# Diffusion Models for Imperceptible and Transferable Adversarial Attack

**Jianqi Chen   Hao Chen   Keyan Chen   Yilan Zhang**
**Zhengxia Zou   Zhenwei Shi**[*]
Beihang University
https://github.com/WindVChen/DiffAttack

## Abstract

Many existing adversarial attacks generate $L_p$-norm perturbations on image RGB space. Despite some achievements in transferability and attack success rate, the crafted adversarial examples are easily perceived by human eyes. Towards visual imperceptibility, some recent works explore unrestricted attacks without $L_p$-norm constraints, yet lacking transferability of attacking black-box models. In this work, we propose a novel imperceptible and transferable attack by leveraging both the generative and discriminative power of diffusion models. Specifically, instead of direct manipulation in pixel space, we craft perturbations in latent space of diffusion models. Combined with well-designed content-preserving structures, we can generate human-insensitive perturbations embedded with semantic clues. For better transferability, we further "deceive" the diffusion model which can be viewed as an additional recognition surrogate, by distracting its attention away from the target regions. To our knowledge, our proposed method, *DiffAttack*, is the first that introduces diffusion models into adversarial attack field. Extensive experiments on various model structures (including CNNs, Transformers, MLPs) and defense methods have demonstrated our superiority over other attack methods.

## 1   Introduction

Recent years have witnessed remarkable performance exhibited by deep neural networks (DNNs) across a range of domains, including autonomous driving [1, 2], medical image analysis [3, 4], remote sensing [5, 6], *etc*. Notwithstanding the indisputable advances, early investigations [7, 8] have elucidated the susceptibility of DNNs to meticulously engineered subversions (hereafter referred to as "adversarial examples"), which may induce grievous mistakes in real-world applications. Moreover, the transferability of these adversarial examples across distinct model architectures [9] poses an even greater hazard to practical implementations. Therefore, it is of the utmost necessity to uncover as many lacunae in machine perception – what may be termed "blind spots" – as can feasibly be achieved, so as to bolster the DNNs' resilience when faced with adversarial challenges.

Compared to white-box attacks [10–12, 7] that the attacker can access the architecture and parameters of the target model, black-box attacks [13–26, 9] can not obtain the target's information and thus are much closer to real-world scenarios. Among black-box directions, we here focus on the transfer-based attacks [15, 19–23, 9] that directly apply the adversarial examples constructed on a surrogate model to fool the target model. By adopting different optimization strategies [27, 23], designing various loss functions [16, 17, 19, 22], leveraging multiple data augmentations [13, 14, 20, 21], *etc*., existing approaches have achieved much success and well improved the attack's transferability.

$L_p$-**norm based methods**. Most of the above methods adopt $L_p$-norm in RGB color space as an indicator of human perception and constrain the amplitude of the adversarial perturbations under a

---

[*]Corresponding Author

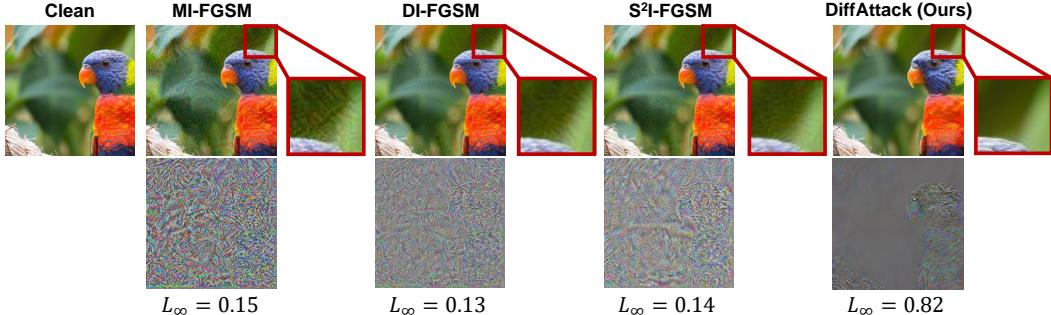| Clean | MI-FGSM | DI-FGSM | S²I-FGSM | DiffAttack (Ours) |
|---|---|---|---|---|
| | $L_\infty = 0.15$ | $L_\infty = 0.13$ | $L_\infty = 0.14$ | $L_\infty = 0.82$ |

Figure 1: **Adversarial perturbations crafted by some attacks.** The second row denotes the difference between the clean image and the adversarial example. Please zoom in for a better view.

specific value. Despite the efforts paid, these pixel-based attacks are still easy to be perceived by human eyes, and $L_p$-norm was recently found unsuitable to measure the perceptual distance between two images [28, 29]. From the examples displayed in Figure 1, the perturbations optimized by $L_p$ loss are noticeable and appear similar to high-frequency noise (indicate overfits on the surrogate model) despite low $L_\infty$ values, which can hinder the transferability to other black-box models [30, 31] and is easy to be defended against by purification defenses [32, 33].

**Towards imperceptible attacks**. Recent works [34, 30, 35, 36, 28, 37] explored new ways to deceive human perception without using the $L_p$-norm constraint (*a.k.a.* unrestricted attacks). By applying perturbations on spaces such as object attribute [36, 30], color mapping matrix [35], image texture [37], *etc.*, the adversarial examples are well imperceptible despite large $L_p$-norm values in RGB space. Furthermore, recent works [30, 37, 35] revealed that the perturbation generated by unrestricted attacks is more focused on relatively large-scale patterns with high-level semantics, instead of manipulating pixel-level intensity, thus benefiting the attack's transferability to other black-box models and even the defended ones. However, these methods' transferability still lags behind the pixel-based ones.

In this work, we propose a novel unrestricted attack based on diffusion models [38–42]. Instead of manipulating pixels directly, we optimize the latent of an off-the-shelf pretrained diffusion model [42]. Besides the basic transferability advantages of high-level perturbations mentioned above, our motivation for introducing the diffusion model into the adversarial attack domain stems primarily from its two beneficial properties. 1) *Good imperceptibility*. The diffusion model tends to sample natural images that conform to human perception, which naturally satisfies the required imperceptibility of the adversarial attack. Recent research on real image editing [43–45] further supports our ideas that perturbations can be applied to high-level semantics without compromising image realism. 2) *Approximation of a strong surrogate*. On the one hand, the denoised process of diffusion models can be taken as a powerful purification defense as pointed out by [32]. On the other hand, diffusion models that are trained on large-scale datasets are inherently powerful in discriminative capabilities [46, 47] and can be approximated as strong surrogate models for transfer-based attacks. If the attacks can fool this "strong model", then we may expect good transferability to other models and defenses.

To effectively utilize the above good properties of diffusion models, our work is carried out in three aspects. *Firstly*, we construct the basic framework following recent real editing approaches [43, 44] that inverse images to noise and then apply modifications in latent space. Compared with methods [48, 49] that change guided text to achieve content editing, we focus on direct operations on latents, which can benefit the attack success a lot. *Secondly*, we propose to deviate the cross-attention maps between text and image pixels, in which way we can transform the diffusion model into a surrogate model that can be practically deceived and attacked. *Finally*, to avoid distorting the initial semantics, specific measures, including self-attention constraint and inversion strength, are considered. We term the proposed unrestricted attack as *DiffAttack*, and our contributions can be summarized as follows:

- As far as we know, we are the first to unveil the potential of diffusion models for crafting adversarial examples with satisfactory imperceptibility and transferability.

- We propose *DiffAttack*, a novel unrestricted attack where the good properties of diffusion models are leveraged by specific designs.

- Extensive experiments on a variety of architectures (*e.g.*, CNNs, Transformers, and MLP-Mixers) and defenses (including recent diffusion purification) have demonstrated the superiorities of our work over the existing transfer-based methods.

## 2 Related Works

**Adversarial Attacks.** Since Szegedy *et al.* [8] found that DNNs can be deceived by imperceptibly small perturbations applied on images, adversarial attacks have long attracted significant attention in the deep learning field [7, 9]. Generally, the existing attacks can be divided into two parts: white-box attack and black-box attack. For white-box scenarios [10–12, 7], the attacker can get access to the model structures and parameters of the target model. Thus, strong adversarial examples can be crafted by directly using the backpropagated gradients. While for black-box scenarios [13–26, 9], there is limited information of the target model and it is closer to the real-world situation. Current approaches are either based on queries [18, 24–26] or on the cross-model transferability of adversarial examples [13, 14, 16, 17, 15, 19–23, 9]. Specifically, the former ones query the black-box model many times. With the queried results, they generate adversarial examples either by gradients approximation [24] or by random search [25, 26, 18]. The transfer-based attacks resort to a surrogate model. By crafting perturbations in the same way as white-box attacks, they expect these adversarial examples can also have a good effect on the target model. In this work, we focus on the transfer-based part.

**Transferable Attacks.** To enhance the generalization of adversarial examples crafted on surrogate models, previous works put a lot of effort into keeping perturbations from getting stuck in a model-specific local optimum that overfits the surrogate model and cannot transfer well to other methods. [50, 51] adopted the straightforward strategy of *model ensembles* to attack as many models as possible by finding an optimum updated direction. [13, 14, 20, 21] proposed to leveraged *data augmentations* to diversify the inputs, which ensure the attack robustness under different scenarios. [16, 17, 19, 22] applied *loss functions* on the feature space which demonstrated good performance on black-box targets. [27, 23] combined momentum into *optimization schedules* to help jump out of local optimum. Despite the much improvement in the transferability, these works mostly conduct attacks with $L_p$-norm constraint on RGB pixel space, resulting in high-frequency noises and patterns (see Figure 1) which, though hold a relatively low value on $L_p$-norm, are easy to be perceived by humans. Different from them, our *DiffAttack* perturbs the image's latent in diffusion models, which is of good imperceptibility together with excellent transferability across various black-box models.

**Unrestricted Attacks.** Since $L_p$-norm in RGB space was found not ideal for measuring the perceptual distance [30, 37, 35], recent researches [34, 30, 35, 36, 28, 37] turn to unconstraint and propose unrestricted but imperceptible attacks. Song *et al.* [34] proposed to sample adversarial examples from scratch through generative adversarial networks [52]. Bhattad *et al.* [37] perturbed images from the perspective of color and texture. Zhao *et al.* [28] adopted CIEDE2000 which can better indicate the perceptual color loss. Qiu *et al.* [36] and Jia *et al.* [30] achieve imperceptibility by modifying the attributes of the images. Yuan *et al.* [35] constructed a color distribution library, which is used to find a successful distribution for adversarial attacks. However, despite their good imperceptibility, these methods generally cannot compete with the aforementioned pixel-based methods in terms of transferability. Our work also falls in this direction but achieves better transferability and imperceptibility, and is the first to explore the strength of diffusion models in crafting unrestricted attacks.

