.
- Compare standard distillation (SD), with other existing losses with similar goals (ARD, RSLAD) and our proposed faithful distillation loss (FD) which *empirically* minimizes KL divergence of teacher and student.

Table 2: *Faithfulness and relative calibration - empirical lower bounds (EMPLB) and faithfulness upper bounds (FAITHUB) on MNIST, F-MNIST and CIFAR-10. Lower is better.*

		EMPLB				FAITHUB			
	ϵ	f_{SD}	f_{ARD}	f_{RSLAD}	f_{FD}	f_{SD}	f_{ARD}	f_{RSLAD}	f_{FD}
MNIST	0.025	0.042 ± 0.102	0.045 ± 0.113	0.039 ± 0.087	0.033 ± 0.079	0.073 ± 0.150	0.061 ± 0.143	0.060 ± 0.118	0.054 ± 0.113
	0.05	0.060 ± 0.129	0.055 ± 0.130	0.049 ± 0.101	0.041 ± 0.091	0.179 ± 0.257	0.097 ± 0.198	0.101 ± 0.175	0.094 ± 0.171
	0.1	0.106 ± 0.186	0.078 ± 0.164	0.072 ± 0.129	0.061 ± 0.117	0.731 ± 0.343	0.262 ± 0.337	0.290 ± 0.319	0.248 ± 0.305
	0.15	0.172 ± 0.246	0.106 ± 0.197	0.101 ± 0.160	0.086 ± 0.145	0.974 ± 0.123	0.636 ± 0.383	0.668 ± 0.355	0.588 ± 0.367
	0.2	0.258 ± 0.302	0.140 ± 0.229	0.138 ± 0.192	0.118 ± 0.174	0.999 ± 0.017	0.905 ± 0.230	0.889 ± 0.229	0.884 ± 0.238
F-MNIST	4/255	0.099 ± 0.159	0.067 ± 0.112	0.053 ± 0.094	0.060 ± 0.099	0.215 ± 0.267	0.128 ± 0.182	0.118 ± 0.170	0.123 ± 0.171
	8/255	0.155 ± 0.212	0.096 ± 0.145	0.081 ± 0.130	0.089 ± 0.131	0.617 ± 0.393	0.330 ± 0.342	0.341 ± 0.341	0.322 ± 0.334
	12/255	0.219 ± 0.265	0.129 ± 0.179	0.112 ± 0.163	0.122 ± 0.165	0.892 ± 0.267	0.678 ± 0.381	0.690 ± 0.370	0.668 ± 0.380
	16/255	0.297 ± 0.310	0.165 ± 0.214	0.148 ± 0.197	0.160 ± 0.199	0.990 ± 0.071	0.906 ± 0.235	0.917 ± 0.214	0.902 ± 0.237
	20/255	0.378 ± 0.343	0.204 ± 0.243	0.190 ± 0.226	0.202 ± 0.229	0.999 ± 0.013	0.986 ± 0.089	0.988 ± 0.079	0.985 ± 0.090
CIFAR-10	4/255	0.282 ± 0.145	0.242 ± 0.139	0.226 ± 0.126	0.187 ± 0.096	0.671 ± 0.191	0.546 ± 0.203	0.497 ± 0.176	0.474 ± 0.182
	8/255	0.388 ± 0.166	0.317 ± 0.155	0.291 ± 0.139	0.245 ± 0.114	0.983 ± 0.033	0.910 ± 0.127	0.873 ± 0.133	0.864 ± 0.153
	12/255	0.489 ± 0.175	0.394 ± 0.169	0.357 ± 0.148	0.302 ± 0.129	1.000 ± 0.002	1.000 ± 0.016	0.990 ± 0.035	0.989 ± 0.038
	16/255	0.583 ± 0.177	0.462 ± 0.178	0.415 ± 0.157	0.359 ± 0.142	1.000 ± 0.000	1.000 ± 0.002	0.999 ± 0.008	0.999 ± 0.009
	20/255	0.666 ± 0.165	0.530 ± 0.176	0.474 ± 0.164	0.416 ± 0.151	1.000 ± 0.000	1.000 ± 0.000	1.000 ± 0.001	1.000 ± 0.001

