

# WASSERSTEIN DISTRIBUTIONAL ROBUSTNESS OF NEURAL NETWORKS

Xingjian Bai, Guangyi He, Yifan Jiang, Jan Obłój

xingjian.bai@sjc.ox.ac.uk, guangyihe2002@outlook.com, yifan.jiang@maths.ox.ac.uk, jan.obloj@maths.ox.ac.uk



## MAIN CONTRIBUTION

1. A unified approach to adversarial attacks and training based on sensitivity analysis for Wasserstein DRO from Bartl et al. (2021).
2. A fast-to-compute first order adversarial attack method for *distributional* threat models. As a special case, this recovers the classical FGSM attack, lending it further theoretical underpinning.
3. Asymptotically certified bounds and out-of-sample performance of adversarial accuracy, applicable to a general threat, including classical pointwise perturbations.

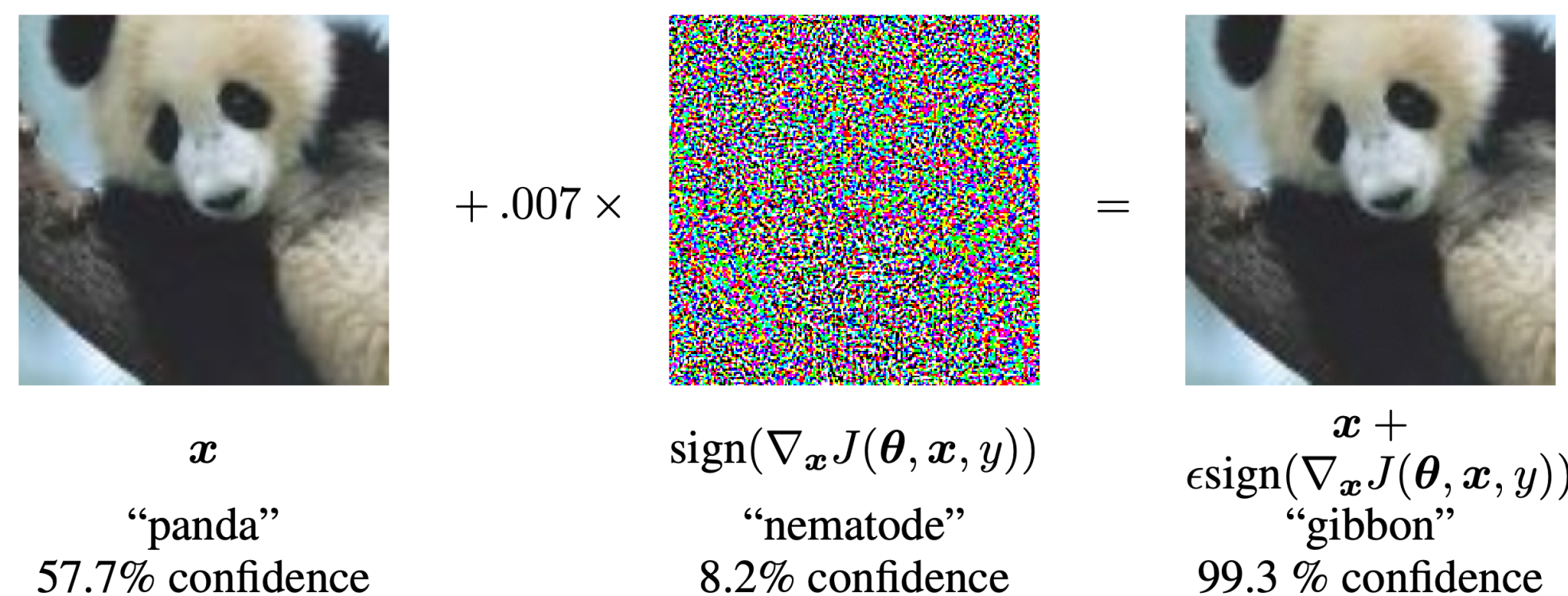
## ADVERSARIAL ATTACK OF NNS

Deep neural networks have achieved a great success in image recognition tasks. However, they are notoriously known to be vulnerable under adversarial attacks, which are small perturbations of the input data that lead to misclassification.

To generate adversarial images, we reverse the training process by maximizing the loss over input data:

$$E_P \left[ \sup_{\|x-x'\|_\infty \leq \delta} J_\theta(x', y) \right] \rightsquigarrow x^* \approx x + \delta \text{sgn}(\nabla_x J_\theta(x, y)),$$

where  $J_\theta(x, y) = L(f_\theta(x), y)$ , and  $P$  is the input distribution.