**Diffusion Models.** Recently, diffusion models [38, 39, 53, 54] have attracted extensive attention and shown their fabulous power. Images are first converted into purely Gaussian noise in the *forward process* and then a U-Net [55] structure is trained to predict the added noise in each timestep of the *reversed process*. Being trained on large numbers of data, the diffusion models [40–42] can either generate high-quality images from randomly sampled noise, or more specific ones that well follow the guidance of text prompt. Due to its significant performance, the diffusion model has also diffused to other areas, such as image inpainting [56, 57], image super-resolution [58], real image editing [44, 43], *etc*. Recent work further showed that the pretrained diffusion models can be taken as strong recognition models [46, 47] and denoisers [32]. Despite the many applications mentioned above, the potential of diffusion models in the adversarial attack field remains unexplored.

## 3 Method

### 3.1 Problem Formulation

Given a clean image $x$ and its corresponding label $y$, attackers aim to craft perturbations that can deviate the decision of a classifier $F_\theta$ ($\theta$ denotes the model's parameters) from correct to wrong:

$$F_\theta(\text{Attack}(\text{x}; \text{G}_\phi)) = F_\theta(x') \neq y \tag{1}$$
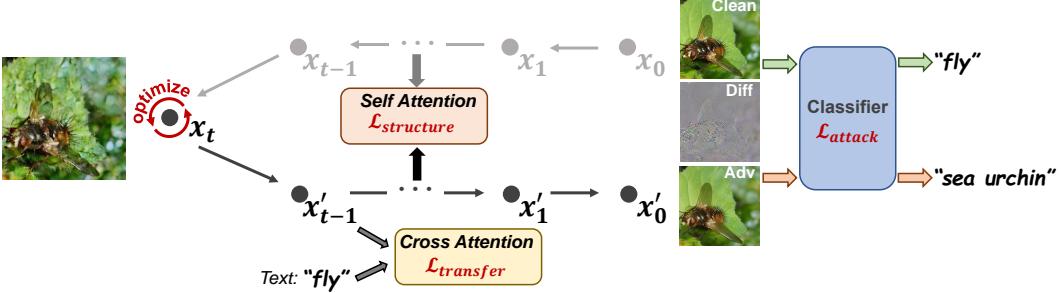
Figure 2: **Framework of *DiffAttack*.** We adopt Stable Diffusion [42] and leverage DDIM Inversion [59] to convert the clean image into the latent space. The latent is optimized to deceive the classifier. The cross-attention maps are leveraged to "deceive" the diffusion model, and we use self-attention maps to preserve the structure. For simplicity, we here do not display the unconditional optimization, whose details can be referred to Section 3.4.

where $\text{Attack}(\cdot)$ is the attack approach and $x'$ denotes the crafted adversarial example. Since the information of $F_\theta$ is inaccessible in black-box scenarios, the adversarial examples are crafted on a surrogate model $G_\phi$.

Different from classic pixel-based attacks [7, 10, 23] that apply $L_p$-norm constraints on perturbations ($\|\epsilon\|_p < c$, where $\epsilon$ denotes the perturbation and $c$ is a hyper-parameter), we impose perturbations in the latent space of the diffusion model and rely on properties of the diffusion model to achieve visually natural and successful attacks. We will describe our design in detail below.

## 3.2 Basic Framework of Destruction and Reconstruction

We display in Figure 2 the whole framework of *DiffAttack*, where we adopt the open-source Stable Diffusion [42] that pretrained on extremely massive text-image pairs. Since adversarial attacks aim to fool the target model by perturbing the initial image, they can be approximated as a special kind of real image editing. Inspired by recent diffusion editing approaches [44, 43, 49], our framework also leverages the DDIM Inversion technology [59] (details can be found in Appendix B), that the clean image is mapped back into the diffusion latent space by reversing the deterministic sampling process:

$$x_t = \text{Inverse}(x_{t-1}) = \underbrace{\text{Inverse} \circ \cdots \circ \text{Inverse}}_{t}(x_0) \qquad (2)$$

where $\text{Inverse}(\cdot)$ denotes the DDIM Inversion operation (For simplicity, we ignore the autoencoder stage in the Stable Diffusion [42]). We apply the inversion for several timesteps from $x_0$ (the initial image) to $x_t$. A high-quality reconstruction of $x_0$ can then be expected if we conduct the deterministic denoising process from $x_t$ [54, 59].

Many of the existing approaches [44, 43, 49] proposed to modify text embeddings for image editing, through which way, the image latent $x_t$ can gradually shift to the target semantic space during the iterative denoising process with the text guidance. However, in our explorations (see Appendix C), we found that the perturbations on the guided text embeddings would be hard to work on the other black-box models, leading to weak transferability. Therefore, different from the editing approaches, we here propose to directly perturb the latent $x_t$:

$$\underset{x_t}{\arg\min} \mathcal{L}_{attack} = -J(x', y; G_\phi), \quad \text{where } x' = x'_0 = \underbrace{\text{Denoise} \circ \cdots \circ \text{Denoise}}_{t}(x_t) \qquad (3)$$

where $J(\cdot)$ is the cross-entropy loss and $\text{Denoise}(\cdot)$ denotes the diffusion denoising process. A possible concern is that such a straightforward approach may lead to unnatural results, however, we can observe in Figure 1 that the difference is almost indistinguishable between the image reconstructed from the perturbed latent and the initial clean one. Furthermore, we can notice that the difference image may embody many high-level semantics clues, instead of the high-frequency noise when using pixel-based attacks. We attribute these to the denoising process of the Stable Diffusion model, together with its autoencoder's strong mapping capability. These semantically rich perturbations can not only enhance the imperceptibility but also benefit a lot the attack's transferability [30, 37, 35].

### 3.3 "Deceive" Strong Diffusion Model

According to the research by [32], the reversed process of the diffusion model is a strong adversarial purification defense. Thus, our perturbed latent has experienced purification before being decoded to the final image, which then ensures the naturalness of crafted adversarial examples and also the robustness towards other purification denoises (see Section 4.2.2).

Apart from evading the denoising part, we here go a further step to enhance our attack's transferability. Given an image and its corresponding caption, we can see from Figure 3 that in the reconstruction process of the inversed latent, the cross-attention maps have displayed a strong relationship between the guided text and the image pixels, which demonstrates the powerful recognition capability of pretrained diffusion models. Such a relationship is also verified by Hertz *et al.* [48]



Figure 3: **Visualization of cross- and self- attention maps.** There is a strong relationship between text and pixels in cross attention, while self attention can well reveal structure.

and its recognition power recently has been leveraged on the downstream tasks [46, 47]. Thus, the diffusion model that trained on massive data can be approximated as a strong recognition model, and our motivation is that, if our crafted attacks can "deceive" this powerful model, we may expect much improvement of the transferability to other black-box models.

Denote $C$ as the caption of the clean image, which we set to the groundtruth category's name (we can also simply use the predicted category of $G_\phi$, thus not relying on true labels). We accumulate the cross-attention maps between image pixels and $C$ in all the denoising steps and get the average. To "deceive" the pretrained diffusion model, we propose to minimize the following formula:

$$\underset{x_t}{\arg\min} \mathcal{L}_{transfer} = \text{Var}(\text{Average}(\text{Cross}(x_t, t, C; \text{SDM}))) \tag{4}$$

where $\text{Var}(\cdot)$ calculates the input's variance, $\text{Cross}(\cdot)$ denotes the accumulation of all the cross-attention maps in the denoising process, and SDM is the Stable Diffusion. The insight is to distract the diffusion model's attention from the labeled objects. By evenly distributing cross attention to each pixel of the image, we can disrupt the original strong semantic relationship, ensuring that our crafted adversarial examples well "deceive" the strong diffusion model.

### 3.4 Preserve Content Structure

As mentioned in Section 3.2, our unrestricted attack can be approximated as an image editing approach, thus the change of the content structure is unavoidable. If the degree of the changes is not under control, the resulting adversarial examples may lose most semantics of the initial clean image (see Figure 5), which losses the significance of the adversarial attacks and is not what we want. Therefore, we here preserve the content structure mainly from two perspectives.

**Self-Attention Control.** Researches by [60, 61] have discovered that the self-similarly-based descriptors can capture structural information while ignoring image appearance. Along with this idea, we can observe from Figure 3 that the self attention in the diffusion models also has that property embedded in it, which is in contrast to cross attention that mainly focuses on high-level semantics. Therefore, we propose to leverage the self-attention maps for structure retention. We set a copy $x_{t(fix)}$ of the inversed latent which is fixed without perturbations. By respectively calculating the self-attention maps (denoted as $S_{t(fix)}$ and $S_t$) of $x_{t(fix)}$ and $x_t$, we force $S_t$ to get close to $S_{t(fix)}$ as follows:

$$\underset{x_t}{\arg\min} \mathcal{L}_{structure} = \|S_t - S_{t(fix)}\|_2^2 \tag{5}$$

Similar to Eq. 4, we here apply the self-attention constraint to all the denoising steps. Since $x_{t(fix)}$ reconstructs the initial clean image well [59, 54], we can in this way preserve the structure.

**Inversion Strength Trade-off.** With DDIM Inversion strength increased, the latent $x_t$ will get closer to pure Gaussian distribution and the perturbations on it may cause serious distortion due to

influence on more denoising steps (see Figure 5). Whereas, a limited inversion cannot provide enough space for attacking, since the latent image prior is too strong. The inversion strength is a trade-off between imperceptibility and the attack success. Recent work [62] has found that the diffusion models tend to add coarse semantic information (*e.g.*, layout) in the early denoising steps while more fine details in the later steps. Thus, we control the inversion at the back of the denoising process for retention of high-level semantics, and reduce the total DDIM sample steps for more editing space.

Besides the above operations, we also adopt the approach of [43] to get a good initial reconstruction by optimizing the unconditional embeddings. Details can be found in the source paper.

In general, the final objective function of *DiffAttack* is as follows, where $\alpha$, $\beta$, and $\gamma$ represent the weight factors of each loss:

$$\arg \min_{x_t} \mathcal{L} = \alpha \mathcal{L}_{attack} + \beta \mathcal{L}_{transfer} + \gamma \mathcal{L}_{structure} \tag{6}$$

## 4 Experiments

### 4.1 Experimental Setup

**Datasets.** Following the previous methods [13, 35, 28], we evaluate the performance of our attack on the development set of ImageNet-Compatible Dataset[2], which consists of 1,000 images with size 299×299×3. Considering that the Stable Diffusion cannot handle the original input size of the ImageNet-Compatible Dataset, we focused on a resized version of 224×224×3 in all the experiments.

