

Certified Training: Small Boxes are All You Need



Mark Müller*



Franziska Eckert*



Marc Fischer



Martin Vechev

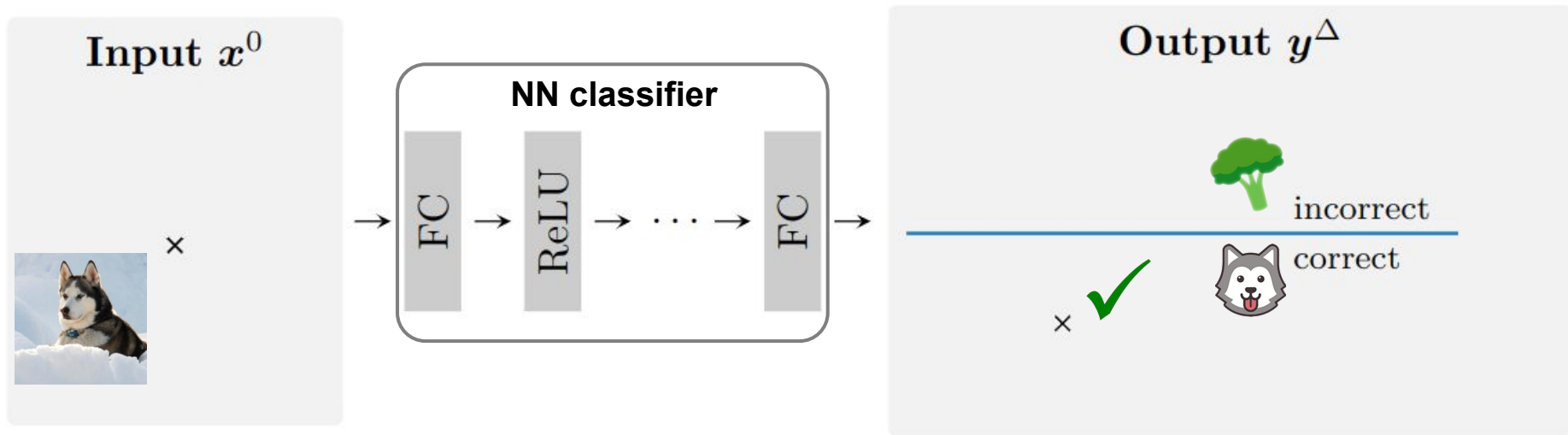
ETH zürich



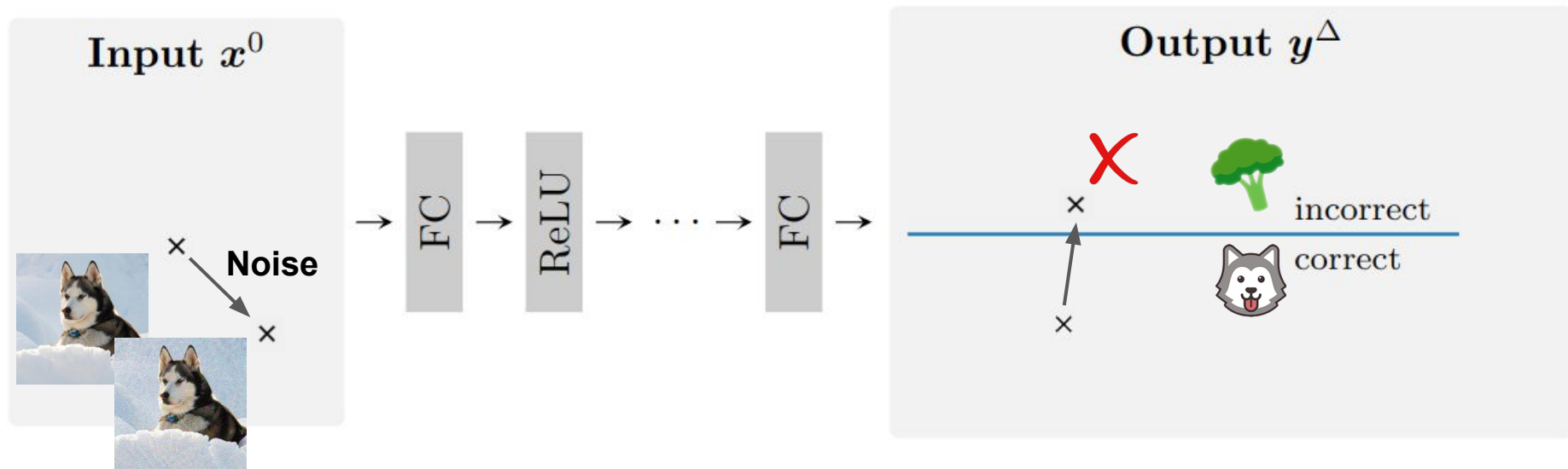
ICLR

 **SRILAB**

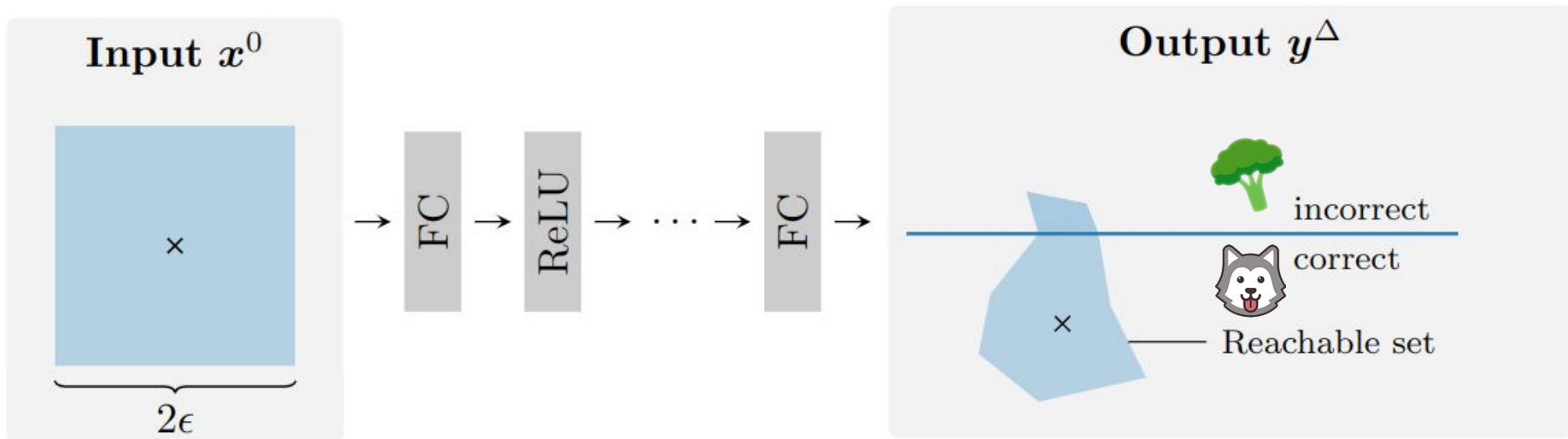
Standard Classification



