
Watch out! Motion is Blurring the Vision of Your Deep Neural Networks

Qing Guo¹Felix Juefei-Xu²Xiaofei Xie^{1*}Lei Ma^{3*}Jian Wang¹Bing Yu³Wei Feng⁴Yang Liu¹

¹Nanyang Technological University, Singapore ²Alibaba Group, USA
³Kyushu University, Japan ⁴Tianjin University, China

Abstract

The state-of-the-art deep neural networks (DNNs) are vulnerable against adversarial examples with additive random-like noise perturbations. While such examples are hardly found in the physical world, the image blurring effect caused by object motion, on the other hand, commonly occurs in practice, making the study of which greatly important especially for the widely adopted real-time image processing tasks (*e.g.*, object detection, tracking). In this paper, we initiate the first step to comprehensively investigate the potential hazards of blur effect for DNN, caused by object motion. We propose a novel adversarial attack method that can generate visually natural motion-blurred adversarial examples, named motion-based adversarial blur attack (ABBA). To this end, we first formulate the kernel-prediction-based attack where an input image is convolved with kernels in a pixel-wise way, and the misclassification capability is achieved by tuning the kernel weights. To generate visually more natural and plausible examples, we further propose the saliency-regularized adversarial kernel prediction, where the salient region serves as a moving object, and the predicted kernel is regularized to achieve naturally visual effects. Besides, the attack is further enhanced by adaptively tuning the translations of object and background. A comprehensive evaluation on the NeurIPS'17 adversarial competition dataset demonstrates the effectiveness of ABBA by considering various kernel sizes, translations, and regions. The in-depth study further confirms that our method shows more effective penetrating capability to the state-of-the-art GAN-based deblurring mechanisms compared with other blurring methods. We release the code to <https://github.com/tsingqguo/ABBA>.

1 Introduction

Deep neural networks (DNN) have been widely applied in various vision perception tasks (*e.g.*, object recognition, segmentation, scene understanding), permeating many aspects of our daily life, such as autonomous driving, robotics, video surveillance, photo taking, *etc.* However, the state-of-the-art DNNs are still vulnerable against adversarial examples. Extensive previous works are proposed (*e.g.*, FGSM[1], BIM [2], MI-FGSM [3], C&W [4]), to mislead the DNN through additive noise perturbations that could be obtained by optimizing the adversarial objectives. To be imperceptible to human, L_p -norm plays an important role in such attacks, confining the perturbation noise to be small. However, the random-noise like perturbation often does not pose imminent threats to the camera systems, which does not usually occur in natural environment. Thus, some recent attempts [2] were made to physically fashion adversarial examples such as by putting up stickers or printed patterns on the physical stop sign, *etc.* Again, these artifacts are often intentionally prepared for adversarial attacks, which do not found ‘naturally’ in the real-world environment either.

*Xiaofei Xie and Lei Ma are corresponding authors (xfxie@ntu.edu.sg, malei@ait.kyushu-u.ac.jp).

While the blurring effect caused object motion commonly occurs in practical image perception systems, the potential hazards of motion blur effect to the DNN are largely untouched so far. Motion blur naturally happens during the exposure time of image capturing. When an object moves at a relatively high speed, all information of the object during the image capture process is integrated, constituting a blur-like image along a relative moving direction. Compared with other kinds of image blur (e.g., defocus blur caused by using unsuitable camera focus), motion blur is directly related to the motion of object and camera, whose effect cannot be easily removed by adjusting the camera’s setting. As a result, motion blur almost coexists with the camera and potentially posts serious effects on DNN perception-based systems. However, up to present, there are limited studies discuss how motion blur affects the DNN perception tasks. It is not even clear whether and what kinds of motion blur can systematically mislead a DNN.

In this paper, we initiate first step to comprehensively investigate blur effects to DNNs from the adversarial attack perspective, where systematic motion blur-based adversarial example discovery would be an important step towards further DNN enhancement. In particular, we propose a new type of adversarial attack, termed motion-based adversarial blur attack (ABBA), which can generate visually natural and plausible motion-blurred adversarial examples. We first formulate the kernel-prediction-based attack where an input image is convolved with kernels in a pixel-wise way, and the misclassification ability is guided by systematically kernel weight tuning. In order to generate more natural motion blurred examples, we also propose the saliency-regularized adversarial kernel prediction, where the salient region serves as a moving object, and the predicted kernel is regularized to achieve visually natural effects. Besides, our method could easily adjust the blur effects of differ exposure time during the image capturing in the real world, by adaptively tuning translations of object and background.