**Models.** We evaluate the transferability of the attacks across a variety of network structures, including CNNs, Transformers, and MLPs. For CNNs, we adopt normally trained models including ConvNeXt [63], ResNet-50 (Res-50) [64], VGG-19 [65], Inception-v3 (Inc-v3) [66], and MobileNet-v2 (Mob-v2) [67]. For Transformers, we consider normally trained ViT-B/16 (ViT-B) [68], Swin-B [69], DeiT-B and DeiT-S [70]. For MLPs, we adopt normally trained Mixer-B/16 (Mix-B) and Mixer-L/16 (Mix-L) [71]. Furthermore, we also consider various defense methods, including DiffPure [32], RS [72], R&P [73], HGD [74], NIPS-r3 [75], NRP [19], and adversarially trained models (Adv-Inc-v3 [76], Inc-v3$_{ens3}$, Inc-v3$_{ens4}$, and IncRes-v2$_{ens}$ [77]).

**Implementation Details.** We leverage DDIM [59] as the sampler of the Stable Diffusion [42]. The number of steps is set to 20 and we apply 5 DDIM Inversion steps of the initial clean image. In the inversion process, the guidance scale is set to 0, while in the denoising process, we set it to 2.5. For optimizing the latent $x_t$, we adopt AdamW [78] with the learning rate set to 1e$^{-2}$ and the iterations set to 30. The weight factors $\alpha$, $\beta$, $\gamma$ in Eq. 6 are set to 10, 10000, 100 respectively. All the experiments are run on a single RTX 3090 GPU.

**Evaluation Metrics.** We adopt top-1 accuracy to evaluate the performance of the attack methods. Besides, we leverage Frechet Inception Distance (FID) [79] as the indicator of the human-imperceptibility of the crafted adversarial examples.

### 4.2 Comparisons

#### 4.2.1 Results on Normally Trained Models

Here, we compared the performance of *DiffAttack* on normally trained models with other transfer-based black-box attacks. We select five pixel-based attacks (MI-FGSM [23], DI-FGSM [20], TI-FGSM [21], PI-FGSM [15], S$^2$I-FGSM [13]) and two unrestricted attacks (PerC-AL [28], NCF [35]). Except that the resolution is changed to 224×224×3, the implementations of these methods follow their original optimal settings. We craft the adversarial examples via Res-50, VGG-19, Mov-v2, Inc-v3, ConvNeXt, and Swin-B. The transferability of different methods is displayed in Table 1.

From the results, we can observe that *DiffAttack* can achieve the best transferability across a variety of model structures, while other unrestricted attacks (PerC-AL and NCF) usually fail to compete with pixel-based attacks. In some architectures such as VGG-19 and Mob-v2, our method can even outperform the second-best method by nearly 10 points (38.2% vs. 49.0%, 40.5% vs. 49.9%). It may be noticed that our method fails to compete with MI-FGSM and PI-FGSM under Inc-v3 structure, however, from the FID results, we have a large advantage over them (about 20 or more points lower).

---

Table 1: **Transferability and imperceptibility comparisons on normally trained models.** We report top-1 accuracy(%) of each method. "S." denotes surrogate models while "T." denotes target models. For the white-box attacks (surrogate model same to target one), we set their background to gray. "AVG(w/o self)" denotes the average accuracy on all the target models except the one that same as the surrogate one. The best result is bolded, and the second-best result is underlined.

| T. \ S. | Attacks | CNNs | | | | | Transformers | | | | MLPs | | AVG(w/o self) | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Res-50 | VGG-19 | Mob-v2 | Inc-v3 | ConvNeXt | ViT-B | Swin-B | DeiT-B | DeiT-S | Mix-B | Mix-L | | |
| | Clean | 92.7 | 88.7 | 86.9 | 80.5 | 97.0 | 93.7 | 95.9 | 94.5 | 94.0 | 82.5 | 76.5 | 89.4 | 57.8 |
| Res-50 | MI-FGSM | 0 | 19.9 | 20.2 | 28.9 | 57.8 | 67.3 | 67.0 | 72.4 | 67.0 | 52.2 | 45.4 | 49.8 | 81.2 |
| | DI-FGSM | 0 | 21.2 | 20.5 | 34.5 | 71.6 | 82.0 | 75.3 | 80.5 | 76.0 | 61.3 | 56.8 | 58.0 | 85.3 |
| | TI-FGSM | 0 | 42.4 | 37.1 | 46.0 | 83.6 | 81.6 | 83.7 | 84.5 | 79.0 | 66.0 | 61.7 | 66.6 | 66.0 |
| | PI-FGSM | 0 | 14.1 | 15.0 | 24.0 | 72.5 | 65.3 | 77.5 | 76.7 | 65.0 | 50.5 | 43.8 | 50.5 | 97.9 |
| | $S^2$I-FGSM | 0 | 9.2 | 6.6 | 18.6 | 44.1 | 63.9 | 52.0 | 65.9 | 59.0 | 45.6 | 44.3 | 40.9 | 79.8 |
| | PerC-AL | 6.5 | 83.1 | 80.2 | 76.4 | 96.0 | 93.9 | 94.8 | 94.4 | 93.0 | 81.6 | 75.1 | 86.8 | 58.2 |
| | NCF | 11.3 | 30.5 | 30.3 | 52.6 | 78.3 | 65.7 | 76.8 | 75.1 | 67.0 | 53.7 | 47.6 | 57.8 | 70.9 |
| | DiffAttack(Ours) | 3.7 | 24.4 | 22.9 | 31.0 | 41.0 | 48.8 | 43.8 | 49.5 | 45.0 | 42.9 | 42.2 | 39.2 | 62.6 |
| VGG-19 | MI-FGSM | 22.7 | 0 | 15.4 | 33.5 | 53.2 | 73.2 | 63.3 | 74.7 | 68.0 | 54.3 | 48.6 | 50.6 | 82.4 |
| | DI-FGSM | 32.2 | 0 | 23.9 | 46.5 | 67.2 | 84.7 | 71.9 | 84.8 | 80.0 | 65.7 | 60.9 | 61.8 | 70.9 |
| | TI-FGSM | 44.5 | 0 | 32.8 | 47.4 | 77.8 | 81.4 | 79.3 | 83.6 | 79.0 | 64.9 | 60.3 | 65.1 | 66.6 |
| | PI-FGSM | 22.7 | 0 | 16.4 | 29.8 | 68.3 | 68.0 | 75.7 | 79.5 | 68.0 | 50.9 | 41.8 | 52.1 | 96.4 |
| | $S^2$I-FGSM | 17.9 | 0 | 11.3 | 31.8 | 49.5 | 74.1 | 57.9 | 76.0 | 68.0 | 52.6 | 50.8 | 49.0 | 82.9 |
| | PerC-AL | 87.5 | 4.6 | 79.0 | 76.1 | 95.1 | 94.2 | 94.0 | 94.3 | 93.0 | 81.3 | 75.1 | 87.0 | 57.9 |
| | NCF | 38.3 | 6.8 | 31.5 | 52.4 | 80.5 | 67.5 | 77.6 | 77.4 | 71.0 | 53.5 | 47.2 | 59.7 | 70.4 |
| | DiffAttack(Ours) | 21.1 | 2.7 | 19.4 | 29.7 | 43.1 | 52.9 | 41.6 | 51.3 | 45.0 | 39.6 | 38.5 | 38.2 | 63.9 |
| Mob-v2 | MI-FGSM | 26.4 | 18.7 | 0 | 31.0 | 62.0 | 69.5 | 65.2 | 71.6 | 63.0 | 46.9 | 44.4 | 49.9 | 76.4 |
| | DI-FGSM | 28.7 | 18.9 | 0 | 33.9 | 73.4 | 79.9 | 71.4 | 79.6 | 75.0 | 57.7 | 57.1 | 57.6 | 78.6 |
| | TI-FGSM | 47.2 | 37.9 | 0 | 45.2 | 83.0 | 79.9 | 80.9 | 81.8 | 76.0 | 61.7 | 58.3 | 65.1 | 65.6 |
| | PI-FGSM | 21.1 | 13.3 | 0 | 27.6 | 74.4 | 65.3 | 77.0 | 77.4 | 66.0 | 49.7 | 41.5 | 51.4 | 98.7 |
| | $S^2$I-FGSM | 21.0 | 13.4 | 0 | 27.2 | 64.3 | 74.1 | 62.6 | 75.2 | 68.0 | 51.4 | 48.3 | 50.5 | 79.4 |
| | PerC-AL | 88.2 | 84.2 | 5.9 | 76.8 | 96.2 | 93.9 | 94.2 | 94.3 | 94.0 | 81.2 | 74.3 | 87.7 | 58.1 |
| | NCF | 36.0 | 29.4 | 7.4 | 51.9 | 77.4 | 67.2 | 76.1 | 76.1 | 68.0 | 54.9 | 48.3 | 58.6 | 69.7 |
| | DiffAttack(Ours) | 23.6 | 23.4 | 1.8 | 31.6 | 50.3 | 51.4 | 45.8 | 53.4 | 46.0 | 38.5 | 40.8 | 40.5 | 62.9 |
| Inc-v3 | MI-FGSM | 42.8 | 36.6 | 34.4 | 0 | 79.8 | 75.3 | 79.4 | 78.6 | 73.0 | 56.5 | 48.9 | 60.6 | 80.5 |
| | DI-FGSM | 61.7 | 57.4 | 51.9 | 0.2 | 89.9 | 84.6 | 86.8 | 86.7 | 82.0 | 68.4 | 62.3 | 73.2 | 67.1 |
| | TI-FGSM | 76.0 | 70.1 | 66.7 | 0.1 | 93.8 | 88.7 | 91.2 | 89.7 | 88.0 | 73.8 | 66.8 | 80.5 | 62.8 |
| | PI-FGSM | 37.9 | 22.4 | 28.4 | 0 | 81.0 | 74.3 | 83.0 | 81.9 | 72.0 | 57.1 | 45.8 | 58.4 | 92.5 |
| | $S^2$I-FGSM | 52.3 | 47.8 | 43.3 | 0 | 86.3 | 80.8 | 84.1 | 83.8 | 78.0 | 63.5 | 57.3 | 67.8 | 72.5 |
| | PerC-AL | 90.8 | 87.0 | 85.8 | 7.7 | 97.5 | 93.6 | 95.1 | 94.2 | 94.0 | 81.5 | 75.3 | 89.4 | 58.4 |
| | NCF | 52.6 | 45.8 | 46.2 | 17.4 | 85.7 | 75.9 | 83.4 | 82.7 | 76.0 | 61.1 | 52.9 | 66.2 | 66.7 |
| | DiffAttack(Ours) | 59.5 | 55.6 | 55.4 | 13.9 | 76.9 | 75.2 | 72.8 | 74.0 | 71.0 | 58.9 | 54.7 | 65.4 | 62.3 |
| ConvNeXt | MI-FGSM | 34.5 | 22.4 | 26.5 | 41.9 | 0 | 63.4 | 18.0 | 56.5 | 56 | 41.4 | 40.2 | 40.1 | 84.5 |
| | DI-FGSM | 33.6 | 24.3 | 29.8 | 46.6 | 0 | 71.0 | 18.8 | 62.2 | 64.0 | 49.6 | 46.7 | 44.6 | 79.6 |
| | TI-FGSM | 50.7 | 37.3 | 41.1 | 51.8 | 0 | 70.9 | 38.8 | 68.6 | 69.0 | 52.3 | 47.2 | 52.7 | 73.5 |
| | PI-FGSM | 23.6 | 14.2 | 17.1 | 22.4 | 0 | 43.0 | 37.2 | 48.7 | 43.0 | 33.2 | 31.7 | 31.4 | 101.8 |
| | $S^2$I-FGSM | 13.6 | 9.6 | 11.9 | 20.2 | 0 | 35.4 | 4.2 | 31.0 | 31.0 | 23.2 | 25.6 | 20.5 | 99.4 |
| | PerC-AL | 89.0 | 84.5 | 84.0 | 77.5 | 88.9 | 92.8 | 90.0 | 92.4 | 92.0 | 79.3 | 74.5 | 85.6 | 57.7 |
| | NCF | 47.1 | 41.4 | 39.2 | 54.7 | 41.4 | 61.6 | 63.9 | 64.8 | 62.0 | 52.2 | 47.8 | 53.5 | 67.0 |
| | DiffAttack(Ours) | 20.9 | 24.8 | 21.8 | 25.8 | 1.9 | 26.7 | 11.4 | 21.6 | 24.0 | 21.7 | 24.0 | 22.2 | 73.3 |
| Swin-B | MI-FGSM | 55.7 | 42.3 | 42.7 | 55.2 | 42.5 | 70.6 | 0.9 | 64.2 | 64.0 | 52.4 | 47.9 | 53.7 | 72.8 |
| | DI-FGSM | 52.7 | 43.0 | 44.5 | 56.4 | 33.9 | 66.6 | 2.7 | 57.2 | 58.0 | 52.4 | 50.8 | 51.5 | 65.7 |
| | TI-FGSM | 71.9 | 61.7 | 56.9 | 60.2 | 66.0 | 76.3 | 1.9 | 72.2 | 72.0 | 61.2 | 56.9 | 65.6 | 65.9 |
| | PI-FGSM | 38.3 | 21.6 | 25.8 | 35.7 | 54.8 | 48.4 | 0.6 | 52.4 | 47.0 | 43.5 | 38.5 | 40.6 | 89.7 |
| | $S^2$I-FGSM | 47.4 | 37.8 | 35.4 | 45.3 | 26.8 | 48.5 | 1.0 | 46.2 | 45.0 | 39.3 | 39.0 | 41.1 | 68.2 |
| | PerC-AL | 92.2 | 87.4 | 85.5 | 78.5 | 94.6 | 94.0 | 6.3 | 94.1 | 93.0 | 81.4 | 75.5 | 87.6 | 57.9 |
| | NCF | 49.5 | 44.9 | 44.9 | 60.5 | 70.1 | 63.7 | 36.9 | 66.0 | 63.0 | 51.7 | 49.1 | 56.3 | 65.5 |
| | DiffAttack(Ours) | 43.5 | 42.1 | 40.7 | 41.4 | 34.0 | 39.0 | 9.9 | 35.0 | 37.0 | 37.7 | 37.4 | 38.8 | 65.5 |