01

Safety in Human Interaction/
Understanding Decision-Making Process

02

Correctness of ML-based Systems

01

Safety in Human Interaction/
Understanding Decision-Making Process

02

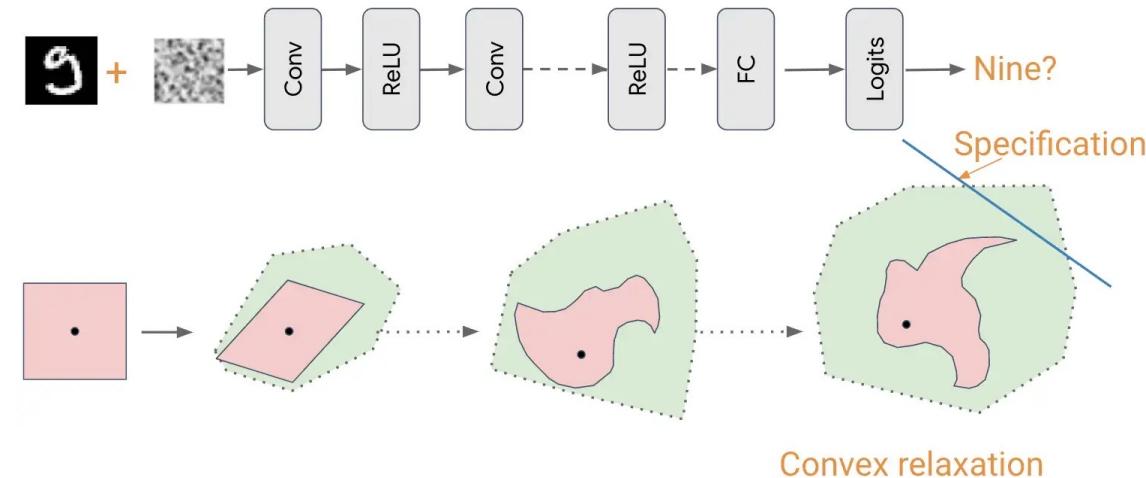
Correctness of ML-based Systems



- 
- 01 Safety in Human Interaction/
Understanding Decision-Making Process
 - 02 Correctness of ML-based Systems
- 

Scalable Adversarial Robustness Verification/Certification

- Complete vs. incomplete methods:
 - *Can it certify a ResNet50? If YES → incomplete.*
- “Sound, but incomplete” methods:
 - Bound propagation (e.g., DiffAI, IBP)
 - Convex relaxation optimization-based methods (e.g. CROWN)
 - **Smoothing-based methods** (e.g. randomized smoothing)



Randomized Smoothing [Cohen et. al, 2019]

- Given a *base classifier* f , we can obtain a *smooth classifier*.

Theorem 1. Let $f : \mathbb{R}^d \rightarrow \mathcal{Y}$ be any deterministic or random function, and let $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$. Let g be defined as in (1). Suppose $c_A \in \mathcal{Y}$ and $\underline{p}_A, \overline{p}_B \in [0, 1]$ satisfy:

$$\mathbb{P}(f(x + \varepsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P}(f(x + \varepsilon) = c) \quad (2)$$

