# Efficiently Training Neural Networks for Verified Robustness

**Alessandro De Palma**
Imperial College London

# Adversarial Examples



"panda"
57.7% confidence

[Goodfellow et al., 2015]

# Adversarial Examples



"panda"
57.7% confidence

$+ .007 \times$

"nematode"
8.2% confidence

[Goodfellow et al., 2015]

# Adversarial Examples



$+.007 \times$

"panda"
57.7% confidence

"nematode"
8.2% confidence

$=$

"gibbon"
99.3 % confidence

[Goodfellow et al., 2015]

# Outline

- **Neural Network Verification**

- Training for Verified Robustness

- NLP?

- Discussion
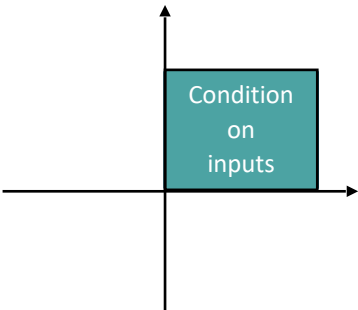
# Neural Network Verification
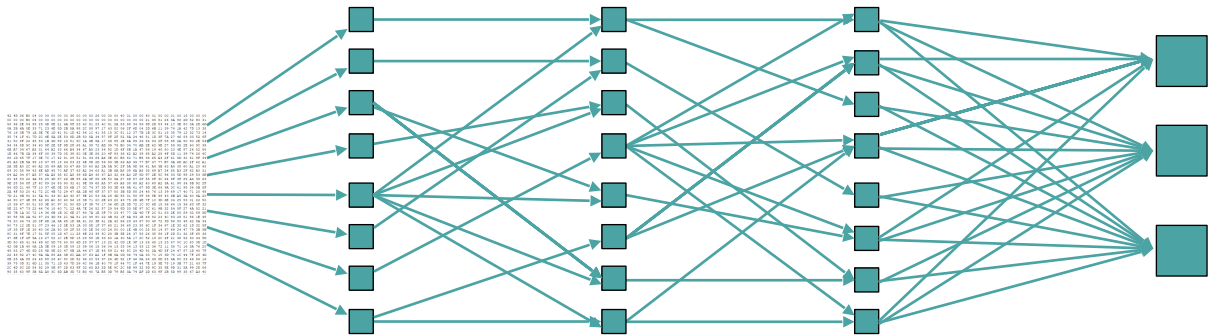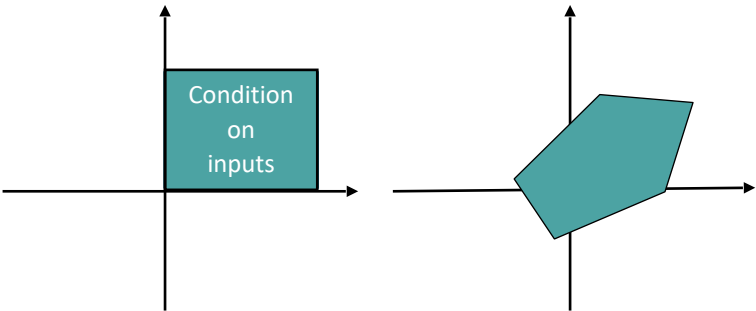


