Adversarial Attacks are Reversible with Natural Supervision







Columbia University, Rutgers University*

Chengzhi Mao, Mia Chiquier, Hao Wang*, Junfeng Yang, Carl Vondrick





Deep Networks are Vulnerable under Adversarial Attacks



Input Image

Adversarial Perturbations

Moosavi-Dezfooli et al.

Attacked Image





[1] Madry et al. Towards deep learning models resistant to adversarial attacks [2] Mao et al. Metric Learning for Adversarial Robustness. [3] Mao et al. Multitask learning strengthens adversarial robustness.

Existing defense focus on the training algorithm

Multiple Tasks [3]



 $\delta = \operatorname{argmax}_{\delta < \epsilon} \mathscr{L}(Y, F_{\theta}(X + \delta))$ Adversarial Training [1] $\theta = \operatorname{argmin}_{\theta}(\mathscr{L}(Y, F_{\theta}(X + \delta)))$



Additional tasks during training secure the model

		None (Baseline)	SemSeg	DepthZ	Edge2D	Normal	Reshading	Key2D	Key3D	DepthE	AutoE	Edge3D	PCurve	Mean
raining Task 1 + Testing Task	SemSeg *	13.36	0.00	44.61	4.42	24.48	9.13	17.51	3.14	11.60	5.61	10.18	11.53	14.22
	DepthZ (10 ⁻²)	11.49	59.00	0.00	40.99	-1.10	-8.02	3.22	27.27	8.64	56.65	-6.43	56.18	23.64
	Edge2D (10 ⁻²)	10.67	7.79	12.27	0.00	10.55	6.83	8.81	9.54	8.98	6.85	6.50	5.41	8.35
	Normal (10 ⁻²)	40.93	14.06	-4.75	3.89	0.00	1.04	3.58	2.43	-2.80	12.71	9.42	-0.70	3.89
	Reshading (10^{-2})	57.89	15.66	0.18	5.05	2.42	0.00	3.36	7.93	-3.68	-5.37	14.80	0.46	4.08
	Key2D (10 ⁻²)	11.72	7.40	7.08	8.69	10.19	7.13	0.00	6.56	9.41	6.27	8.23	9.88	8.08
	Key3D (10 ⁻²)	49.70	37.72	0.18	-2.19	7.73	15.14	11.81	0.00	-3.04	34.47	-7.48	-6.39	8.80
	DepthE (10 ⁻³)	4.85	27.15	29.95	32.96	11.87	-17.11	24.24	23.07	0.00	23.57	31.19	39.44	22.63
	AutoE (10 ⁻²)	59.31	2.60	-1.62	1.75	-5.08	-0.17	-0.03	-2.37	1.80	0.00	-1.94	-3.68	-0.87
	Edge3D (10 ⁻²)	15.90	8.36	3.58	-2.71	3.38	4.45	1.96	-6.32	3.12	20.69	0.00	6.98	4.35
F	PCurve (10 ⁻⁴)	11.47	22.22	22.38	9.74	19.51	16.14	22.39	9.04	2.81	19.91	9.31	0.00	15.34

Training Task 2

(a) Performance Under Attack



Natural Images contain intrinsic structure

















Natural Images contain intrinsic structure

















Natural Images contain intrinsic structure









Adversarial Attack also damages the incidental structure of image











If the image is natural, the backward path should not If the image is attacked, the backward path can fix the

Algorithm



Contrastive Loss

$$\mathcal{L}_s(\mathbf{x}) = -\mathbb{E}_{i,j} \left[\mathbf{y}_{ij}^{(s)} \log \frac{\exp(\cos(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_k \exp(\cos(\mathbf{z}_i, \mathbf{z}_k)/\tau)} \right]$$



Contrastive Loss for Self-supervision



$$\mathcal{L}_s(\mathbf{x}) = -\mathbb{E}_{i,j} \left[\mathbf{y}_{ij}^{(s)} \log \frac{\exp(\cos(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_k \exp(\cos(\mathbf{z}_i, \mathbf{z}_k)/\tau)} \right]$$