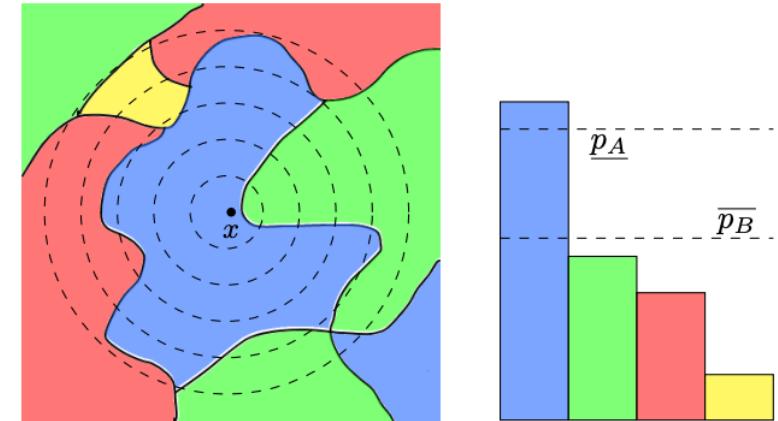
Then $g(x + \delta) = c_A$ for all $\|\delta\|_2 < R$, where

$$R = \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)) \quad (3)$$

Inverse Gaussian
CDF

$$g(x) = \arg \max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \varepsilon) = c)$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$



[Left] Decision boundaries of the base classifier f and level sets of the noise distribution [Right] output distribution given noise.

- Proof based on Neyman-Pearson lemma [Neyman & Pearson, 1933], as tight as possible for ℓ_2 .
- Other works extend this to other ℓ_p norms [Yang et al., 2020].
- Can we simplify the analysis?

Cohen, Jeremy, Elan Rosenfeld, and Zico Kolter. "Certified adversarial robustness via randomized smoothing." international conference on machine learning. PMLR, 2019.

Neyman, Jerzy, and Egon Sharpe Pearson. "IX. On the problem of the most efficient tests of statistical hypotheses." Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character 231.694-706 (1933): 289-337.

Yang, Greg, et al. "Randomized smoothing of all shapes and sizes." International Conference on Machine Learning. PMLR, 2020.

A Lipschitz view of Randomized Smoothing (ANCER, 2022)

Proposition 1. Consider a differentiable function $g : \mathbb{R}^n \rightarrow \mathbb{R}$. If $\sup_x \|\nabla g(x)\|_* \leq L$ where $\|\cdot\|_*$ has a dual norm $\|z\| = \max_x z^\top x$ s.t. $\|x\|_* \leq 1$, then g is L -Lipschitz under norm $\|\cdot\|_*$, that is $|g(x) - g(y)| \leq L\|x - y\|$.

Given the previous proposition, we formalize $\|\cdot\|$ certification as follows:

Theorem 1. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^K$, g^i be L -Lipschitz continuous under norm $\|\cdot\|_*$ $\forall i \in \{1, \dots, K\}$, and $c_A = \arg \max_i g^i(x)$. Then, we have $\arg \max_i g^i(x + \delta) = c_A$ for all δ satisfying:

$$\|\delta\| \leq \frac{1}{2L} \left(g^{c_A}(x) - \max_c g^{c \neq c_A}(x) \right).$$

- **Remark:** randomized smoothing is an instance of **Theorem 1** where the smooth classifier enjoys an analytical form for L by design.
 - Choose a smoothing distribution and compute an analytic Lipschitz constant under the dual norm for $g \rightarrow$ gives you a certificate under the norm by **Theorem 1**.
 - Can recover the certificates from [Cohen et. al, 2019] for ℓ_2 and from [Yang et al., 2020] for several ℓ_p ones.

Anisotropic Certification (ANCER, 2022)

- Lipschitz analysis allows us to obtain **anisotropic** certificates
- **Ellipsoid certificates:** $g_\Sigma(x) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \Sigma)} [f(x + \epsilon)]$ $\|\delta\|_{\Sigma,2} = \sqrt{\delta^\top \Sigma^{-1} \delta}$

Proposition 2. $\Phi^{-1}(g_\Sigma(x))$ is 1-Lipschitz (i.e. $L = 1$) under the $\|\cdot\|_{\Sigma^{-1},2}$ norm.

