

# Expressive Losses for Verified Robustness via Convex Combinations

Alessandro De Palma *Inria*  | PSL 

Rudy Bunel, Krishnamurthy Dvijotham,  
M. Pawan Kumar, Robert Stanforth 

Alessio Lomuscio 

# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

+ .007 ×



“nematode”  
8.2% confidence

# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

+ .007 ×



“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

[\[Goodfellow et al., 2015\]](#)

# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

+ .007 ×



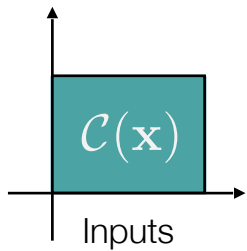
“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

[Goodfellow et al., 2015]



# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

+ .007 ×



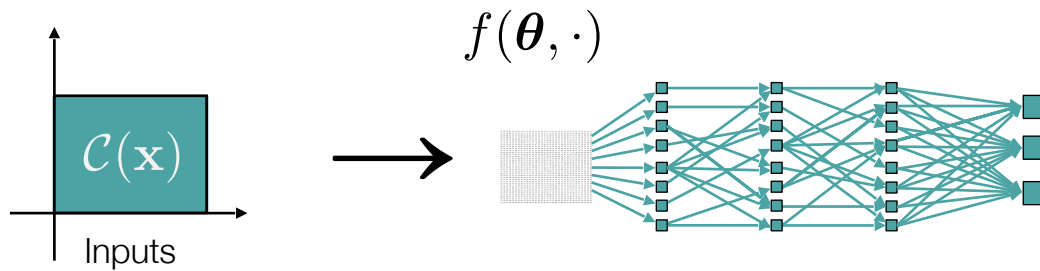
“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

[Goodfellow et al., 2015]



# (Verified) Adversarial Robustness



“panda”  
57.7% confidence

+ .007 ×



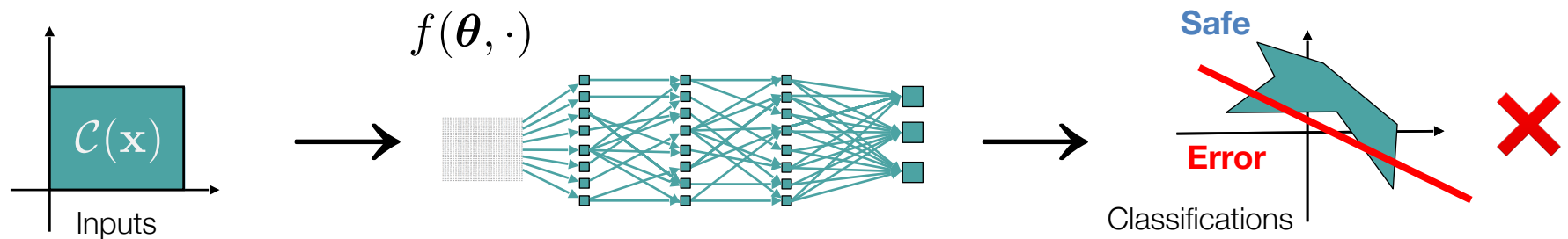
“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

[Goodfellow et al., 2015]



# (Verified) Adversarial Robustness



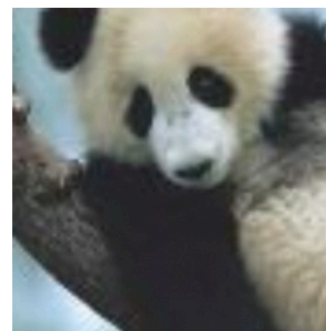
“panda”  
57.7% confidence

+ .007 ×



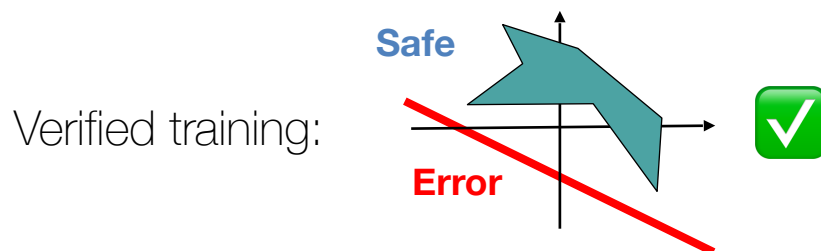
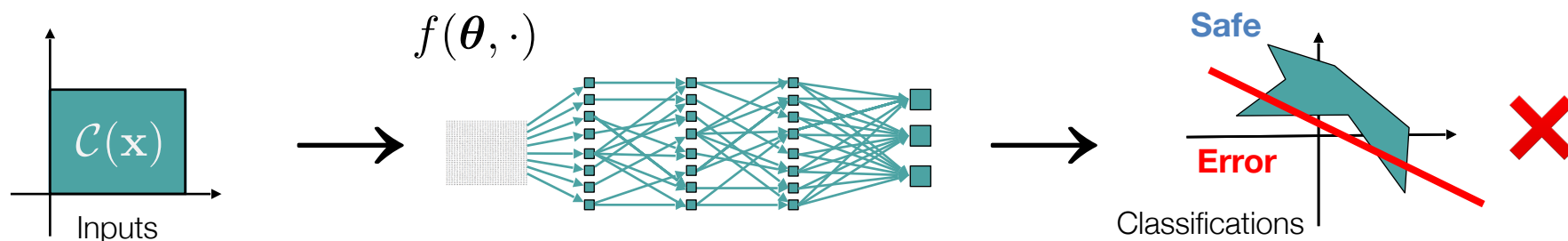
“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