Condition on inputs

*Image Deformations*

# Neural Network Verification
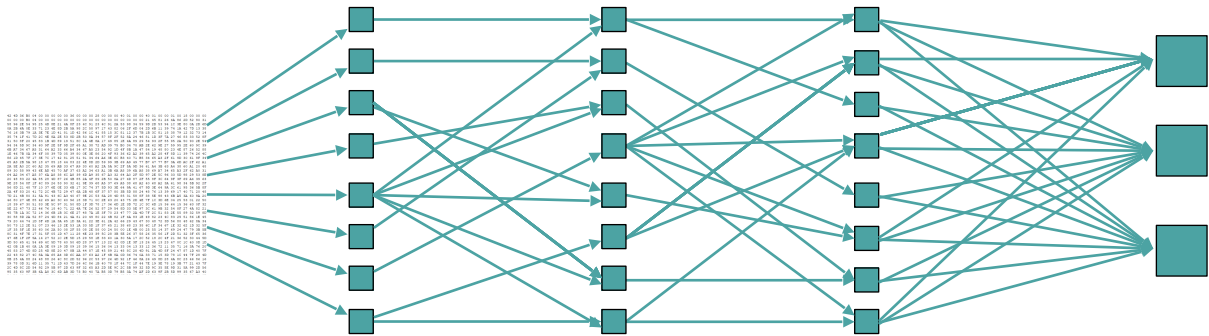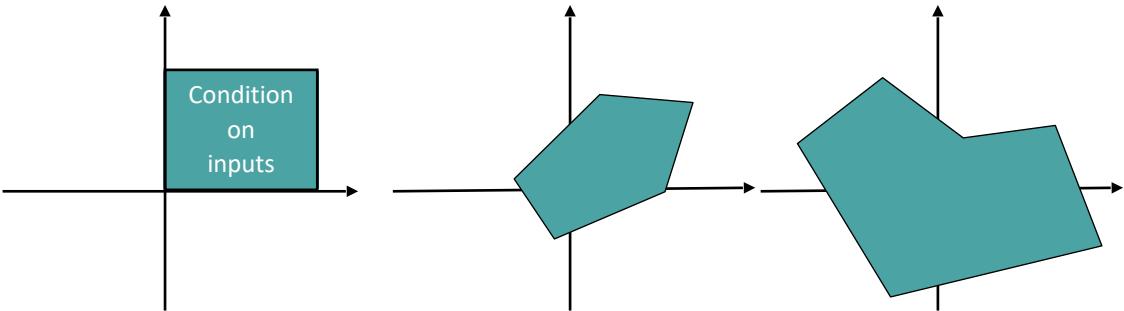


Image Deformations

# Neural Network Verification



Condition on inputs

*Image Deformations*