Contrastive Loss for Self-supervision

Our Method

Cannon

Defended DNN Wong et al. 2020

Algorithm 1 Self-supervised Reverse Attack

- 1: Input: Potentially attacked image x, step size η , number of iterations K, a classifier F, reverse attack bound ϵ_v , and self-supervised loss function \mathcal{L}_s .
- 2: **Output:** Class prediction \hat{y}
- 3: Inference:
- 4: $\mathbf{x}' \leftarrow \mathbf{x} + \mathbf{n}$, where **n** is the initial random noise
- 5: for k = 1, ..., K do
- 6: $\mathbf{x}' \leftarrow \mathbf{x}' \eta \nabla_{\mathbf{x}'} \mathcal{L}_s(\mathbf{x}')$
- 7: $\mathbf{x}' \leftarrow \Pi_{(\mathbf{x},\epsilon_v)}\mathbf{x}'$, which projects the image back into the bounded region.
- 8: **end for**
- 9: Predict the final output by $\hat{y} = F(\mathbf{x}')$

Our Method

Cannon

Defended DNN Wong et al. 2020

Algorithm 1 Self-supervised Reverse Attack

- 1: Input: Potentially attacked image x, step size η , number of iterations K, a classifier F, reverse attack bound ϵ_v , and self-supervised loss function \mathcal{L}_s .
- 2: **Output:** Class prediction \hat{y}
- 3: Inference:
- 4: $\mathbf{x}' \leftarrow \mathbf{x} + \mathbf{n}$, where **n** is the initial random noise
- 5: for k = 1, ..., K do
- 6: $\mathbf{x}' \leftarrow \mathbf{x}' \eta \nabla_{\mathbf{x}'} \mathcal{L}_s(\mathbf{x}')$
- $\mathbf{x}' \leftarrow \Pi_{(\mathbf{x},\epsilon_v)} \mathbf{x}'$, which projects the image back into 7: the bounded region.
- 8: **end for**
- 9: Predict the final output by $\hat{y} = F(\mathbf{x}')$

Adaptive attack: Lose-lose situation for attacker

If the attack ignores our defense, our defense improves accuracy

If the attacker adapts to our defense, it will hurt their attack and increase our robust accuracy.

This makes sense which means the attacker help us find a attack that also produces a good prior

- ---- RO (No Defense)
- RO (Our Defense)
- ---- Semi SL (No Defense)
- Semi SL (Our Defense)
- Optimal λ

Theoretical Guarantee

Our approach can improve the classification accuracy upper bound.

Lemma 1. The standard classifier under adversarial attack is equivalent to predicting with $P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}_a^{(s)})$, and our approach is equivalent to predicting with $P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$.

Theorem 1. Assume the base classifier operates better than chance and instances in the dataset are uniformly distributed over n categories. Let the prediction accuracy bounds be $P(\mathbf{Y}|\mathbf{Y}_{a}^{(s)}, \mathbf{X}_{a}) \in [b_{1}, c_{1}]$ and $P(\mathbf{Y}|\mathbf{Y}^{(s)}, \mathbf{X}_{a}) \in [b_{2}, c_{2}]$. If the conditional mutual information $I(\mathbf{Y}; \mathbf{Y}^{(s)}|\mathbf{X}_{a}) > 0$, we have $b_{2} \ge b_{1}$ and $c_{2} > c_{1}$, which means our approach strictly improves the bound for classification accuracy.

Theoretical Guarantee

Lemma 1. The standard classifier under adversarial attack is equivalent to predicting with $P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}_a^{(s)})$, and our approach is equivalent to predicting with $P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$. *Proof.* For the standard classifier under attack, we know that $P(\mathbf{Y}^{(s)} = \mathbf{y}_a^{(s)} | \mathbf{X} = \mathbf{x}_a) = 1$. Thus we know the standard classifier under adversarial attack is equivalent to

$$P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a) = \sum_{\mathbf{Y}^{(s)}} P(\mathbf{Y}^{(s)}|\mathbf{X} = \mathbf{x}_a) P(\mathbf{Y}|\mathbf{Y}^{(s)}, \mathbf{X} = x_a)$$
$$= P(\mathbf{Y}|\mathbf{Y}^{(s)} = \mathbf{y}_a^{(s)}, \mathbf{X} = \mathbf{x}_a).$$

Our algorithm finds a new input image $\mathbf{x}_{\max}^{(n)}$ that

$$\operatorname{argmax}_{\mathbf{x}^{(n)}} P(\mathbf{X}^{(n)} = \mathbf{x}^{(n)} | \mathbf{X} = \mathbf{x}_a) P(\mathbf{Y}^{(s)} = \mathbf{y}^{(s)} | \mathbf{X}^{(n)} = \mathbf{x}^{(n)})$$
$$= \operatorname{argmax}_{\mathbf{x}^{(n)}} P(\mathbf{X}^{(x)} = \mathbf{x}^{(n)} | \mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)}).$$