[Goodfellow et al., 2015]





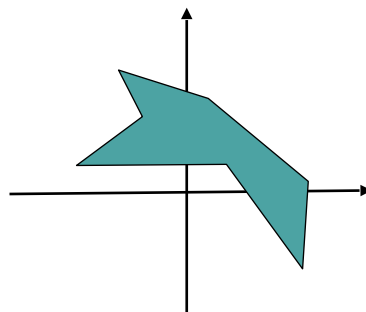
# Robust Training

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

# Robust Training

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

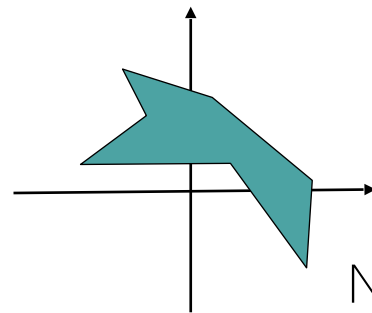
$$\mathcal{L}^*(f(\boldsymbol{\theta}, \mathbf{x}), \mathbf{y})$$



# Robust Training

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

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



NP-Complete

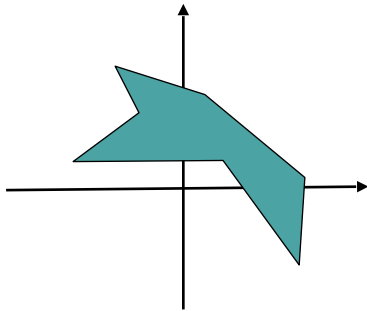
[Katz et al., 2017]

# Adversarial Training

Lower bound  $\rightarrow$  adversarial training

[Madry et al., 2018](#), [Wong et al., 2020](#)

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

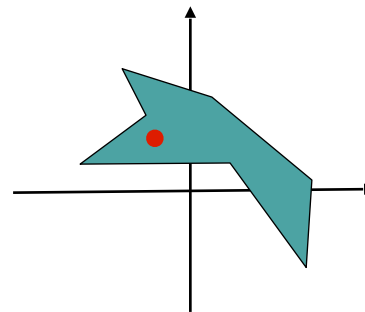
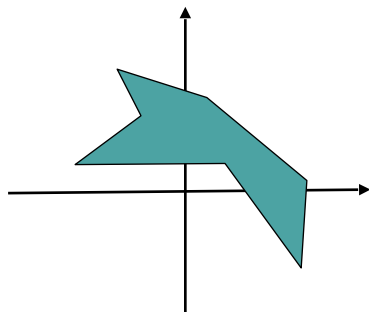


# Adversarial Training

Lower bound  $\rightarrow$  adversarial training

[Madry et al., 2018](#), [Wong et al., 2020](#)

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

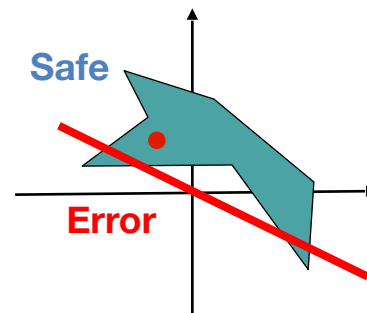
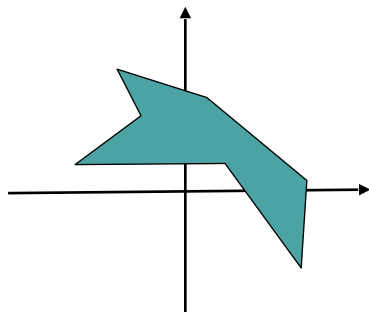


# Adversarial Training

Lower bound  $\rightarrow$  adversarial training

[\[Madry et al., 2018, Wong et al., 2020\]](#)

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

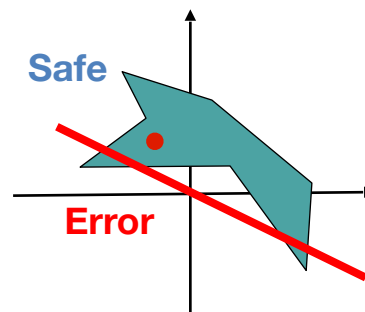
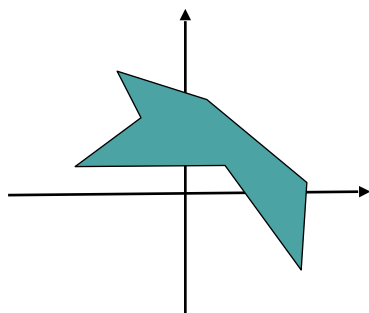