# Neural Network Verification



Condition on inputs

*Image Deformations*

# Neural Network Verification
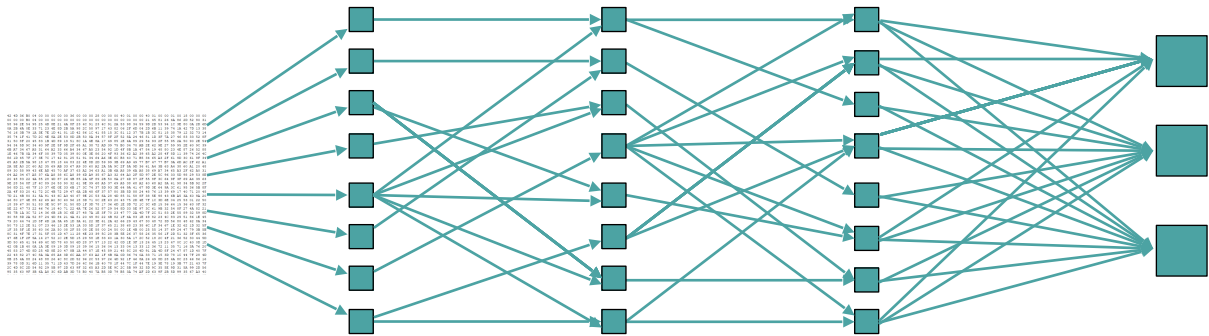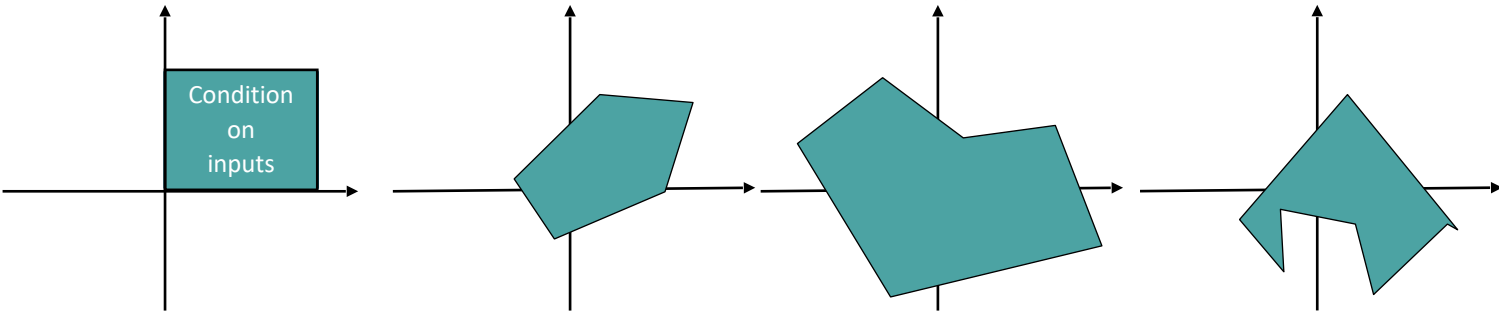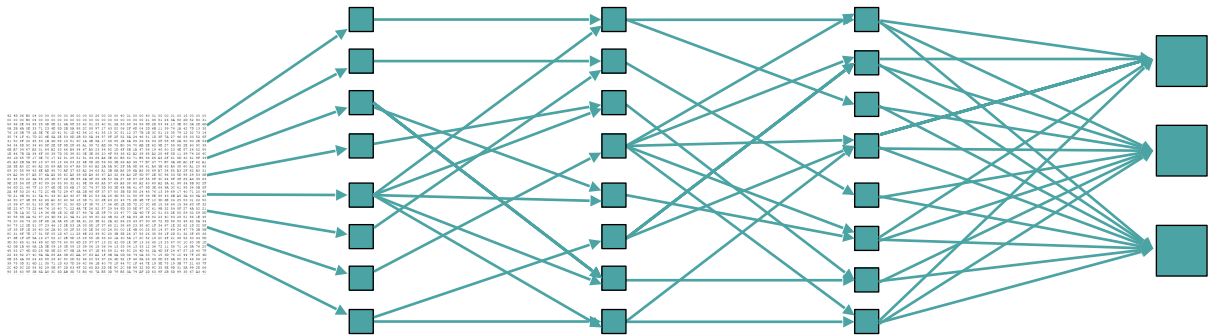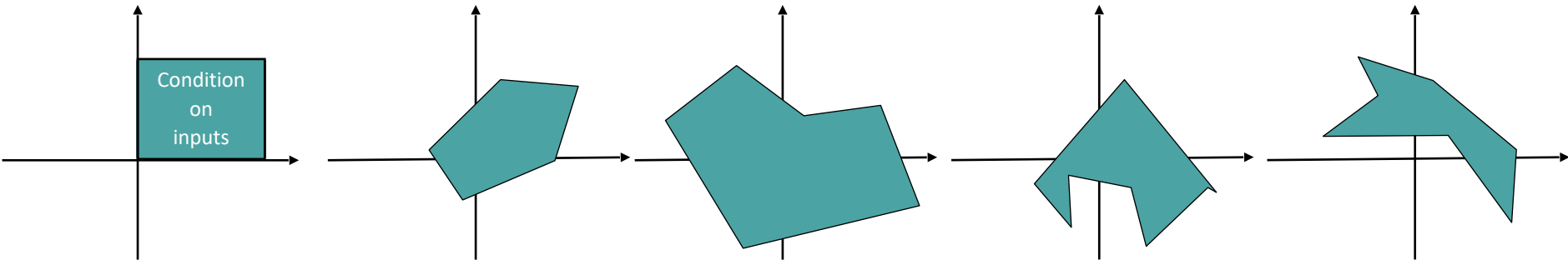


