

Towards Transferable Unrestricted Adversarial Examples with Minimum Changes

Fangcheng Liu¹, Chao Zhang¹ *, Hongyang Zhang² *

¹Key Laboratory of Machine Perception (MOE), School of EECS, Peking University

²David R. Cheriton School of Computer Science, University of Waterloo & Vector Institute

{lfc@stu., c.zhang@}{pku.edu.cn hongyang.zhang@waterloo.ca}

Abstract

Transfer-based adversarial example is one of the most important classes of black-box attacks. However, there is a trade-off between transferability and imperceptibility of the adversarial perturbation. Prior work in this direction often requires a fixed but large ℓ_p -norm perturbation budget to reach a good transfer success rate, leading to perceptible adversarial perturbations. On the other hand, most of the current unrestricted adversarial attacks that aim to generate semantic-preserving perturbations suffer from weaker transferability to the target model. In this work, we propose a *geometry-aware framework* to generate transferable adversarial examples with minimum changes. Analogous to model selection in statistical machine learning, we leverage a validation model to select the optimal perturbation budget for each image under both the ℓ_∞ -norm and unrestricted threat models. Extensive experiments verify the effectiveness of our framework on balancing imperceptibility and transferability of the crafted adversarial examples. The methodology is the foundation of our entry to the *CVPR’21 Security AI Challenger: Unrestricted Adversarial Attacks on ImageNet*, in which we ranked 1st place out of 1,559 teams and surpassed the runner-up submissions by 4.59% and 23.91% in terms of final score and average image quality level, respectively. Code is available at <https://github.com/Equationliu/GA-Attack>.

1 Introduction

Though deep neural networks have exhibited impressive performance in various fields (He et al., 2016; Dosovitskiy et al., 2021), they are vulnerable to adversarial examples (Szegedy et al., 2014; Goodfellow et al., 2015; Shao et al., 2021; Bhojanapalli et al., 2021; Bai et al., 2021), where test inputs that have been modified slightly strategically cause misclassification. Adversarial examples have posed serious threats to various security-critical applications, such as autonomous driving (Bojarski et al., 2016) and face recognition (Parkhi et al., 2015). Most positive results on adversarial attacks have focused on white-box settings (Athalye et al., 2018; Tramer et al., 2020), where the attacker can get full access to the defense models. However, the problem becomes more challenging when it comes to the black-box setting, where the attacker has no information about the model architecture, hyper-parameters, and even the outputs of the black-box model. In this setting, adversarial examples

*Corresponding author

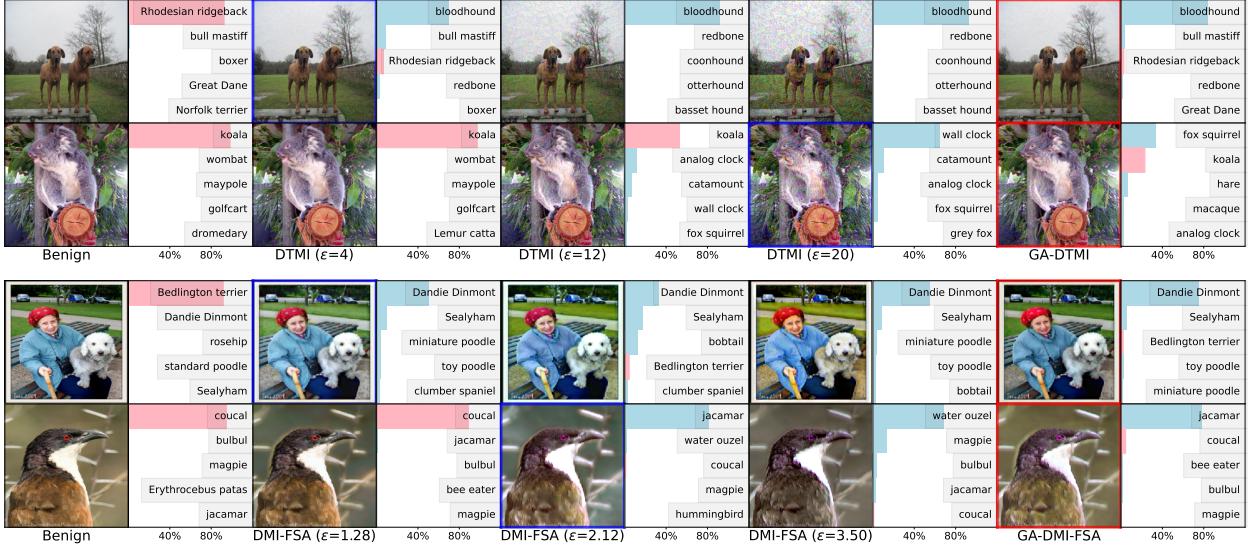


Figure 1: Comparison between our method and two baselines under both the ℓ_∞ -norm (top) and unrestricted (bottom) threat models using various perturbation radii. In the even columns, we present the top-5 confidence bars of the target model for the images in the left. The ground-truth label is marked by pink and other labels are marked by blue. In each row, the misclassified adversarial example with minimum perturbation radius is highlighted by a blue bounding box, indicating that *the perturbation budgets required for distinct images are different*. Note that the “human-imperceptible” constraint is violated when the ℓ_∞ -norm perturbation radius is too large the fifth and seventh columns in the top sub-figure. However, our Geometry-Aware (GA) framework generates transferable unrestricted adversarial examples (highlighted by red bounding boxes) with lower budgets and smaller changes when compared to the benign images.

are typically generated via *transfer-based* methods (Szegedy et al., 2014; Papernot et al., 2016, 2017), e.g., attacking an ensemble of accessible source models and hoping that the same adversarial examples are able to fool the unknown target/test model (Liu et al., 2017; Tramer et al., 2018).

Despite a large amount of work on transfer-based attacks, many fundamental questions remain unresolved. For example, existing transfer-based attacks (Dong et al., 2018; Xie et al., 2019b; Dong et al., 2019) that search for adversarial examples in a fixed-radius ℓ_∞ -norm ball often require a high perturbation budget to reach a satisfactory transfer success rate. However, such perturbations might be perceptible to humans (see Figs. 1 and 7). On the other hand, unrestricted adversarial attacks that aim to generate minimum human-imperceptible perturbations (Xiao et al., 2018; Wong et al., 2019; Laidlaw et al., 2021) suffer from weaker transferability to the target model. This is in part due to the difference between the decision boundaries of the source and target models. Given the trade-off between transferability and imperceptibility, one of the long-standing questions is generating transferable adversarial examples by minimum changes of natural examples.

1.1 Our methodology and results

In this work, we propose a novel geometry-aware framework, into which fixed-budget methods can be integrated to generate transferable unrestricted adversarial examples with minimum changes.

Our intuition is that the smallest perturbation budgets w.r.t. distinct images should be different (see Fig. 1) and should depend on their geometrical relationship with the decision boundary of the target model (see Fig. 2). Unfortunately, finding transferable minimum-budget adversarial perturbations is an intractable optimization problem (see Eq (4.1)) as the target model is unknown. We approximately solve this problem by discretizing the continuous space of perturbation radius into a finite set and choosing the minimum perturbation budget that is able to fool the test model. The main challenge here is to evaluate whether a given perturbation can transfer well to the unknown target model (Cheng et al., 2019b; Katzir and Elovici, 2021).

To overcome this challenge, we split all accessible white-box source models into training and validation sets, where adversarial perturbations are crafted only on the training set. We use the validation set to select the smallest perturbation radius for each input that suffices to fool the validation model with a certain confidence level through an *early-stopping* mechanism. When the training (or validation) set consists of multiple models, we use their average ensemble (Liu et al., 2017). Experimentally, our method yields a significant performance boost on the trade-off between transferability and imperceptibility. As shown in Fig. 4, the transfer success rate of our method GA-DTMI-FGSM surpasses the baseline DTMI-FGSM (see Eq (3.1)) by up to 16% in absolute value under the same average perturbation reward (see Eq (3.4)). Besides, our method GA-DMI-FSA is able to generate semantic-preserving yet transferable adversarial examples under the unrestricted threat model (see Eq (3.2), Figs. 1, 6, 7 and 9).