# Adversarial Training

Lower bound  $\rightarrow$  adversarial training

[\[Madry et al., 2018, Wong et al., 2020\]](#)

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



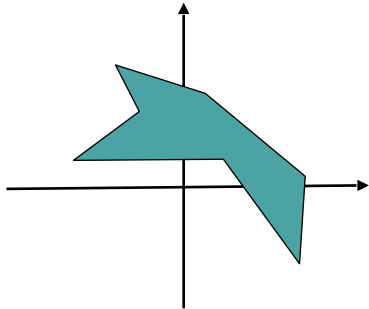
formal guarantees?

# Verified Training

Upper bound  $\rightarrow$  certified training

[\[Gowal et al., 2018, Zhang et al., 2020, Shi et al., 2021\]](#)

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



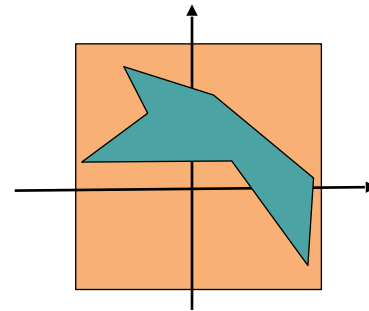
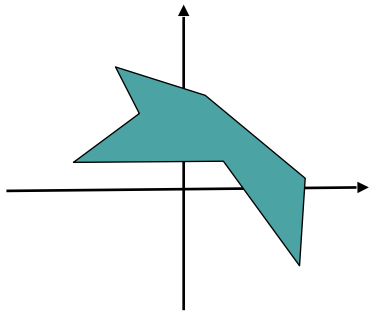


# Verified Training

Upper bound  $\rightarrow$  certified training

[\[Gowal et al., 2018, Zhang et al., 2020, Shi et al., 2021\]](#)

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

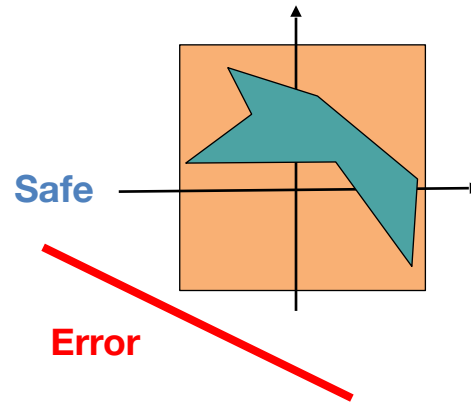
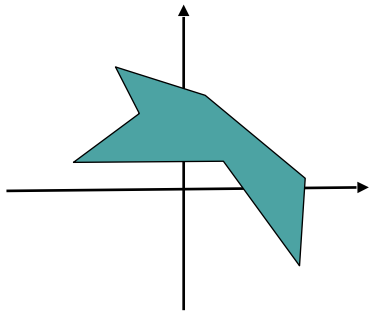


# Verified Training

Upper bound  $\rightarrow$  certified training

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

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

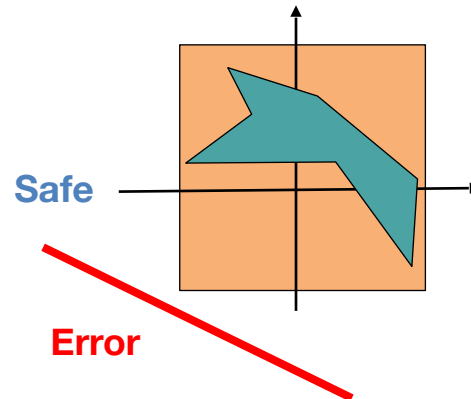
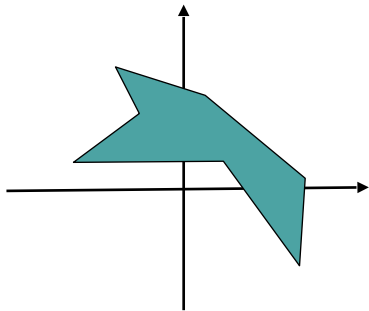


# Verified Training

Upper bound  $\rightarrow$  certified training

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

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

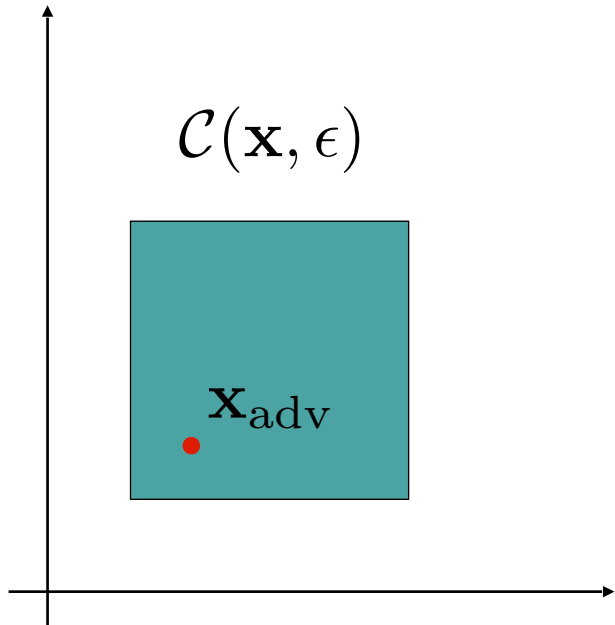


large cost in standard performance

# Hybrid Training: SABR

[Müller et al., 2023]