A demonstration of Fast Gradient Sign Method (FGSM).

## DISTRIBUTIONAL THREAT MODELS

Let  $d((x, y), (x', y')) = \|x - x'\|_s + \infty \mathbb{1}\{y \neq y'\}$  and  $\mathcal{W}_p$  be the Wasserstein distance given by

$$\mathcal{W}_p(P, Q) = \inf \left\{ E[d(X, Y)^p]^{1/p} : X \sim P, Y \sim Q \right\}.$$

We propose a novel *distributional* threat model

$$\sup_{Q: \mathcal{W}_p(P, Q) \leq \delta} E_Q[J_\theta(x, y)].$$

In *distributional* threat models, the attacker has a greater flexibility and can perturb images close to the decision boundary only slightly while spending more of the attack budget on images farther away from the boundary.

The classical *pointwise* threat model can be covered as an extreme case  $p = s = \infty$

$$E_P \left[ \sup_{\|x-x'\|_\infty \leq \delta} J_\theta(x', y) \right] = \sup_{Q: \mathcal{W}_\infty(P, Q) \leq \delta} E_Q[J_\theta(x, y)].$$

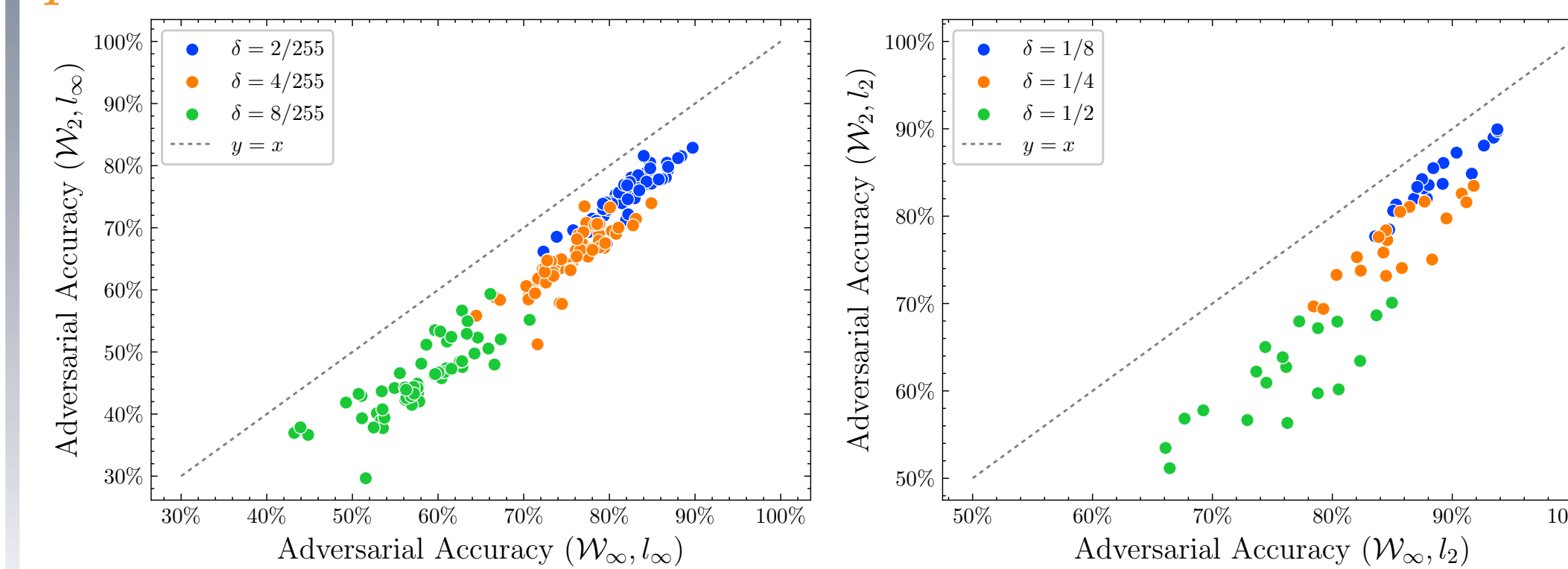
## WASSERSTEIN DISTRIBUTIONAL ADVERSARIAL ATTACKS

Our proposed PGD-type attack is:

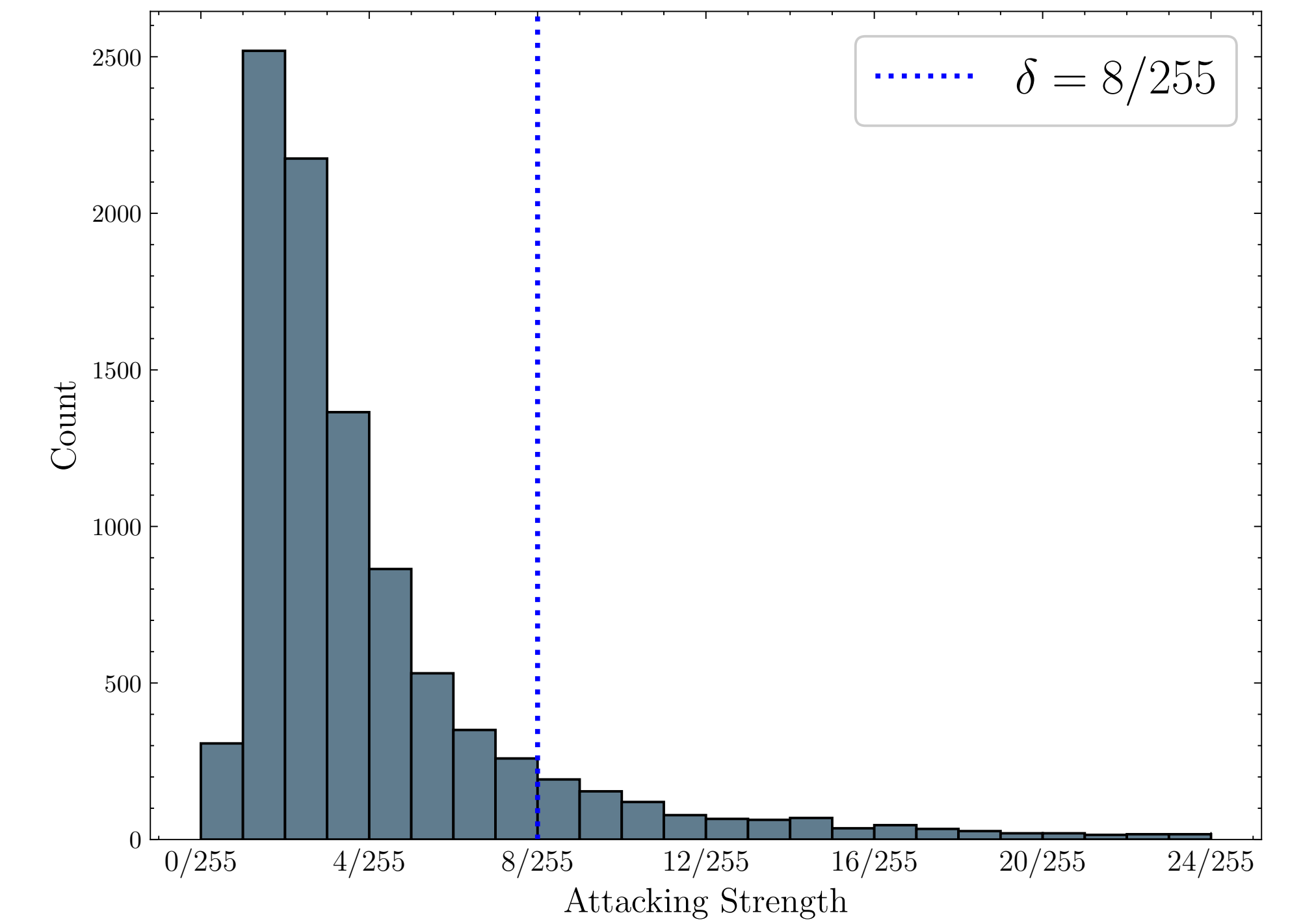
$$x^{n+1} = \text{pj}(x^n + \alpha h(\nabla_x J_\theta(x^n, y)) \|\Upsilon^{-1} \nabla_x J_\theta(x^n, y)\|_*^{q-1}),$$

where  $\alpha$  is the step size and  $\text{pj}$  is a projection.

The classical FGSM attack is retrieved by taking  $p = s = \infty$ .



Shortfall of WD-adversarial accuracy on CIFAR-10.



Histogram of *distributional* attacking strength on CIFAR-10.