For the imperceptibility of the crafted adversarial examples, we can observe that PerC-AL has always achieved the best FID result. However, it can hardly deceive other black-box models and achieves the worst transferability where the accuracy value of AVG(w/o self) is very close to the clean image. Therefore, we choose to ignore the PerC-AL results here, and our method achieves the best performance. We turn the color of the best performance in Table 1 to red for a better view.

In Figure 4, we visualize the adversarial examples crafted by different attack approaches. We can see that our attack is much more imperceptible compared with MI-FGSM, DI-FGSM, TI-FGSM, PI-FGSM, and $S^2$I-FGSM, where there is high-frequency noise that can be perceived easily. Compared with NCF, *DiffAttack* is more natural in color space. For PerC-AL, although the attack can hard to be perceived, its transferability is the worst as mentioned above. Thus, our method's superiority is well verified. More visualizations can be found in **Supplementary Material** (Appendix F).

### 4.2.2  Results on Defense Approaches

To further verify the robustness of each attack method, we evaluate the performance of the crafted adversarial examples on defense approaches. Following [35, 13], we consider both input preprocessing defenses [72–75, 19] and adversarially trained models [76, 77] (see Section 4.1). We further consider the recent DiffPure defense [32] to better demonstrate our superiority. We take Inc-v3 as an example surrogate model and all the adversarial examples are crafted from it. For NRP and DiffPure, we set

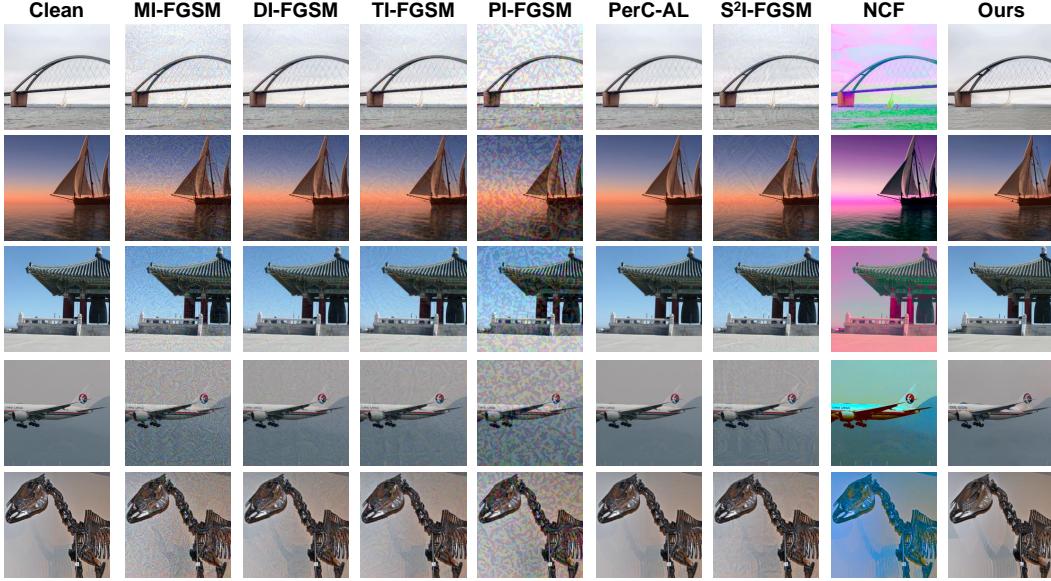| Clean | MI-FGSM | DI-FGSM | TI-FGSM | PI-FGSM | PerC-AL | S²I-FGSM | NCF | Ours |

Figure 4: **Visual comparisons among different attacks.** Please zoom in for a better view.

Table 2: **Robustness on defense approaches.** We report top-1 accuracy(%) of each method. "A." denotes attack methods while "D." denotes defense approaches. "Inc-v3$_{normal}$" denotes the accuracy on normally trained Inc-v3. For NRP and DiffPure, we display the accuracy differences after purification. The best result is bolded, and the second-best result is underlined.

| D. / A. | HGD | R&P | NIP-r3 | RS | Adv-Inc-v3 | Inc-v3$_{ens3}$ | Inc-v3$_{ens4}$ | IncRes-v2$_{ens}$ | Inc-v3$_{normal}$ | NRP | DiffPure |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MI-FGSM | 77.9 | 76.8 | _65.2_ | 62.7 | 51.1 | 49.8 | 53.7 | 70.5 | 0 | +5.9 | +38.4 |
| DI-FGSM | 80.5 | 83.8 | _79.4_ | 68.7 | 64.2 | 58.5 | 61.5 | 74.9 | 0.2 | +22.9 | +52.4 |
| TI-FGSM | 84.7 | 86.1 | 87.0 | 69.4 | 66.1 | 62.4 | 64.5 | 76.8 | 0.1 | +25.8 | +55.5 |
| PI-FGSM | 73.4 | 68.6 | **57.2** | _37.1_ | **42.3** | _45.0_ | _44.6_ | 62.0 | 0 | +7.7 | +21.5 |
| S²I-FGSM | 72.5 | 76.5 | 73.3 | 65.0 | 51.8 | 47.0 | 52.2 | 67.7 | 0 | +3.2 | +47.0 |
| PerC-AL | 95.6 | 94.4 | 96.7 | 74.2 | 80.8 | 76.7 | 75.4 | 88.6 | 7.7 | +56.8 | +55.9 |
| NCF | _71.1_ | _66.4_ | 74.6 | **29.0** | 48.8 | 47.2 | 49.0 | _60.5_ | 17.4 | +11.0 | +14.8 |
| DiffAttack(Ours) | **62.0** | **65.5** | 70.0 | 52.8 | _46.0_ | **43.8** | **43.1** | **58.3** | 13.9 | **+2.3** | **+13.9** |

the target model as Inc-v3 itself, thus better revealing the robustness. For other defenses, the target models are the same as the official papers. We display the results in Table 2.

From the results, we can see that our method can achieve good robustness and outperform other methods when some defenses are applied. For the adversarial purification defenses, it can be seen that the attack success of our attack has the least change compared with other ones, which does verify the robustness of *DiffAttack* and the effectiveness of our designs in Section 3.3.

### 4.3 Ablation Studies

In Table 3, we ablate the designs mentioned in Section 3.3. The adversarial examples are crafted on Inc-v3. We can observe that with the loss in Eq. 4 added, the attack success improves, verifying our design's effectiveness. It can also be noted that prompt guidance is important for transferability, which we attribute to the fact

Table 3: **Transferability.**

| Prompt Guidance | Diffusion Deception | AVG (w/o self) |
|---|---|---|
| ✗ | ✗ | 70.0 |
| ✓ | ✗ | 66.5 |
| ✓ | ✓ | 65.4 |

Table 4: **Imperceptibility.**

| Steps Inversed | Self-Attn Control | Initial Recon | FID |
|---|---|---|---|
| 10 | ✗ | ✗ | 97.9 |
| 5 | ✗ | ✗ | 66.7 |
| 5 | ✗ | ✓ | 63.5 |
| 5 | ✓ | ✓ | 62.3 |



| Clean | inverse 10 steps w/o self-attn control w/o initial recon | inverse 5 steps w/o self-attn control w/o initial recon | inverse 5 steps w/ self-attn control w/o initial recon | inverse 5 steps w/ self-attn control w/ initial recon (final) |

Figure 5: **Visualization of imperceptibility ablations.** Please zoom in for a better view.

that prompts can help guide the attack on the target objects. Table 4 verifies the effectiveness of our designs for structure retention. With the inversion strength and self attention controlled, the FID

result gradually improves. We also visualize the structure ablation in Figure 5, which can display the visual improvement obviously. It can be seen that the control of inversion strength can help a lot preserve the structure, and the usage of self-attention maps can ensure better texture. More ablation experiments on parameter settings can be found in **Supplementary Material** (Appendix F).