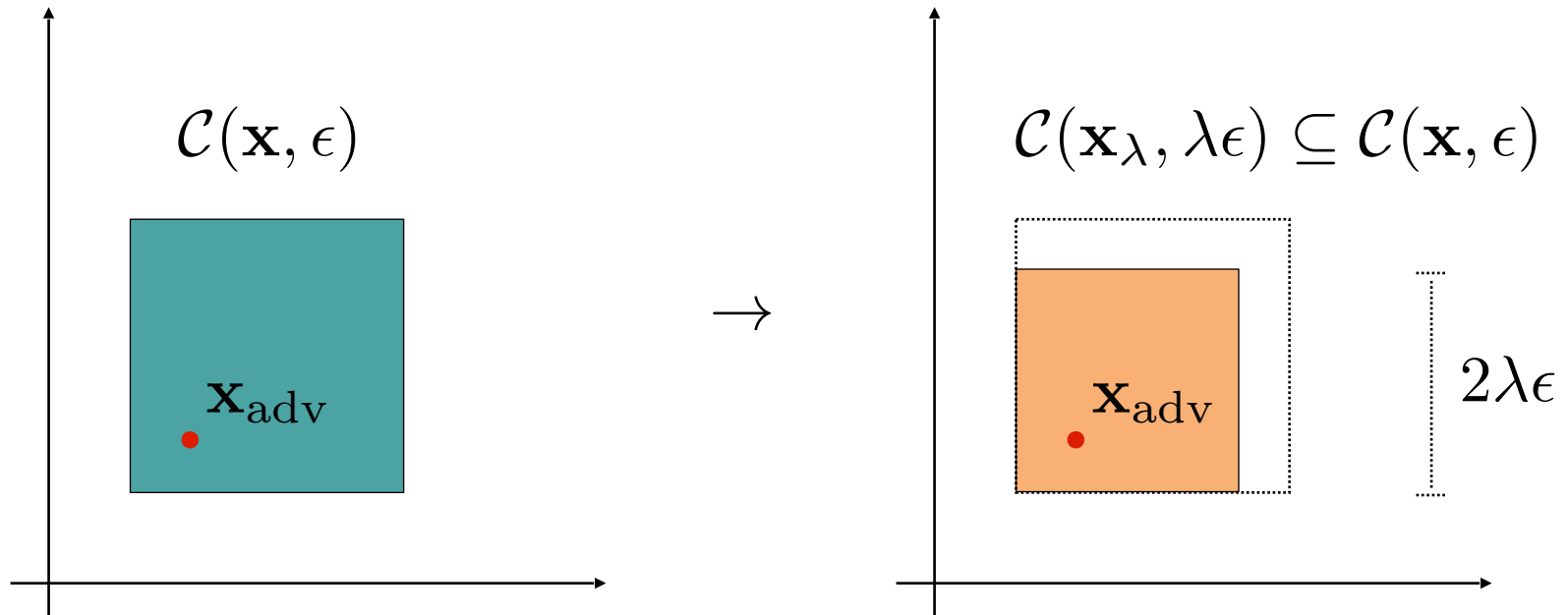
Compute over-approximation over a parametrized subset of the input domain that includes an adversarial attack.



# Hybrid Training: SABR

[Müller et al., 2023]

