

Better Safe Than Sorry: Preventing Delusive Adversaries with Adversarial Training

Lue Tao^{1,2}, Lei Feng³, Jinfeng Yi⁴, Sheng-Jun Huang^{1,2}, Songcan Chen^{1,2*}

¹Nanjing University of Aeronautics and Astronautics

²MIT Key Laboratory of Pattern Analysis and Machine Intelligence

³Chongqing University

⁴JD AI Research

*s.chen@nuaa.edu.cn



Adversarial Examples

Adversarial Examples: worst-case data at **test** time

Clean training data \mathcal{D}



Clean test data \mathcal{D}



Accuracy : 90+%

Perturbations

Worst test data \mathcal{D}^*



Accuracy: 0%

[1] Biggio, et al. Evasion attacks against machine learning at test time. ECML-KDD, 2013.

[2] Szegedy, et al. Intriguing properties of neural networks. ICLR, 2014.

What if the training data can be perturbed?

Clean training data \mathcal{D}



Perturbations

Worst training data $\hat{\mathcal{D}}$



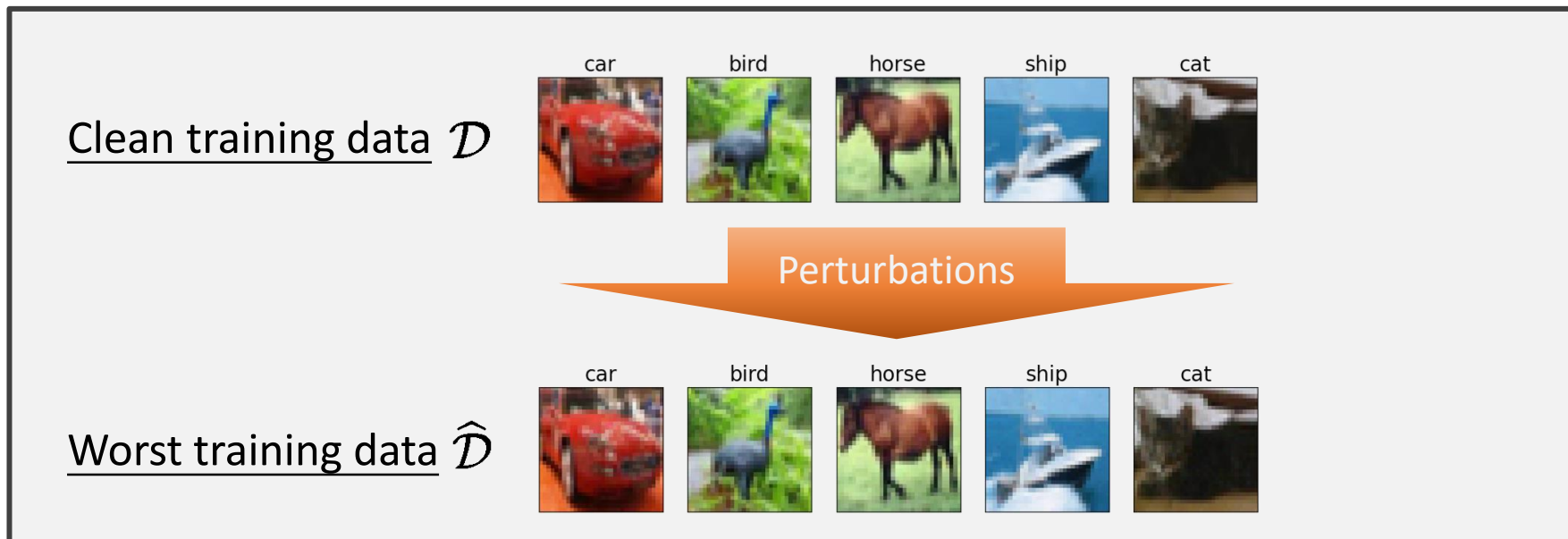
Clean test data \mathcal{D}



Accuracy : 0%

Delusive Attacks

Delusive Attacks: worst-case data at **training** time

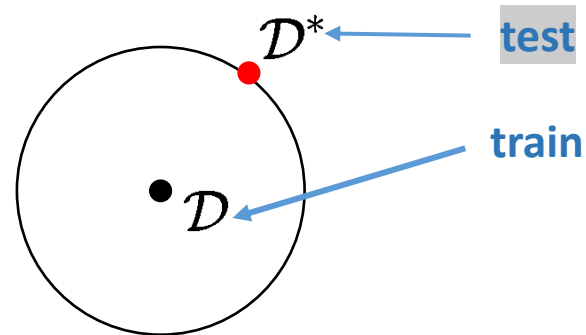


[3] Newsome, et al. Paragraph: Thwarting Signature Learning by Training Maliciously. Recent advances in intrusion detection, 2006.