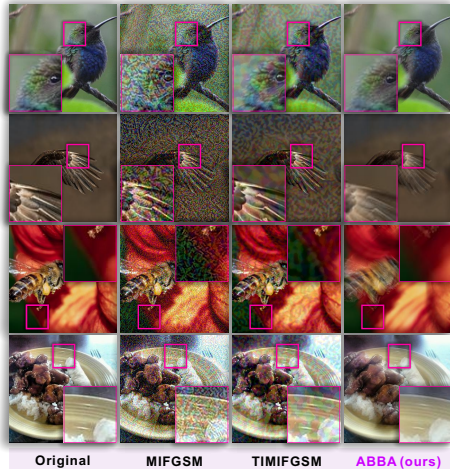


Figure 1: Four adversarial examples of MIFGSM [3], TIMIFGSM [5] and our ABBA. MIFGSM and TIMIFGSM produces apparent noise on all four cases. Our ABBA generates visually natural motion blur. All adversarial images fool the Inception-v3 model.

We perform comprehensive evaluation on the effectiveness of the proposed ABBA, benchmarked against various noise-based attacks on both attack success rates and transferability. The main contributions of this work can be summarized as follows. ❶ To the best of our knowledge, we make the very first attempt to investigate kernel-based adversarial attack. ❷ We propose a motion-based adversarial blur attack as a new type of attack mode, to be added to the adversarial attack family. ❸ In order to produce more visually plausible blur attack, we introduce a saliency regularizer that forces consistent blur patterns within the boundary of the objects (or background in some cases). ❹ Compared with state-of-the-art (SOTA) additive-noise based adversarial attacks and common blur techniques, our proposed method achieves better attack success rate and transferability. ❺ Furthermore, our proposed method has demonstrated higher penetration capability against the SOTA GAN-based deblur mechanism, compared to normal image motion blur.

Related Work. Since the discovery of adversarial examples to attack a deep neural network (DNN) both theoretically [1] and physically [2], there has been extensive research towards developing adversarial attack and defense mechanisms. The basic iterative method (BIM) [2], the C&W method [4], the fast gradient signed method (FGSM) [1], and the momentum iterative fast gradient sign method (MI-FGSM) [3], etc., are a few popular ones among early adopters in the research community. Building upon these ideas, researchers have been continuously pushing the envelope in many ways. For example, serious attempts have been made to integrate momentum term into the iterative process for the attacks [3]. By doing so, the momentum can help stabilize the update directions, begetting more transferable adversarial examples and posing more threats to adversarially trained defense mechanisms. More recently, [5] proposes to optimize the noise perturbation over an ensemble of translated images, making the generated adversarial examples more robust against white-box models being attacked while achieving better transferability. The mainstream adversarial attack is an additive noise pattern that is learnable given the model parameters under a white-box setting. Perhaps the prevalence is partially due to the fact that the adversarial noise with the ‘addition’ operation is

relatively straightforward to optimize for. Of course, there are many other ways to alter a benign image beyond the addition operation that are all potential candidates for coming up with new types of adversarial attack modes. One caveat of additive noise attack is the lack of balance between being visually plausible and imperceptible while having high attack success rate. Usually, it has to compromise one for the other. Researchers are looking beyond additive noise attack to seek novel attack modes that strike a better balance between visual plausibility and performance, *e.g.*, multiplicative attack [6], deformation attack [7, 8, 9], and semantic manipulation attack [10].

We are proposing a new type of motion-based adversarial blur attack that can generate visually natural and plausible motion-blurred adversarial examples, inspired from kernel prediction such as [11, 12, 13] and motion blur generation [14]. One desired property of the proposed method is the immunity and robustness against the SOTA deblurring techniques (*e.g.*, [15, 16]). In an effort to understand black-box DNN better, an image saliency paradigm is proposed [17] to learn where an algorithm looks at, by discovering which parts of an image most affect output score when perturbed in terms of Gaussian blur, replacing the region with constant value and injecting noise. The localization of the blur region is performed through adaptive iteration whilst ours is saliency regularized that leads to a visually plausible motion-blurred image. The major difference is that their Gaussian blur kernel is fixed while ours is learnable to maximally jeopardize the image recognition DNNs.