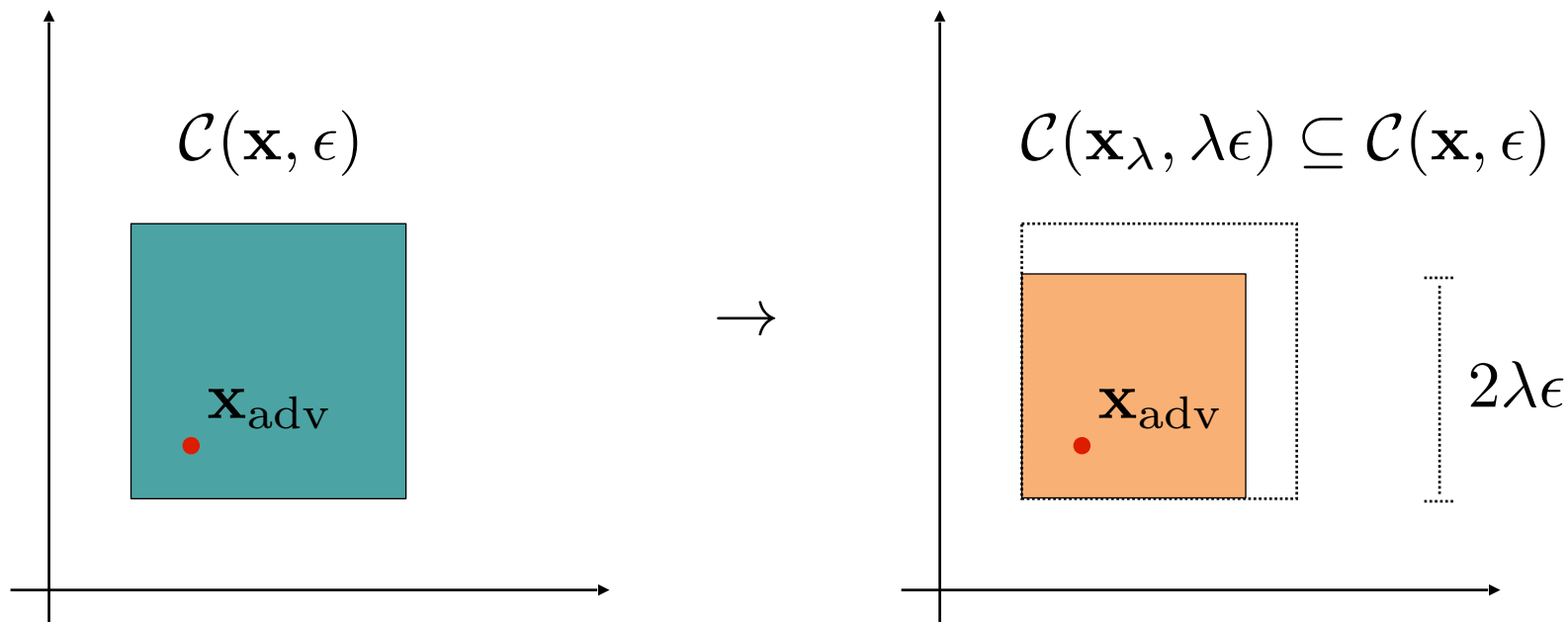
Compute over-approximation over a parametrized subset of the input domain that includes an adversarial attack.



# Hybrid Training: SABR

[Müller et al., 2023]