Since Φ^{-1} is a strictly increasing function, by combining Proposition 2 with Theorem 1, we have:

Corollary 1. Let $c_A = \arg \max_i g_\Sigma^i(x)$, then $\arg \max_i g_\Sigma^i(x + \delta) = c_A$ for all δ satisfying:

$$\|\delta\|_{\Sigma,2} \leq \frac{1}{2} \left(\Phi^{-1}(g_\Sigma^{c_A}(x)) - \Phi^{-1} \left(\max_c g_\Sigma^{c \neq c_A}(x) \right) \right).$$

- **Generalized cross-polytope certificates:** $g_\Lambda(x) = \mathbb{E}_{\epsilon \sim \mathcal{U}[-1,1]^n} [f(x + \Lambda\epsilon)]$

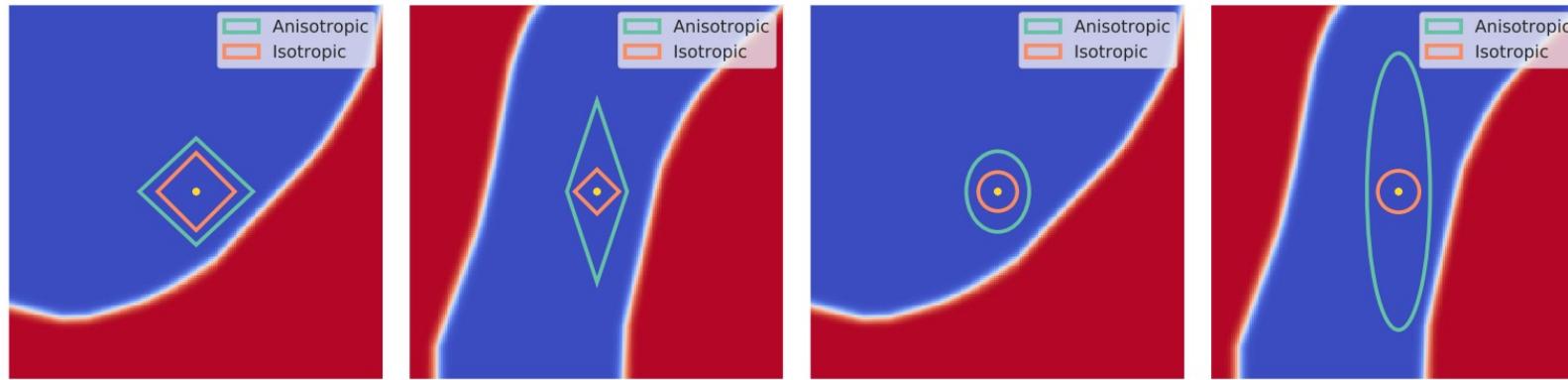
Proposition 3. The classifier g_Λ is $1/2$ -Lipschitz (i.e. $L = 1/2$) under the $\|\Lambda x\|_\infty$ norm.

Similar to Corollary 1, by combining Proposition 3 with Theorem 1, we have that:

Corollary 2. Let $c_A = \arg \max_i g_\Lambda^i(x)$, then $\arg \max_i g_\Lambda^i(x + \delta) = c_A$ for all δ satisfying:

$$\|\delta\|_{\Lambda,1} = \|\Lambda^{-1}\delta\|_1 \leq \left(g_\Lambda^{c_A}(x) - \max_c g_\Lambda^{c \neq c_A}(x) \right).$$

Anisotropic Certification (ANCER, 2022)



Examples of points in a toy 2D dataset, as well as potential isotropic and anisotropic certificates.



Visualization of natural CIFAR-10 images (top) and modified with an imperceptible change that is **not** inside the optimal isotropic certified region but **is** inside the optimal anisotropic one.

- Data-dependent certification requires memorization technique from [Alfarra et al., 2022], imposing a linear memory cost on the number of samples.