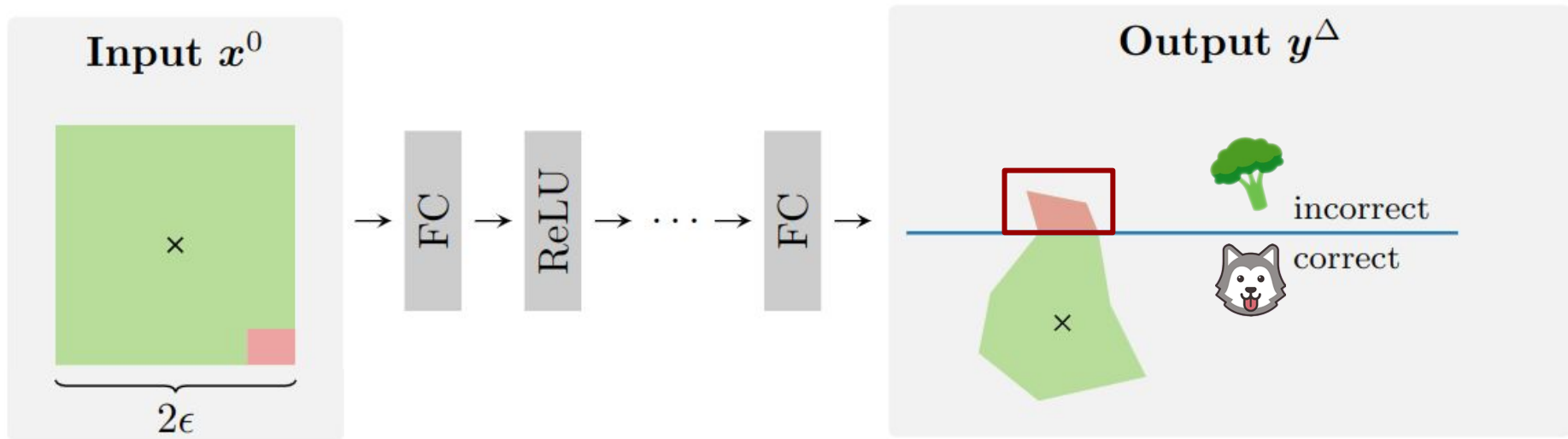
Adversarial Examples



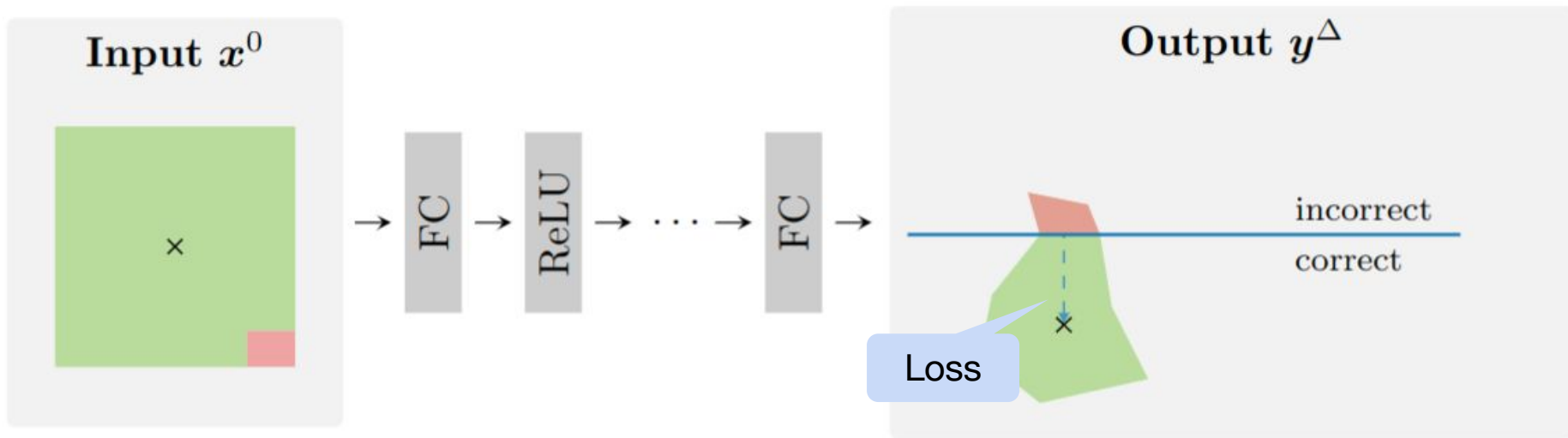
Exact Propagation



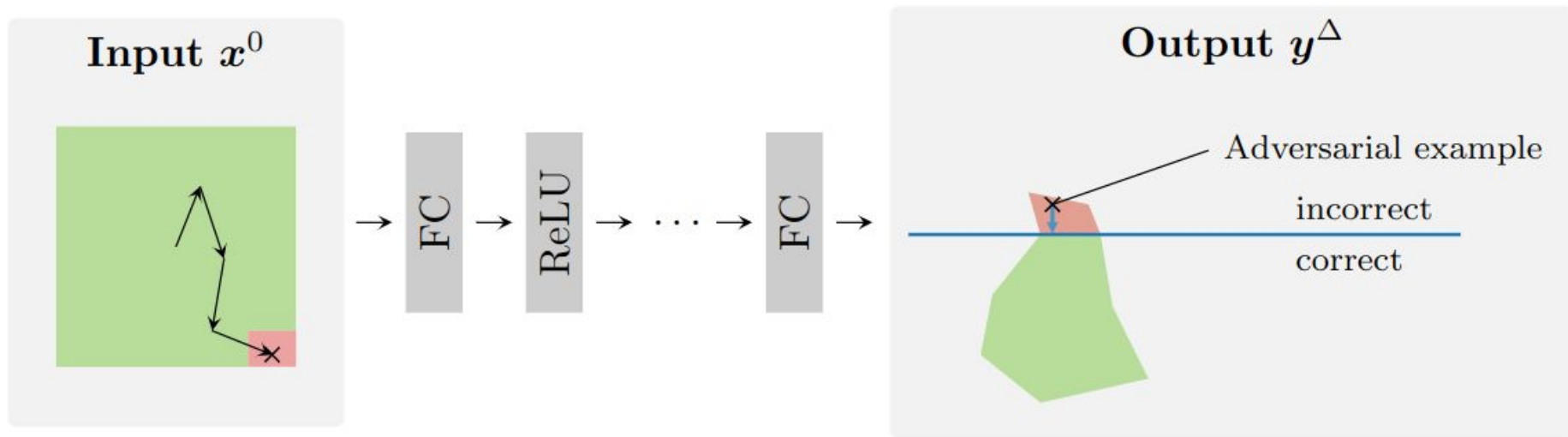
Exact Propagation



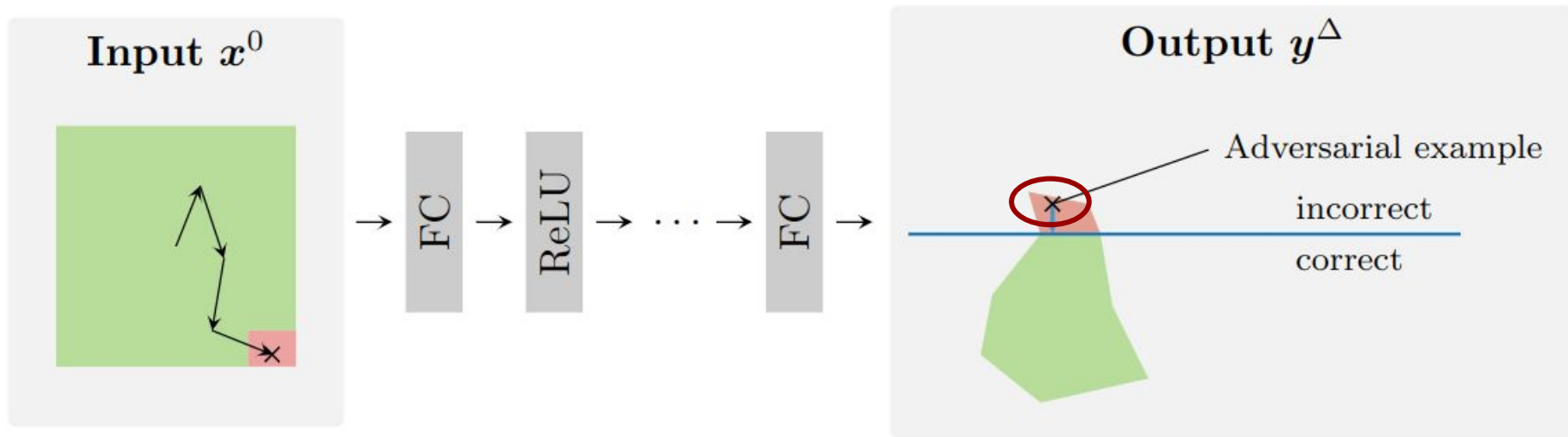
Standard Training



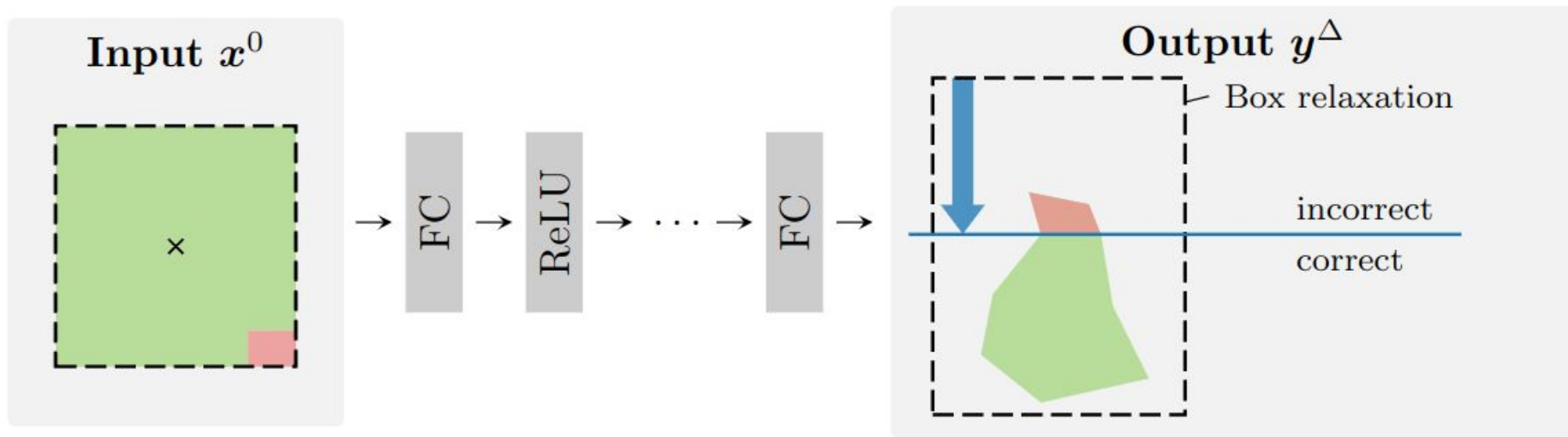
Adversarial Training (PGD)