2 Methodology

2.1 Background: Additive-Perturbation-Based Attack

Let \mathbf{X}^{real} be a real example, *e.g.*, images in the ImageNet dataset, and y denotes its ground truth label. A classifier denoted as $f(\mathbf{X}) : \mathcal{X} \rightarrow \mathcal{Y}$ predicts the label of \mathbf{X}^{real} . An attack method aims to generate an adversarial example denoted as \mathbf{X}^{adv} that can fool the classifier to predict an incorrect label with imperceptible perturbation. Existing attack methods mainly focus on the additive adversarial perturbation that is added to the real example to get $\mathbf{X}^{\text{adv}} = g(\mathbf{X}^{\text{real}}, \delta) = \mathbf{X}^{\text{real}} + \delta$, where δ is generated by maximizing a loss function $J(\mathbf{X}^{\text{adv}}, y)$ with a constrained term:

$$\arg \max_{\delta} J(\mathbf{X}^{\text{real}} + \delta, y) \quad \text{subject to} \quad \|\delta\|_p \leq \epsilon_a, \quad (1)$$

For example, gradient descent is widely employed by many methods to generate adversarial examples, *e.g.*, FGSM, BIM, MIFGSM, DIM and TIMIFGSM.

2.2 ABBA_{pixel}: Kernel-Prediction-Based Adversarial Attack

Besides ‘+’, there are various techniques that can perform advanced image transformation for different objectives, *e.g.*, Gaussian filter for image denoising, Laplacian filter for image sharpening, and guided filter for edge-preserving smooth [18], which are all kernel-based techniques, processing each pixel of the image with a hand-crafted or guided kernel. In general, compared with the addition, kernel-based operation can handle more complex image processing tasks via different kinds of kernels. More recently, several works [19, 12, 13] found that the kernel weights can be carefully predicted for advanced tasks with high performance, *e.g.*, high quality noise-free rendering and video frame interpolation. Inspired by these works, in this paper, we propose the kernel-prediction-based attack. Specifically, we process each pixel (*i.e.*, $\mathbf{X}_p^{\text{real}}$) of a real example \mathbf{X}^{real} with a kernel \mathbf{k}_p ,

$$\mathbf{X}_p^{\text{adv}} = g(\mathbf{X}_p^{\text{real}}, \mathbf{k}_p, \mathcal{N}(p)) = \sum_{q \in \mathcal{N}(p)} \mathbf{X}_q^{\text{real}} k_{pq}, \quad (2)$$

where p denotes the p th pixel in \mathbf{X}^{adv} and \mathbf{X}^{real} , *i.e.*, $\mathbf{X}_p^{\text{adv}}$ and $\mathbf{X}_p^{\text{real}}$, respectively. $\mathcal{N}(p)$ is a set of N pixels in \mathbf{X}^{real} and $p \in \mathcal{N}(p)$. The kernel \mathbf{k}_p also has the size of N and determines the weights of N pixels in $\mathcal{N}(p)$. In general, we have $\sum_{q \in \mathcal{N}(p)} k_{pq} = 1$ to ensure the generated image lies within the respective neighborhood of the input image, where a softmax activation is often adopted for this requirement [19]. To better understand Eq. (2), we discuss with two simple cases: 1) when we let $\mathcal{N}(p)$ be a neighborhood of the pixel p and the kernel of each pixel is a fixed Gaussian kernel, \mathbf{X}^{adv} is the Gaussian-blurred \mathbf{X}^{real} . Similarly, we can obtain defocus-blurred image with disk kernels. 2) when we set $\max(\mathbf{k}_p) = k_{pp}$ and $\forall q, k_{pq} \neq 0$, the perturbation of \mathbf{X}^{adv} becomes more imperceptible as length of \mathbf{k}_p decreases. To achieve high attack success rate, we need to optimize kernels of all pixels independently, *i.e.*, $\mathcal{K} = \{\mathbf{k}_p | \forall p \in \mathbf{X}^{\text{real}}\}$, according to the loss function of AI-related tasks, *e.g.*, image classification, and constrained terms

$$\arg \max_{\mathcal{K}} J(\{\sum_{q \in \mathcal{N}(p)} \mathbf{X}_q^{\text{real}} k_{pq}\}, y) \quad \text{subject to} \quad \forall p, \|\mathbf{k}_p\|_0 \leq \epsilon, \max(\mathbf{k}_p) = k_{pp}, \sum_{q \in \mathcal{N}(p)} k_{pq} = 1, \quad (3)$$