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



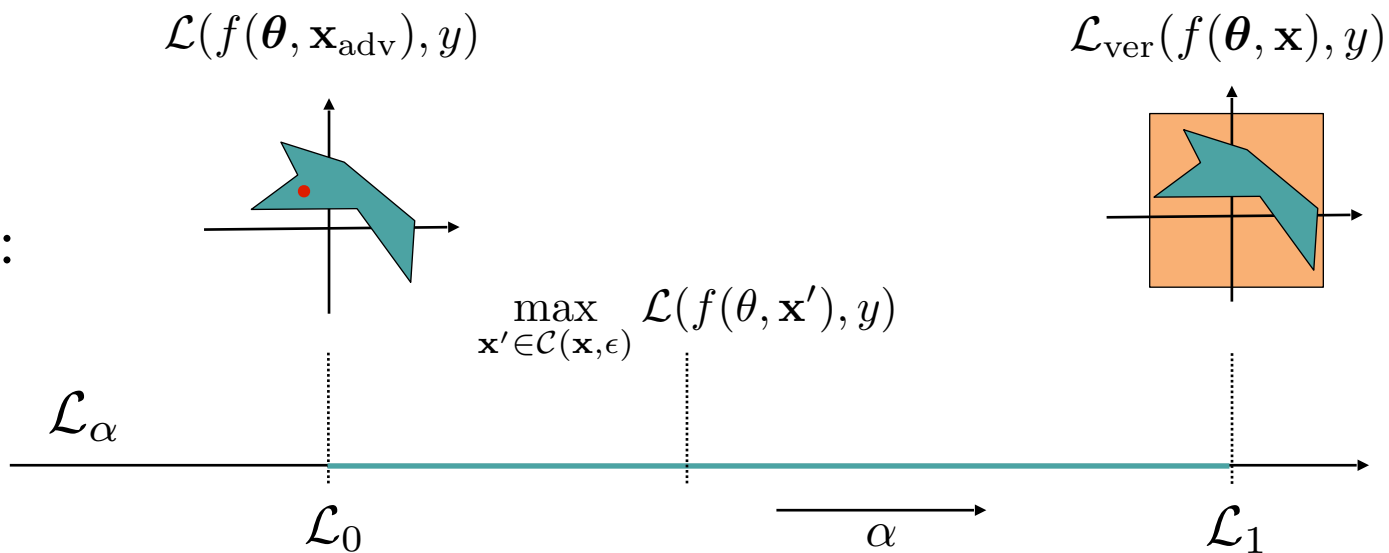
state-of-the-art results

# Loss Expressivity

$$\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) :$$

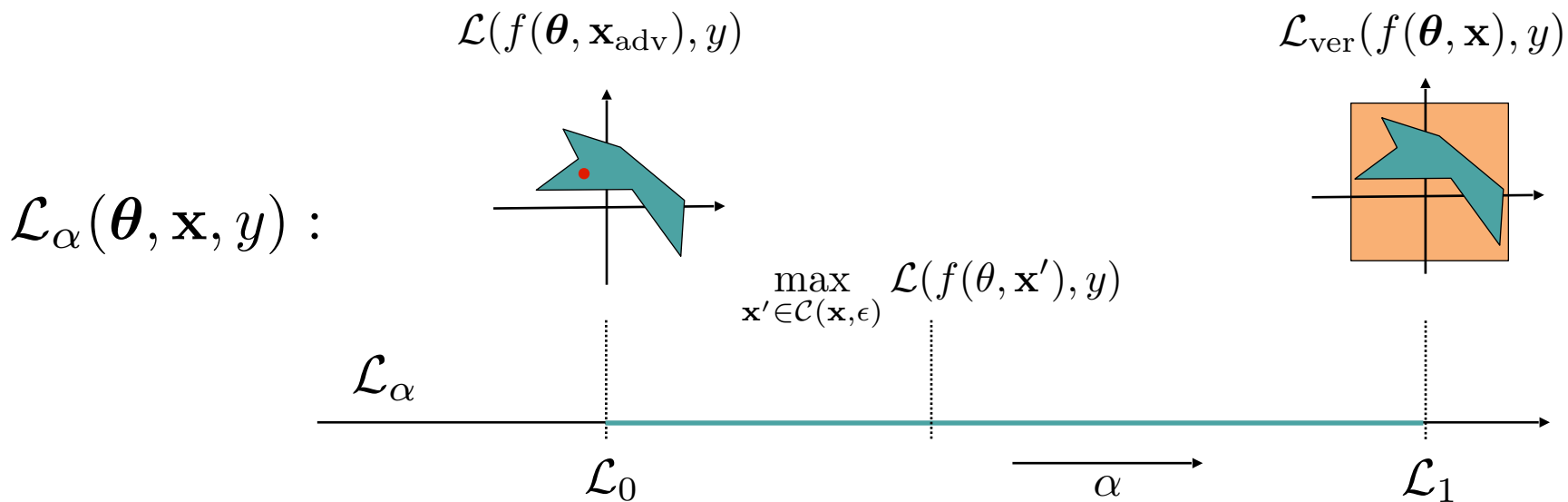
# 💡 Loss Expressivity

$\mathcal{L}_\alpha(\theta, \mathbf{x}, y) :$



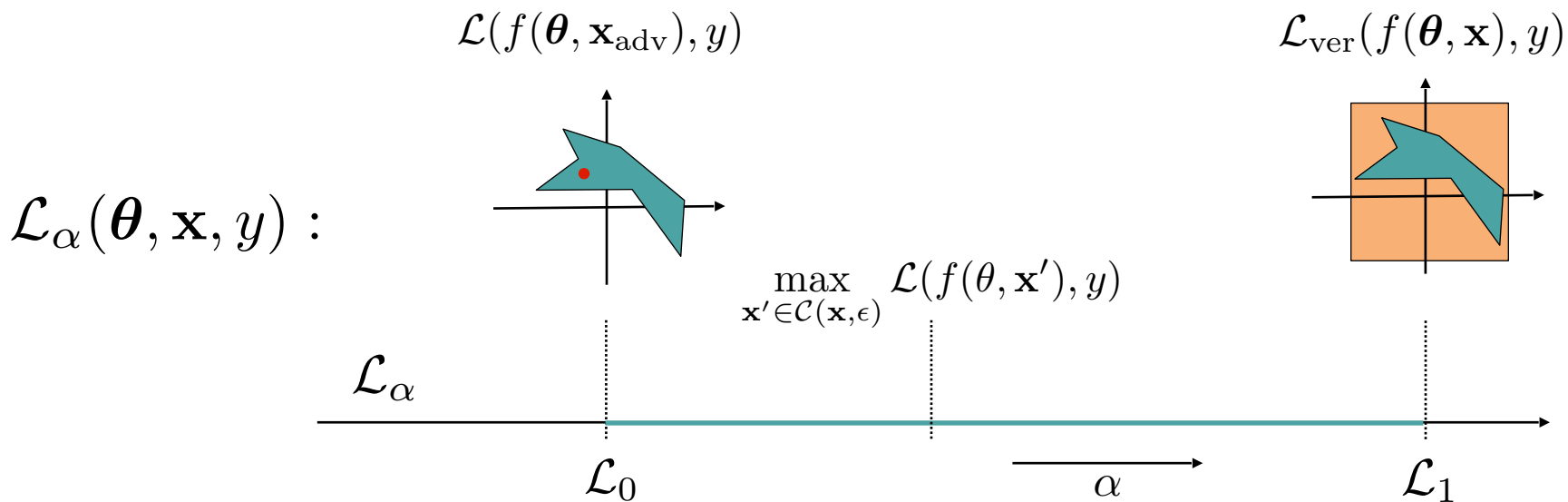


# 💡 Loss Expressivity



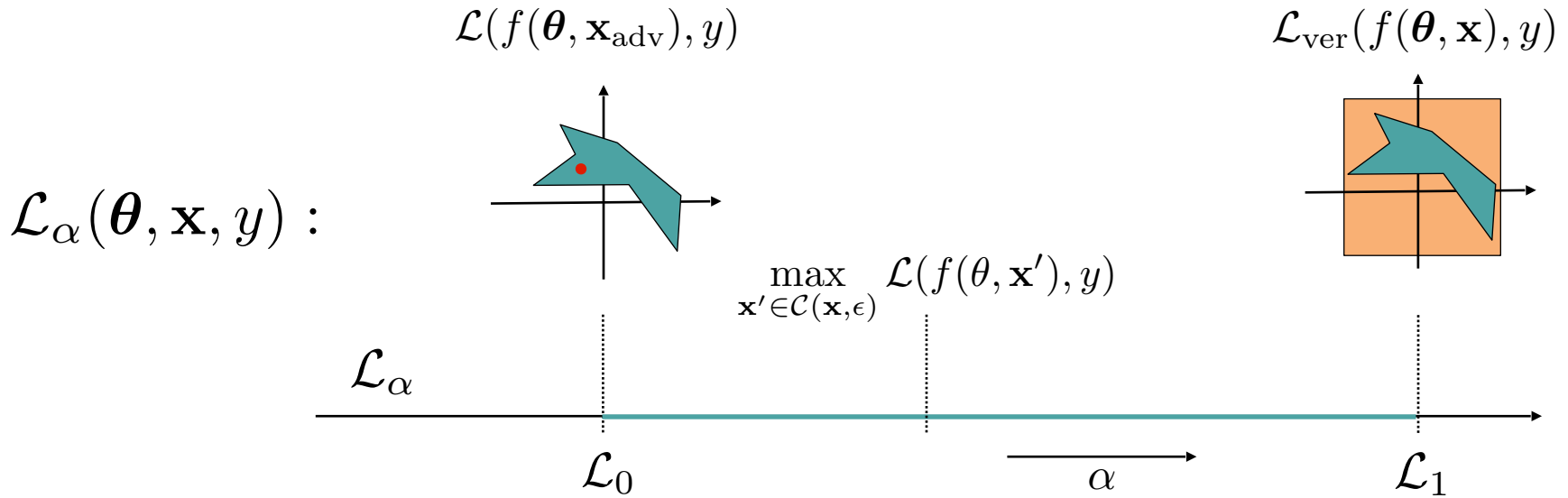
- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \forall \alpha \in [0, 1];$

# 💡 Loss Expressivity



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

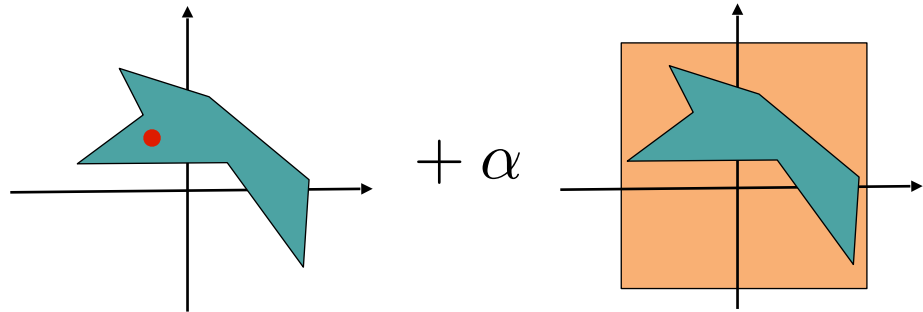
# 💡 Loss Expressivity



- $\mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y) \leq \mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y) \leq \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y) \quad \forall \alpha \in [0, 1];$
- $\mathcal{L}_\alpha(\boldsymbol{\theta}, \mathbf{x}, y)$  continuous and monotonically increasing for  $\alpha \in [0, 1];$
- $\mathcal{L}_0(\boldsymbol{\theta}, \mathbf{x}, y) = \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y);$       •  $\mathcal{L}_1(\boldsymbol{\theta}, \mathbf{x}, y) = \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y).$

# Expressivity via Convex Combinations

## CC-IBP

$$\mathcal{L}(- [(1 - \alpha) \text{ (teal shape with red dot)} + \alpha \text{ (teal shape inside orange square)}], y)$$


# Expressivity via Convex Combinations

## CC-IBP

$$\mathcal{L}\left(- \left[ (1 - \alpha) \begin{array}{c} \uparrow \\ \text{Teal polygon with red dot} \\ \downarrow \end{array} + \alpha \begin{array}{c} \uparrow \\ \text{Teal polygon inside orange square} \\ \downarrow \end{array} \right], y\right)$$