*Image Deformations*

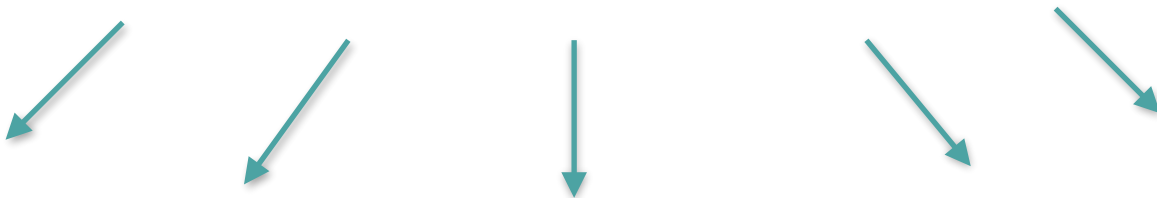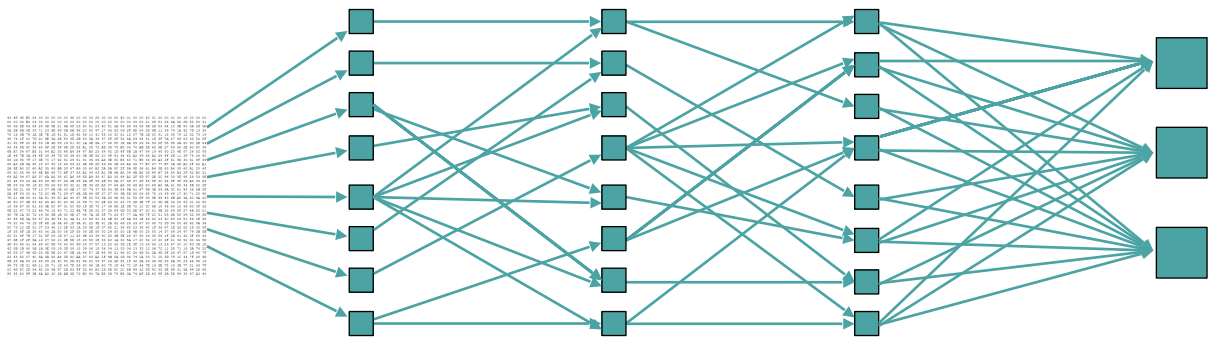# Neural Network Verification



Image Deformations

Condition on inputs

Safe

Error

Classifications

**Property is satisfied**

# Neural Network Verification



*Image Deformations*

Condition on inputs

Safe

Error

*Classifications*

NP-HARD

**Property is satisfied**

# Branch and Bound: Bounding



*Image Deformations*

# Branch and Bound: Bounding



Image Deformations

# Branch and Bound: Bounding



*Image Deformations*

# Branch and Bound: Bounding



Image Deformations

# Branch and Bound: Bounding



*Image Deformations*

# Branch and Bound: Bounding



*Image Deformations*

Verify <u>subset</u> of properties

# Branch and Bound: Bounding



Safe

Error

Condition on inputs

*Image Deformations*

*Classifications*

Verify <u>subset</u> of properties

**Property is not satisfied**

# Convex Relaxation: Planet

# Convex Relaxation: Planet

$$\sigma\left(\hat{\mathbf{x}}_k\right)$$

# Convex Relaxation: Planet

$$\sigma(\hat{\mathbf{x}}_k)$$

# Convex Relaxation: Planet

$$\sigma\left(\hat{\mathbf{x}}_k\right) \qquad\qquad \mathrm{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \hat{\mathbf{u}}_k)$$

# Convex Relaxation: Planet

$$\sigma\left(\hat{\mathbf{x}}_k\right) \qquad \text{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \hat{\mathbf{u}}_k)$$



[Ehlers 2017, Wong and Kolter, 2018; Zhang et al., 2018; Dvijotham et al. 2018; Singh et al. 2018; Bunel et al., 2020, Xu et al. 2021, Wang et al. 2021]