1.2 Summary of our contributions

- We propose a Geometry-Aware (GA) framework, where fixed-budget attacking methods can be integrated, to generate transferable unrestricted adversarial examples with approximately minimum changes. To the best of our knowledge, we are the first to explore transfer-based black-box attacks with adaptive perturbation budgets.
- Under ℓ_∞ -norm setting, our geometry-aware framework improves the imperceptibility of the crafted adversarial examples by a large margin without the decrease of transfer success rate (see Fig. 4). By applying our method GA-DTMI-FGSM to the *CVPR’21 Security AI Challenger: Unrestricted Adversarial Attacks on ImageNet* (Chen et al., 2021), we ranked 1st place out of 1,559 teams and surpassed the runner-up submissions by 4.59% and 23.91% in terms of final score and average image quality level, respectively.
- Under unrestricted setting, we propose a transfer-based unrestricted attack (see Eq (3.2)) by combining the *white-box* feature space attack (Xu et al., 2021) with the transfer-based ℓ_∞ -norm attacks to generate semantic-preserving yet transferable adversarial examples (see Figs. 1 and 6). Moreover, the crafted adversarial examples can even transfer to adversarially trained robust models (see Tab. 3, Figs. 7 and 9).

2 Related Work

ℓ_p -norm adversarial examples. Existing gradient-based white-box attacks either search for adversarial examples in a fixed ℓ_p -norm ball (Madry et al., 2018; Kurakin et al., 2019), or optimize

the perturbation for each image independently to get a *minimum-norm* solution such as DeepFool (Moosavi-Dezfooli et al., 2016), CW attack (Carlini and Wagner, 2017), and fast adaptive boundary attack (Croce and Hein, 2020). However, white-box assumption usually does not hold in real-world scenarios. In query-based black-box threat model, attackers utilize output logits (Chen et al., 2017; Andriushchenko et al., 2020) or predicted label (Brendel et al., 2018; Cheng et al., 2019a) of the target model to generate adversarial examples. But these query-based attacks typically suffer from high query complexity, making it easy to be detected (Willmott et al., 2021). Transfer-based black-box attacks can pose serious threats in practice as they need no information about the defense models. Dong et al. (2018) boosted transferability by integrating momentum into gradient-based methods. Liu et al. (2017) found that attacking a group of substituted source models simultaneously can improve transferability. Besides, transferability benefits from input transformations such as input diversity (Xie et al., 2019b) and translation-invariant method (Dong et al., 2019).

Unrestricted adversarial examples. The ℓ_p -norm distance is not an ideal perceptual similarity metric (Johnson et al., 2016; Isola et al., 2017), which oversimplifies the diversity of real-world perturbations. Unrestricted adversarial examples have received significant attention in recent years (Brown et al., 2018). Most of the current unrestricted attacks aim to generate imperceptible adversarial examples under *white-box* setting, such as geometric transformations (Xiao et al., 2018; Alaifari et al., 2019; Engstrom et al., 2019) and distance metrics beyond ℓ_p norm (Wong et al., 2019; Laidlaw et al., 2021). Color-based attacks (Hosseini and Poovendran, 2018; Laidlaw and Feizi, 2019; Zhao et al., 2020a,b; Shamsabadi et al., 2020; Bhattad et al., 2020) were also proposed to generate large but imperceptible perturbations, however, the modified color can sometimes be unnatural. Instead of optimizing in the input space directly, generative approaches (Song et al., 2018; Gowal et al., 2020; Qiu et al., 2020; Wong and Kolter, 2021) search for adversarial embeddings in the latent space. Style transfer (Prabhu et al., 2018; Bhattad et al., 2020) is inherently an unrestricted attack as it preserves the semantic of the content image. However, constructing transferable unrestricted adversarial examples under black-box settings is still less explored. In this work, we will fill this gap by combining the white-box feature space attack (Xu et al., 2021) with the transfer-based ℓ_∞ -norm attacks to generate semantic-preserving yet transferable adversarial examples, which is directly an extension of our previous technical report (Liu et al., 2021a) that mainly focus on ℓ_∞ -norm setting.

Adversarial defenses. There have been long-standing arms races between defenders and attackers. *Adversarial training* (Goodfellow et al., 2015) is one of the most promising defense methods. Many variants of adversarial training framework were proposed, e.g., ensemble adversarial training (Liu et al., 2017) for transfer-based attacks, PGD-based adversarial training (Madry et al., 2018), and TRADES (Zhang et al., 2019) with a new robust loss based on the trade-off between robustness and accuracy. Geometry-aware instance-reweighted adversarial training (Zhang et al., 2021), which is proven falling into gradient masking (Hitaj et al., 2021), shares similar insights with us that the importance of distinct inputs in adversarial training should be different. Laidlaw et al. (2021) integrate adversarial training with Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018), aiming to improve robustness against perturbations that were unseen during training. Unlike empirical defenses, Certified defenses (Cohen et al., 2019; Salman et al., 2019; Zhang et al., 2020; Leino et al., 2021) could provide robustness guarantee under a certain ℓ_p -norm budget.

3 Preliminaries

Notation. A deep neural network classifier can be described as a function $f(\mathbf{x}; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathbb{R}^C$, parameterized by weights $\boldsymbol{\theta}$, which maps a vector $\mathbf{x} \in \mathcal{X}$ to its output logits. Given an input \mathbf{x} of class $y \in \{1, 2, \dots, C\}$, the predicted label of $f(\mathbf{x}; \boldsymbol{\theta})$ is $\hat{f}(\mathbf{x}) := \arg \max_j f_j(\mathbf{x}; \boldsymbol{\theta})$, where $f_j(\mathbf{x}; \boldsymbol{\theta})$ represents the j -th entry of $f(\mathbf{x}; \boldsymbol{\theta})$. We use $L(f(\mathbf{x}; \boldsymbol{\theta}), y)$ to represent the cross-entropy loss and denote the ε -neighborhood of \mathbf{x} by $\mathbb{B}(\mathbf{x}, \varepsilon) := \{\mathbf{x}' \in \mathcal{X} : \mathcal{D}(\mathbf{x}, \mathbf{x}') \leq \varepsilon\}$, where \mathcal{D} is a distance metric that describes the change between the adversarial example \mathbf{x}' and the nature example \mathbf{x} . We denote the black-box test model by g , and split the set of accessible source models into the set of training models f and the set of validation models h .

3.1 Transfer-based ℓ_∞ -norm attacks

Existing transfer-based attacks typically search for adversarial examples in a fixed-radius ℓ_p -norm ball, i.e., $\mathcal{D}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x}' - \mathbf{x}\|_p \leq \varepsilon$. Various methods were proposed to boost transferability of the generated adversarial examples, such as input Diversity Iterative Fast Gradient Sign Method (DI-FGSM) (Xie et al., 2019b), Momentum-based Iterative (MI-FGSM) method (Dong et al., 2018) and Translation-invariant Iterative (TI-FGSM) method (Dong et al., 2019). We formulate a strong ℓ_∞ -norm baseline DTMI-FGSM by combining all these techniques, i.e.,

$$\begin{aligned} \mathbf{m}_{t+1} &= \gamma \cdot \mathbf{m}_t + \frac{\mathbf{W} * \nabla_{\mathbf{x}_t} L(f(T(\mathbf{x}_t, p); \boldsymbol{\theta}), y)}{\|\mathbf{W} * \nabla_{\mathbf{x}_t} L(f(T(\mathbf{x}_t, p); \boldsymbol{\theta}), y)\|_1}, \\ \mathbf{x}_{t+1} &= \Pi_{\mathbb{B}(\mathbf{x}, \varepsilon)}(\mathbf{x}_t + \alpha \cdot \text{sign}(\mathbf{m}_{t+1})), \end{aligned} \quad (3.1)$$

where $\mathbf{m}_0 = \mathbf{0}$, \mathbf{W} is a pre-defined kernel with a convolution operation $*$, α is the step size, Π is the projection operator, and γ is the decay factor for the momentum term. $T(\mathbf{x}_t, p)$ represents the input transformation on \mathbf{x}_t with probability p . When $\gamma = 0$, DTMI-FGSM attack degenerates to the DTI-FGSM attack. When $p = 0$, DTMI-FGSM attack degenerates to the DMI-FGSM attack.