## CERTIFIED ADVERSARIAL ACCURACY BOUNDS

- $S$  the set of images equipped with their labels generated by  $f_\theta$ , i.e.,

$$S = \left\{ (x, y) : \arg \max_{1 \leq i \leq m} f_\theta(x)_i = \{y\} \right\}.$$

- $A$  the clean accuracy given by  $A = E_P[\mathbb{1}_S]$ .
- $A_\delta$  the adversarial accuracy given by  $A_\delta = \inf_{Q: \mathcal{W}_p(P, Q) \leq \delta} E_Q[\mathbb{1}_S]$ .
- $W_\delta$  the loss condition on the misclassified image, given by

$$W(\delta) = \sup_{Q \in \mathcal{B}_\delta(P)} E_Q[J_\theta(x, y) | S^c].$$

- Upper bound:  $\mathcal{R}_\delta \leq \mathcal{R}_\delta^u := Q_\delta(S)/A$ .

- An asymptotic lower bound:

$$\mathcal{R}_\delta \geq \frac{W(0) - V(\delta)}{W(0) - V(0)} + o(\delta) = \tilde{\mathcal{R}}_\delta^l + o(\delta) = \bar{\mathcal{R}}_\delta^l + o(\delta),$$

where we utilize the first order approximations of  $V(\delta)$  and write  $\tilde{\mathcal{R}}_\delta^l = \frac{W(0) - E_{Q_\delta}[J_\theta(x, y)]}{W(0) - V(0)}$  and  $\bar{\mathcal{R}}_\delta^l = \frac{W(0) - V(0) - \delta \Upsilon}{W(0) - V(0)}$ .

## WASSERSTEIN DRO SENSITIVITY

We write  $V(\delta) = \sup_{Q: \mathcal{W}_p(P, Q) \leq \delta} E_Q[J_\theta(x, y)]$  and assume  $J_\theta$  is Lipschitz under  $d$ .

The following result follows readily from Bartl et al. and its proof.

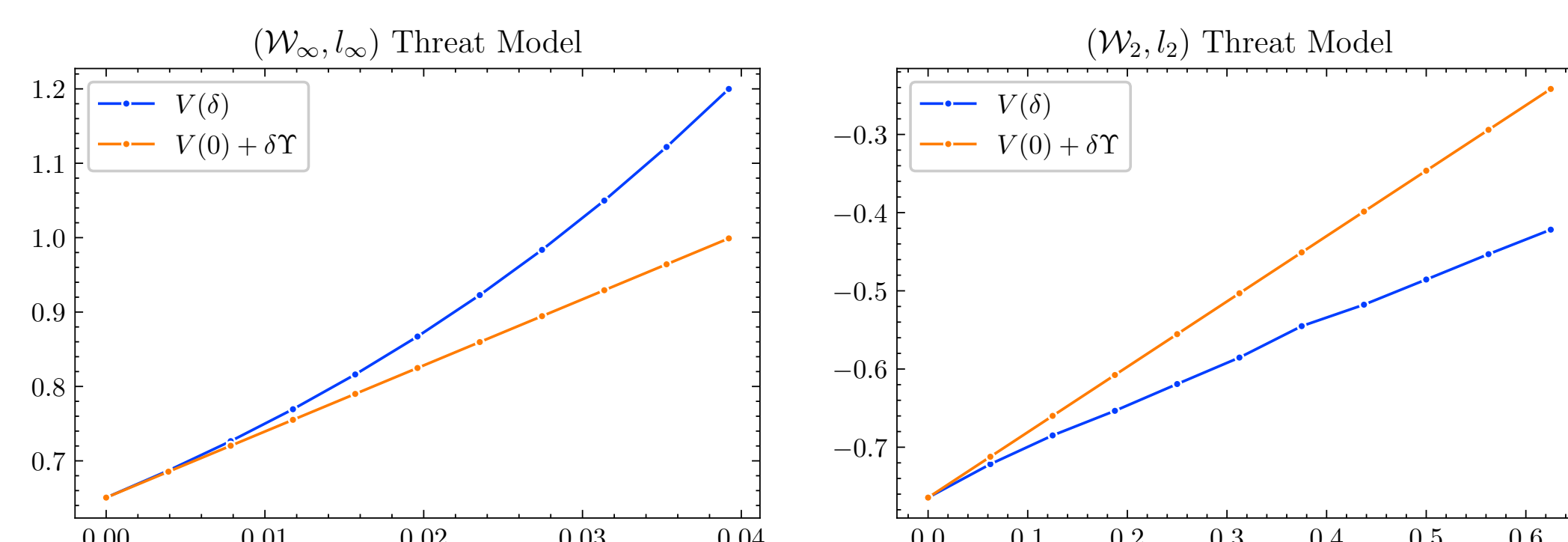
- $V(\delta) = V(0) + \delta \Upsilon + o(\delta)$ , where

$$\Upsilon = \left( E_P \left[ \|\nabla_x J_\theta(x, y)\|_*^q \right] \right)^{1/q}.$$