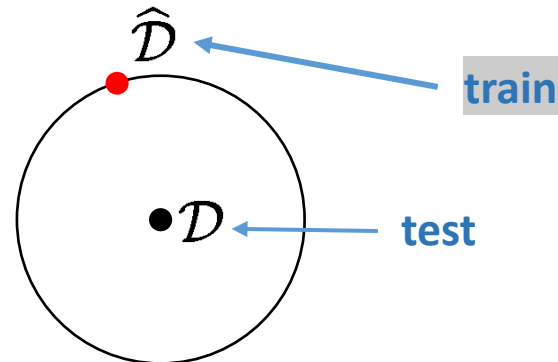
[4] Feng, et al. Learning to Confuse: Generating Training Time Adversarial Data with Auto-Encoder. NeurIPS, 2019.

Our Perspective: Twins of Evil

Adversarial Examples: worst-case **test** data

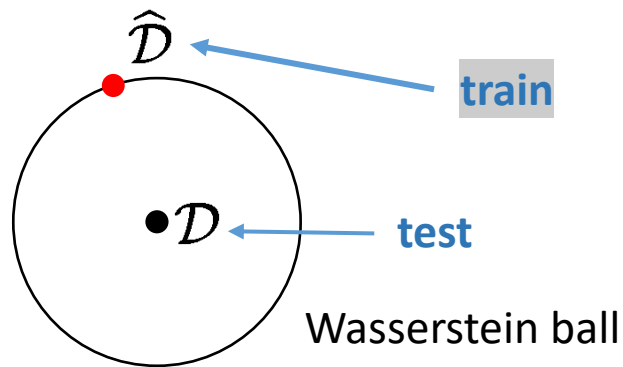


Delusive Attacks: worst-case **training** data



Our Contribution

[Contribution 1] Formulation of delusive attacks



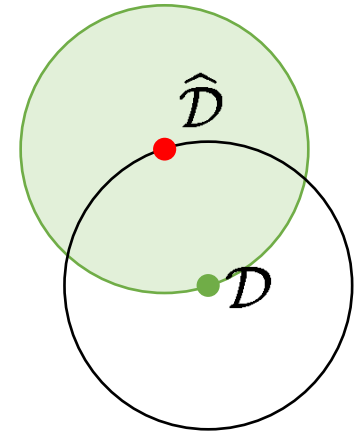
$$\begin{aligned} \max_{\hat{\mathcal{D}} \in \mathcal{B}_{W_\infty}(\mathcal{D}, \epsilon)} \quad & \mathbb{E}_{(\mathbf{x}, y) \sim \hat{\mathcal{D}}} [\ell(f_{\hat{\mathcal{D}}}(\mathbf{x}), y)], \\ \text{s.t.} \quad & f_{\hat{\mathcal{D}}} = \arg \min_f \mathbb{E}_{(\mathbf{x}, y) \sim \hat{\mathcal{D}}} [\ell(f(\mathbf{x}), y)]. \end{aligned}$$

Our Contribution

[Contribution 2] The principled defense

Theorem 1. For any data distribution \mathcal{D} and any delusive distribution $\hat{\mathcal{D}}$ such that $\hat{\mathcal{D}} \in \mathcal{B}_{W_\infty}(\mathcal{D}, \epsilon)$ generated by a delusive adversary, we have

$$\mathcal{R}_{\text{nat}}(f, \mathcal{D}) \leq \max_{\mathcal{D}' \in \mathcal{B}_{W_\infty}(\hat{\mathcal{D}}, \epsilon)} \mathcal{R}_{\text{nat}}(f, \mathcal{D}') = \mathcal{R}_{\text{adv}}(f, \hat{\mathcal{D}}).$$



Take-aways

1. Minimizing the adversarial risk on the perturbed data \Leftrightarrow Minimizing an upper bound of natural risk on the original data
2. Adversarial Training: A principled defense against delusive attacks

Our Contribution

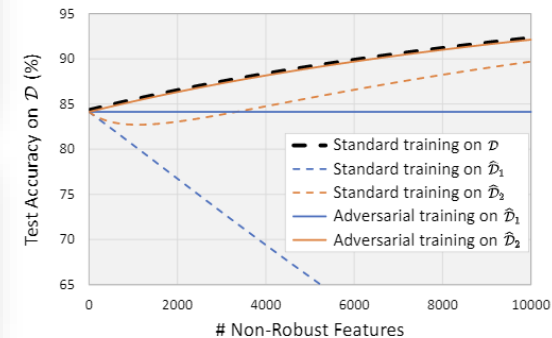
[Contribution 3] Internal Mechanisms

Theorem 2. Let $f_{\mathcal{D}}$, $f_{\hat{\mathcal{D}}_1}$, and $f_{\hat{\mathcal{D}}_2}$ be the Bayes optimal classifiers for the mixture-Gaussian distributions \mathcal{D} , $\hat{\mathcal{D}}_1$, and $\hat{\mathcal{D}}_2$, defined in Eqs. (5), (6), and (7), respectively. For any $\eta > 0$, we have

$$\mathcal{R}_{\text{nat}}(f_{\mathcal{D}}, \mathcal{D}) < \mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_2}, \mathcal{D}) < \mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_1}, \mathcal{D}).$$