# Branch and Bound: Branching



Condition on inputs

*Image Deformations*

# Branch and Bound: Branching



Condition on inputs

Image Deformations

# Branch and Bound: Branching



Condition on inputs

*Image Deformations*

# Branch and Bound: Branching

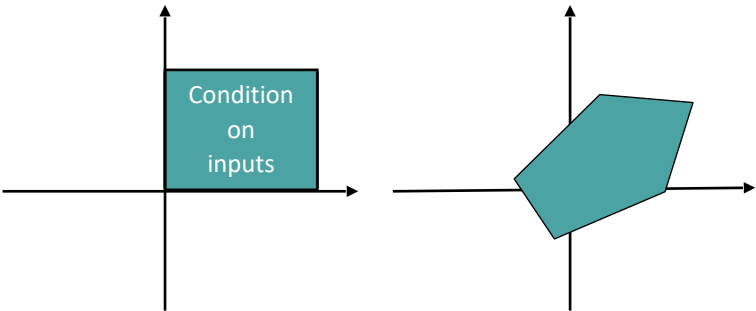Condition
on
inputs

*Image Deformations*
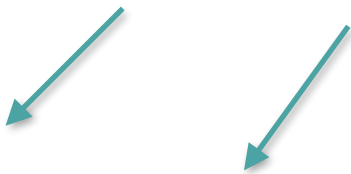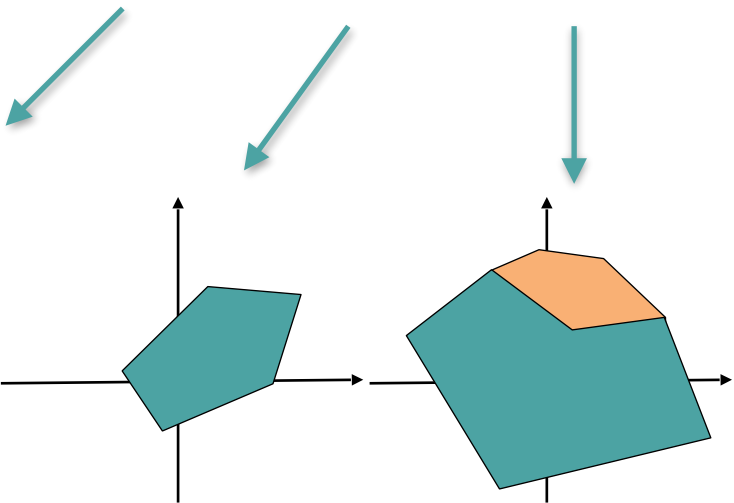
# Branch and Bound: Branching



Condition on inputs

*Image Deformations*

# Branch and Bound: Branching



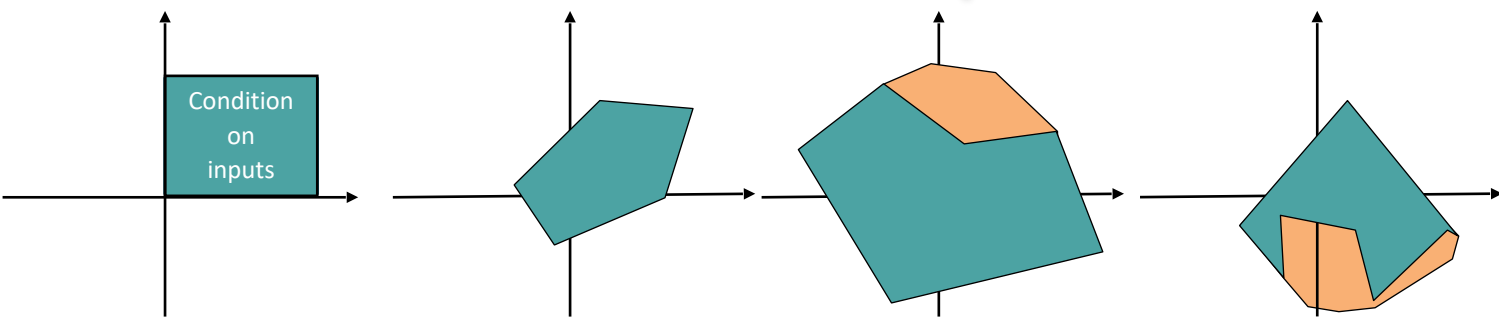*Image Deformations*

Verify <u>all</u> properties via iterative branching

# Branch and Bound: Branching



Condition on inputs

*Image Deformations*

Safe

Error

*Classifications*

Verify <u>all</u> properties via iterative branching

**Property is satisfied**

# Activation Splitting

$$\mathrm{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \hat{\mathbf{u}}_k)$$



[Bunel et al. 2020, De Palma et al. 2021, Mueller et al. 2022]

# Activation Splitting

$$\mathrm{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \hat{\mathbf{u}}_k)$$



$\longrightarrow$

[Bunel et al. 2020, De Palma et al. 2021, Mueller et al. 2022]

# Activation Splitting

$\text{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \hat{\mathbf{u}}_k)$

$\text{Conv}(\sigma(\hat{\mathbf{x}}_k), \mathbf{0}, \hat{\mathbf{u}}_k)$

$\bigcup$

$\text{Conv}(\sigma(\hat{\mathbf{x}}_k), \hat{\mathbf{l}}_k, \mathbf{0})$

$\longrightarrow$

[Bunel et al. 2020, De Palma et al. 2021, Mueller et al. 2022]

# Outline

- Neural Network Verification

- **Training for Verified Robustness**

- NLP?