Adversarial Training (PGD)



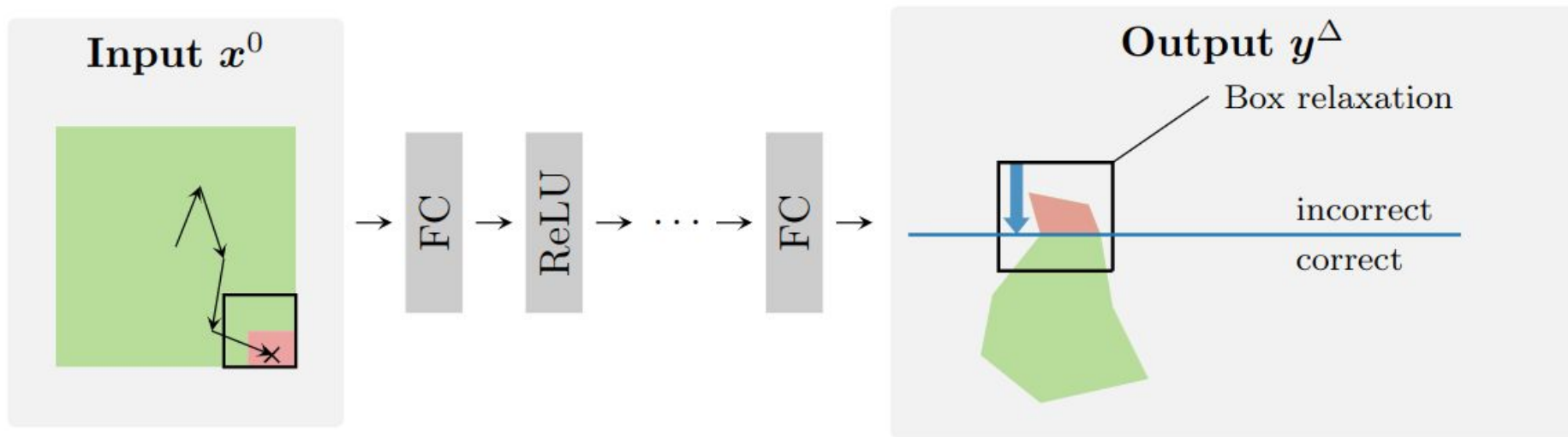
Certified Training (IBP)



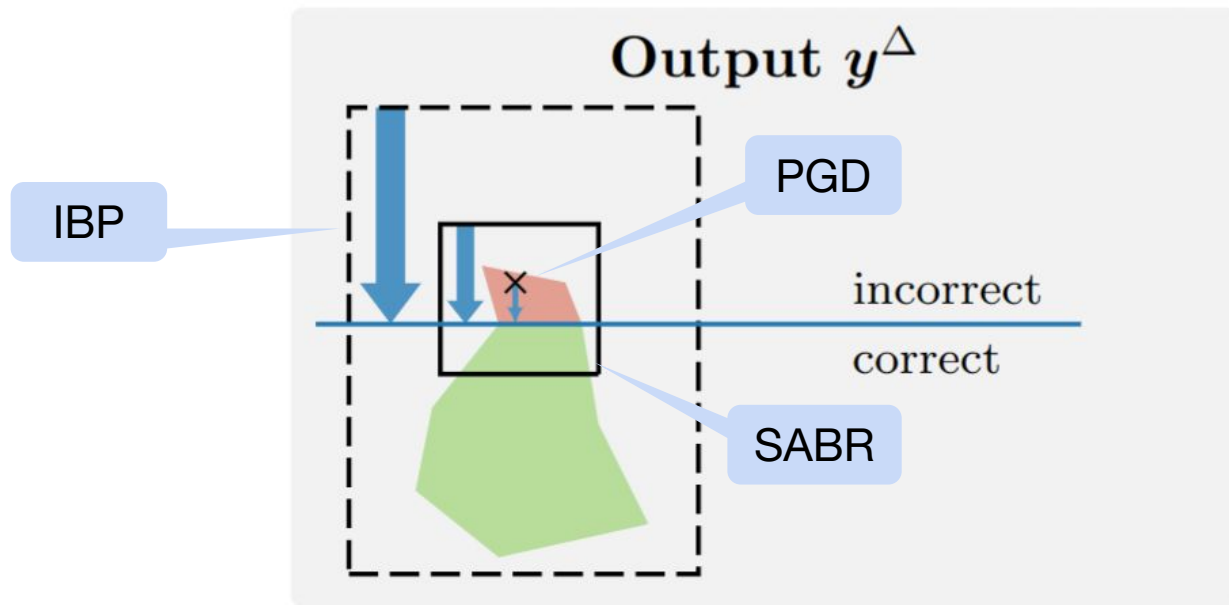
Gowal, et al. "On the effectiveness of interval bound propagation for training verifiably robust models." arXiv 2018

Mirmann et al. "Differentiable abstract interpretation for provably robust neural networks." ICML 2018

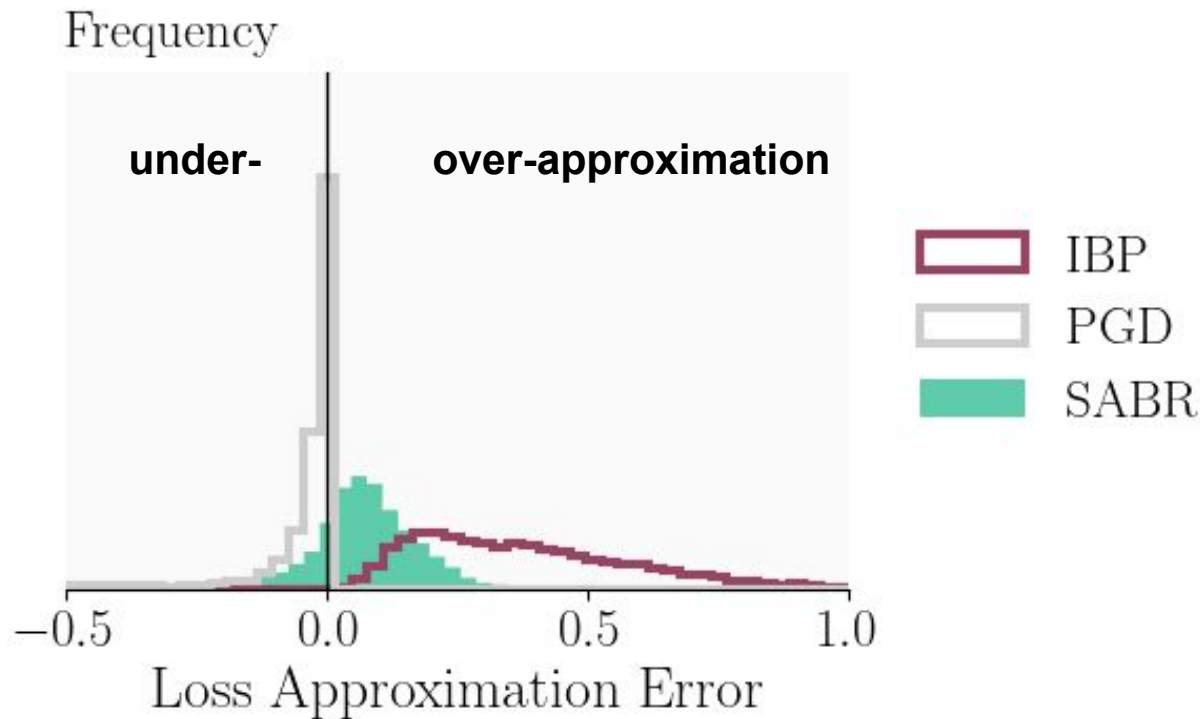
SABR – This Work



Regularisation Comparison



Worst-Case Loss Approximation Precision



Box Abstraction Size Growth

Growth rate $\kappa = \frac{\mathbb{E}[\text{Output Box Size}]}{\text{Input Box Size}} = \frac{\text{Diagram of Output Box}}{\text{Diagram of Input Box}}$

The diagram shows a fraction where the numerator is a large gray square with a bracket to its right labeled δ_{out} , and the denominator is a small blue square with a bracket to its right labeled δ_{in} .