3.2 Transfer-based unrestricted attack

Inspired by prior work (Huang and Belongie, 2017) in style transfer, Xu et al. (2021) tries to find stylized adversarial examples by assuming that the image pairs from the same class share consistent content and differ mainly in their styles. Here we propose to generate semantic-preserving yet transferable unrestricted adversarial examples by combining the Feature Space Attack (FSA) (Xu et al., 2021) with transfer-based ℓ_∞ -norm attacks (Dong et al., 2018; Xie et al., 2019b). Given an encoder ϕ , we extract the style features of input \mathbf{x} as channel-wise mean $\boldsymbol{\mu}(\phi(\mathbf{x})) \in \mathbb{R}^C$ and channel-wise standard deviation $\boldsymbol{\sigma}(\phi(\mathbf{x})) \in \mathbb{R}^C$. Specifically,

$$\mu_c = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W \phi_c(\mathbf{x})_{hw}, \quad \sigma_c = \sqrt{\frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (\phi_c(\mathbf{x})_{hw} - \mu_c)^2},$$

where $\phi(\mathbf{x}) \in \mathbb{R}^{C \times H \times W}$ represents the latent embedding. Xu et al. (2021) adds adversarial perturbations on $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ before projecting $\phi(\mathbf{x})$ back to the input space \mathcal{X} with a pre-trained¹

¹We use the official pre-trained shallowest decoder: <https://github.com/qiulingxu/FeatureSpaceAttack>.

decoder ϕ^{-1} , namely,

$$\begin{aligned}\tilde{\phi}(\mathbf{x}) &= (e^{\tau^\sigma} \cdot (\phi(\mathbf{x}) - \boldsymbol{\mu}) + e^{\tau^\mu} \cdot \boldsymbol{\mu},) \\ \mathbf{x}' &= \phi^{-1}(\tilde{\phi}(\mathbf{x})), \quad \|\boldsymbol{\tau}^\mu\|_\infty \leq \ln \varepsilon, \|\boldsymbol{\tau}^\sigma\|_\infty \leq \ln \varepsilon,\end{aligned}\tag{3.2}$$

where $\tilde{\phi}(\mathbf{x})$ enlarges or shrinks the mean $\boldsymbol{\mu}$ and the standard deviation $\boldsymbol{\sigma}$ of the embedding $\phi(\mathbf{x})$ by a factor of e^{τ^μ} and e^{τ^σ} , respectively. In this way, the distance metric $\mathcal{D}(\mathbf{x}, \mathbf{x}') = \max(e^{\|\boldsymbol{\tau}^\mu\|_\infty}, e^{\|\boldsymbol{\tau}^\sigma\|_\infty}) \leq \varepsilon$. In order to preserve the semantic of the unrestricted adversarial example \mathbf{x}' , an additional content loss was added during the attacking process, i.e.,

$$\min_{\boldsymbol{\tau}^\mu, \boldsymbol{\tau}^\sigma} \mathcal{L}(\mathbf{x}', y) = \lambda \cdot \mathcal{L}_{\text{top-5}}(f(\mathbf{x}'; \boldsymbol{\theta}), y) + \|\phi(\mathbf{x}') - \tilde{\phi}(\mathbf{x})\|_2,$$

where λ balance the trade-off between adversarial loss and the content loss. Following Xu et al. (2021), we set $\lambda = 128$ and use the top-5 margin loss for adversarial attack. With all above, the unrestricted attack (see Eq (3.2)) can be solved by conventional ℓ_∞ -norm attack on parameters $\boldsymbol{\tau}^\mu$ and $\boldsymbol{\tau}^\sigma$. Moreover, the same techniques in Sec. 3.1 such as input diversity $T(\mathbf{x}', p)$ (only for the margin loss) and momentum-based method can be integrated to improve transferability, i.e.,

$$\begin{aligned}\mathbf{m}_{t+1} &= \gamma \cdot \mathbf{m}_t + \frac{\nabla_{\boldsymbol{\tau}_t} \mathcal{L}(T(\mathbf{x}', p), y)}{\|\nabla_{\boldsymbol{\tau}_t} \mathcal{L}(T(\mathbf{x}', p), y)\|_1}, \\ \boldsymbol{\tau}_{t+1} &= \Pi_{\mathbb{B}(\mathbf{x}, \varepsilon)}(\boldsymbol{\tau}_t - \alpha \cdot \text{sign}(\mathbf{m}_{t+1})),\end{aligned}\tag{3.3}$$

where $\boldsymbol{\tau}$ represents $\boldsymbol{\tau}^\mu$ or $\boldsymbol{\tau}^\sigma$. When $\gamma = 0$, the DMI-FSA attack degenerates to the DI-FSA attack.

3.3 Evaluation metric for transfer-based attack

Different from transferability, the imperceptibility of adversarial examples is hard to evaluate due to the lack of precise quantization of human perception (Wang et al., 2004). Sharif et al. (2018) found that the ℓ_p -norm distance is not an ideal perceptual similarity metric and suggest setting adaptive perturbation budget for every sample to ensure that the attacks' output would be imperceptible. Therefore, we choose transfer success rate and the perturbation budget under certain distance metric \mathcal{D} as our main evaluation metrics. Consider a dataset $\hat{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and the corresponding adversarial examples $\hat{S}_{\text{adv}} = \{(\mathbf{x}'_i, y_i)\}_{i=1}^n$ that are crafted on the training model f . Let $n_0 = \sum_{i=1}^n \mathbb{1}\{\hat{g}(\mathbf{x}'_i) \neq y_i\}$ be the number of misclassified adversarial examples on the test model g . We define the average total score as follows:

$$\begin{aligned}S_{\text{total}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\hat{g}(\mathbf{x}'_i) \neq y_i\} \cdot \mathcal{F}_{\text{reward}}(\mathcal{D}(\mathbf{x}_i, \mathbf{x}'_i)) \\ &= \frac{n_0}{n} \cdot \frac{1}{n_0} \sum_{i=1}^n \mathbb{1}\{\hat{g}(\mathbf{x}'_i) \neq y_i\} \cdot \mathcal{F}_{\text{reward}}(\mathcal{D}(\mathbf{x}_i, \mathbf{x}'_i)) \\ &\stackrel{\text{def}}{=} \frac{n_0}{n} \cdot S_{APR},\end{aligned}\tag{3.4}$$

where S_{APR} is the Average Perturbation Reward of adversarial examples that are misclassified by test model g and the reward function $\mathcal{F}_{\text{reward}}$ is a decreasing function w.r.t. the metric $\mathcal{D}(\mathbf{x}, \mathbf{x}')$.

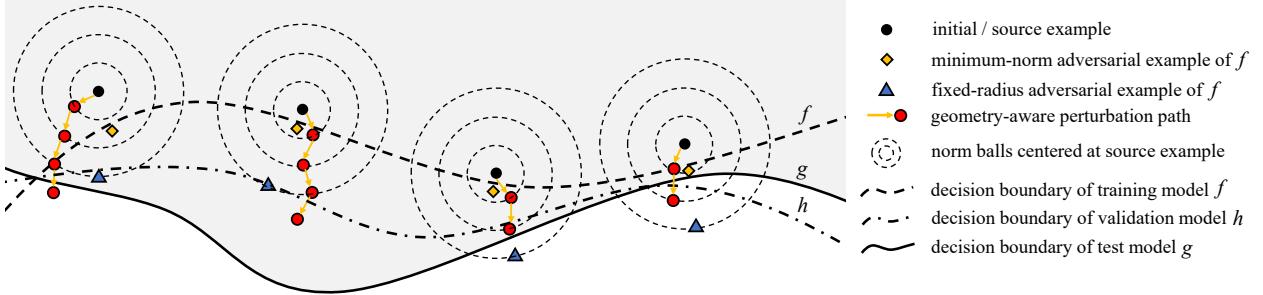


Figure 2: **Our geometry-aware framework.** Existing fixed-budget methods typically overlook the importance of geometrical distances from inputs to the decision boundary of the test model g . In contrast, our geometry-aware framework aims to find geometry-aware minimal-change perturbation via a validation model h . The goal of the validation model is to prevent an attack algorithm overfitting f by forcing the solution to cross the decision boundary of h with a certain margin. Our framework consists of multiple sub-procedures with adaptive perturbation budgets. In each sub-procedure, we start from the solution of the last sub-procedure (the red solid points) and re-run the attack algorithm on the training model f . The procedure stops if the output probability of the true class on the validation model h is smaller than a certain threshold η .

4 Methodology: Geometry-Aware Framework

Eq (3.4) factorizes the average total score as the product of transfer success rate and average perturbation reward, which motivates us to find the adversarial example with minimum changes under metric \mathcal{D} , i.e.,²

$$\min_{\mathbf{x}'} \mathcal{D}(\mathbf{x}, \mathbf{x}') \quad \text{s.t. } \hat{g}(\mathbf{x}') \neq y. \quad (4.1)$$

However, direct optimization of problem (4.1) is intractable, in part due to the lack of information about test model g . We approximately solve this problem by discretizing the continuous space of perturbation radius into a discrete set and choosing the minimum perturbation budget such that the attack is able to fool the test model g . However, the challenge is that it is typically difficult to decide whether a given perturbation radius can also fool the test model (Cheng et al., 2019b; Katzir and Elovici, 2021). This problem is also known as model selection problem, and a classic approach to tackle this problem is to use a validation model h to select the right perturbation radius. More specifically, we split all accessible source models into training model set and validation model set. With validation model h , we are able to generate transferable adversarial examples with smaller changes. To approximately solve problem (4.1), we first divide the attack in the ball $\mathbb{B}(\mathbf{x}, \varepsilon)$ into K sub-procedures. In the k -th sub-procedure, we re-run a fixed-radius attack algorithm \mathcal{A} such as DMI-FSA (see Eq (3.3)) under the perturbation budget

$$\varepsilon_k = \frac{k}{K} \times \varepsilon, \quad k = 1, 2, \dots, K.$$

We start each sub-procedure from the solution of the last sub-procedure to accelerate the convergence. To obtain a minimum-radius solution, we perform an early-stopping mechanism at the end of each

²In contrast to existing white-box minimum-norm attack (Carlini and Wagner, 2017), g is an unknown model here.

sub-procedure if the probability of the true class on the validation model h is smaller than a certain threshold η , i.e.,

$$\mathbb{P}(\hat{h}(\mathbf{x}) = y) = \frac{\exp(h_y(\mathbf{x}; \boldsymbol{\theta}))}{\sum_j \exp(h_j(\mathbf{x}; \boldsymbol{\theta}))} < \eta. \quad (4.2)$$

Our Geometry-Aware (GA) framework is summarized in Algorithm 1 and illustrated in Fig. 2.