## 5 Further Trials and Outlook

Besides the designs in Section 3, we have also conducted some other trials during the exploration of adversarial attacks with diffusion models. Although these trials achieve little success, considering the possible help for future research, we give a discussion here.

As mentioned in Section 3.4, for some specific images, the adversarial examples crafted by *DiffAttack* may distort a lot compared with the original ones. For better control of the changes, we try to generate "pseudo" masks with the cross attention. With these masks, we can then filter out the background regions and only perturb the foreground objects, thus achieving better human-imperception. However, we found that although the results could more easily evade the human eyes, their transferability dropped a lot. We infer this may be because background information is also beneficial for image recognition. More details about the implementation and experiments of the trial can be found in Appendix D. In practice, we will weaken the inversion strength for overly distorted images.

Moreover, we also explore further improving the transferability of *DiffAttack*. For image classification, the classifier will output each category's confidence, and *top1* is usually taken as the final decision. We here try to also make use of the following 4 categories in *top5* for better transferability. Specifically, different from Section 3.3 that we only set the guided text to *top1* category's name, we here set the text to a stack of the 5 categories' names in *top5* (sort by confidence from largest to smallest). Then, we optimize $x_t$ to reduce the intensity of cross attention between pixels and the first category text and increase that between pixels and the other four categories text. The motivation is that, the confidence denotes, to some extent, the amount of related information of the category in the image, thus it may be much easier to deceive the classifier to the nearest category on the decision plane. However, from the experiments in Appendix E, this trial fails to improve the transferability and even hurts it. We attribute this to the limitation of the search space. More details can be found in Appendix E.

In addition to the trials described above, here we provide future directions about diffusion adversarial attacks. One avenue for future research is to take diffusion models as a novel input augmentation. Recently, there are many works [13, 14, 20] that enhance the attack's transferability by applying differentiable augmentations on the input image, in which way, the crafted adversarial examples gain robustness under different scenarios. In line with these approaches, we can also take diffusion models as novel augmentations. By directly adding noise (or applying DDIM Inversion), we first convert the input image into the latent space, then we conduct the diffusion denoising process to reconstruct images. This reconstruction process, with small differences from the input image every time, can be seen as an augmentation when we leverage stochastic sampling in each step (the way like DDPM [39] but not deterministic DDIM [59]). Therefore, we may expect good transferability in this way.

Moreover, as the adversarial example crafted by diffusion models have many semantic clues embedded in it, it is also interesting and worth exploring whether the accuracy of clean images can be improved if we merge these examples in the training dataset and whether such an adversarial training can enhance the robustness of the classifier without sacrificing the clean image accuracy compared with previous attacks [80]. We leave these to future work.

## 6 Conclusion

In this work, we explore the potential of diffusion models in crafting adversarial examples and propose a powerful transfer-based unrestricted attack. By leveraging the properties of diffusion models, our approach achieves both imperceptibility and transferability. Experiments across extensive black-box models and defenses have demonstrated our method's superiority. Furthermore, we also comprehensively discussed the possible future work with diffusion models. We hope that our work can pave the way for further research on diffusion-based adversarial attacks.

**Negative Societal Impacts.** Since images crafted by *DiffAttack* are natural from human eyes but can lead to wrong decisions across various black-box models and defenses, some could maliciously use our methods to undermine real-world applications, inevitably raising more concerns about AI safety.

# References

[1] Shuo Feng, Haowei Sun, Xintao Yan, Haojie Zhu, Zhengxia Zou, Shengyin Shen, and Henry X Liu. Dense reinforcement learning for safety validation of autonomous vehicles. *Nature*, 615(7953):620–627, 2023. 1

[2] Zhengxia Zou, Rusheng Zhang, Shengyin Shen, Gaurav Pandey, Punarjay Chakravarty, Armin Parchami, and Henry X Liu. Real-time full-stack traffic scene perception for autonomous driving with roadside cameras. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 890–896. IEEE, 2022. 1

[3] Yilan Zhang, Fengying Xie, Jianqi Chen, and Jie Liu. Tformer: A throughout fusion transformer for multi-modal skin lesion diagnosis. *Computers in Biology and Medicine*, page 106712, 2023. 1

[4] Yilan Zhang, Fengying Xie, Xuedong Song, Yushan Zheng, Jie Liu, and Juncheng Wang. Dermoscopic image retrieval based on rotation-invariance deep hashing. *Medical Image Analysis*, 77:102301, 2022. 1

[5] Jianqi Chen, Keyan Chen, Hao Chen, Wenyuan Li, Zhengxia Zou, and Zhenwei Shi. Contrastive learning for fine-grained ship classification in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–16, 2022. 1

[6] Jianqi Chen, Keyan Chen, Hao Chen, Zhengxia Zou, and Zhenwei Shi. A degraded reconstruction enhancement-based method for tiny ship detection in remote sensing images with a new large-scale dataset. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–14, 2022. 1

[7] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 1, 3, 4

[8] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 1, 3

[9] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016. 1, 3

[10] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. 1, 3, 4

[11] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*, pages 39–57. Ieee, 2017.

[12] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2574–2582, 2016. 1, 3

[13] Yuyang Long, Qilong Zhang, Boheng Zeng, Lianli Gao, Xianglong Liu, Jian Zhang, and Jingkuan Song. Frequency domain model augmentation for adversarial attack. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IV*, pages 549–566. Springer, 2022. 1, 3, 6, 7, 9

[14] Junyoung Byun, Seungju Cho, Myung-Joon Kwon, Hee-Seon Kim, and Changick Kim. Improving the transferability of targeted adversarial examples through object-based diverse input. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15244–15253, 2022. 1, 3, 9

[15] Lianli Gao, Qilong Zhang, Jingkuan Song, Xianglong Liu, and Heng Tao Shen. Patch-wise attack for fooling deep neural network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 307–322. Springer, 2020. 1, 3, 6

[16] Yantao Lu, Yunhan Jia, Jianyu Wang, Bai Li, Weiheng Chai, Lawrence Carin, and Senem Velipasalar. Enhancing cross-task black-box transferability of adversarial examples with dispersion reduction. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 940–949, 2020. 1, 3

[17] Weibin Wu, Yuxin Su, Xixian Chen, Shenglin Zhao, Irwin King, Michael R Lyu, and Yu-Wing Tai. Boosting the transferability of adversarial samples via attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1161–1170, 2020. 1, 3

[18] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII*, pages 484–501. Springer, 2020. 3

[19] Muhammad Muzammal Naseer, Salman H Khan, Muhammad Haris Khan, Fahad Shahbaz Khan, and Fatih Porikli. Cross-domain transferability of adversarial perturbations. *Advances in Neural Information Processing Systems*, 32, 2019. 1, 3, 6, 7

[20] Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L Yuille. Improving transferability of adversarial examples with input diversity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2730–2739, 2019. 1, 3, 6, 9

[21] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4312–4321, 2019. 1, 3, 6

[22] Nathan Inkawhich, Wei Wen, Hai Helen Li, and Yiran Chen. Feature space perturbations yield more transferable adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7066–7074, 2019. 1, 3

[23] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9185–9193, 2018. 1, 3, 4, 6

[24] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 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*, pages 15–26, 2017. 3

[25] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *International Conference on Learning Representations*, 2018. 3

[26] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple black-box adversarial perturbations for deep networks. *arXiv preprint arXiv:1612.06299*, 2016. 1, 3

[27] Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. In *International Conference on Learning Representations*, 2020. 1, 3

[28] Zhengyu Zhao, Zhuoran Liu, and Martha Larson. Towards large yet imperceptible adversarial image perturbations with perceptual color distance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1039–1048, 2020. 2, 3, 6

[29] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 694–711. Springer, 2016. 2

[30] Shuai Jia, Bangjie Yin, Taiping Yao, Shouhong Ding, Chunhua Shen, Xiaokang Yang, and Chao Ma. Adv-attribute: Inconspicuous and transferable adversarial attack on face recognition. In *Advances in Neural Information Processing Systems*, 2022. 2, 3, 4

[31] Yash Sharma, Gavin Weiguang Ding, and Marcus A Brubaker. On the effectiveness of low frequency perturbations. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 3389–3396, 2019. 2

[32] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar. Diffusion models for adversarial purification. In *International Conference on Machine Learning*, pages 16805–16827. PMLR, 2022. 2, 3, 5, 6, 7

[33] Changhao Shi, Chester Holtz, and Gal Mishne. Online adversarial purification based on self-supervised learning. In *International Conference on Learning Representations*, 2021. 2

[34] Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial examples with generative models. *Advances in Neural Information Processing Systems*, 31, 2018. 2, 3

[35] Shengming Yuan, Qilong Zhang, Lianli Gao, Yaya Cheng, and Jingkuan Song. Natural color fool: Towards boosting black-box unrestricted attacks. In *Advances in Neural Information Processing Systems*, 2022. 2, 3, 4, 6, 7

[36] Haonan Qiu, Chaowei Xiao, Lei Yang, Xinchen Yan, Honglak Lee, and Bo Li. Semanticadv: Generating adversarial examples via attribute-conditioned image editing. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 19–37. Springer, 2020. 2, 3

11

[37] Anand Bhattad, Min Jin Chong, Kaizhao Liang, Bo Li, and DA Forsyth. Unrestricted adversarial examples via semantic manipulation. In *International Conference on Learning Representations*, 2020. 2, 3, 4

[38] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 2, 3, 15

[39] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 9, 15

[40] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 3

[41] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

[42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 3, 4, 6

[43] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022. 2, 3, 4, 6, 15, 16

[44] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 2, 3, 4, 15, 16

[45] Gihyun Kwon and Jong Chul Ye. Diffusion-based image translation using disentangled style and content representation. *arXiv preprint arXiv:2209.15264*, 2022. 2

[46] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. *arXiv preprint arXiv:2303.04803*, 2023. 2, 3, 5

[47] Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero-shot classifiers. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023. 2, 3, 5