Box Abstraction Size Growth

Growth rate $\kappa = \frac{\mathbb{E}[\text{Output Box Size}]}{\text{Input Box Size}} = \frac{\text{Diagram: A large gray square with a bracket to its right labeled } \delta_{\text{out}}}{\text{Diagram: A small blue square with a bracket to its right labeled } \delta_{\text{in}}}$

Linear layers: κ is independent of input box scale:

$$\kappa \sim [10, 100]$$

Diagram: A large gray square = W Diagram: A small blue square + b

Box Abstraction Size Growth

Growth rate $\kappa = \frac{\mathbb{E}[\text{Output Box Size}]}{\text{Input Box Size}} = \frac{\text{gray box} \delta_{\text{out}}}{\text{blue box} \delta_{\text{in}}}$

Linear layers: κ is independent of input box scale:

$$\kappa \sim [10, 100]$$

$$\text{gray box} = W \text{ blue box} + b$$

ReLU layers: κ depends on box scale and box centre:

$$\kappa \sim [0, 1]$$

$$\text{gray box} = \max(\text{blue box}, 0)$$

Box Abstraction Size Growth

Growth rate $\kappa = \frac{\mathbb{E}[\text{Output Box Size}]}{\text{Input Box Size}} = \frac{\text{Diagram: A large gray square with a bracket to its right labeled } \delta_{\text{out}}}{\text{Diagram: A small blue square with a bracket to its right labeled } \delta_{\text{in}}}$

Linear layers: κ is independent of input box scale:

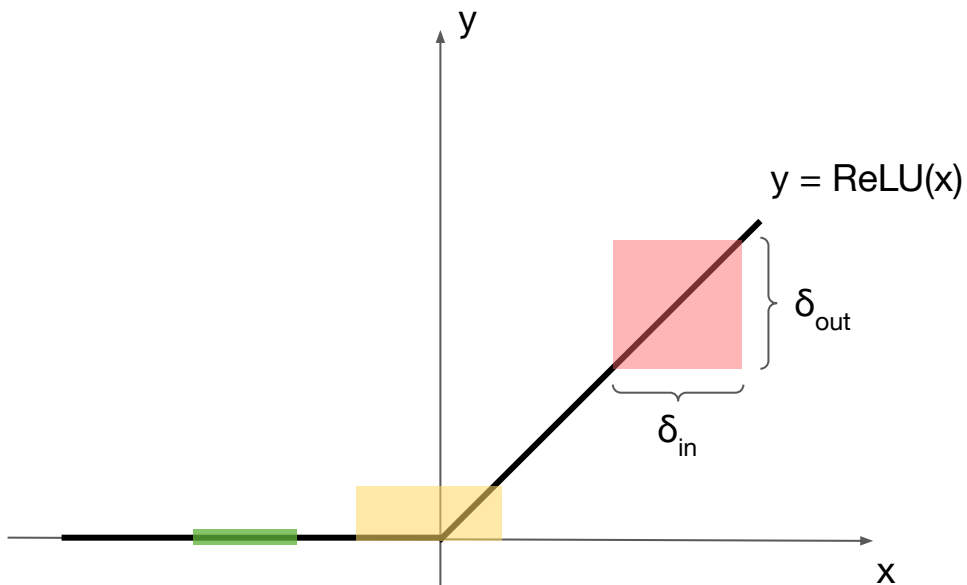
$$\kappa \sim [10, 100]$$

Diagram: A large gray square = W Diagram: A small blue square + b

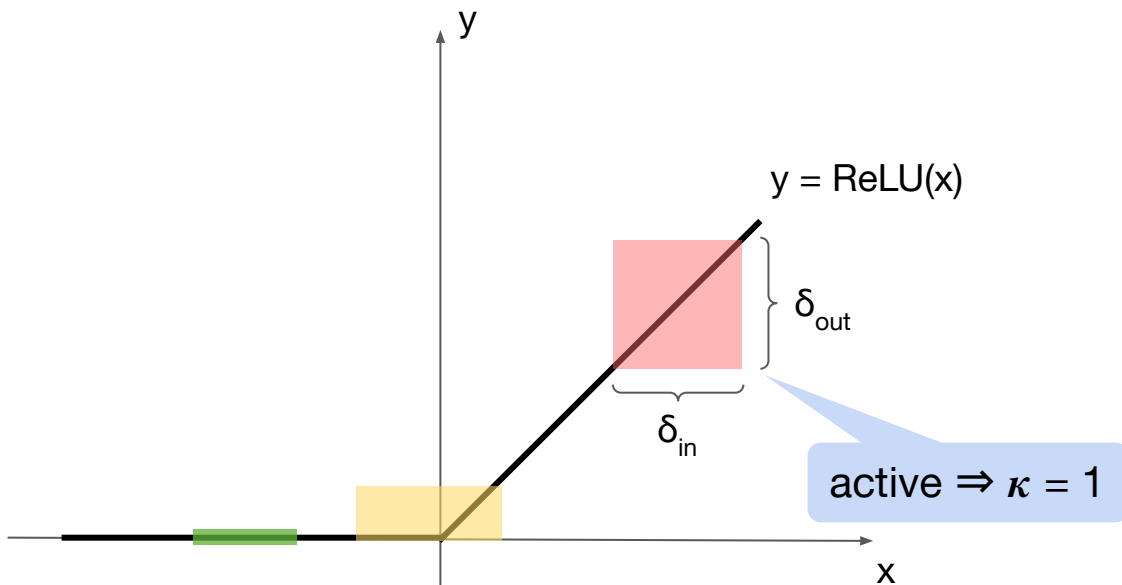
ReLU layers: κ depends on box scale and box centre: Diagram: A large gray square = $\max(\text{Diagram: A small blue square}, 0)$