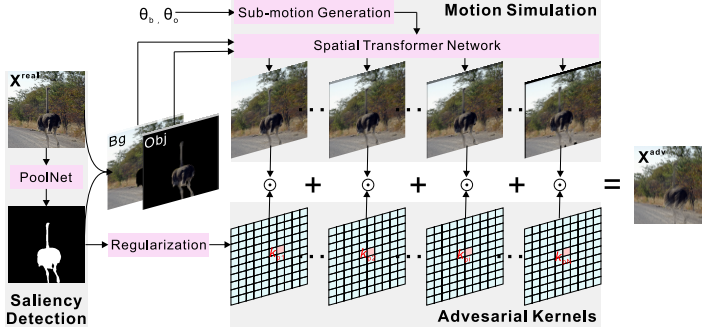


Figure 3: Pipeline of our motion-based adversarial blur attack, *i.e.*, Eq. (4). First, We use the PoolNet [20] to extract the salient object in the image and obtain the object and background regions. Then, the translation parameters θ_o and θ_b are divided to N parts to simulate the motion process and generate N images with the spatial transformer network [21]. Finally, we get the adversarial example by adding these images with adversarial kernels as weights for each pixel. The adversarial kernels and translation parameters are tuned to realize effective attack by optimizing Eq. (5). \odot denotes the element-wise product.

where $\|\mathbf{k}_p\|_0$ represents the number of valid kernel elements (*i.e.*, $\{k_{pq} \neq 0\}$) and $\epsilon \in [1, N]$ controls the upper bound of $\|\mathbf{k}_p\|_0$. When $\epsilon = 1$, we have $\mathbf{X}^{\text{adv}} = \mathbf{X}^{\text{real}}$, and when $\epsilon = N$, the perturbation would be the most serious case. We can calculate the gradient of the loss function with respect to all kernels, to realize the gradient-based attack. As a result, the attack method can be integrated into any gradient-based additive-perturbation attack methods, *e.g.*, FGSM, BIM, MIFGSM. Since the kernel-prediction-based adversarial attack tunes each pixel’s kernel independently, it can achieve a significantly high attack success rate, but generating unnatural images that are easily perceptible. For an easier reference, we name this method $\text{ABBA}_{\text{pixel}}$. As shown in Fig. 2, $\text{ABBA}_{\text{pixel}}$ distorts the original inputs and produces additive noise-like results. To reach the balance between high attack success rate and natural visual effect, we propose to regularize the kernels to produce visually natural motion blur via the guidance of a visual saliency map.



Figure 2: From left to right: original image, adversarial examples generated by the kernel-prediction-based attack ($\text{ABBA}_{\text{pixel}}$), and motion-based adversarial blur attack (ABBA).

2.3 ABBA: Motion-Based Adversarial Blur Attack

Motion blur is a frequently occurring effect during image capture. We first introduce how to generate visually natural motion-blurred adversarial examples by regularizing the kernel-prediction-based attack. Then, we describe the workflow of our proposed attack in Fig. 3 for better understanding.

2.3.1 Saliency-Regularized Adversarial Kernel Prediction

Motion blur is often generated during the exposure time by integrating the light from a moving object. To synthesize the motion blur, we need to know where the object is and specify how it moves. To this end, given an image, we first use the SOTA saliency detection method, *i.e.*, PoolNet [20], to extract the salient object \mathbf{S} from \mathbf{X}^{real} and assume it is moving at the time when capturing the image. The saliency map \mathbf{S} is a binary image, indicating the salient object region (*i.e.*, $\mathbf{X}^{\text{real}} \odot \mathbf{S}$) and background region (*i.e.*, $\mathbf{X}^{\text{real}} \odot (1 - \mathbf{S})$), as shown in Fig. 3. Then, we specify translation transformations to the object and background, respectively. We denote them as $T(\mathbf{X}^{\text{real}} \odot \mathbf{S}, \theta_o)$ and $T(\mathbf{X}^{\text{real}} \odot (1 - \mathbf{S}), \theta_b)$ that are simplified as \mathbf{X}^{S, θ_o} and $\mathbf{X}^{1-S, \theta_b}$, where θ_o and θ_b are the translation parameters² for the object and background, respectively.