Theorem 3. Let $f_{\hat{\mathcal{D}}_1, \text{rob}}$ and $f_{\hat{\mathcal{D}}_2, \text{rob}}$ be the optimal linear ℓ_{∞} robust classifiers for the delusive distributions $\hat{\mathcal{D}}_1$ and $\hat{\mathcal{D}}_2$, defined in Eqs. (6) and (7), respectively. For any $0 < \eta < 1/3$, we have

$$\mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_1}, \mathcal{D}) > \mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_1, \text{rob}}, \mathcal{D}) \quad \text{and} \quad \mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_2}, \mathcal{D}) > \mathcal{R}_{\text{nat}}(f_{\hat{\mathcal{D}}_2, \text{rob}}, \mathcal{D}).$$



Take-aways

1. Adversarial training works under delusive attacks by mitigating model reliance on non-robust features
2. Adversarial perturbations are more harmful than hypocritical perturbations

Our Contribution

[Contribution 4] Empirical evidences

➤ Practical delusive attacks

- Adversarial perturbations (P1),
- Hypocritical perturbations (P2),
- Universal adversarial perturbations (P3),
- Universal hypocritical perturbations (P4),
- Universal random perturbations (P5),
- DeepConfuse (L2C) [4]

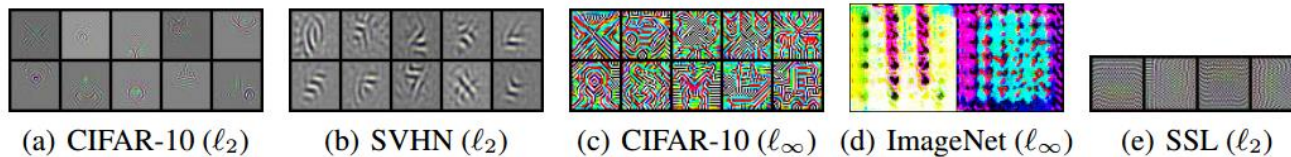
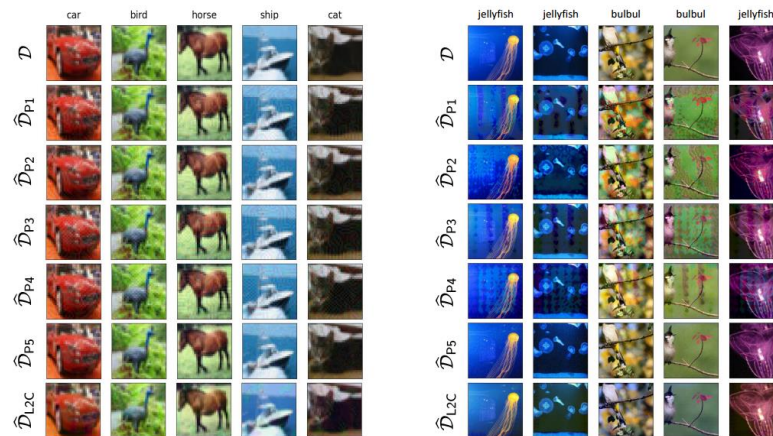


Figure 3: Universal perturbations for the P3 and P4 attacks across different datasets and threat models. Perturbations are rescaled to lie in the $[0, 1]$ range for display. The resulting inputs are nearly indistinguishable from the originals to a human observer (see Appendix B Figures 10, 11, and 12).

[4] Feng, et al. Learning to Confuse: Generating Training Time Adversarial Data with Auto-Encoder. NeurIPS, 2019.

Our Contribution

[Contribution 4] Empirical evidences

➤ Six delusive attacks

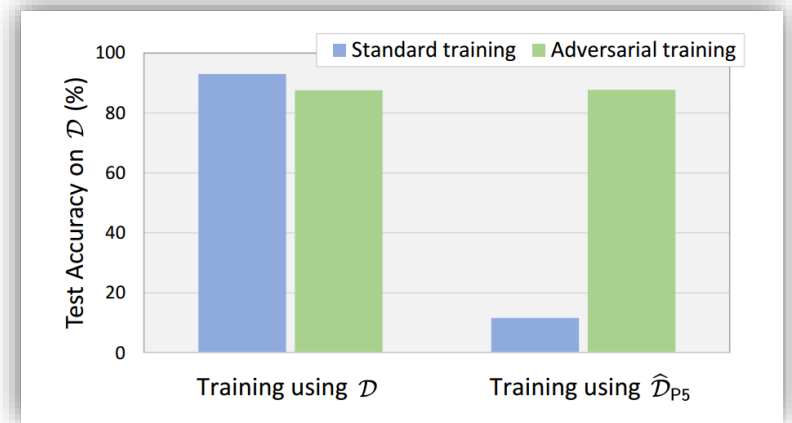
Adversarial perturbations (P1), Hypocritical perturbations (P2), Universal adversarial perturbations (P3), Universal hypocritical perturbations (P4), Universal random perturbations (P5), and DeepConfuse (L2C)

➤ Three datasets

CIFAR-10, SVHN, and a subset of ImageNet

➤ Three tasks

supervised learning, self-supervised learning, and overcoming simplicity bias



Take-aways

1. The defense withstands all the attacks on all the datasets/tasks.
2. Both theoretical and empirical results vote for adversarial training.

Thanks!