$$\kappa \sim [0, 1]$$

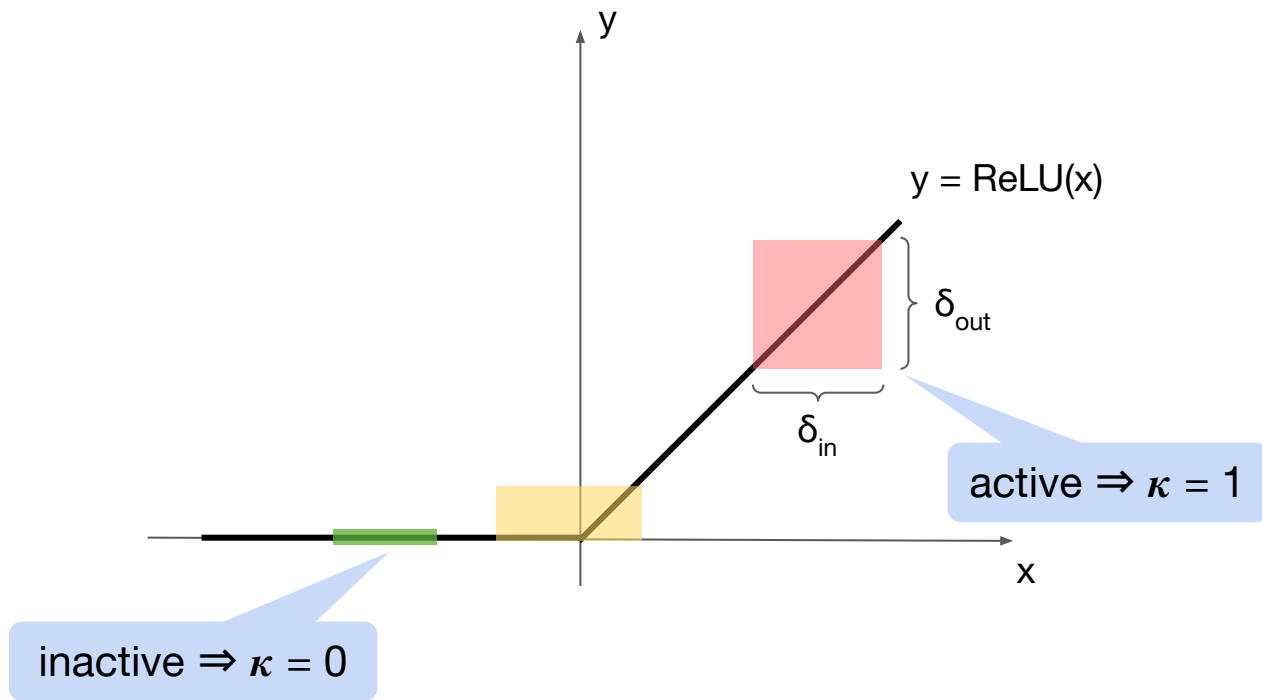
Box Abstraction Size Growth – ReLUs



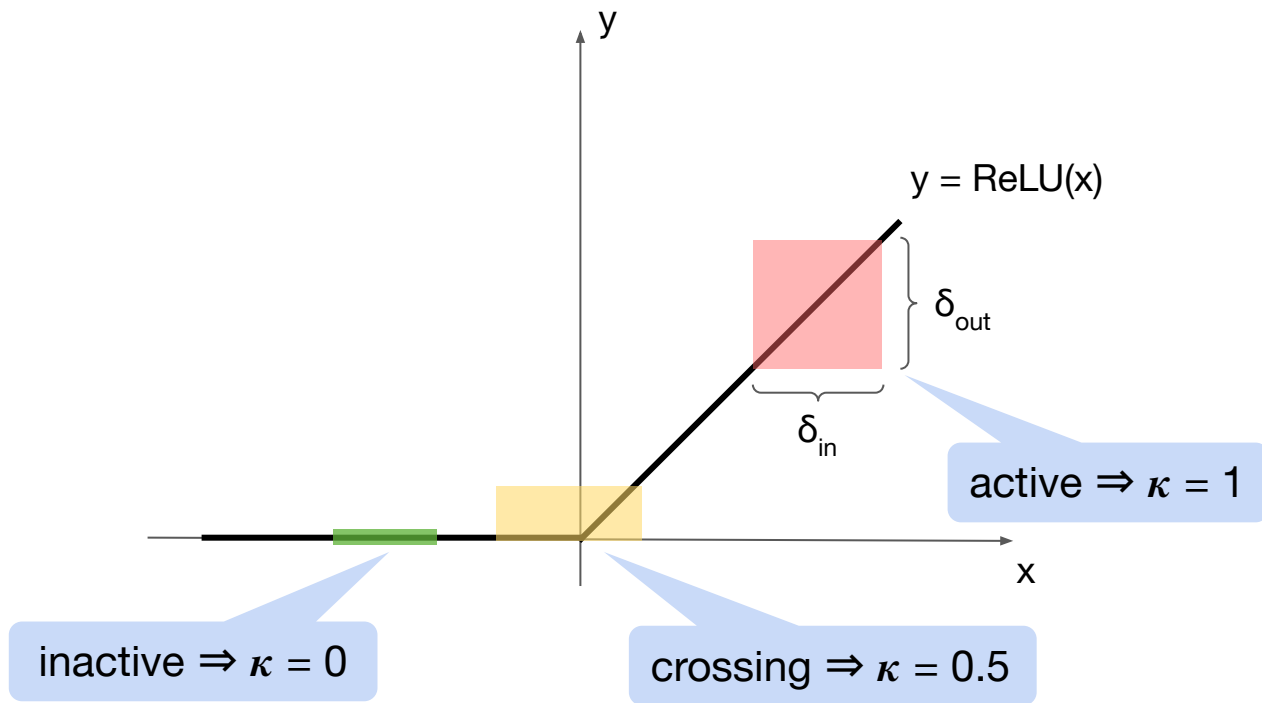
Box Abstraction Size Growth – ReLUs



Box Abstraction Size Growth – ReLUs



Box Abstraction Size Growth – ReLUs



Box Abstraction Size Growth – ReLUs

For input box sizes $\varepsilon \rightarrow 0$

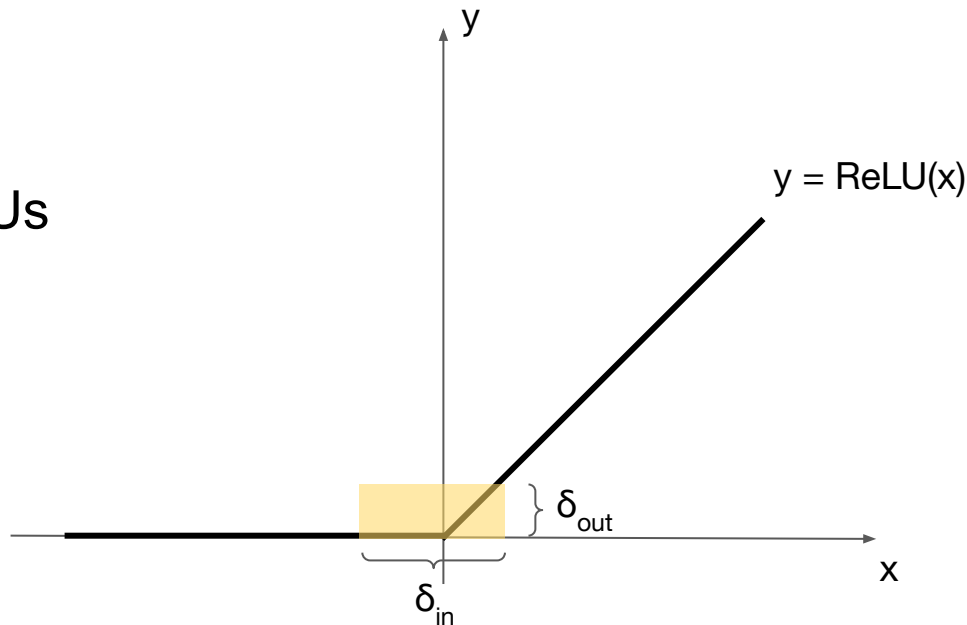
κ = Portion of active ReLUs

For input box sizes $\varepsilon \rightarrow \infty$

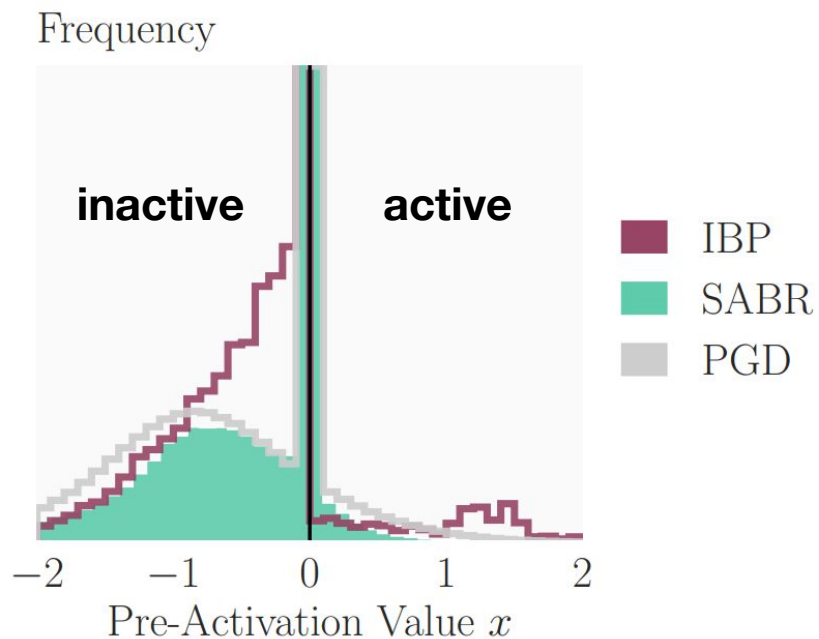
$\kappa = 0.5$

In-between:

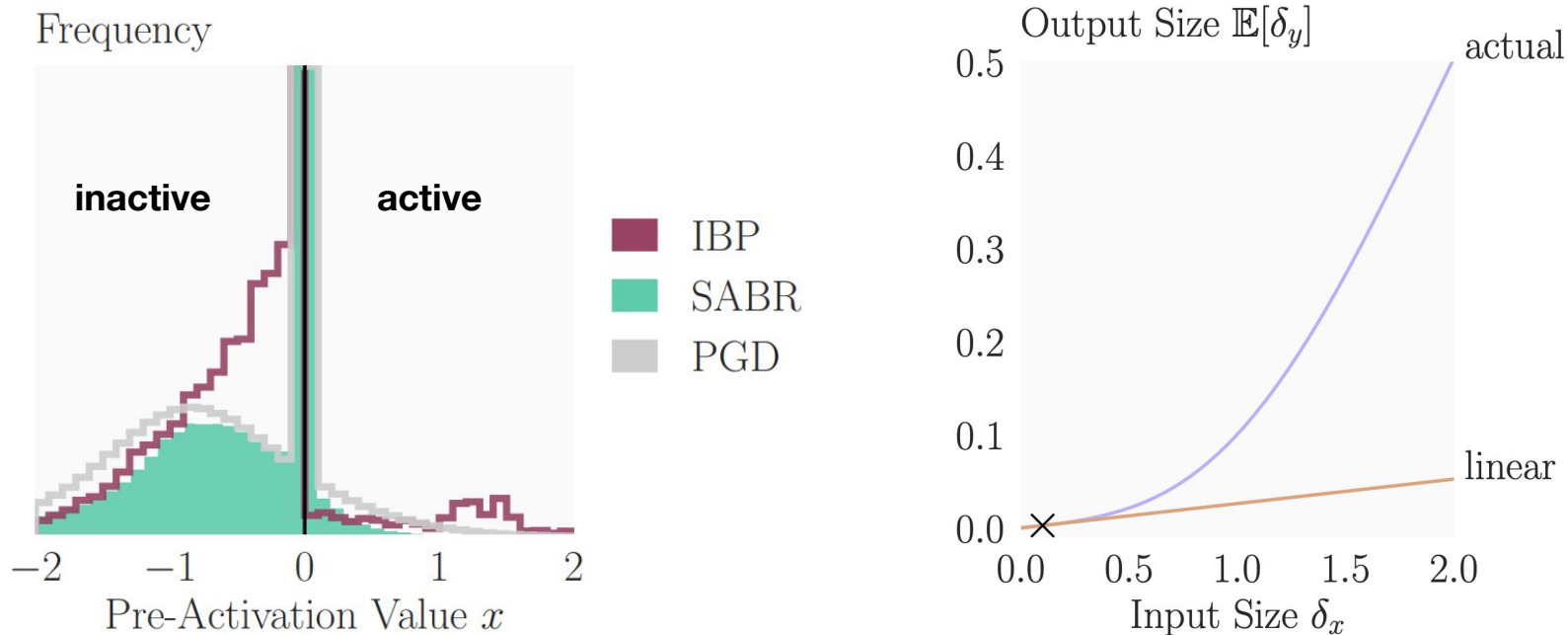
κ depends on box positions



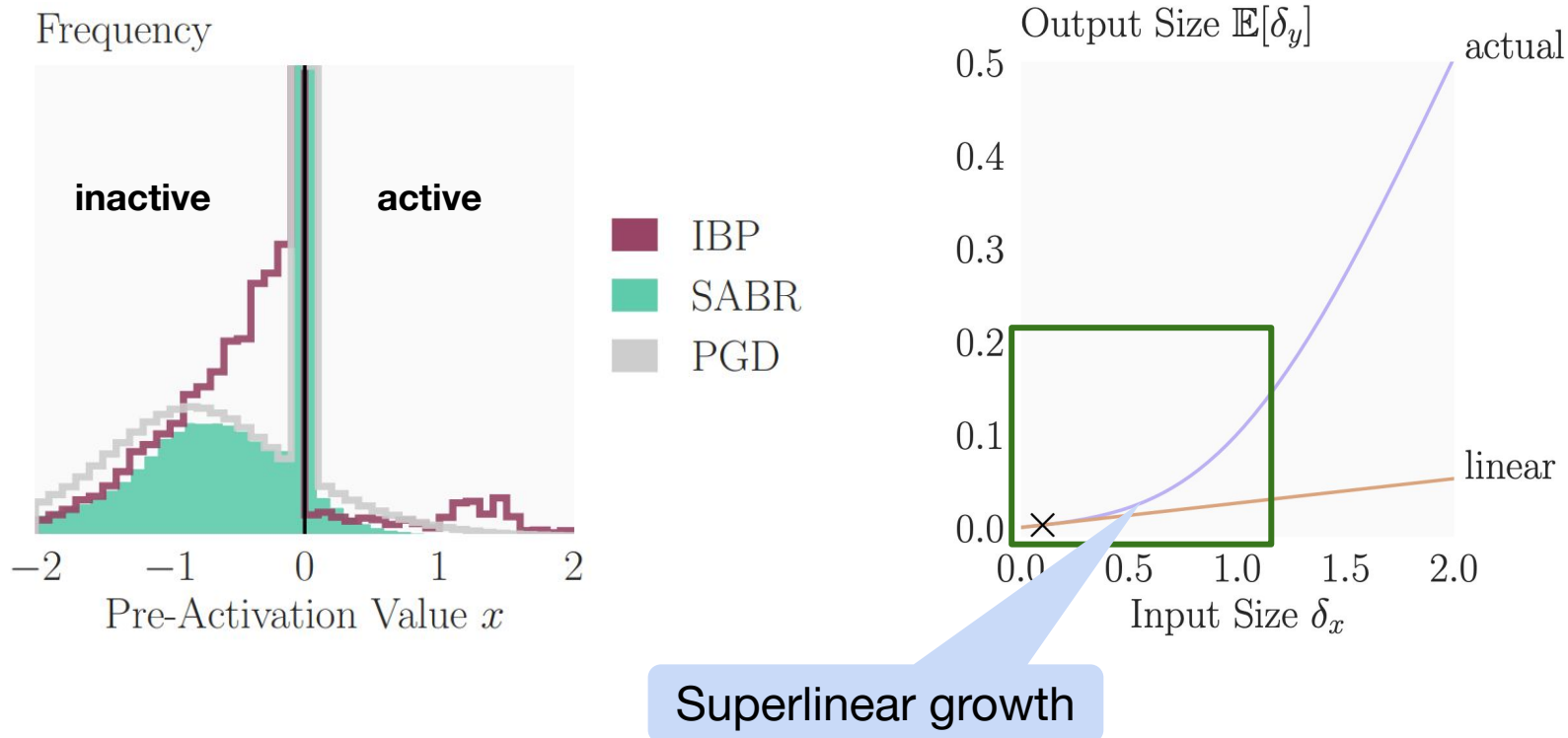
Box Abstraction Size Growth – ReLUs



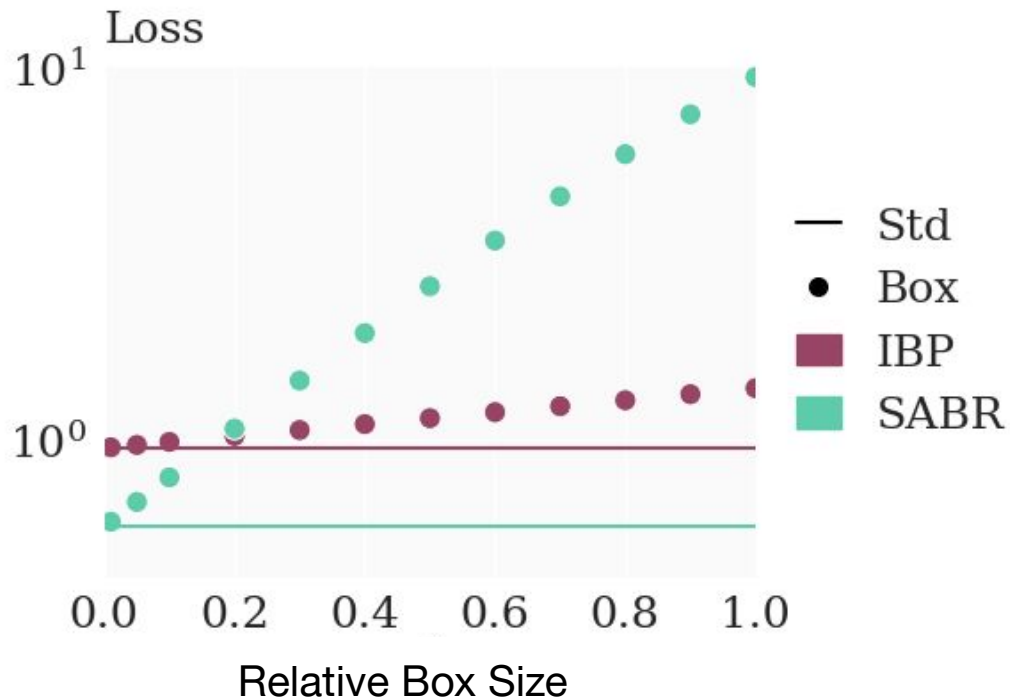
Box Abstraction Size Growth – ReLUs



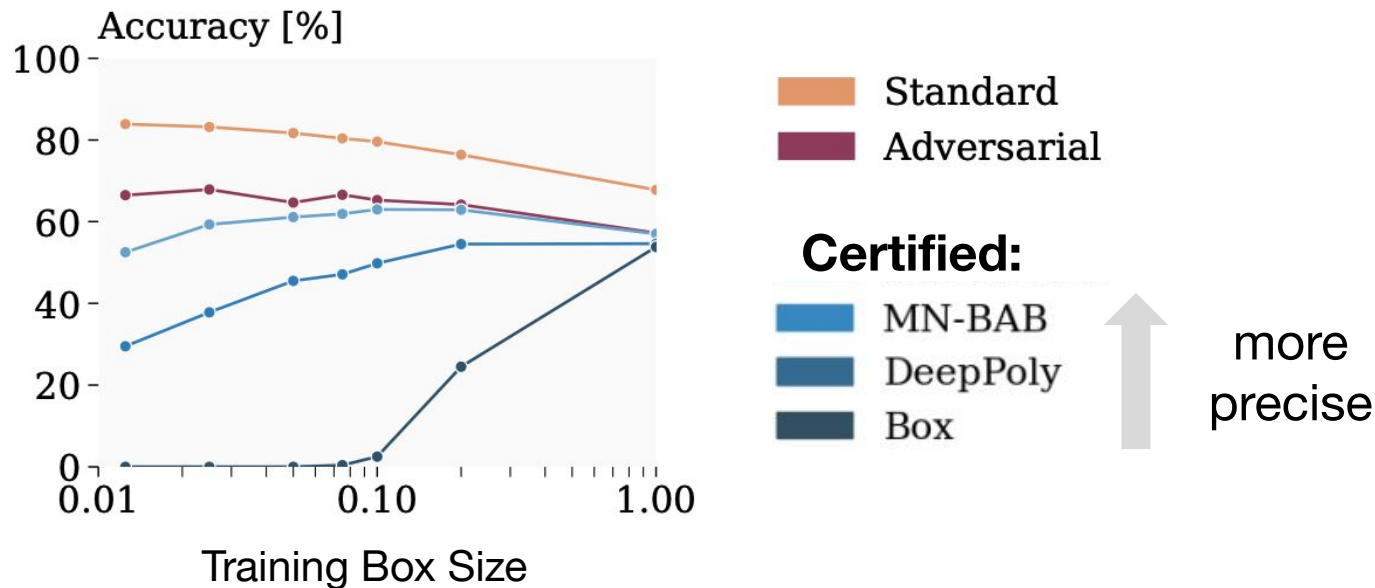
Box Abstraction Size Growth – ReLUs



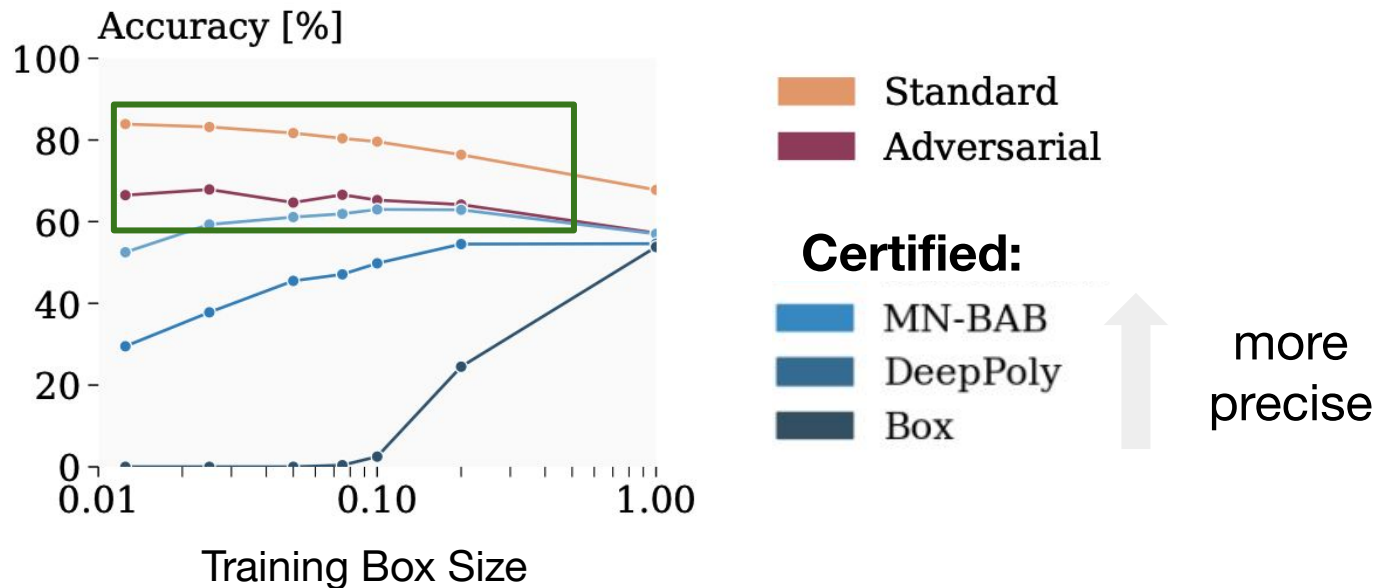