Our algorithm first estimate $\mathbf{x}_{\max}^{(n)}$ with adversarial image \mathbf{x}_a and self-supervised label $\mathbf{y}^{(s)}$. We then predict the label \mathbf{Y} using our new image $\mathbf{x}_{\max}^{(n)}$. Thus, our approach in fact estimates $P(\mathbf{Y}|\mathbf{X}^{(n)} = \mathbf{x}_{\max}^{(n)})P(\mathbf{X}^{(n)} = \mathbf{x}_{\max}^{(n)}|\mathbf{X} =$ $\mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$. Note the following holds:

$$P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$$
(9)

$$= \sum_{\mathbf{x}^{(n)}} P(\mathbf{Y}|\mathbf{x}^{(n)}) P(\mathbf{x}^{(n)}|\mathbf{X}=\mathbf{x}_a, \mathbf{Y}^{(s)}=\mathbf{y}^{(s)})$$
(10)

$$\approx P(\mathbf{Y}|\mathbf{X}^{(n)} = \mathbf{x}_{\max}^{(n)})P(\mathbf{X}^{(n)} = \mathbf{x}_{\max}^{(n)}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$$
(11)

Thus our approach is equivalent to estimating $P(\mathbf{Y}|\mathbf{X} = \mathbf{x}_a, \mathbf{Y}^{(s)} = \mathbf{y}^{(s)})$.

Theoretical Guarantee

Theorem 1. Assume the base classifier operates better than chance and instances in the dataset are uniformly distributed over n categories. Let the prediction accuracy bounds be $P(\mathbf{Y}|\mathbf{Y}_a^{(s)},\mathbf{X}_a) \in [b_1,c_1]$ and $P(\mathbf{Y}|\mathbf{Y}^{(s)}, \mathbf{X}_a) \in [b_2, c_2]$. If the conditional mutual information $I(\mathbf{Y}; \mathbf{Y}^{(s)} | \mathbf{X}_a) > 0$, we have $b_2 \ge b_1$ and $c_2 > c_1$, which means our approach strictly improves the bound for classification accuracy.

We here use Fano's Inequality to connect entropy to our prediction accuracy

Venn diagram for our defense

Adversarial Atta Classificatio

Attack

nsic Signal from Reversed Supervision Task

Our defense still work when adaptive attack know our defense

Defense-Awar Attacker Kept out of Defense-Awar Attack - Supervision Task Attack

The state-of-the-art Defense

	CIFAR10	CIFAR100	SVHN	ImageNet		
STOA Defense	60.2%	27.6%	53.6%	27.81%		
Adding Ours	+7.5%	+5.5%	+11.8%	+3.1%		

Experiment Result

Our defense generalize to larger adversarial attack budget

improves the robustness.

Figure 7: The trade-off between adversarial robust accuracy vs. clean accuracy on CIFAR-10, CIFAR-100, and SVHN under the L_{∞} norm. We increase the noise budget ϵ_v from small to large, which causes the clean accuracy to drop from right to left. Our method produces a better reversal of the adversarial perturbation than just adding random noise to reverse it.

Our approach has better trade-off between robustness and clean accuracy

Figure 6: The adversarial robust accuracy vs. perturbation budget curves on CIFAR-10, CIFAR-100, SVHN, and ImageNet, under the L_{∞} norm. The red line is applying our inference algorithm to the baseline models [50, 61, 65]. Using our inference algorithm significantly

Steps for Reversal

Visualization for feature space before and after reversal

Visualization

Desk

Wall Clock

Church

Trailer Truck

Image Attack

Reverse Attack

Desk

Oscilloscope

Analog

Castle

Wall Clock

Church

Harvester Trailer Truck

Malamute Eskimo Dog Malamute

Image Attack

Eagle

Cannon

Beach Wagon

Worm Fence

Cello

African Grey

Traffic Light

Reverse

Attack

Chainlink

Fence

Violin

Cello

Conclusion and future work

- 1. We propose an inference-time defense for adversarial attacks by restoring the intrinsic structure of the input.
- 2. Our method achieves the state-of-the-art defense on the auto-attack leaderboard
- 3. Our test time defense method is generic, which is compatible with existing language, and beyond.

Code: <u>https://github.com/cvlab-columbia/SelfSupDefense</u> Paper: https://arxiv.org/pdf/2103.14222.pdf

training-time defense method, and can be used in other modalities such as sound,