

1 We thank the reviewers for the positive feedback: new state-of-the-art results
 2 (R1,2&3), first to explore cross-domain transferability (R1), high significance
 3 to the community (R3), very well written and clear presentation (R1,2&3).
 4 **Code:** Code will be made public. Fig.(1, 2, 3) best viewed in zoom.

5 **R1,2&3: Significance of Relativistic Cross-Entropy (RCE):** Adversarial per-
 6 turbations are crafted via loss function gradients. An effective loss helps in
 7 adversary generation by back-propagating *stronger* gradients. Below, we show
 8 that $\mathcal{RC}\mathcal{E}$ ensures this requisite and thus leads to better performance than $\mathcal{C}\mathcal{E}$.

9 **Notation:** classifier \mathcal{F} , clean sample \mathbf{x} , adversarial example \mathbf{x}' , output scores $a = \mathcal{F}(\mathbf{x})$, $a' = \mathcal{F}(\mathbf{x}')$.

10 **Gradient Perspective:** Let $\mathcal{C}\mathcal{E}(a', y) = -\log(e^{a'_y} / \sum_k e^{a'_k})$ be the CE loss for input \mathbf{x}' . For clarity, we define
 11 $p'_y = e^{a'_y} / \sum_k e^{a'_k}$. The derivative of p'_y w.r.t a'_i is $\partial p'_y / \partial a'_i = p'_y (\mathbb{1}[i=y] - p'_i)$. From chain rule, $\partial \mathcal{C}\mathcal{E} / \partial a'_i = p'_i - \mathbb{1}[i=y]$
 12 (Eq. 1). For relativistic loss, $\mathcal{RC}\mathcal{E}(a', a, y) = -\log(e^{a'_y - a_y} / \sum_k e^{a'_k - a_k})$, we define $r_y = (e^{a'_y - a_y} / \sum_k e^{a'_k - a_k})$. The
 13 derivative of r_y w.r.t a'_i is $\partial r_y / \partial a'_i = r_i (\mathbb{1}[i=y] - r_y)$. From chain rule, $\partial \mathcal{RC}\mathcal{E} / \partial a'_i = r_i - \mathbb{1}[i=y]$ (Eq. 2).

14 In light of above relations, $\mathcal{RC}\mathcal{E}$ has three important properties: (a) Comparing (Eq. 2) with (Eq. 1) shows that $\mathcal{RC}\mathcal{E}$ gra-
 15 dient is a function of 'difference' ($a'_y - a_y$) as opposed to only scores a'_y in $\mathcal{C}\mathcal{E}$ loss. Thus it measures the relative change
 16 in prediction as an explicit objective during optimization. (b) $\mathcal{RC}\mathcal{E}$ loss back-propagates larger gradients compared to
 17 $\mathcal{C}\mathcal{E}$, resulting in efficient training and stronger adversaries (see Fig. 1 for empirical evidence). **Sketch Proof:** We can
 18 factorize the denominator in (Eq. 2) as follows: $\partial \mathcal{RC}\mathcal{E} / \partial a'_i = (e^{a'_y - a_y} / (e^{a'_y - a_y} + \sum_{k \neq y} e^{a'_k - a_k})) - \mathbb{1}[i=y]$. Consider
 19 the fact that maximization of $\mathcal{RC}\mathcal{E}$ is only possible when $e^{(a'_y - a_y)}$ decreases and $\sum_{k \neq y} e^{(a'_k - a_k)}$ increases. Generally,
 20 $a_y \gg a_{k \neq y}$ for the score generated by a pre-trained model and $a'_y \ll a'_{k \neq y}$. Thus, $\partial \mathcal{RC}\mathcal{E} / \partial a'_i > \partial \mathcal{C}\mathcal{E} / \partial a'_i$ since
 21 $e^{(a'_y - a_y)} < e^{(a'_y)}$ and $\sum_{k \neq y} e^{(a'_k - a_k)} > \sum_{k \neq y} e^{(a'_k)}$. In simple words, the gradient strength of $\mathcal{RC}\mathcal{E}$ is higher than $\mathcal{C}\mathcal{E}$.

22 (c) In case \mathbf{x} is misclassified by $\mathcal{F}(\cdot)$, the gradient strength of $\mathcal{RC}\mathcal{E}$ is still higher than $\mathcal{C}\mathcal{E}$ (here noise update with the
 23 $\mathcal{C}\mathcal{E}$ loss will be weaker since adversary's goal is already achieved i.e., \mathbf{x} is misclassified). We will add it in final version.

24 **Evaluation:** We further validate (see Fig. 2) the significance of $\mathcal{RC}\mathcal{E}$ compared to $\mathcal{C}\mathcal{E}$ in terms of three criterion
 25 (accuracy, logits difference and transfer to unseen classes). For the test on unseen classes, we divide ImageNet into two
 26 mutually exclusive sets (500 classes each), named IN1 and IN2. VGG16 is trained on IN1 & IN2 from scratch.