Full Network Loss Growth



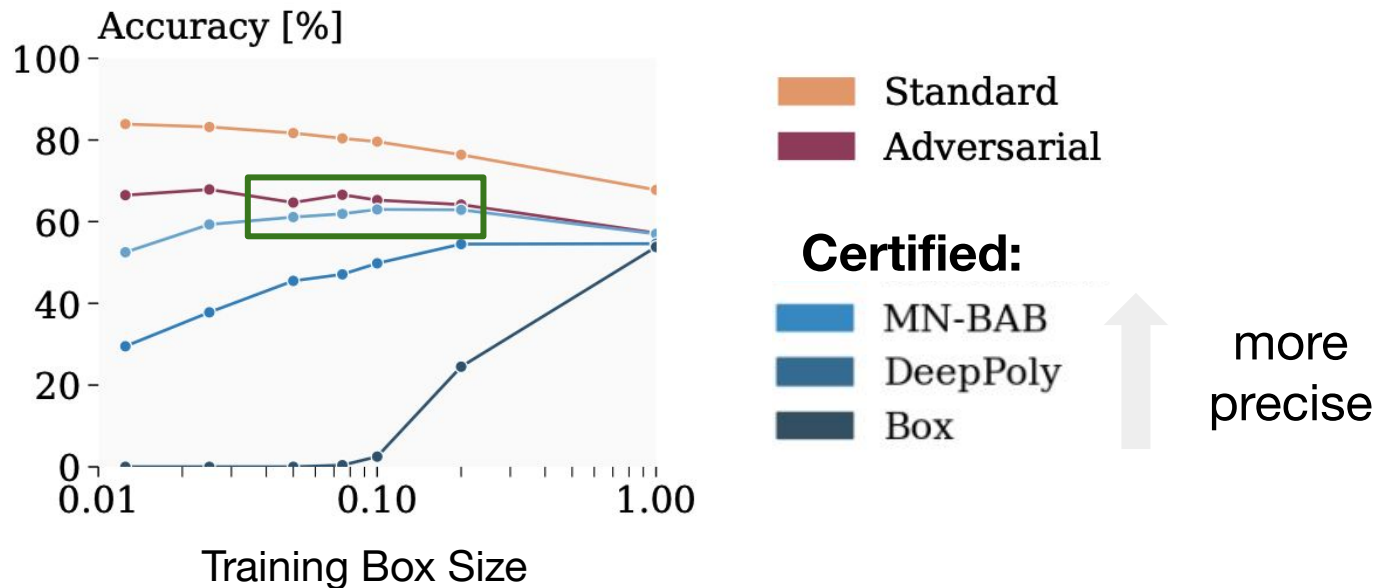
Impact of Verification Method



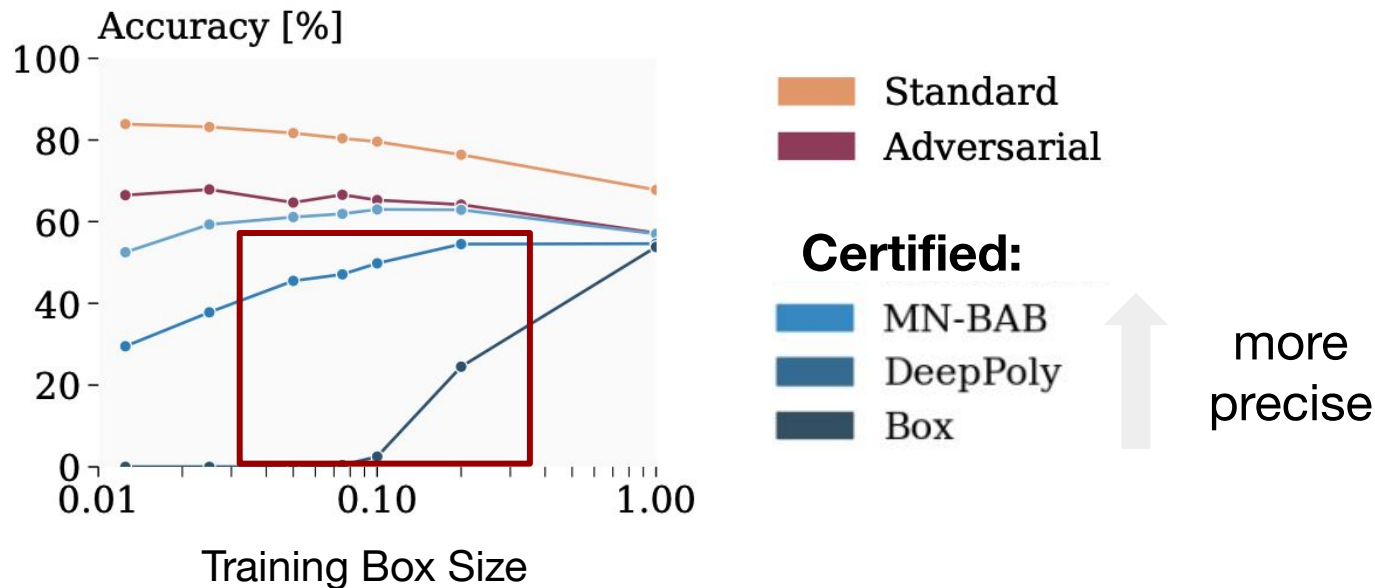
Impact of Verification Method



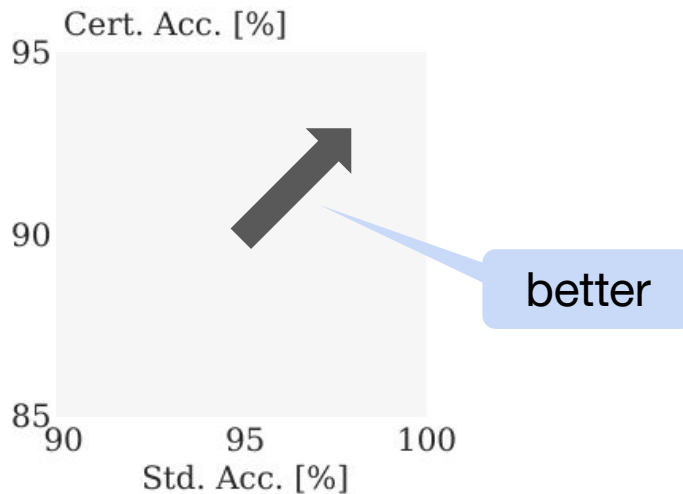
Impact of Verification Method



Impact of Verification Method

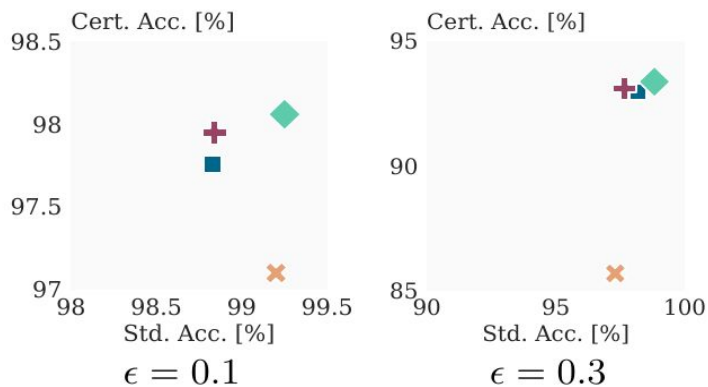


Empirical Results