[48] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2, 5, 16

[49] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. *arXiv preprint arXiv:2302.03027*, 2023. 2, 4, 16

[50] Yifeng Xiong, Jiadong Lin, Min Zhang, John E Hopcroft, and Kun He. Stochastic variance reduced ensemble adversarial attack for boosting the adversarial transferability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14983–14992, 2022. 3

[51] Xiaosen Wang and Kun He. Enhancing the transferability of adversarial attacks through variance tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1924–1933, 2021. 3

[52] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 3

[53] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 3

[54] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 3, 4, 5

[55] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3

[56] Wenbo Li, Xin Yu, Kun Zhou, Yibing Song, Zhe Lin, and Jiaya Jia. Sdm: Spatial diffusion model for large hole image inpainting. *arXiv preprint arXiv:2212.02963*, 2022. 3

[57] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. *arXiv preprint arXiv:2212.05034*, 2022. 3

[58] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3

[59] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 4, 5, 6, 9, 15

[60] Narek Tumanyan, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Splicing vit features for semantic appearance transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10748–10757, 2022. 5

[61] Eli Shechtman and Michal Irani. Matching local self-similarities across images and videos. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. 5

[62] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021. 6

[63] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 6

[64] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6

[65] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6

[66] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 6

[67] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 6

[68] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 6

[69] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 6

[70] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 6

[71] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021. 6

[72] Jinyuan Jia, Xiaoyu Cao, Binghui Wang, and Neil Zhenqiang Gong. Certified robustness for top-k predictions against adversarial perturbations via randomized smoothing. In *International Conference on Learning Representations*, 2020. 6, 7

[73] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations*, 2018. 6

[74] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1778–1787, 2018. 6

[75] Anil Thomas and Oguz Elibol. Defense against adversarial attacks-3rd place. https://github.com/anlthms/nips-2017/blob/master/poster/defense.pdf, 2017. 6, 7

[76] Alexey Kurakin, Ian Goodfellow, Samy Bengio, Yinpeng Dong, Fangzhou Liao, Ming Liang, Tianyu Pang, Jun Zhu, Xiaolin Hu, Cihang Xie, et al. Adversarial attacks and defences competition. In *The NIPS'17 Competition: Building Intelligent Systems*, pages 195–231. Springer, 2018. 6, 7

[77] F Tramèr, D Boneh, A Kurakin, I Goodfellow, N Papernot, and P McDaniel. Ensemble adversarial training: Attacks and defenses. In *6th International Conference on Learning Representations, ICLR 2018-Conference Track Proceedings*, 2018. 6, 7

[78] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 6

[79] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6

[80] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2017. 9

[81] Zhuotun Zhu, Lingxi Xie, and Alan Yuille. Object recognition with and without objects. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3609–3615, 2017. 17

# A  Overview

In this supplement, we will first provide a detailed review of DDIM Inversion in Appendix B. Then, we conduct an analysis on the effect of perturbations on the text embeddings in Appendix C. In Appendix D and Appendix E, considering the possible help for future research, we display our further trials (with little success) on improving the imperceptibility and transferability of the attacks. Finally, more quantitive studies and visualizations are shown in Appendix F.

# B  Detailed Formulation of DDIM Inversion

Diffusion Denoising Probabilistic Models (DDPMs) [39] are a class of generative models that sample images by gradually denoising an initial Gaussian noise. There is a *forward process* and a *reversed process* in DDPMs. The *forward process* is to gradually add Gaussian noise to the original image $x_0$ and thus produces a series of noisy latents $x_1, x_2, \cdots, x_T$:

$$q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{1-\beta_t}x_{t-1}, \beta_t\mathbf{I}) \tag{7}$$

where $\beta_t \in (0,1)$. When $T$ is large enough, the last latent $x_T$ will approximately follow an isotropic Gaussian distribution.

Instead of iteratively calculating the intermediate latents to get $x_t$, a good property of the *forward process* is that we can directly sample $x_t$ from $x_0$:

$$q(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)\mathbf{I}) \tag{8}$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, \ \ \epsilon \sim \mathcal{N}(0,\mathbf{I}) \tag{9}$$

where $\alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{s=0}^{t}\alpha_s$.

The *reversed process* is to draw a new sample from the distribution $q(x_0)$. Starting from $x_T \sim \mathcal{N}(0,\mathbf{I})$, we can get a new sample by iteratively sampling the posteriors $q(x_{t-1}|x_t)$. Since $q(x_{t-1}|x_t)$ is intractable due to the unknown data distribution $q(x_0)$, a neural network $p_\theta$ is trained to approximate that by predicting the mean and covariance of $q(x_{t-1}|x_t)$, which is shown to also be Gaussian distributions [38]:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t,t), \Sigma_\theta(x_t,t)) \tag{10}$$

Since $\mu_\theta(x_t,t) = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t,t)\right)$, Ho *et al.* [39] simplified the objective function by only predicting the noise $\epsilon_\theta(x_t,t)$:

$$\min_\theta \mathcal{L}(\theta) = \mathbb{E}_{x_0,\epsilon \sim N(0,I),t}\|\epsilon - \epsilon_\theta(x_t,t)\|_2^2 \tag{11}$$

After we get the trained $\epsilon_\theta(x_t,t)$, we can conduct a sampling as follows:

$$x_{t-1} = \mu_\theta(x_t,t) + \sigma_t z, \ \ z \sim N(0,I). \tag{12}$$

Since the classic DDPMs are essentially a Markov chain and they require a large timesteps $T$ to achieve good performance. To accelerate DDPMs sampling process, Song *et al.* [59] generalize DDPMs from a particular Markovian process to non-Markovian processes:

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}}\left(\frac{x_t - \sqrt{1-\bar{\alpha}_t}\epsilon_\theta(x_t)}{\sqrt{\bar{\alpha}_t}}\right) + \sqrt{1-\bar{\alpha}_{t-1}-\sigma_t^2}\cdot\epsilon_\theta(x_t) + \sigma_t z, \ \ z \sim N(0,I) \tag{13}$$

By setting $\sigma_t = 0$, we then get a deterministic sampling process (from $x_T$ to $x_0$), which is the DDIM's principle.

Since the deterministic process of DDIM can be further taken as Euler integration for solving ordinary differential equations (ODEs)[59], we can map a real image back to its corresponding latent by reversing the process. This operation, named DDIM Inversion, paves the way for later editing of real images [44, 43]. By rewriting Eq. 13, the denoising process of DDIM is as follows:

$$x_{t-1} - x_t = \sqrt{\bar{\alpha}_{t-1}}\left[\left(\sqrt{1/\bar{\alpha}_t} - \sqrt{1/\bar{\alpha}_{t-1}}\right)x_t + \left(\sqrt{1/\bar{\alpha}_{t-1}-1} - \sqrt{1/\bar{\alpha}_t-1}\right)\epsilon_\theta(x_t)\right] \tag{14}$$

We can then encode the real image into the latent space by reversing the above formulation:

$$x_{t+1} - x_t = \sqrt{\bar{\alpha}_{t+1}}\left[\left(\sqrt{1/\bar{\alpha}_t} - \sqrt{1/\bar{\alpha}_{t+1}}\right)x_t + \left(\sqrt{1/\bar{\alpha}_{t+1}-1} - \sqrt{1/\bar{\alpha}_t-1}\right)\epsilon_\theta(x_t)\right] \tag{15}$$

## C  Perturbation on Guided Text Embeddings

As mentioned in Section 3.2 in the main paper, we choose to perturb the latent $x_t$ but not the guided text $C$, which is different from the mainstream image editing approaches [44, 43, 49]. The reason is that text perturbation will be hard to transfer to other black-box models. In the following, we display the details of text perturbation designs and some necessary experiments and analyses.

### C.1  Design Details

Here we first define two text prompts: $C_1$, $C_2$, which are the first and second most possible category predicted by the classifier. We leverage $C_1$ for the optimization of unconditional embeddings mentioned in Section 3.4 in the main paper. Then, we replace $C_1$ with $C_2$ which follows [43, 48] and can expect the changes of object semantics in the image. For the loss functions, we remove $\mathcal{L}_{transfer}$ in Eq. 6 in the main paper, and modify $\mathcal{L}_{attack}$ as follows:

$$\arg\min_{C_2} \mathcal{L}_{attack} = J(x', C_2; G_\phi) \tag{16}$$

The equation above is similar to the objective function of targeted attacks, and the insight is to trick the classifier into predicting the nearest wrong label. Other implementation details are the same as Section 4.1 in the main paper.

### C.2  Experiments and Analysis

In this subsection, we compare the results between the text perturbation and the latent perturbation. From Table 5, we can observe that although the text perturbation has a slightly higher attack success in a white-box way (0.5 point accuracy lower on Inc-v3), the attack itself is hard to work on the other black-box models, thus not competitive with the latent perturbations. We attribute this phenomenon to the fact that the text perturbation is more high-level than the latent perturbation, due to text semantics. Therefore, it will tend to generate more realistic results (lower FID in Table 5), but has limited control over the local area, while the latent perturbation does the opposite.

Table 5: **Comparisons of perturbations on the latent and text.** "S." denotes surrogate models while "T." denotes target models. For the white-box attacks (surrogate model same to target one), we set their background to gray. "AVG(w/o self)" denotes the average accuracy on all the target models except the one that same as the surrogate one. The best result is bolded.

| T. / S. | CNNs | | | | | Transformers | | | | MLPs | | AVG(w/o self) | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Res-50 | VGG-19 | Mob-v2 | Inc-v3 | ConvNeXt | ViT-B | Swin-B | DeiT-B | DeiT-S | Mix-B | Mix-L | | |
| Clean | 92.7 | 88.7 | 86.9 | 80.5 | 97.0 | 93.7 | 95.9 | 94.5 | 94.0 | 82.5 | 76.5 | 89.4 | 57.8 |
| Text Perturbation | 79.6 | 73.3 | 74.7 | **13.4** | 91.3 | 85.9 | 86.9 | 87.9 | 86.0 | 71.6 | 63.8 | 80.1 | **58.8** |
| Latent Perturbation | **59.5** | **55.6** | **55.4** | 13.9 | **76.9** | **75.2** | **72.8** | **74.0** | **71.0** | **58.9** | **54.7** | **65.4** | 62.3 |

## D  "Pseudo" Mask for Better Imperceptibility

We further conduct a trial to improve the imperceptibility as mentioned in Section 5 in the main paper. Here, we display the details and experimental results.