- Discussion

# Robust Loss

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x},\mathbf{y})\in\mathcal{D}} \left[ \max_{\mathbf{x}'\in\mathcal{C}(\mathbf{x})} \mathcal{L}(f(\boldsymbol{\theta},\mathbf{x}'),\mathbf{y}) \right]$$

# Robust Loss

$$\min_{\boldsymbol{\theta}} \ \mathbb{E}_{(\mathbf{x},\mathbf{y})\in\mathcal{D}} \left[ \max_{\mathbf{x}'\in\mathcal{C}(\mathbf{x})} \mathcal{L}(f(\boldsymbol{\theta},\mathbf{x}'),\mathbf{y}) \right]$$

$$\mathcal{L}^*(f(\boldsymbol{\theta},\mathbf{x}),y)$$

# Adversarial Training

Lower bound $\to$ adversarial training

[Madry et al. 2018]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

# Adversarial Training

Lower bound $\rightarrow$ adversarial training

[Madry et al. 2018]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \geq \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y)$$

# Adversarial Training

Lower bound $\rightarrow$ adversarial training

[Madry et al. 2018]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \qquad \geq \qquad \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y)$$



formal guarantees?

# Verified Training

Upper bound $\rightarrow$ certified training

[Wong and Kolter 2018, Gowal et al. 2018, Zhang et al. 2020, Shi et al. 2021]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

# Verified Training

Upper bound $\rightarrow$ certified training

[Wong and Kolter 2018, Gowal et al. 2018, Zhang et al. 2020, Shi et al. 2021]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \qquad \leq \qquad \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

# Verified Training

Upper bound $\rightarrow$ certified training

[Wong and Kolter 2018, Gowal et al. 2018, Zhang et al. 2020, Shi et al. 2021]

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \quad \leq \quad \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$$



verification via cheap incomplete verifiers

# Verified Training

# Verified Training

IBP



CROWN

# Loss Expressivity

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

# Loss Expressivity

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \qquad \leq$$

# Loss Expressivity

$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), y) \qquad \leq \qquad ?$$

# Hybrid Training Methods: SABR

Compute over-approximation over a parametrized subset of the input domain that includes an adversarial attack.



[Mueller et al., 2023]

# Hybrid Training Methods: SABR

Compute over-approximation over a parametrized subset of the input domain that includes an adversarial attack.



$\mathcal{C}(\mathbf{x}, \epsilon)$

$\mathbf{x}_{\text{adv}}$

$\rightarrow$

$\mathcal{C}(\mathbf{x}_\lambda, \lambda\epsilon) \subseteq \mathcal{C}(\mathbf{x}, \epsilon)$

$\mathbf{x}_{\text{adv}}$

$2\lambda\epsilon$

[Mueller et al., 2023]

# Hybrid Training Methods: SABR

Compute over-approximation over a parametrized subset of the input domain that includes an adversarial attack.



$$\mathcal{C}(\mathbf{x}, \epsilon)$$

$$\rightarrow$$

$$\mathcal{C}(\mathbf{x}_\lambda, \lambda\epsilon) \subseteq \mathcal{C}(\mathbf{x}, \epsilon)$$

$$2\lambda\epsilon$$

$\mathbf{x}_{\text{adv}}$

verification via BaB

[Mueller et al., 2023]

# Expressive Losses for Verified Robustness via Convex Combinations

**Alessandro De Palma**
Imperial College London[1]
adepalma@ic.ac.uk

**Rudy Bunel**
Google DeepMind

**Krishnamurthy (Dj) Dvijotham**
Google DeepMind

**M. Pawan Kumar**
Google DeepMind

**Robert Stanforth**
Google DeepMind

**Alessio Lomuscio**
Imperial College London

https://arxiv.org/abs/2305.13991

# Loss Expressivity

A parametrized family of losses $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is *expressive* if:

# Loss Expressivity

A parametrized family of losses $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is *expressive* if:

- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \ \forall \ \alpha \in [0, 1]$;