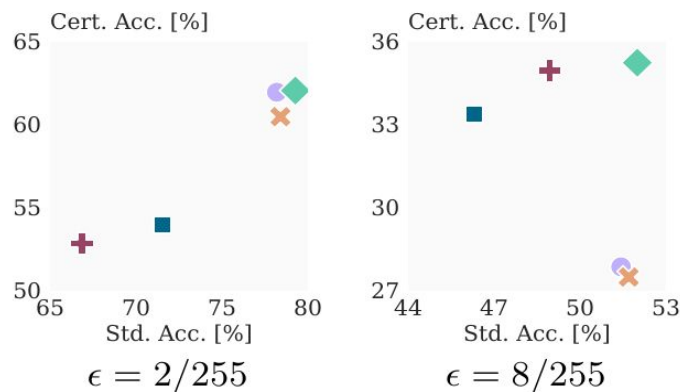


Empirical Results

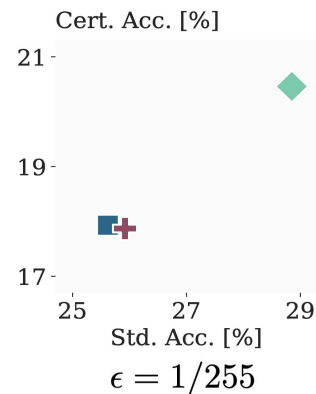
MNIST



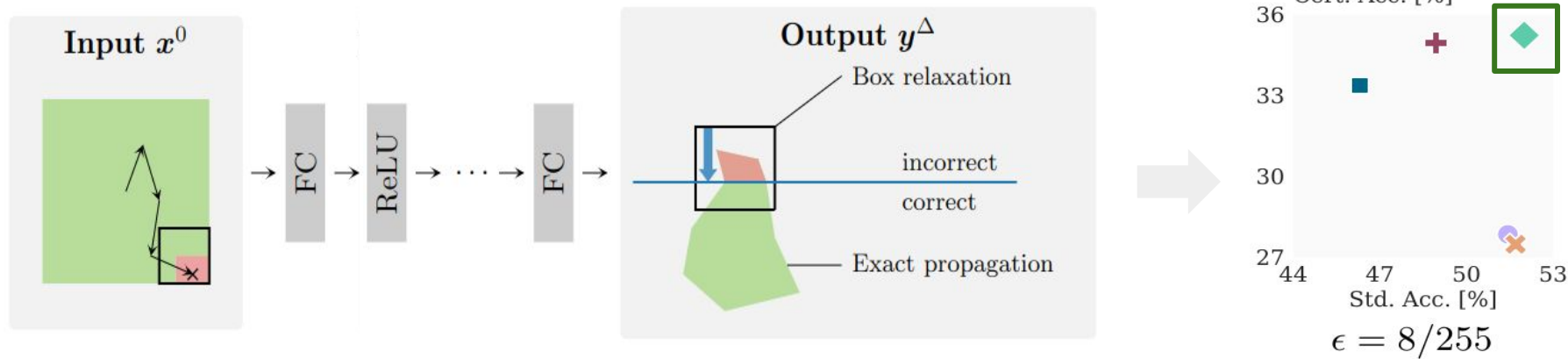
CIFAR-10



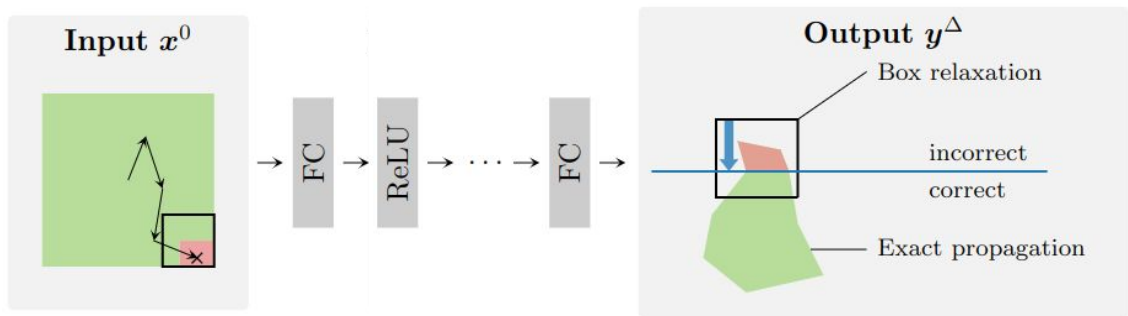
TinyImageNet



Conclusion



Thank You For Your Attention!



Paper & Code:

<https://www.sri.inf.ethz.ch/publications/mueller2022sabr>

<https://github.com/eth-sri/SABR>