### D.1  Design Details

As mentioned in Section 3.3 in the main paper, there is a strong relationship in the cross-attention maps between the text prompt and the image pixels. Thus, we can make use of this property to generate the true label's "pseudo" mask:

$$P = \text{Average}(\text{Cross}(x_t, t, C; \text{SDM})) \tag{17}$$

$$M_{soft} = \text{Up}\left(\frac{P}{\text{Max}(P)}\right) \tag{18}$$

$$(\text{Optional}) \quad M_{hard} = \begin{cases} 1, & M_{soft} > 0.5 \\ 0, & M_{soft} \leq 0.5 \end{cases} \tag{19}$$

Table 6: **Comparisons of different mask types and upsampling strategies.** For the white-box attacks (surrogate model same to target one), we set their background to gray. "AVG(w/o self)" denotes the average accuracy on all the target models except the one that same as the surrogate one. The best result is bolded.

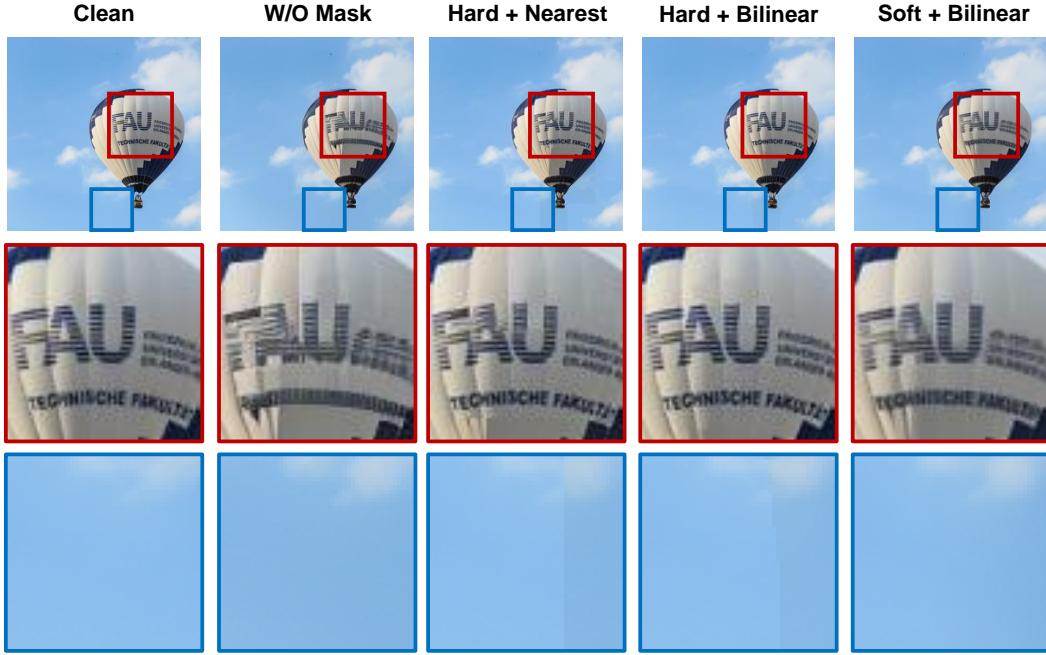| Mask Types | Upsampling Strategy | CNNs | | | | | Transformers | | | | MLPs | | AVG(w/o self) | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Res-50 | VGG-19 | Mob-v2 | Inc-v3 | ConvNeXt | ViT-B | Swin-B | DeiT-B | DeiT-S | Mix-B | Mix-L | | |
| None | None | **59.5** | **55.6** | **55.4** | **13.9** | 76.9 | 75.2 | 72.8 | **74.0** | 71.0 | 58.9 | 54.7 | **65.4** | 62.3 |
| hard | nearest | 71.4 | 67.9 | 64.8 | 17.8 | 85.2 | 80.9 | 82.9 | 80.3 | 81.0 | 68.2 | 60.2 | 74.2 | 59.1 |
| hard | bilinear | 68.8 | 66.8 | 65.6 | 18.3 | 84.0 | 79.0 | 81.4 | 79.9 | 79.5 | 66.0 | 61.8 | 73.3 | 59.3 |
| soft | bilinear | 73.9 | 69.4 | 66.9 | 18.4 | 88.2 | 82.7 | 86.5 | 84.8 | 82.0 | 68.5 | 62.0 | 76.5 | **58.8** |



Figure 6: **Visualization of the adversarial example crafted by leveraging the mask.** The second and third rows denote the scaled-up regions in the first row.

where $\mathrm{Up}(\cdot)$ is an upsampling operation to resize the cross-attention map (due to the existing downsamplings in the encoder of the Autoencoder and U-Net). $\mathrm{Max}(\cdot)$ is to extract the maximum value and normalize the cross-attention maps $P$. Since $P \geq 0$, the normalized $M_{soft} \in [0, 1]$. Eq. 19 is optional to get a hard mask. With the mask, we then filter out the background area and only apply perturbations on the foreground (area covered by true objects). The Eq. 3 in the main paper is then changed as follows:

$$\arg\min_{x_t} \mathcal{L}_{attack} = -J(x' \times M + x \times (1 - M), y; G_\phi) \tag{20}$$

The optimization details are the same as the implementation details in Section 4.1 in the main paper.

### D.2   Experiments and Analysis

Here we conduct experiments to see the impact of different upsampling strategies and different mask types. In Table 6, we display the performance when the mask is applied. It can be perceived from the results that there is an obvious trade-off between transferability and imperceptibility. The use of masks lowers the FID value, yet also lowers the attack success by a large margin. We infer that it is because the recognition of an image is not only related to its foreground but also its background [81]. Thus the attack success rate will drop when the mask is applied. We also visualize the adversarial example crafted by leveraging the mask in Figure 6, from which we can see that the applied mask can better preserve words on hot air balloon skin, and the hard mask tends to generate blocky artifacts compared with soft-mask.

Table 7: **Explorations on the effect of different categories as text prompts.** For the white-box attacks (surrogate model same to target one), we set their background to gray. "AVG(w/o self)" denotes the average accuracy on all the target models except the one that same as the surrogate one. The best result is bolded.

| top-N | CNNs | | | | | Transformers | | | | MLPs | | AVG(w/o self) | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Res-50 | VGG-19 | Mob-v2 | Inc-v3 | ConvNeXt | ViT-B | Swin-B | DeiT-B | DeiT-S | Mix-B | Mix-L | | |
| 1 | **59.5** | **55.6** | **55.4** | 13.9 | **76.9** | 75.2 | 72.8 | **74.0** | 71.0 | **58.9** | 54.7 | **65.4** | 62.3 |
| 2 | 67.0 | 63.7 | 61.2 | 9.1 | 81.7 | 77.1 | 77.5 | 79.1 | 76.4 | 64.0 | 56.8 | 70.4 | **60.9** |
| 5 | 64.4 | 61.0 | 59.2 | **5.9** | 81.4 | 78.8 | 78.0 | 77.5 | 75.9 | 61.0 | **54.0** | 69.1 | 62.8 |

# E   Trial on Further Improving Transferability

We also try to further improve the transferability by leveraging more guided text prompts as mentioned in Section 5 in the main paper, yet failed. We provide details here, considering possible help for future research.

## E.1   Design Details

In Eq. 4 in the main paper, $C =$ "{True Label / Category $1_{st}$}" that the guided text can be either the true label or the *top1* predicted category. We here extend the text to leverage more categories:

$$C_{ext} = \text{``\{Category } 1_{st}\}, \{\text{Category } 2_{nd}\}, \cdots, \{\text{Category } N_{th}\}\text{''} \tag{21}$$

where {Category $N_{th}$} denotes the name of the $N_{th}$ most possible category predicted by the classifier. Then, Eq. 3 in the main paper is modified to:

$$\mathcal{L}_{attack} = -J(x', \text{Category } 1_{st}; G_\phi) + \underbrace{(J(x', \text{Category } 2_{nd}; G_\phi) + \cdots + J(x', \text{Category } N_{th}; G_\phi))}_{N-1}$$
$$\tag{22}$$

By minimizing the above equation, the adversarial examples are crafted to lead the classification results towards the most error-prone categories, which may have benefits on the transferability (but failed through our experiments). We further add an extra loss to force the perturbed $x_t$ to have lower cross-attention intensity with {Category $1_{st}$} and higher with other text prompts, which we expect can help deceive the diffusion models:

$$P_i = \text{Average}(\text{Cross}(x_t, t, C_i; \text{SDM})) \tag{23}$$

$$\arg\max_{x_t} \mathcal{L}_{ext} = \text{Average}(\underbrace{P_2 + \cdots + P_N}_{N-1}) - P_1 \tag{24}$$

where $C_i$ denotes Category $i_{th}$, and $P_i$ denotes the cross attention between image pixels and $C_i$. Average($\cdot$) here represents the averaging operation in pixel space. We then add $\mathcal{L}_{ext}$ to Eq. 6 in the main paper with a weight factor set to 100.

## E.2   Experiments and Analysis

We here analyze the impact of different numbers of categories leveraged as text prompts. From Table 7, leveraging more guided category texts failed to improve the attack's transferability, and even damage the performance. We infer that it is because the search space of the attack is limited when we force the adversarial examples to be classified as some specific categories. When we set the category number from 2 to 5, we can observe a slight increase in the attack success, while when we set it to 1, we have no constraint on the predicted category, and thus gain a large increase in the attack success.

# F   More Quantitive Studies and Visualizations

As a supplement to the experiments in Section 4, we here reveal more experimental results about the parameter settings and also display more visualized comparisons.
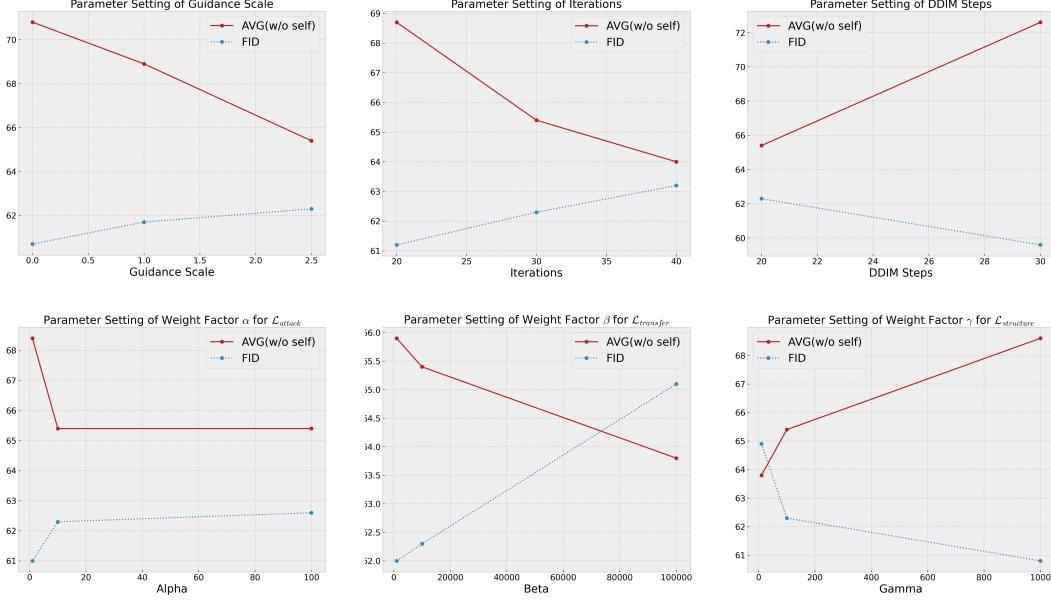
Figure 7: **The effect of different parameter settings.** We conduct a quantitative study on the parameter settings of the guidance scale, iterations, DDIM steps, and weight factors of each loss. "AVG(w/o self)" denotes the average accuracy on all the target models except the one that same as the surrogate one.