# Loss Expressivity

A parametrized family of losses $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is *expressive* if:

- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \ \forall \ \alpha \in [0, 1];$

- $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is monotonically increasing with $\alpha$;

# Loss Expressivity

A parametrized family of losses $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is *expressive* if:

- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \; \forall \; \alpha \in [0, 1]$;

- $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is monotonically increasing with $\alpha$;

- $\mathcal{L}_0(\boldsymbol{\theta}, \mathbf{x}, y) = \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y)$;

# Loss Expressivity

A parametrized family of losses $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is *expressive* if:

- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \; \forall \; \alpha \in [0, 1]$;

- $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$ is monotonically increasing with $\alpha$;

- $\mathcal{L}_0(\boldsymbol{\theta}, \mathbf{x}, y) = \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y)$;

- $\mathcal{L}_1(\boldsymbol{\theta}, \mathbf{x}, y) = \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$.

# Expressivity via Convex Combinations

**CC-IBP**

$$\mathcal{L}(-\,[(1-\alpha)\ \text{}\ +\alpha\ \text{}\ ]\,,y)$$

# Expressivity via Convex Combinations

**CC-IBP**

$$\mathcal{L}(-[(1-\alpha) \quad  \quad + \alpha \quad  \quad ], y)$$

**MTL-IBP**

$$(1-\alpha)\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\mathrm{adv}}), y) + \alpha \; \mathcal{L}_{\mathrm{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

# Loss Sensitivity

Sensitivity of CC-IBP and MTL-IBP to the convex combination coefficient $\alpha$ on the first 1000 CIFAR-10 test images.



(a) CC-IBP, $\epsilon = {}^2/_{255}$.    (b) MTL-IBP, $\epsilon = {}^2/_{255}$.    (c) CC-IBP, $\epsilon = {}^8/_{255}$.    (d) MTL-IBP, $\epsilon = {}^8/_{255}$.

# Experimental Results

Performance of different verified training algorithms under $\ell_\infty$ norm perturbations on the CIFAR-10 dataset.

| Dataset | $\epsilon$ | Method | Standard acc. [%] | Verified rob. acc. [%] | Training time [s] |
|---------|-----------|--------|-------------------|------------------------|-------------------|
| CIFAR-10 | $\frac{2}{255}$ | CC-IBP | <u>80.09</u> | **63.78** | $1.77 \times 10^4$ |
| | | MTL-IBP | **80.11** | <u>63.24</u> | $1.76 \times 10^4$ |
| | | STAPS | 79.76 | 62.98 | $1.41 \times 10^5$ |
| | | SABR | 79.24 | 62.84 | $2.56 \times 10^4$ |
| | | SORTNET | 67.72 | 56.94 | $4.04 \times 10^4$ |
| | | IBP-R | 78.19 | 61.97 | $9.34 \times 10^3$ |
| | | CROWN-IBP | 71.52 | 53.97 | $9.13 \times 10^4$ |
| | $\frac{8}{255}$ | CC-IBP | <u>53.71</u> | <u>35.27</u> | $1.72 \times 10^4$ |
| | | MTL-IBP | <u>53.35</u> | <u>35.44</u> | $1.70 \times 10^4$ |
| | | STAPS | 52.82 | 34.65 | $2.70 \times 10^4$ |
| | | SABR | 52.38 | 35.13 | $2.64 \times 10^4$ |
| | | SORTNET | **54.84** | **40.39** | $4.04 \times 10^4$ |
| | | IBP-R | 52.74 | 27.55 | $5.89 \times 10^3$ |
| | | IBP | 48.94 | 34.97 | $9.51 \times 10^3$ |

# Experimental Results

Performance of different verified training algorithms under $\ell_\infty$ norm perturbations on the TinyImageNet and downscaled ($64 \times 64$) ImageNet datasets.