Algorithm 1 Geometry-Aware Framework

Require:

Benign input \mathbf{x} with label y ; training models f ; validation model h ; number of sub-procedures K ; maximum perturbation size ε and threshold η ; attack algorithm \mathcal{A} ;

Ensure:

Transfer-based unrestricted adversarial example \mathbf{x}' with approximately minimum change;

```

1:  $\mathbf{x}_0 = \mathbf{x}$ ;
2: for  $k = 1, 2, \dots, K$  do
3:    $\mathbf{x}_k = \mathcal{A}(\mathbf{x}, \mathbf{x}_{k-1}, f, \frac{k\varepsilon}{K})$ ;            $\triangleright$  fixed-budget attack
4:   conf  $\leftarrow \frac{\exp(h_y(\mathbf{x}_k; \boldsymbol{\theta}))}{\sum_j \exp(h_j(\mathbf{x}_k; \boldsymbol{\theta}))}$ ;
5:   if conf  $< \eta$  then
6:     Return  $\mathbf{x}_k$ ;            $\triangleright$  early-stopping in Eq (4.2)
7:   end if
8: end for
9: Return  $\mathbf{x}_K$ ;

```

5 Experiments

5.1 Experimental setup

Datasets & Networks. Similar to Xie et al. (2019b), we randomly select 1,000 images from ILSVRC 2012 validation set (Deng et al., 2009), which are almost correctly classified by all the attacking models. All these images are resized to $229 \times 229 \times 3$ beforehand. We consider eight normally trained models and two ensemble adversarially trained models (Tramer et al., 2018). The weights of all these models are publicly available at Wightman (2019). More details about the networks are summarised in Tab. 1. The transferability between these models under ℓ_∞ -norm setting is summarised in Fig. 3. It is much easier for the generated adversarial examples to transfer from vision transformers to convolutional neural networks (CNNs), which is consistent with the empirical observation in Shao et al. (2021). Surprisingly, the robustness of naturally trained vision transformers under transfer attack is even on par with two ensemble adversarially trained CNNs.

Implementation details. Given the maximum perturbation size ε and number of sub-procedures K (5 as default) in our geometry-aware framework, we set the step size $\alpha = \frac{1.25 \times \varepsilon_k}{T}$ in the k -th sub-procedure, where the number of iterations T is set to 10 in the ℓ_∞ setting and 50 in the

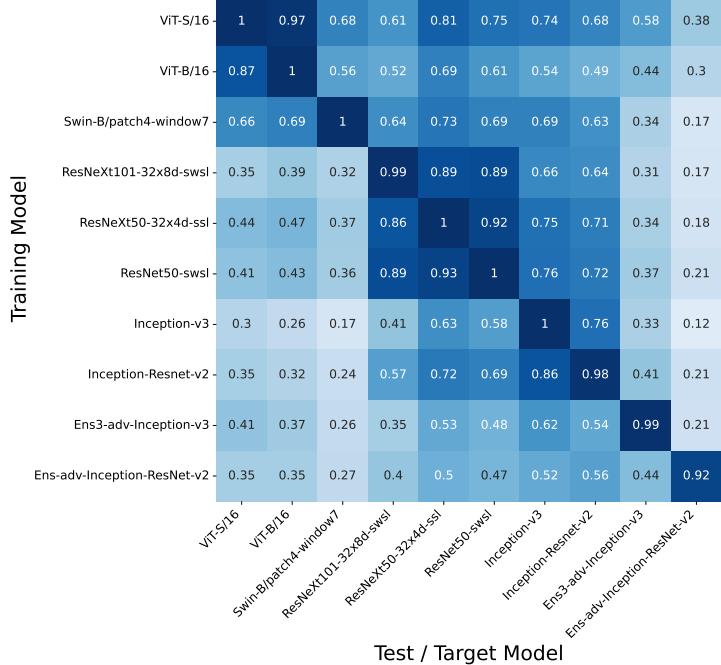


Figure 3: **Transferability between different networks under DTMI-FGSM attack.** The rows stand for source models and the columns stand for target models. Adversarial examples transfer well between models with similar architectures.

Table 1: **An overview of all considered networks for generating adversarial examples.** Top-1_{ImageNet} represents the accuracy on the ILSVRC 2012 validation set while Top-1₁₀₀₀ represents the accuracy on the randomly selected 1000 images.

Training	Index	Model Name	Top-1 _{ImageNet}	Top-1 ₁₀₀₀	Index	Model Name	Top-1 _{ImageNet}	Top-1 ₁₀₀₀
Normal	0	ViT-S/16	76.01%	99.8%	1	ViT-B/16	81.08%	99.1%
	2	Swin-B/patch4-window7	84.23%	99.4%	3	ResNeXt101-32x8d-sws1	83.62%	99.9%
	4	ResNeXt50-32x4d-ssl	78.90%	99.7%	5	ResNet50-sws1	79.97%	99.5%
	6	Inception-v3	76.94%	100.0%	7	Inception-ResNet-v2	79.85%	99.9%
Ensemble	8	Ens3-adv-Inception-v3	76.49%	100.0 %	9	Ens-adv-Inception-ResNet-v2	78.98%	99.9%

unrestricted setting. ε is set to 20 in the ℓ_∞ -norm setting and 3.5 in the unrestricted setting. When running a fixed-radius baseline at perturbation budget ε_k , we set the number of iterations $N = \frac{T}{2} \left(1 + \frac{K\varepsilon_k}{\varepsilon} \right)$ with step size α to keep the same total perturbation budget (the sum of step size across all iterations) as our geometry-aware framework for fair comparison. The reward function $\mathcal{F}_{\text{reward}}(\varepsilon_0)$ is set to $1/\varepsilon_0$ as smaller perturbation radius exhibits significantly higher image quality. For the momentum term, we set the decay factor $\mu = 1$ as in Dong et al. (2018). For DI-FGSM (Xie et al., 2019b), we set the transformation probability to $p = 0.7$. The input is first randomly resized to be an $r \times r \times 3$ image with $r \in [(1 - \gamma)s, (1 + \gamma)s]$, and then padded to size $(1 + \gamma)s \times (1 + \gamma)s \times 3$. The transformed input is then resized to $s \times s \times 3$ for different input size s of various models, i.e., 224, 299 and 384. We set $\gamma = 0.1$ as default. For TI-FGSM (Dong et al., 2019), we use Gaussian kernel with kernel size 5×5 .

5.2 Balancing transfer success rate and average perturbation reward

Implementation details. Benefiting from the adaptive choice of perturbation budgets, our geometry-aware framework can generate transferable unrestricted adversarial examples with smaller changes. In this experiment, the training model f and validation model h are an ensemble of models $\{2, 3, 5\}$ and $\{1, 4, 6\}$ in Tab. 1, respectively. The test model is Inception-ResNet-v2. The optimal threshold η (see Eq (4.2)) is searched from a finite set ranging from 0.001 to 0.9 by querying³ the test model g to achieve the best average total score S_{total} . For each η , we execute our method and compute the average perturbation reward S_{APR} (the x -axis of each red point in Fig. 4). Then the corresponding fixed-radius baseline is run at the same x -axis. We conduct experiments on two threat models. For the ℓ_∞ -norm setting, we combine our Geometry-Aware (GA) framework with DI-FGSM, DTI-FGSM, and DTM-FGSM, named GA-DTI-FGSM, GA-DTI-FGSM, and GA-DTM-FGSM, respectively; For the unrestricted setting, we combine our GA framework with DMI-FSA and DI-FSA, named GA-DMI-FSA and GA-DI-FSA, respectively.

