

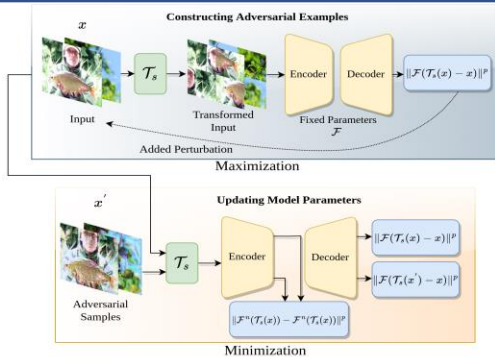


Introduction

In the black-box setting, adversarial examples are typically created using surrogate models trained on the target model's large labeled distribution. Such attacks are thus limited by the availability of a pretrained surrogate model and label space information. Our work focuses on a stronger threat model on how to learn an effective surrogate model from the limited **unlabelled** data and then how to generate self-supervised transferable adversarial examples.

With limited samples, training surrogate models in a supervised fashion (**conventional**) causes severe overfitting, leading to poor adversarial transferability. We propose a **Self-Supervised Adversarial Training** method to find highly transferable patterns by learning over flatter loss surfaces

Self-supervised Adversarial Training



We train an autoencoder-based surrogate model via self-supervised adversarial pixel restoration to learn generalizable representations from limited data samples (≤ 20) or cross-domain samples.

In the **maximization** step, adversarial examples are generated by fooling model's reconstruction ability; in the **minimization** step, model parameters are updated using restoration objectives between adversarial and clean samples.

Self-supervised Adversarial Objectives

Maximization Objective: Adversarial examples can be generated by maximizing a self-supervised objective based on reconstructing a transformed (rotated/shuffled) image.

$$\underset{x'}{\text{maximize}} \quad \mathcal{L}_{max} = \|\mathcal{F}(T_s(x')) - x\|^p \quad x' = x + \alpha \times \nabla_x \mathcal{L}_{max}$$

x' is the adversarial image, T_s represents pixel transformation (e.g., rotation or jigsaw shuffle)

Minimization Objective: The surrogate model is trained by minimizing the reconstruction error of the adversarial and clean output with respect to the original image. Furthermore, the model's feature space is regulated by enforcing alignment between clean and adversarial features at the encoder stage.

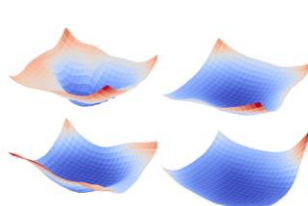
$$\mathcal{L}_{out} = \|\mathcal{F}(T_s(x')) - x\|^p + \|\mathcal{F}(T_s(x)) - x\|^p \quad \text{and} \quad \mathcal{L}_{feature} = \|\mathcal{F}^n(T_s(x')) - \mathcal{F}^n(T_s(x))\|^p$$

$$\mathcal{L}_{min} = \mathcal{L}_{out} + \lambda \mathcal{L}_{feature}$$

Training Settings & Unsupervised Attack

Limited Samples: Surrogate model are trained only on the few data samples (≤ 20) on which adversarial examples need to be crafted.

Cross-Domain Samples: Due to abundance of unlabelled data, we scale our self-supervised adversarial training to large-scale datasets and then test the cross-domain transferability of our method.



The first row shows loss landscapes of surrogate models trained by reducing the reconstruction objective (**top**: rotation, **bottom**: jigsaw). The second row shows the smoother loss surfaces obtained by using our training method. This has significant effect on finding generalizable adversarial examples with better transferability.

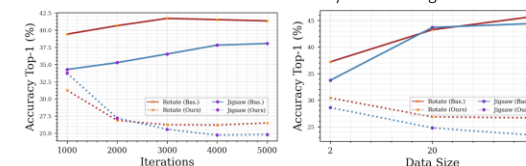
Crafting Adversarial examples: The attack objective for surrogate models trained in a self-supervised manner is based on maximizing the reconstruction error between clean and adversarial samples. Adversarial transferability is evaluated on a selected set of 5000 images from ImageNet validation set, with perturbation budget of 0.1.

Quantitative Results

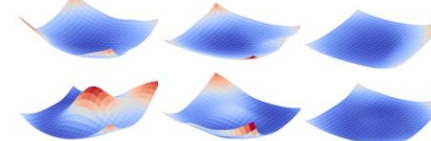
Limited Samples

Transformation	Method	VGG-19	Inc-V3	Res152	Dense121	SeNet	WRN	MNet-V2	Average
Jigsaw	Baseline (Ours)	31.54	50.28	46.24	42.38	59.06	51.24	25.24	43.71
	Ours	16.82	25.54	31.18	22.64	38.06	25.76	13.70	24.81 (-18.9)
Rotation	Baseline (Ours)	31.14	48.14	47.40	41.26	58.20	50.72	26.00	43.27
	Ours	19.02	25.76	33.60	25.60	38.92	29.78	15.38	26.87 (-16.4)
Prototypical	Baseline (Ours)	18.74	33.68	34.72	26.06	42.36	33.14	16.34	29.29
	Ours	17.02	21.48	28.66	21.06	35.04	23.56	13.06	22.84 (-6.45)

Our attack boosts adversarial transferability across ImageNet models.

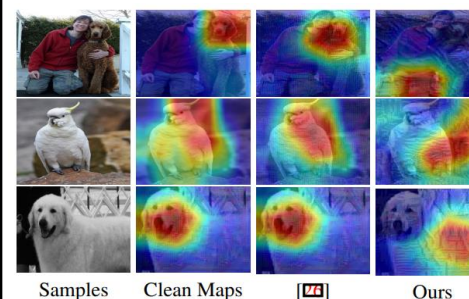


The performance of our method improves with more training iterations and data size in contrast to the baseline.



Loss landscape of surrogate models with increasing robustness strength (α) from left to right. The first row shows the loss surface on the clean samples, while as the second row plots the loss surface with respect to adversarial samples. It becomes harder to maximize the reconstruction error or flip decisions on the excessively smooth loss surface during attack.

Attention to salient regions



Cross-Domain Samples

Transformation	Method (→)	CoCo	Paintings	Comics
Rotation	CoCo	28.56	23.31	17.75
	Ours	27.83	17.75	24.19
Jigsaw	CoCo	43.93	31.28	31.28
	Ours	44.07	33.42	41.54

Our method provides favorable results on cross-domain transferability averaged across ImageNet models.

Transformation (↓)	Dataset (→)	CoCo	Paintings	Comics
Rotation	No Attack	39.7	19.3	17.2
	Ours	19.3	14.6	11.9
Jigsaw	No Attack	39.7	24.1	14
	Ours	38	20.8	13.3

Transferability to object detector (DETR) based on mAP is evaluated on CoCo validation set.

Transformation (↓)	Dataset (→)	CoCo	Paintings	Comics
Rotation	No Attack	61.8	53.2	52.6
	Ours	48.9	46.9	47.81
Jigsaw	No Attack	61.8	53.9	48.5
	Ours	58.29	51.65	51.65

Transferability to object segmentation (DINO) based on Jaccard index is evaluated on DAVIS validation set.

Conclusion

- We show the benefits of unsupervised adversarial training to learn transferable adversarial perturbations.
- Our adversarial training method reduces overfitting during training and can exploit very few data samples to learn meaningful adversarial features while it can also scale to large unsupervised datasets.
- Our unsupervised attack is task independent and allows cross-domain attacks (e.g., learning surrogate on comics and transferring its perturbations to models trained on natural images).