| Dataset | $\epsilon$ | Method | Standard acc. [%] | Verified rob. acc. [%] | Training time [s] |
|---|---|---|---|---|---|
| TinyImageNet | $\frac{1}{255}$ | CC-IBP | <u>32.71</u> | <u>23.10</u> | $6.58 \times 10^4$ |
| | | MTL-IBP | **32.76** | **24.14** | $6.56 \times 10^4$ |
| | | STAPS | 28.98 | 22.16 | $3.06 \times 10^5$ |
| | | SABR | 28.85 | 20.46 | $2.07 \times 10^5$ |
| | | SORTNET | 25.69 | 18.18 | $1.56 \times 10^5$ |
| | | IBP | 25.92 | 17.87 | $3.53 \times 10^4$ |
| ImageNet64 | $\frac{1}{255}$ | CC-IBP | <u>19.62</u> | <u>11.87</u> | $3.26 \times 10^5$ |
| | | MTL-IBP | **20.15** | **12.13** | $3.52 \times 10^5$ |
| | | SORTNET | 14.79 | 9.54 | $6.58 \times 10^5$ |
| | | CROWN-IBP | 16.23 | 8.73 | / |

# Future Work

- Theoretical understanding of relative method performance;

# Future Work

- Theoretical understanding of relative method performance;

- Would more network capacity help?

# Future Work

- Theoretical understanding of relative method performance;

- Would more network capacity help?

- Examine applications to different data domains.

# Outline

- Neural Network Verification

- Training for Verified Robustness

- **NLP?**

- Discussion

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

- More complexity than PLNN;

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

- More complexity than PLNN;

- LLMs: very large, pre-trained;

## ChatGPT

Software

ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched on November 30, 2022. It is notable for enabling users to refine and steer a conversation towards a desired length, format, style, level of detail, and language used. Wikipedia

**Initial release date:** 30 November 2022

**Platform:** Cloud computing

**Programming language:** Python

**Developer:** OpenAI, Microsoft Corporation

**Engine:** GPT-3.5; GPT-4

**License:** Proprietary

**Stable release:** May 24, 2023; 40 days ago

People also search for                        View 5+ more

GPT-4          Bard          GPT-3          Generative pre-train...

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

- More complexity than PLNN;

- LLMs: very large, pre-trained;

- Interaction: black-box, prompt-based (zero) few-shot learning;

## ChatGPT
Software

ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched on November 30, 2022. It is notable for enabling users to refine and steer a conversation towards a desired length, format, style, level of detail, and language used. Wikipedia

**Initial release date:** 30 November 2022

**Platform:** Cloud computing

**Programming language:** Python

**Developer:** OpenAI, Microsoft Corporation

**Engine:** GPT-3.5; GPT-4

**License:** Proprietary

**Stable release:** May 24, 2023; 40 days ago

People also search for                    View 5+ more
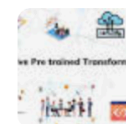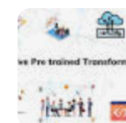
GPT-4        Bard        GPT-3        Generative pre-train...

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

- More complexity than PLNN;

- LLMs: very large, pre-trained;

- Interaction: black-box, prompt-based (zero) few-shot learning;

- Hallucinations?

## ChatGPT

Software

ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched on November 30, 2022. It is notable for enabling users to refine and steer a conversation towards a desired length, format, style, level of detail, and language used. Wikipedia

**Initial release date:** 30 November 2022

**Platform:** Cloud computing

**Programming language:** Python

**Developer:** OpenAI, Microsoft Corporation

**Engine:** GPT-3.5; GPT-4

**License:** Proprietary

**Stable release:** May 24, 2023; 40 days ago

People also search for          View 5+ more

GPT-4          Bard          GPT-3          Generative pre-train...

# NLP Challenges

- Threat model: robustness to synonym-based perturbations;

- More complexity than PLNN;

- LLMs: very large, pre-trained;

- Interaction: black-box, prompt-based (zero) few-shot learning;

- Hallucinations?

- Textual inputs are discrete.

## ChatGPT

Software

ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched on November 30, 2022. It is notable for enabling users to refine and steer a conversation towards a desired length, format, style, level of detail, and language used. Wikipedia

**Initial release date:** 30 November 2022

**Platform:** Cloud computing

**Programming language:** Python

**Developer:** OpenAI, Microsoft Corporation

**Engine:** GPT-3.5; GPT-4

**License:** Proprietary

**Stable release:** May 24, 2023; 40 days ago

People also search for                    View 5+ more
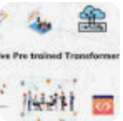
GPT-4          Bard          GPT-3          Generative pre-train...

# Discussion