## MTL-IBP

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

# Expressivity via Convex Combinations

## CC-IBP

$$\mathcal{L}\left(- \left[ (1 - \alpha) \begin{array}{c} \uparrow \\ \text{Teal polygon with red dot} \\ \downarrow \end{array} + \alpha \begin{array}{c} \uparrow \\ \text{Teal polygon inside orange square} \\ \downarrow \end{array} \right], y\right)$$

## MTL-IBP

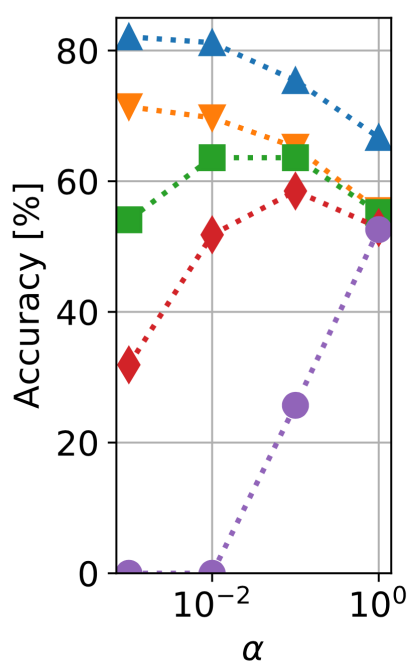
$$(1 - \alpha) \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y) + \alpha \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)$$

## Exp-IBP

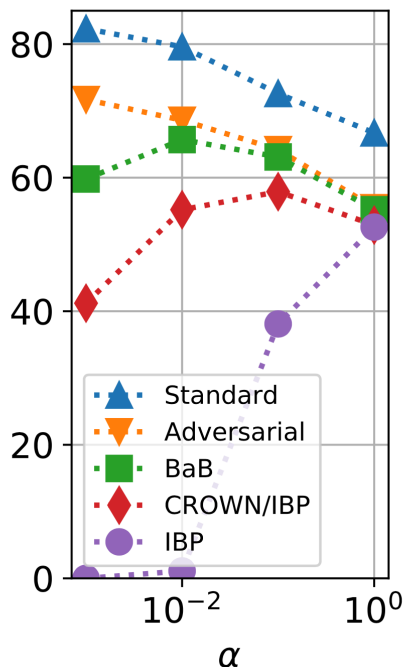
$$\mathcal{L}_{\alpha, \text{Exp}}(\boldsymbol{\theta}, \mathbf{x}, y) := \mathcal{L}(f(\boldsymbol{\theta}, \mathbf{x}_{\text{adv}}), y)^{(1-\alpha)} \mathcal{L}_{\text{ver}}(f(\boldsymbol{\theta}, \mathbf{x}), y)^\alpha$$

# Loss Sensitivity

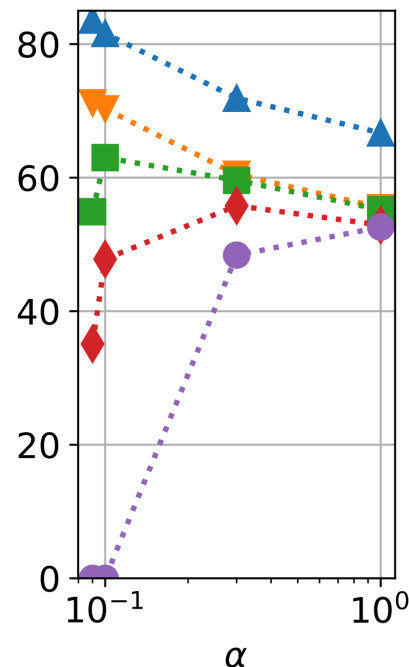
Sensitivity of CC-IBP, MTL-IBP and Exp-IBP to  $\alpha$  for robustness to  $\ell_\infty$  perturbations on CIFAR-10.



(a) CC-IBP,  
 $\epsilon = 2/255$ .



(b) MTL-IBP,  
 $\epsilon = 2/255$ .



(c) Exp-IBP,  
 $\epsilon = 2/255$ .

# 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. [%]
TinyImageNet	$\frac{1}{255}$	CC-IBP	38.61	26.39
		MTL-IBP	37.56	26.09
		EXP-IBP	38.71	26.18
		STAPS	28.98	22.16
		SABR	28.97	21.36
		SORTNET	25.69	18.18
		IBP	25.40	19.92
		CROWN-IBP	25.62	17.93
		ImageNet64	$\frac{1}{255}$	CC-IBP
MTL-IBP	20.15			12.13
EXP-IBP	22.73			13.30
SORTNET	14.79			9.54
CROWN-IBP	16.23			8.73
IBP	15.96			6.13



# Conclusions

- Expressivity  $\rightarrow$  state-of-the-art certified training;
- Expressivity easily obtained via convex combinations;
- Verified accuracy still comes at great cost in standard performance.

## Code and models

<https://github.com/alessandrodepalma/expressive-losses>



[alessandro.de-palma@inria.fr](mailto:alessandro.de-palma@inria.fr)