Since motion blur is the integration of all light during the object moving process, we divide the motion represented by θ_o and θ_b into N sub-motions (corresponding to the kernel size N in Eq. (2)) to simulate blur generation. The sub-motions are represented by $\{i\Delta\theta_o, |i \in [1, N]\}$ and $\{i\Delta\theta_b, |i \in [1, N]\}$, where $\Delta\theta_o = \theta_o/N$ and $\Delta\theta_b = \theta_b/N$. Then, we redefine Eq. (2) as

$$\mathbf{X}_p^{\text{adv}} = g(\mathbf{X}_p^{\text{real}}, \mathbf{S}, \mathbf{k}_p, \mathcal{N}(p)) = \sum_{q=\mathcal{N}(p, i), i \in [1, N]} (\mathbf{X}_q^{S, i\Delta\theta_o} + \mathbf{X}_q^{1-S, i\Delta\theta_b}) k_{pq}, \quad (4)$$

where $\mathcal{N}(p)$ is a set of the p th pixel in all translated examples. $\mathbf{X}^{S, i\Delta\theta_o}$ and $\mathbf{X}^{1-S, i\Delta\theta_b}$ denote the object and background images translated $i\Delta\theta_o$ pixels. Compared with the attack in Sec. 2.2, the perturbation amplitude is affected by the kernel and translation parameters. The objective function is

² θ_o and θ_b are 2D vectors in the range of $[0, 1]$ and represent the rates of x-axis and y-axis shifting distance of the object&background regions w.r.t. image size.

defined as

$$\begin{aligned} & \arg \max_{\mathcal{K}, \theta_o, \theta_b} J(\{ \sum_{\substack{q \in \mathcal{N}(p, i) \\ i \in [1, N]}} (\mathbf{X}_q^{S, i \Delta \theta_o} + \mathbf{X}_q^{1-S, i \Delta \theta_b}) k_{pq} \}, y) \quad (5) \\ & \text{subject to } \forall p, \|\mathbf{k}_p\|_0 \leq \epsilon, \max(\mathbf{k}_p) = k_{pp}, \sum_{q \in \mathcal{N}(p)} k_{pq} = 1 \\ & \forall p, q, \mathbf{k}_p = \mathbf{k}_q, \text{ if } \mathbf{S}(p) = \mathbf{S}(q), \|\theta_o\|_\infty \leq \epsilon_\theta, \|\theta_b\|_\infty \leq \epsilon_\theta. \end{aligned}$$

where $\epsilon_\theta \in [0, 1]$ controls the maximum translations of the object/background. Here, we use the spatial transformer network [21] to perform translation according to θ_o and θ_b , enabling the gradient propagate to all kernels. There are two main differences about the constrained terms *c.f.* Eq. (3): (1) The translation parameters are added to guide the generation of the adversarial example; (2) The kernels are set to be the same within the same region, which is needed to generate visually natural motion blur, since pixels in the object region usually have the same motion. As shown in Fig. 2, by incorporating saliency and motion regularization, the ABBA’s adversarial example looks visually more natural than the one by $\text{ABBA}_{\text{pixel}}$.

2.3.2 Attacking Algorithm

We summarize the workflow of our attacking algorithm in the following steps: 1) Calculate the saliency map of an image, *i.e.*, \mathbf{X}^{real} , via PoolNet and obtain \mathbf{S} . 2) Initialize $\theta_{o,t} = \theta_{b,t} = [0, 0]$ and set each kernel $\mathbf{k}_{p,t}$ of \mathcal{K}_t by $\{k_{pp,t} = 1, k_{pq,t} = 0 | \forall q \in \mathcal{N}(p), q \neq p\}$ where $t = 0$ denotes the first iteration. 3) Calculate $\mathbf{X}_t^{\text{adv}}$ via Eq. (4), which is also visualized in Fig. 3 for better understanding. 4) Calculate the gradient of $\mathbf{X}_t^{\text{adv}}$ with respect to the objective function and obtain $\nabla_{\mathbf{X}_t^{\text{adv}}} J(\mathbf{X}_t^{\text{adv}}, y)$. 5) Propagate the gradient through the spatial transformer network and obtain the gradients of \mathcal{K}_t , $\theta_{o,t}$, and $\theta_{b,t}$, *i.e.*, $\nabla_{\mathcal{K}_t} J(\mathbf{X}_t^{\text{adv}}, y)$, $\nabla_{\theta_{o,t}} J(\mathbf{X}_t^{\text{adv}}, y)$ and $\nabla_{\theta_{b,t}} J(\mathbf{X}_t^{\text{adv}}, y)$. 6) Update \mathcal{K}_t , $\theta_{o,t}$, and $\theta_{b,t}$ with a step size. 7) Update $t = t + 1$ and go to the Step 3) for further optimization until it reaches the maximum iteration or $\mathbf{X}_t^{\text{adv}}$ fools the DNN. We will detail our settings in Sec. 3.1.