**Settings of Guidance Scale.** From Figure 7, it can be observed that with the guidance scale increased, the transferability improves while the imperceptibility deteriorates. We infer this is because larger guidance scales will tend to change the latent more and thus potentially generate more perturbations. Since there is a large gap in the attack success between the guidance scale set to 1.0 and 2.5, but a slight change of the FID value, we set the guidance scale to 2.5 finally.

**Settings of Iterations.** We can notice from Figure 7 that more iterations will sacrifice image quality for the attack success. As more iterations will consume longer optimization time, we here set the number of iterations to 30, which strikes a balance between time-consuming, image quality, and attack robustness.

**Settings of DDIM Steps.** In Figure 7, we keep the DDIM Inversion steps the same (5 inversion steps), to see the effect of different DDIM full sample steps. We do not show here the results for the step number set to 10 because the image quality is rather poor and the structure is completely changed. From the results, we can see that the step number does impact a lot both the transferability and the imperceptibility. Here we set the number of DDIM sample steps to 20, which can produce perceptually invisible adversarial samples with stronger attack robustness.

**Settings of Weight Factor for Loss.** We also conduct quantitative studies on the weight factor settings in Eq. 6 in the main paper. From Figure 7, it can be noticed that our designs of $\mathcal{L}_{transfer}$ and $\mathcal{L}_{structure}$ do make sense for improving the attack's transferability and preserving the content structure. For $\mathcal{L}_{attack}$, we can see from the results that there is a negligible performance improvement when $\alpha$ is increased to a certain extent, thus we set $\alpha$ to 10. For $\mathcal{L}_{transfer}$ and $\mathcal{L}_{structure}$, to balance both the transferability and the imperceptibility, we set them to 10000 and 100 respectively.

**More Visualizations.** We display more visual comparisons in Figure 8 and Figure 9, from which it can be observed that the adversarial examples crafted by our attack are human-imperceptible and hard to be perceived.
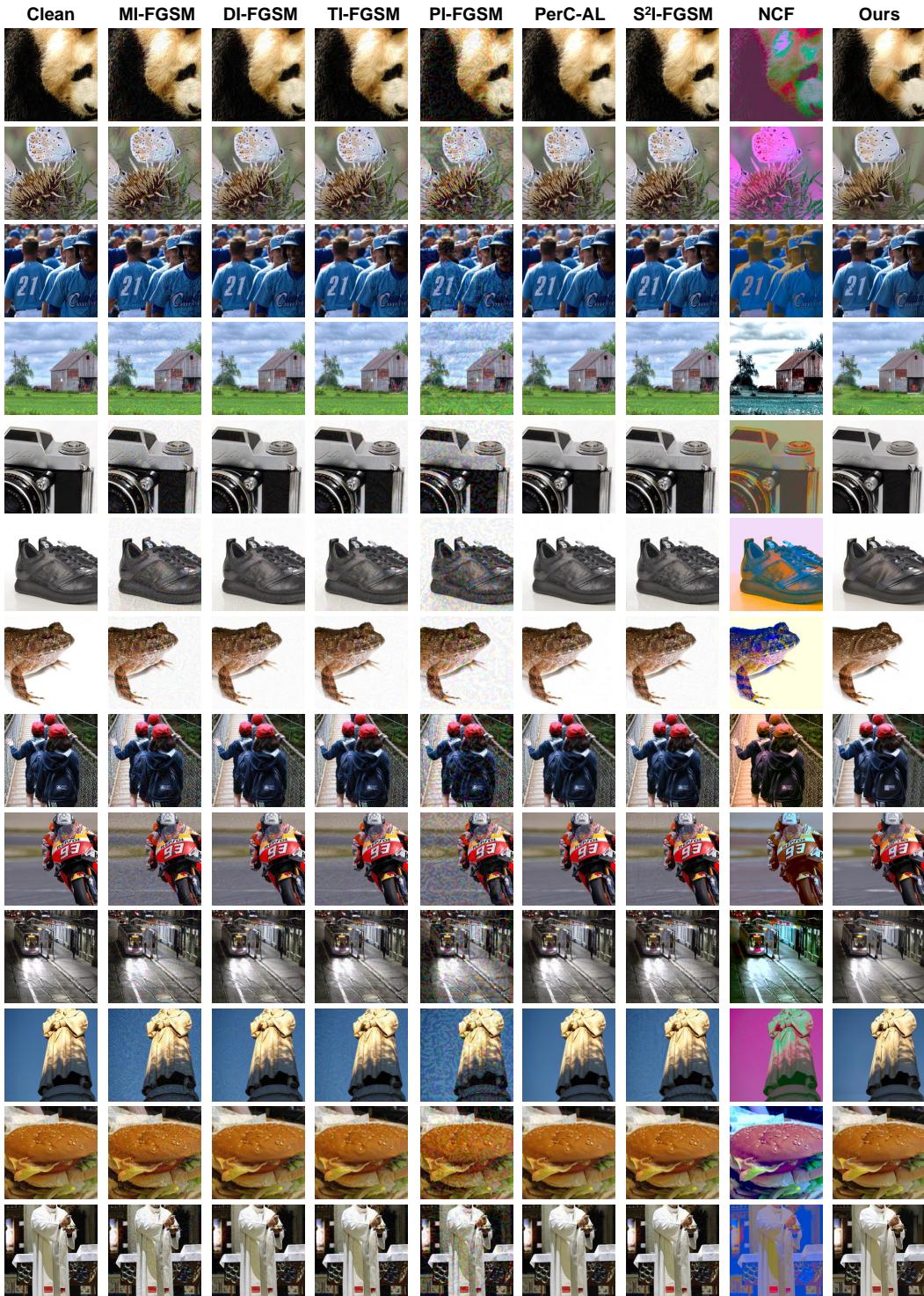
Figure 8: **Supplement visualization of adversarial examples crafted by different attacks.** Please zoom in for a better view.
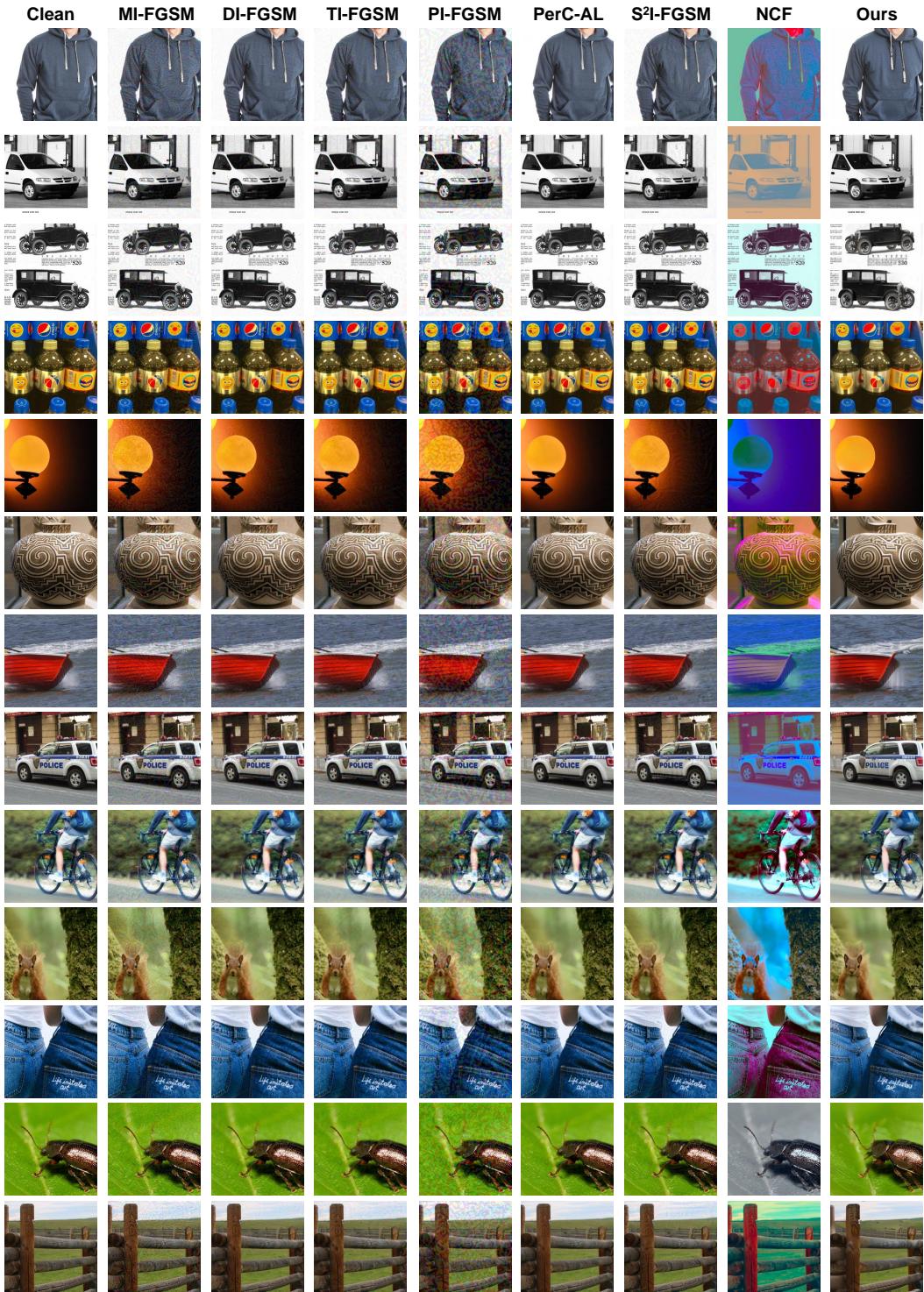
Figure 9: **Supplement visualization of adversarial examples crafted by different attacks.** Please zoom in for a better view.