Anisotropic Certification (ANCER, 2022)

Certificate volume

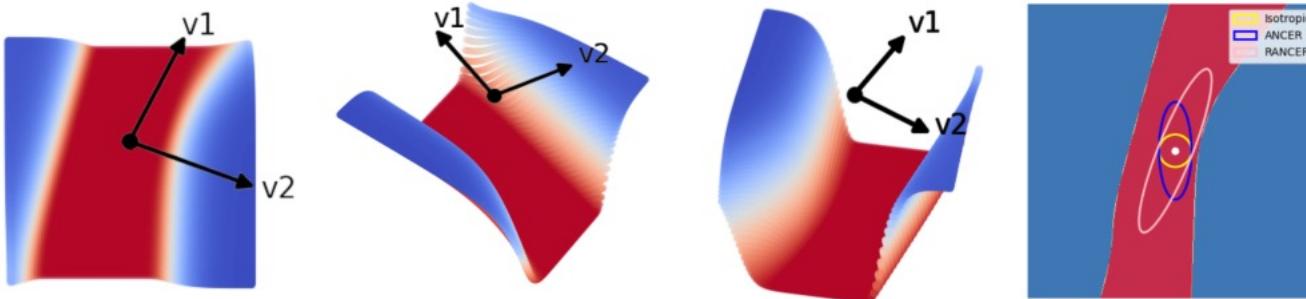


CIFAR-10	Certification	Accuracy @ ℓ_2 radius (%)							$\ell_2 ACR$	$\ell_2^\Sigma ACR\tilde{R}$
		0.0	0.25	0.5	1.0	1.5	2.0	2.5		
COHEN Cohen et al. (2019)	Fixed σ	86	71	51	27	14	6	2	0.722	0.722
	Isotropic DD	82	76	62	39	24	14	8	1.117	1.117
	ANCER	86	85	77	53	31	17	10	1.449	1.772
SMOOTHADV Salman et al. (2019a)	Fixed σ	82	72	55	32	19	9	5	0.834	0.834
	Isotropic DD	82	75	63	40	25	15	7	1.011	1.011
	ANCER	83	81	73	48	30	17	8	1.224	1.573
MACER Zhai et al. (2019)	Fixed σ	87	76	59	37	24	14	9	0.970	0.970
	Isotropic DD	88	80	66	40	17	9	6	1.007	1.007
	ANCER	84	80	67	34	15	11	9	1.136	1.481
ImageNet		Accuracy @ ℓ_2 radius (%)							$\ell_2 ACR$	$\ell_2^\Sigma ACR\tilde{R}$
COHEN Cohen et al. (2019)	Certification	0.0	0.5	1.0	1.5	2.0	2.5	3.0		
	Fixed σ	70	56	41	31	19	14	12	1.098	1.098
	Isotropic DD	71	59	46	36	24	19	15	1.234	1.234
SMOOTHADV Salman et al. (2019a)	ANCER	70	70	62	61	42	36	29	1.810	1.981
	Fixed σ	65	59	44	38	26	20	18	1.287	1.287
	Isotropic DD	66	62	53	41	32	24	20	1.428	1.428
	ANCER	66	66	62	58	44	37	32	1.807	1.965
CIFAR-10		Accuracy @ ℓ_1 radius (%)							$\ell_1 ACR$	$\ell_1^\Sigma ACR\tilde{R}$
RS4A Yang et al. (2020)	Certification	0.0	0.25	0.5	0.75	1.0	1.5	2.0		
	Fixed σ	92	83	75	71	46	0	0	0.775	0.775
	Isotropic DD	92	89	82	76	58	6	2	0.946	0.946
	ANCER	92	90	84	80	63	6	2	0.980	1.104
ImageNet		Accuracy @ ℓ_1 radius (%)							$\ell_1 ACR$	$\ell_1^\Sigma ACR\tilde{R}$
RS4A Yang et al. (2020)	Certification	0.0	0.5	1.0	1.5	2.0	2.5	3.0		
	Fixed σ	78	73	67	63	0	0	0	0.683	0.683
	Isotropic DD	79	76	70	65	46	0	0	0.729	0.729
	ANCER	78	76	70	66	48	0	0	0.730	1.513