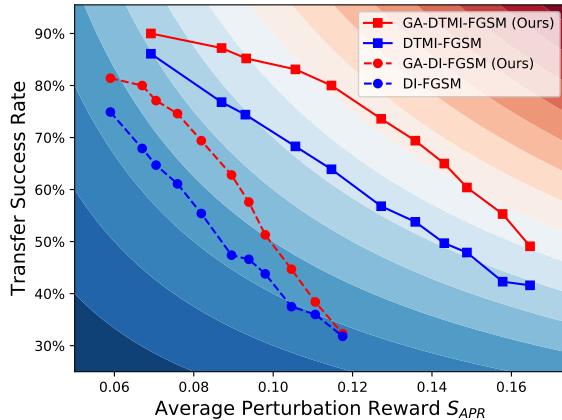


Figure 4: **Contour of average total score S_{total} (see Eq (3.4)).** Fixing transfer success rate as 80%, our approach GA-DTM-FGSM surpasses the baseline DTM-FGSM (see Eq (3.1)) by up to 43.35% in terms of average perturbation reward.

Table 2: **Comparison of our method with baselines.** We report the results when both our approach and baselines achieve highest average total score S_{total} . TSR: Transfer Success Rate.

Method	TSR (\uparrow)	S_{APR} (\uparrow)	S_{total} (\uparrow)
DI-FGSM	61.1%	0.0759	4.64%
GA-DI-FGSM (ours)	69.4%	0.0819	5.68%
DTI-FGSM	57.3%	0.1101	5.68%
GA-DTI-FGSM (ours)	67.9%	0.1176	7.98%
DTMI-FGSM	63.9%	0.1147	7.33%
GA-DTM-FGSM (ours)	69.4%	0.1358	9.42%
DI-FSA	48.3%	0.5328	25.73%
GA-DI-FSA (ours)	50.4%	0.5541	27.93%
DMI-FSA	51.3%	0.5616	28.81%
GA-DMI-FSA (ours)	58.3%	0.5355	31.32%

Experimental results We present the contour plot of average total score in Fig. 4, where the improvement of our method upon baselines depends on the choice of hyper-parameter η (leading to different S_{APR}). Fixing S_{APR} as 0.115, our approach GA-DTM-FGSM surpasses the baseline DTM-FGSM by up to 16.1% in terms of transfer success rate; As shown in Tab. 2, our approach yields a significant performance boost on the average total score S_{total} across various threat models, especially in the ℓ_∞ -norm setting where both the transfer success rate and S_{APR} are improved.

³In contrast to conventional query-based attacks that need the logits or predicted label on the target model, we query whether an adversarial example transfers to the target model successfully. Besides, we have strong prior information on η which depends on the similarity between f, h and g , making the query complexity rather limited.

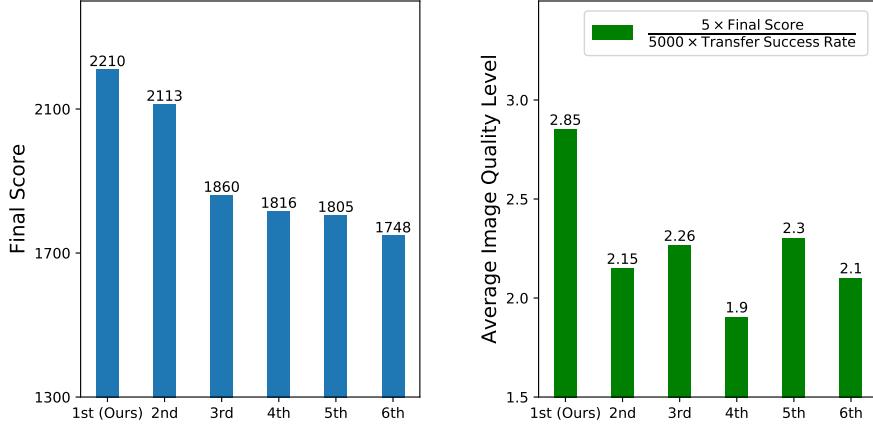


Figure 5: **Top-6 results** (out of 1,559 teams) in the CVPR’21 Security AI Challenger: Unrestricted Adversarial Attacks on ImageNet. The final scores were manually scored by multiple human referees.

5.3 Case study: CVPR’21 Unrestricted Adversarial Attacks on ImageNet

In the *CVPR’21 Security AI Challenger: Unrestricted Adversarial Attacks on ImageNet* (Chen et al., 2021), contestants were asked to submit adversarial examples without any access to the defense models. The dataset is a subset of ILSVRC 2012 validation set (Deng et al., 2009), which consists of 5,000 images with 5 images per class. The final score of each submission was *manually* scored from two aspects: 1) image semantic and 2) quality. If the semantic of the submitted image changes, then $S_s = 0$, otherwise $S_s = 1$. The image quality S_q (equivalent to our reward function $\mathcal{F}_{\text{reward}}$) was quantified with five levels $S_q \in \{1, 2, 3, 4, 5\}$ by multiple human referees. The final score is given by $\sum_i \mathbb{1}\{\hat{g}(\mathbf{x}'_i) \neq y_i\} \times S_s(\mathbf{x}'_i) \times \frac{S_q(\|\mathbf{x}'_i - \mathbf{x}_i\|)}{5}$.

We apply our method GA-DTMI-FGSM ($\eta = 0.01$) to the competition, where our entry ranked 1st place out of 1,559 teams. We report the final score and average image quality level (equivalent to our average perturbation reward) in Fig. 5. It shows that our method outperforms other approaches by a large margin. In particular, we surpass the runner-up submissions by 4.59% and 23.91% in terms of final score and average image quality level, respectively.

5.4 Transferable unrestricted adversarial examples beyond ℓ_p norm

Most of current defenses can be easily broken by unseen attacks in a white-box manner. Adversarial training against multiple ℓ_p -norm attacks (Tramer and Boneh, 2019) solved this issue partially, however, at the cost of robustness against single ℓ_p -norm attack. Laidlaw et al. (2021) integrated adversarial training with Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018), aiming to improve robustness against perturbations that were unseen during training. However, the proposed attack (Laidlaw et al., 2021), similar to other unrestricted attacks (Laidlaw and Feizi, 2019; Engstrom et al., 2019), suffers from weaker transferability to the target model. In practice, attackers typically have no information about the defense models (thus transfer attack is required) and the defenders do not have the ground truth to make pixel-level comparison (perturbation can be large as long as the generated adversarial examples are semantic-preserving). Therefore, we propose to benchmark classification models on ImageNet under *transfer-based unrestricted attacks*.



Figure 6: **Adversarial examples of our method GA-DMI-FSA**, which are misclassified by *all* the models in Tab. 3. Top: benign examples \mathbf{x} . Middle: unrestricted adversarial examples \mathbf{x}' . Bottom: normalized adversarial perturbations $\mathbf{x}' - \mathbf{x}$.

Table 3: **Benchmarking classification on Imagenet under transfer-based unrestricted attacks.** PGD₄₀^{*} indicates the white-box PGD attack (Madry et al., 2018) with 40 iterations ($\varepsilon = 4/255$). We denote the adversarial examples generated by GA-DTMI-FGSM and GA-DMI-FSA in Table 2 as GA_{DTMI-FGSM} and GA_{DMI-FSA}, respectively. FSA represents the feature space attack (Xu et al., 2021). ReColor represents the color-based attack (Laidlaw and Feizi, 2019).

Defenses	Clean	PGD ₄₀ [*]	ReColor	FSA	GA _{DTMI-FGSM}	GA _{DMI-FSA}
Inception-ResNet-v2 _{Ens-adv} (Tramer et al., 2018)	99.9%	10.0%	93.2%	90.0%	87.1%	43.5%
FastAT (Wong et al., 2020)	66.9%	35.7%	62.8%	59.3%	64.7%	28.5%
FreeAT (Shafahi et al., 2019)	77.3%	40.6%	71.6%	68.4%	74.0%	35.1%
Resnet152-Base (Xie et al., 2019a)	67.6%	39.0%	64.1%	61.2%	65.1%	37.4%
Resnext101-DenoiseAll (Xie et al., 2019a)	80.3%	52.2%	77.0%	73.8%	78.7%	47.8%
Resnet152-Denoise (Xie et al., 2019a)	72.2%	41.8%	68.2%	65.2%	70.7%	40.3%
RVT-Tiny (Mao et al., 2021)	96.7%	0.0%	78.9%	81.3%	44.5%	33.6%
DeepAugment+AugMix (Hendrycks et al., 2021)	96.1%	0.0%	82.8%	89.3%	58.9%	63.2%
Efficientnet-l2-ns (Xie et al., 2020)	99.2%	0.0%	95.7%	97.0%	81.3%	86.0%
Swin-L/patch4-window-12 (Liu et al., 2021b)	99.1%	0.0%	88.4%	90.7%	66.6%	61.8%

Implementation details For the ReColor attack (Laidlaw and Feizi, 2019), we set $\varepsilon = 1.0$ and iterations $T = 100$ which achieves 89.7% attack success rate on training model f (the same as Sec. 5.2) and 9.1% transfer success rate on the test model Inception-ResNet-v2. For FSA attack (Xu et al., 2021), we set $\varepsilon = 3.0$ and $T = 500$ which achieves 54.2%⁴ attack success rate and 12.4% transfer success rate on the same training and test models. Besides six adversarially trained and two high-performance classification models, we select two state-of-the-art models on the ImageNet-R dataset (Hendrycks et al., 2021). Note that it’s hard to compare the robustness across different attacks (columns in Tab. 3) due to the lack of a consistent perceptual metric.