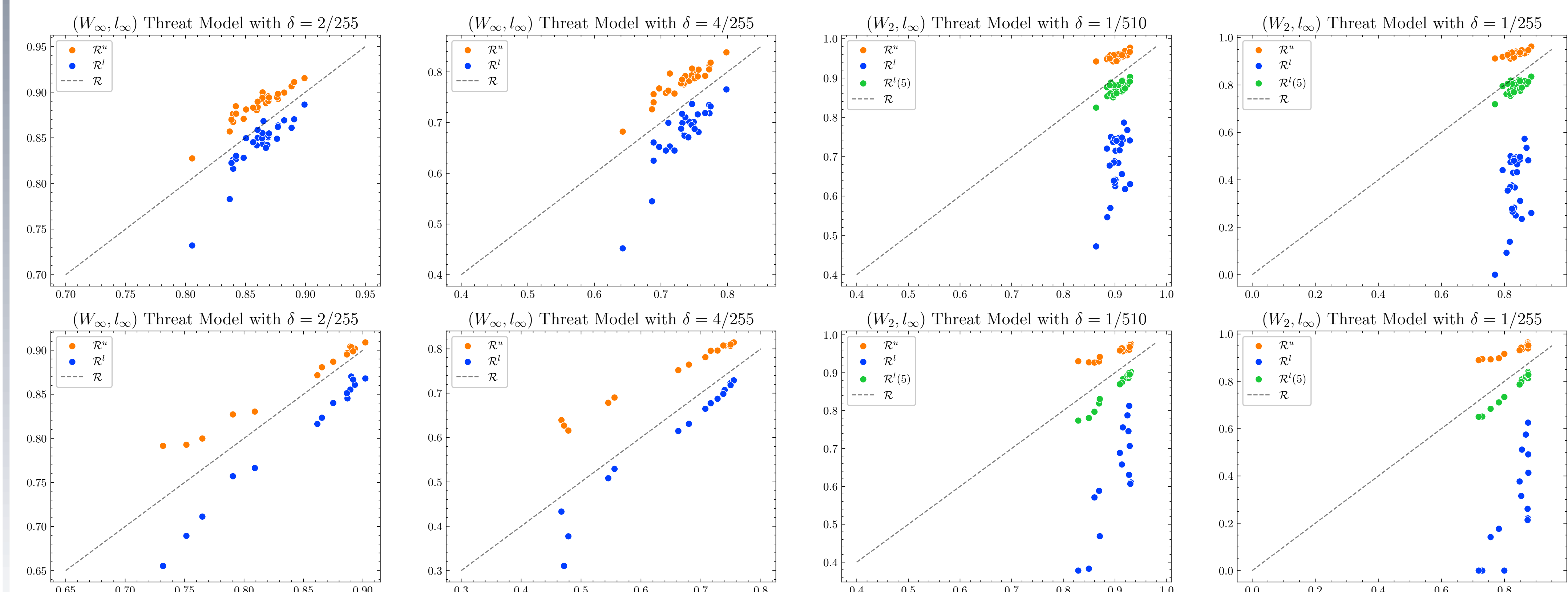
- $V(\delta) = E_{Q_\delta}[J_\theta(x, y)] + o(\delta)$ , where  $Q_\delta$  is the pushforward of  $P$  under the map

$$x \mapsto x + \delta h(\nabla_x J_\theta(x, y)) \|\Upsilon^{-1} \nabla_x J_\theta(x, y)\|_*^{q-1},$$

and  $h$  is uniquely determined by  $\langle h(x), x \rangle = \|x\|_*$ .



Performance of the first order approximation for the W-DRO value on CIFAR-10.



Certified bounds  $\mathcal{R}^u$  &  $\mathcal{R}^l$  versus  $\mathcal{R}$  CIFAR-100 (top) and ImageNet (bottom).