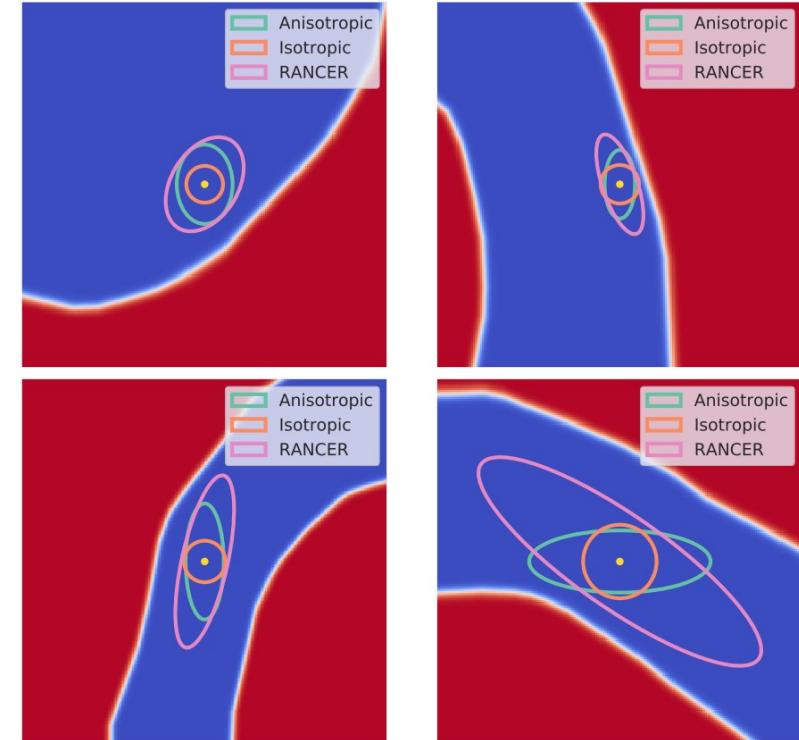
Comparison of top-1 certified accuracy, average certified radius and average certified proxy radius for different certification regimes.

RANCER: Non-Axis Aligned (RANCER, 2023)

- Anisotropic certification so far has been axis-aligned
- RANCER extends ANCER to optimize a full covariance matrix for Gaussian noise, resulting in a non-axis aligned certificate
- To achieve it, we use the eigen decomposition of the Hessian of the loss function at the certified point – corresponding to the “unsafe” directions of maximal change – which should be **smoothed more**
- Once rotated with that basis, optimized in the same way as (ANCER, 2022)



Safe and unsafe directions of the Hessian of the loss function at the point, eventually leading to the pink certified ellipse in the rightmost figure once region is optimized.



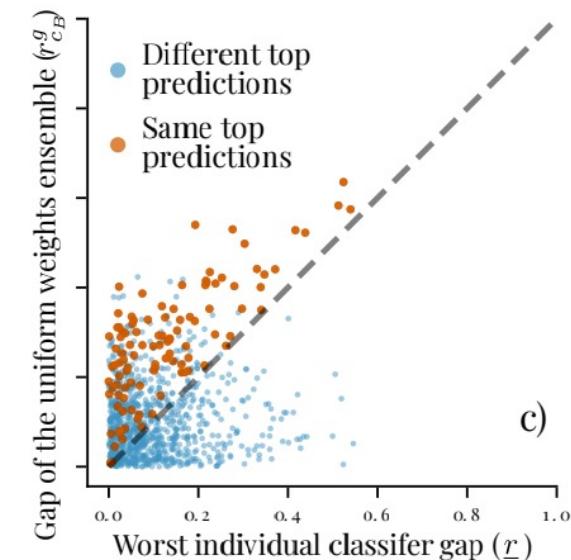
Examples of points in a toy 2D dataset.