⁴Note that our method GA-DMI-FSA achieves 95.5% attack success rate and 58.3% transfer success rate, indicating that the input diversity and momentum techniques in Eq (3.3) boost both the attacking ability and transferability.

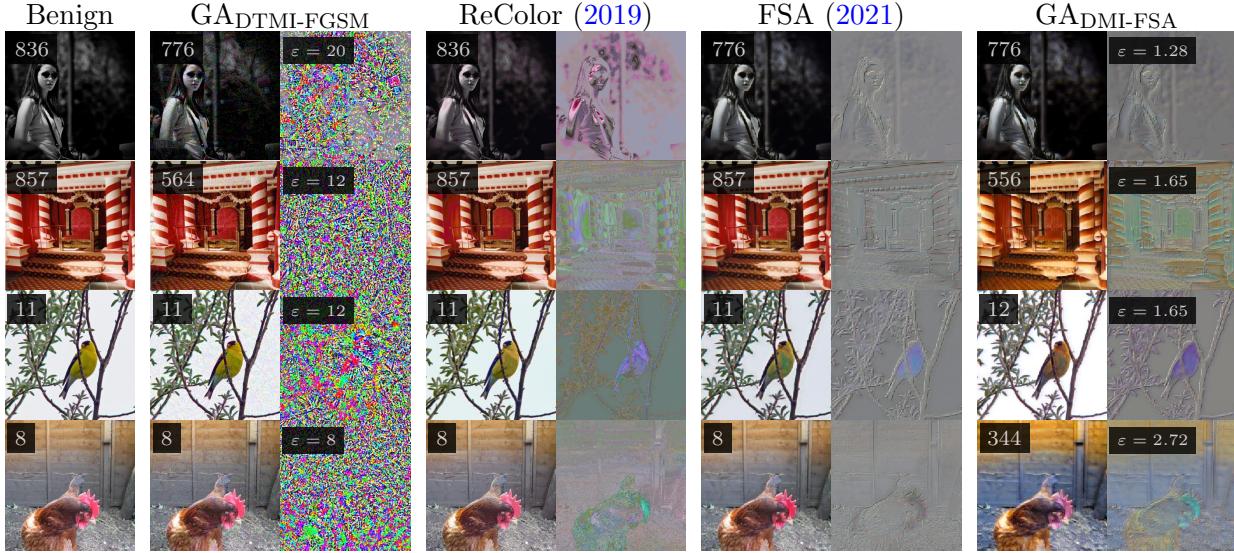


Figure 7: **Visualization of transfer attack results on Resnext101-DenoiseAll** ([Xie et al., 2019a](#)). For each image, we print its predicted label on model Resnext101-DenoiseAll in the upper left corner. For each transfer-based adversarial example, we present the perturbation on its right. For each perturbation crafted via our geometry-aware framework, we print its perturbation budget in the upper left corner. Although the transfer-based ℓ_∞ -norm attack GA-DTMI-FGSM is able to fool the defense test model to a certain extent, the generated perturbations can be “human-perceptible” (the first and third rows of GA-DTMI-FGSM). Besides, the other two unrestricted attacks suffer from weaker transferability when compared to our method GA-DMI-FSA, which adjusts the images’ color and texture that ImageNet-trained CNNs might be biased to ([Geirhos et al., 2018](#)).

Experimental results From Tab. 3, we can conclude the following three observations: 1) ℓ_∞ -norm based transfer attack (GA-DTMI-FGSM) can hardly break ℓ_∞ -norm adversarially trained models while the unrestricted attack (GA-DMI-FSA) reduces the accuracy of these models by a large margin (see also in Fig. 7). 2) DeepAugment ([Hendrycks et al., 2021](#)), which utilizes semantic-preserving augmentations during training, exhibits non-trivial robustness against GA-DMI-FSA attack. 3) Efficientnet-l2-ns ([Xie et al., 2020](#)) performs well under all the transfer-based attacks and enjoys 86% accuracy against GA-DMI-FSA attack, showing that the distribution of our generated adversarial examples is not too far from the natural examples’. Note that Efficientnet-l2-ns is also the best-performing model on the ImageNet-V2 dataset ([Recht et al., 2019; Taori et al., 2020](#)).

We visualize part of the transfer attack results on adversarially trained model Resnext101-DenoiseAll ([Xie et al., 2019a](#)) in Figs. 7 and 9, where our method GA-DMI-FSA is able to generate semantic-preserving yet transferable adversarial examples.

5.5 Ablation studies

The optimal η depends on the combination of f, h, g . Another natural question is how to choose the training model f and validation model h . As illustrated in Fig. 2, the effectiveness of our geometry-aware framework relies on the relationship between the decision boundaries of training

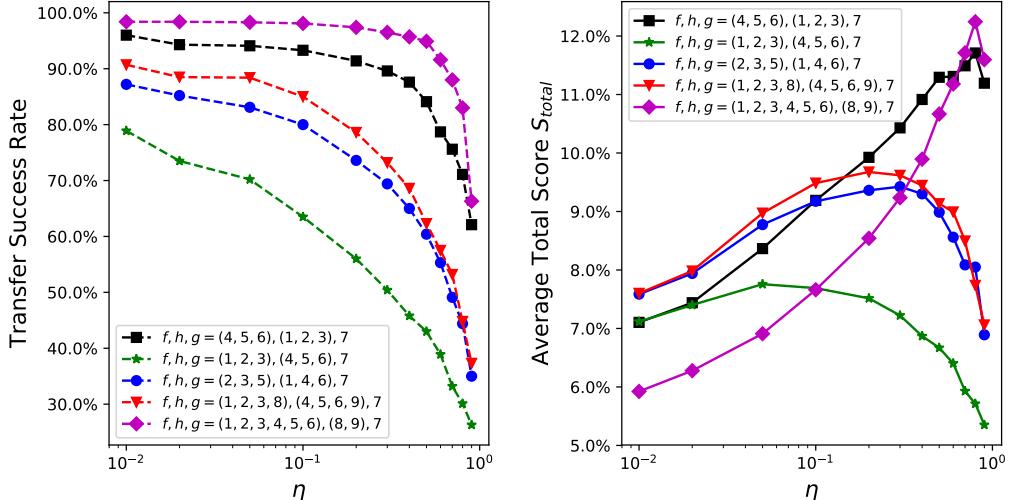


Figure 8: **Comparison between various kinds of partition for training, validation, and test models.** The index in the legend corresponds to the model index in Tab. 1. As the threshold η for early-stopping increases, the generated adversarial examples have higher confidence (probability of the true class) on the validation model h , leading to a lower transfer success rate (**left**). Besides, the optimal η that yields the maximum average total score is dependent on the partition (**right**).

model f , validation model h , and test model g . Thus we compared various kinds of partition of training, validation, and test models in Fig. 8, from which we can conclude the following two observations: 1) The role of training model f and validation model h is not *symmetrical*. The transfer success rate changed a lot when we exchanged the role of models $\{4, 5, 6\}$ and $\{1, 2, 3\}$ for training and validation. 2) Ensembling more models yields better performance. When additional models were added into the training and validation sets, the average total score improved significantly.

6 Conclusion

In this work, we propose a geometry-aware framework, where fixed-radius methods can be integrated to generate transferable unrestricted adversarial examples with minimum changes. Under ℓ_∞ -norm setting, our framework improves the imperceptibility of the crafted adversarial examples by a large margin without the decrease of transfer success rate. Besides, we propose a transfer-based unrestricted attack by combining the white-box feature space attack with transfer-based ℓ_∞ -norm attacks to generate semantic-preserving yet transferable adversarial examples.

References

- ALAIFARI, R., ALBERTI, G. S. and GAUKSSON, T. (2019). ADef: an iterative algorithm to construct adversarial deformations. In *ICLR*.