27 **R1: 1) RCE Justification:** See R1,2&3 above. 2)

28 **Relation with Style-Transfer:** We visualize the
 29 intermediate feature space of cross-domain per-
 30 turbed images and compare it with original and
 31 stylized images (Fig. 3). We note that the feature
 32 space of perturbed images is fairly shifted from
 33 the original and stylized images. This shows that
 34 although some of the generated patterns resemble
 35 "style" of a specific domain (e.g., in Fig. 3 main
 36 paper), the overall behaviour of our proposed ap-
 37 proach is distinct from style transfer. This is potentially due to the existence of "non-robust features" defined as 'features
 38 that are highly predictive but brittle and incomprehensible to humans' [A1]. Since, our generated perturbations are
 39 bounded (as opposed to unbounded style transfer), the attacker is likely to focus on the non-robust features. We will add
 40 further qualitative examples on other domains in final version (Fig. 6 in supp. material). 3) **Notations:** Will update in
 41 the final draft. 4) **On Adversarial Training Defenses:** Our main draft already includes evaluations with adversarial
 42 training (Tab.5&6 in paper). 5) **On the Existence of Universal Adversarial Function (UAF):** Earlier works [A2,A3]
 43 show that universal adversarial perturbations exist due to overlap in decision space of different classification models.
 44 Our work empirically shows that the same holds true even across different domains. This possibly happens due to the
 45 overlap between latent low-dimensional manifolds across different domains.

46 **R2: Theoretical Result:** See R1,2&3 earlier. **Typo:** We thank R2 & fix it.

47 **R3: 1) Use of Instance-Agnostic:** We used this term to differentiate the
 48 one-time training feature of our attack as opposed to instance-specific
 49 attacks. However, we acknowledge R3's point and will replace this term
 50 with *domain-agnostic* for clarity. 2) **Comparison with [1,19]:** [1] trains
 51 conditional generators to learn original data manifold and searches the latent
 52 space conditioned on the human recognizable target class that is misclassified by a target classifier. Different to [1],
 53 our approach learns to add adversarial noise to the original samples. [19] produces adversarial images by employing a
 54 separate discriminator alongside classifier. Different to [19], we train a generator to first produce unbounded adversaries
 55 and then project them to nearby original images. We thank R3 and will add further discussion in final version.

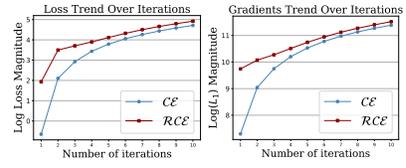


Figure 1: Loss and gradients trend for CE and RCE loss. Results are reported with VGG16 network on 100 random images for MIFGSM attack.

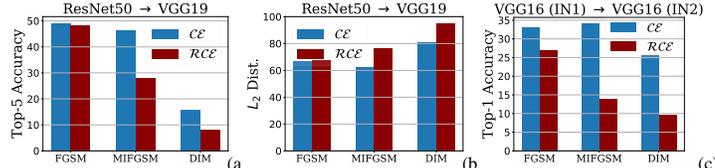


Figure 2: (a) shows Top-5 accuracy of adversaries (lower is better), (b) shows normalized l_2 difference b/w logits of adversarial and benign examples (higher is better) while (c) shows transferability to unseen classes. In each case $\mathcal{RC}\mathcal{E}$ perform significantly better than $\mathcal{C}\mathcal{E}$.

46 **R2: Theoretical Result:** See R1,2&3 earlier. **Typo:** We thank R2 & fix it.

47 **R3: 1) Use of Instance-Agnostic:** We used this term to differentiate the
 48 one-time training feature of our attack as opposed to instance-specific
 49 attacks. However, we acknowledge R3's point and will replace this term
 50 with *domain-agnostic* for clarity. 2) **Comparison with [1,19]:** [1] trains
 51 conditional generators to learn original data manifold and searches the latent
 52 space conditioned on the human recognizable target class that is misclassified by a target classifier. Different to [1],
 53 our approach learns to add adversarial noise to the original samples. [19] produces adversarial images by employing a
 54 separate discriminator alongside classifier. Different to [19], we train a generator to first produce unbounded adversaries
 55 and then project them to nearby original images. We thank R3 and will add further discussion in final version.

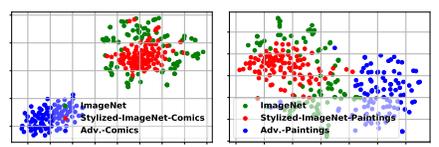


Figure 3: t-SNE visualization for features of 100 images and their corresponding stylized and perturbed versions. VGG16 is used to extract features.

[A1] Ilyas, Andrew, et al. "Adversarial examples are not bugs, they are features." arXiv (2019). [A2] Tramèr, Florian, et al. "The space of transferable adversarial examples." arXiv (2017). [A3] Dezfouli, Seyyed, et al. "Analysis of universal adversarial perturbations." arXiv (2017).