Certifying Ensembles with \mathcal{S} -Lipschitzness (SEnsemb, 2023)

- A generalization of Lipschitzness to \mathcal{S} -Lipschitzness (avoiding symmetry requirements) enables tight analysis on the theoretical robustness of ensembles through \mathcal{S} -certificates.

Definition 2 (\mathcal{S} -Lipschitz function). A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is \mathcal{S} -Lipschitz for a bounded set $\mathcal{S} \subset \mathbb{R}^d$ if it holds that:
$$-\rho_{\mathcal{S}}(x - y) \leq f(y) - f(x) \leq \rho_{\mathcal{S}}(y - x), \forall x, y \in \mathbb{R}^d,$$
with $\rho_{\mathcal{S}}(\delta) = \sup_{c \in \mathcal{S}} c^\top \delta$. If \mathcal{S} is convex, then $\rho_{\mathcal{S}}$ corresponds to its support function.

- **Limitation of ensembling robust classifiers:**
 - One can gain robustness by ensembling... But the more robust the individual classifiers, the lower the possible improvement.
 - One can be worse off ensembling than using a single classifier... Even potentially losing all robustness (radius collapses to 0).
 - One will never be worse off if all classifiers predict the same class – but then why ensemble?
 - Usually end up losing more robustness than gaining → **ensembling robust classifiers reduces robustness.**



01

Safety in Human Interaction/
Understanding Decision-Making Process

02

Correctness of ML-based Systems



Certifiable ML – a network bounding problem

Safety/Robustness

(SEnsemb, 2023) Certifying Ensembles: A General Certification Theory with S-Lipschitzness

A Petrov, F Eiras, A Sanyal, PHS Torr, A Bibi (ICML)

(RANCER, 2023) RANCER: Non-Axis Aligned Anisotropic Certification With Randomized Smoothing

T Rumezhak, F Eiras, PHS Torr, A Bibi (WACV)

(ANCER, 2022) AnCer: Anisotropic certification via sample-wise volume maximization

F Eiras, M Alfarra, MP Kumar, PHS Torr, PK Dokania, B Ghanem, A Bibi (TMLR)

Reliability/Correctness

(PINN, 2023) Provably Correct Physics-Informed Neural Networks

F Eiras, A Bibi, R Bunel, KD Dvijotham, PHS Torr, MP Kumar (Pre-print)

(FKD, 2023) Faithful Knowledge Distillation

T Lamb, R Bunel, KD Dvijotham, PHS Torr, MP Kumar, F Eiras (Pre-print)