- ANDRIUSHCHENKO, M., CROCE, F., FLAMMARION, N. and HEIN, M. (2020). Square attack: A query-efficient black-box adversarial attack via random search. In *ECCV*.
- ATHALYE, A., CARLINI, N. and WAGNER, D. (2018). Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, vol. 1.
- BAI, Y., MEI, J., YUILLE, A. and XIE, C. (2021). Are transformers more robust than cnns? In *NeurIPS*.
- BHATTAD, A., CHONG, M. J., LIANG, K., LI, B. and FORSYTH, D. A. (2020). Unrestricted adversarial examples via semantic manipulation. In *ICLR*.
- BHOJANAPALLI, S., CHAKRABARTI, A., GLASNER, D., LI, D., UNTERTHINER, T. and VEIT, A. (2021). Understanding robustness of transformers for image classification. In *ICCV*.
- BOJARSKI, M., DEL TESTA, D., DWORAKOWSKI, D., FIRNER, B., FLEPP, B., GOYAL, P., JACKEL, L. D., MONFORT, M., MULLER, U., ZHANG, J. ET AL. (2016). End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316* .
- BRENDEL, W., RAUBER, J. and BETHGE, M. (2018). Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *ICLR*.
- BROWN, T. B., CARLINI, N., ZHANG, C., OLSSON, C., CHRISTIANO, P. and GOODFELLOW, I. (2018). Unrestricted adversarial examples. *arXiv preprint arXiv:1809.08352* .
- CARLINI, N. and WAGNER, D. (2017). Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*. IEEE.
- CHEN, P.-Y., ZHANG, H., SHARMA, Y., YI, J. and HSIEH, C.-J. (2017). Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*.
- CHEN, Y., MAO, X., HE, Y., XUE, H., LI, C., DONG, Y., FU, Q.-A., YANG, X., XIANG, W., PANG, T. ET AL. (2021). Unrestricted adversarial attacks on imagenet competition. *arXiv preprint arXiv:2110.09903* .
- CHENG, M., LE, T., CHEN, P.-Y., ZHANG, H., YI, J. and HSIEH, C.-J. (2019a). Query-efficient hard-label black-box attack: An optimization-based approach. In *ICLR*.
- CHENG, S., DONG, Y., PANG, T., SU, H. and ZHU, J. (2019b). Improving black-box adversarial attacks with a transfer-based prior. In *NeurIPS*, vol. 32.
- COHEN, J., ROSENFELD, E. and KOLTER, Z. (2019). Certified adversarial robustness via randomized smoothing. In *ICML*. PMLR.
- CROCE, F. and HEIN, M. (2020). Minimally distorted adversarial examples with a fast adaptive boundary attack. In *ICML*. PMLR.
- DENG, J., DONG, W., SOCHER, R., LI, L.-J., LI, K. and FEI-FEI, L. (2009). Imagenet: A large-scale hierarchical image database. In *CVPR*. IEEE.

- DONG, Y., LIAO, F., PANG, T., SU, H., ZHU, J., HU, X. and LI, J. (2018). Boosting adversarial attacks with momentum. In *CVPR*.
- DONG, Y., PANG, T., SU, H. and ZHU, J. (2019). Evading defenses to transferable adversarial examples by translation-invariant attacks. In *CVPR*.
- DOSOVITSKIY, A., BEYER, L., KOLESNIKOV, A., WEISSENBORN, D., ZHAI, X., UNTERTHINER, T., DEHGHANI, M., MINDERER, M., HEIGOLD, G., GELLY, S., USZKOREIT, J. and HOULSBY, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*.
- ENGSTROM, L., TRAN, B., TSIPRAS, D., SCHMIDT, L. and MADRY, A. (2019). Exploring the landscape of spatial robustness. In *ICML*. PMLR.
- GEIRHOS, R., RUBISCH, P., MICHAELIS, C., BETHGE, M., WICHMANN, F. A. and BRENDL, W. (2018). Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231* .
- GOODFELLOW, I. J., SHLENS, J. and SZEGEDY, C. (2015). Explaining and harnessing adversarial examples. In *ICLR*.
- GOWAL, S., QIN, C., HUANG, P.-S., CEMGIL, T., DVIJOTHAM, K., MANN, T. and KOHLI, P. (2020). Achieving robustness in the wild via adversarial mixing with disentangled representations. In *CVPR*.
- HE, K., ZHANG, X., REN, S. and SUN, J. (2016). Deep residual learning for image recognition. In *CVPR*.
- HENDRYCKS, D., BASART, S., MU, N., KADAVATH, S., WANG, F., DORUNDO, E., DESAI, R., ZHU, T., PARAJULI, S., GUO, M., SONG, D., STEINHARDT, J. and GILMER, J. (2021). The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV* .
- HITAJ, D., PAGNOTTA, G., MASI, I. and MANCINI, L. V. (2021). Evaluating the robustness of geometry-aware instance-reweighted adversarial training. *arXiv preprint arXiv:2103.01914* .
- HOSSEINI, H. and POOVENDRAN, R. (2018). Semantic adversarial examples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*.
- HUANG, X. and BELONGIE, S. (2017). Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*.
- ISOLA, P., ZHU, J.-Y., ZHOU, T. and EFROS, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *CVPR*.
- JOHNSON, J., ALAHI, A. and FEI-FEI, L. (2016). Perceptual losses for real-time style transfer and super-resolution. In *ECCV*. Springer.
- KATZIR, Z. and ELOVICI, Y. (2021). Who's afraid of adversarial transferability? *arXiv preprint arXiv:2105.00433* .

- KURAKIN, A., GOODFELLOW, I. J. and BENGIO, S. (2019). Adversarial examples in the physical world. *ICLR 2017 - Workshop Track Proceedings* 1–14.
- LAIDLAW, C. and FEIZI, S. (2019). Functional adversarial attacks. In *NeurIPS*.
- LAIDLAW, C., SINGLA, S. and FEIZI, S. (2021). Perceptual adversarial robustness: Defense against unseen threat models. In *ICLR*.
- LEINO, K., WANG, Z. and FREDRIKSON, M. (2021). Globally-robust neural networks. In *ICML*.
- LIU, F., ZHANG, C. and ZHANG, H. (2021a). Towards transferable adversarial perturbations with minimum norm. In *ICML 2021 Workshop on Adversarial Machine Learning*.
- LIU, Y., CHEN, X., LIU, C. and SONG, D. (2017). Delving into transferable adversarial examples and black-box attacks. In *ICLR*.
- LIU, Z., LIN, Y., CAO, Y., HU, H., WEI, Y., ZHANG, Z., LIN, S. and GUO, B. (2021b). Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*
- MADRY, A., MAKELOV, A., SCHMIDT, L., TSIPRAS, D. and VLADU, A. (2018). Towards deep learning models resistant to adversarial attacks. In *ICLR*.
- MAO, X., QI, G., CHEN, Y., LI, X., YE, S., HE, Y. and XUE, H. (2021). Rethinking the design principles of robust vision transformer. *arXiv preprint arXiv:2105.07926* .
- MOOSAVI-DEZFOOLI, S.-M., FAWZI, A. and FROSSARD, P. (2016). Deepfool: a simple and accurate method to fool deep neural networks. In *CVPR*.
- PAPERNOT, N., McDANIEL, P. and GOODFELLOW, I. (2016). Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*
- PAPERNOT, N., McDANIEL, P., GOODFELLOW, I., JHA, S., CELIK, Z. B. and SWAMI, A. (2017). Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*.
- PARKHI, O. M., VEDALDI, A. and ZISSERMAN, A. (2015). Deep face recognition. In *BMVC*. BMVA Press.
- PRABHU, V. U., WHALEY, J. and FRANCISCO, S. (2018). Art-attack! on style transfers with textures, label categories and adversarial examples.
- QIU, H., XIAO, C., YANG, L., YAN, X., LEE, H. and LI, B. (2020). Semanticadv: Generating adversarial examples via attribute-conditioned image editing. In *ECCV*. Springer.
- RECHT, B., ROELOFS, R., SCHMIDT, L. and SHANKAR, V. (2019). Do imagenet classifiers generalize to imagenet? In *ICML*. PMLR.

- SALMAN, H., YANG, G., LI, J., ZHANG, P., ZHANG, H., RAZENSHTEYN, I. and BUBECK, S. (2019). Provably robust deep learning via adversarially trained smoothed classifiers. In *NeurIPS*. Curran Associates Inc., Red Hook, NY, USA.
- SHAFABI, A., NAJIBI, M., GHIAZI, M. A., XU, Z., DICKERSON, J., STUDER, C., DAVIS, L. S., TAYLOR, G. and GOLDSTEIN, T. (2019). Adversarial training for free! In *NeurIPS*, vol. 32. Curran Associates, Inc.
- SHAMSABADI, A. S., SANCHEZ-MATILLA, R. and CAVALLARO, A. (2020). Colorfool: Semantic adversarial colorization. In *CVPR*.
- SHAO, R., SHI, Z., YI, J., CHEN, P.-Y. and HSIEH, C.-J. (2021). On the adversarial robustness of visual transformers. *arXiv preprint arXiv:2103.15670* .
- SHARIF, M., BAUER, L. and REITER, M. K. (2018). On the suitability of lp-norms for creating and preventing adversarial examples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*.
- SONG, Y., SHU, R., KUSHMAN, N. and ERMON, S. (2018). Constructing unrestricted adversarial examples with generative models. In *NeurIPS*, vol. 31. Curran Associates, Inc.
- SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I. and FERGUS, R. (2014). Intriguing properties of neural networks. In *ICLR*.
- TAORI, R., DAVE, A., SHANKAR, V., CARLINI, N., RECHT, B. and SCHMIDT, L. (2020). Measuring robustness to natural distribution shifts in image classification. In *NeurIPS*.
- TRAMER, F. and BONEH, D. (2019). Adversarial training and robustness for multiple perturbations. In *NeurIPS*, vol. 32. Curran Associates, Inc.
- TRAMER, F., CARLINI, N., BRENDL, W. and MADRY, A. (2020). On adaptive attacks to adversarial example defenses. In *NeurIPS*, vol. 33. Curran Associates, Inc.
- TRAMER, F., KURAKIN, A., PAPERNOT, N., GOODFELLOW, I., BONEH, D. and McDANIEL, P. (2018). Ensemble adversarial training: Attacks and defenses. In *ICLR*.
- WANG, Z., BOVIK, A. C., SHEIKH, H. R. and SIMONCELLI, E. P. (2004). Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13** 600–612.
- WIGHTMAN, R. (2019). Pytorch image models. <https://github.com/rwightman/pytorch-image-models>.
- WILLMOTT, D., SAHU, A. K., SHEIKHOLESLAMI, F., CONDESSA, F. and KOLTER, Z. (2021). You only query once: Effective black box adversarial attacks with minimal repeated queries. *arXiv preprint arXiv:2102.00029* .
- WONG, E. and KOLTER, J. Z. (2021). Learning perturbation sets for robust machine learning. In *ICLR*.

- WONG, E., RICE, L. and KOLTER, J. Z. (2020). Fast is better than free: Revisiting adversarial training. In *ICLR*.
- WONG, E., SCHMIDT, F. and KOLTER, Z. (2019). Wasserstein adversarial examples via projected sinkhorn iterations. In *ICML*. PMLR.
- XIAO, C., ZHU, J.-Y., LI, B., HE, W., LIU, M. and SONG, D. (2018). Spatially transformed adversarial examples. In *ICLR*.
- XIE, C., WU, Y., MAATEN, L. v. d., YUILLE, A. L. and HE, K. (2019a). Feature denoising for improving adversarial robustness. In *CVPR*.
- XIE, C., ZHANG, Z., ZHOU, Y., BAI, S., WANG, J., REN, Z. and YUILLE, A. (2019b). Improving transferability of adversarial examples with input diversity. In *CVPR*.
- XIE, Q., LUONG, M.-T., HOVY, E. and LE, Q. V. (2020). Self-training with noisy student improves imagenet classification. In *CVPR*.
- XU, Q., TAO, G., CHENG, S. and ZHANG, X. (2021). Towards feature space adversarial attack by style perturbation. *AAAI* **35** 10523–10531.
- ZHANG, H., CHEN, H., XIAO, C., GOWAL, S., STANFORTH, R., LI, B., BONING, D. and HSIEH, C.-J. (2020). Towards stable and efficient training of verifiably robust neural networks. In *ICLR*.
- ZHANG, H., YU, Y., JIAO, J., XING, E. P., GHAOUI, L. E. and JORDAN, M. I. (2019). Theoretically principled trade-off between robustness and accuracy. In *ICML*.
- ZHANG, J., ZHU, J., NIU, G., HAN, B., SUGIYAMA, M. and KANKANHALLI, M. (2021). Geometry-aware instance-reweighted adversarial training. In *ICLR*.
- ZHANG, R., ISOLA, P., EFROS, A. A., SHECHTMAN, E. and WANG, O. (2018). The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*.
- ZHAO, Z., LIU, Z. and LARSON, M. (2020a). Towards large yet imperceptible adversarial image perturbations with perceptual color distance. In *CVPR*.
- ZHAO, Z., LIU, Z. and LARSON, M. A. (2020b). Adversarial color enhancement: Generating unrestricted adversarial images by optimizing a color filter. In *BMVC*.

A Appendix

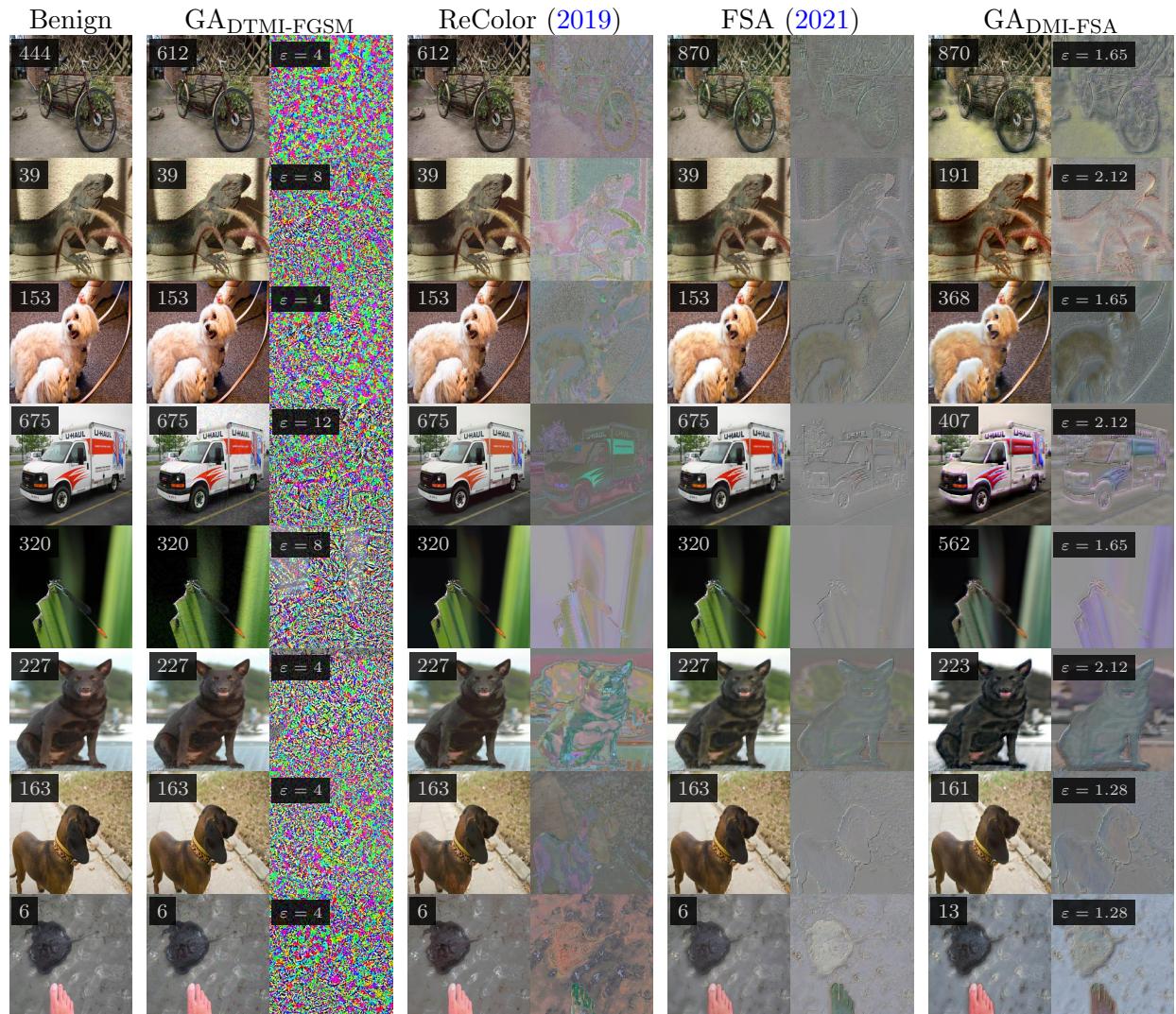


Figure 9: **Visualization of transfer attack on Resnext101-DenoiseAll** (Xie et al., 2019a).