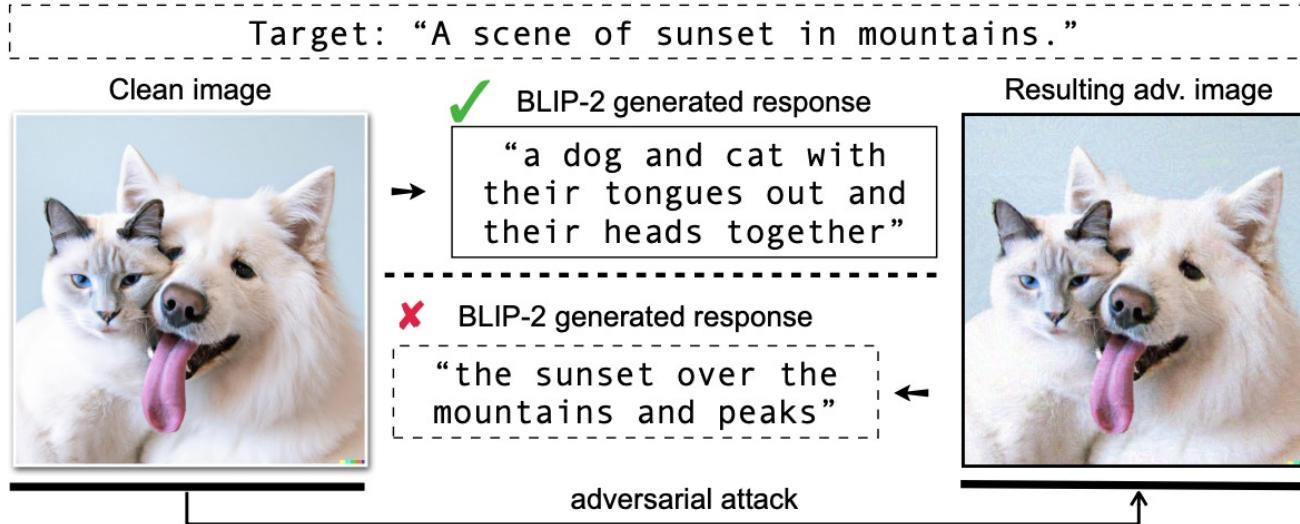
Thank you

eiras@robots.ox.ac.uk

Twitter: @fgirbal

LinkedIn: franciscogirbaleiras

Adversarial Robustness of Vision-Language Models



A safety problem



From Zhao, Yunqing, et al. "On Evaluating Adversarial Robustness of Large Vision-Language Models." arXiv preprint arXiv:

arXiv:2305.16934v1 [cs.CV] 26 May 2023

On Evaluating Adversarial Robustness of Large Vision-Language Models

Yunqing Zhao^{*1}, Tianyu Pang^{*†2}, Chao Du^{†2}, Xiao Yang³, Chongxuan Li⁴, Ngai-Man Cheung^{†1}, Min Lin²

¹Singapore University of Technology and Design
²Sea AI Lab, Singapore
³Tsinghua University
⁴Renmin University of China
yunqing_zhao@mymail.sutd.edu.sg, ngaiman_cheung@sutd.edu.sg, {tianyupang, duchao, linmin}@sea.com, yangxiao19@mails.tsinghua.edu.cn, chongxuanli@ruc.edu.cn

Abstract

Large vision-language models (VLMs) such as GPT-4 have achieved unprecedented performance in response generation, especially with visual inputs, enabling more creative and adaptable interaction than large language models such as ChatGPT. Nonetheless, multimodal generation exacerbates safety concerns, since adversaries may successfully evade the entire system by subtly manipulating the most vulnerable modality (e.g., vision). To this end, we propose evaluating the robustness of open-source large VLMs in the most realistic and high-risk setting, where adversaries have only *black-box* system access and seek to deceive the model into returning the *targeted* responses. In particular, we first craft targeted adversarial examples against pretrained models such as CLIP and BLIP, and then transfer these adversarial examples to other VLMs such as MiniGPT4, LLaVA, UniDiffuser, BLIP-2, and Img2Prompt. In addition, we observe that black-box queries on these VLMs can further improve the effectiveness of targeted evasion, resulting in a surprisingly high success rate for generating targeted responses. Our findings provide a quantitative understanding regarding the adversarial vulnerability of large VLMs and call for a more thorough examination of their potential security flaws before deployment in practice. Code is at <https://github.com/yunqing-me/AttackVLM>.

1 Introduction

Large vision-language models (VLMs) have enjoyed tremendous success and demonstrated promising capabilities in text-to-image generation [52, 65, 69], image-grounded text generation (e.g., image captioning or visual question-answering) [2, 14, 40, 83], and joint generation [5, 30, 95] due to an increase in the amount of data, computational resources, and number of model parameters. Notably, after being finetuned with instructions and aligned with human feedback, GPT-4 [55] is capable of conversing with human users and, in particular, supports visual inputs.

Along the trend of multimodal learning, vision-and-language models (ViLMs) are made available