2.4 $\text{ABBA}_{\text{physical}}$: Towards Real-World Adversarial Blur Attack

As introduced in Sec. 2.3.2, ABBA takes a real example as the input and produces an adversarial blur example, the estimated kernels (\mathcal{K}^*), and translation parameters (θ_o^* and θ_b^*). Then, it posts an interesting problem whether we could use the estimated translation parameters to guide camera or object moving, in generating a real-world adversarial blur examples. There are three main challenges : 1) We cannot control kernels’ values in the real world. The optimized kernels are to let the adversarial examples fool deep models and may not exist in the real world. 2) It is difficult to precisely control the object or camera’s moving without a high-precision robot arm. 3) To transfer the image translation to camera translation, we need know the object depth and camera intrinsic parameters.

To alleviate these challenges, we conduct the following modifications of our ABBA: 1) we fix the kernels \mathcal{K} to be average kernels (*i.e.*, we set each kernel’s elements as $\frac{1}{N}$ where N is the kernel size) that let the generated adversarial blur follow the real-world motion blur³. 2) We force the object and background to share the translation parameters. As a result, ABBA simulates the blur generation through camera moving, where object and background have the same motion thus we can produce blurred examples by moving the camera. 3) In the real world, we could use a RGB-D camera to get an object’s depth and leverage a calibration software to obtain the intrinsic parameters of the camera.

Based on these, we can generate a real-world adversarial blur example by: 1) capturing a picture containing an object in a real-world scene. 2) using our ABBA to calculate θ_o^* or θ_b^* . 3) calculating the camera translations according to object depth and camera intrinsic parameters. 4) moving camera according to the camera translations and take a real blur picture during the moving process as the output. In particular, we will validate $\text{ABBA}_{\text{physical}}$ in Sec. 3.4 through the AirSim simulator [24] within which we can control a simulated camera precisely and obtain the depth map. We also conduct a experiment with a mobile phone to primarily verify our method in the real world.

3 Experimental Results

In this section, we conduct comprehensive experiments to demonstrate the effectiveness of our proposed method, by investigating the following five research questions: 1) Is the transferability of

³We use the average kernels to simulate the real motion blur generation process, which has been used to construct training dataset for deblurring methods [22, 23].

our proposed method across models comparable or even better than SOTA attack methods? 2) Could our method produce visually natural examples? 3) Could we generate adversarial blurred examples that possibly occur in the real world? 4) As a blurring method, could SOTA deblurring methods easily defend our attack? 5) How is attacking success rate affected by hyper-parameters, *e.g.*, ϵ , ϵ_θ , motion direction, and blur region?

3.1 Experimental Settings

Dataset and Models. We use NeurIPS’17 adversarial competition dataset [25], compatible with ImageNet, for all the experiments. To validate our method’s transferability, we consider four widely used models, *i.e.*, Inception v3 (Inc-v3) [26], Inception v4 (Inc-v4), Inception ResNet v2 (IncRes-v2) [27], and Xception [28]. We further compare on four defense models: Inc-v3_{ens3}, Inc-v3_{ens4}, and IncRes-v2_{ens} from [29] and high-level representation guided denoiser (HGD) [30] with the highest ranking in NeurIPS’17 defense competition. We report the results of Inc-v3 in the paper, and put more results of other models (*e.g.*, Inc-v4, IncRes-v2, Xception), in the supplementary material.

Baselines. We consider two kinds of baselines. The first scope is SOTA additive-perturbation-based attacks, *e.g.*, FGSM [1], MIFGSM [3], DIM [31], TIFGSM, TIMIFGSM, and TIDIM [5]⁴, and interpretation-based noise [17]. The second kind contains three blur-based methods including the interpretation-based blur [4], Gaussian blur [32] and Defocus blur. For the transferability comparison in Sec. 3.2, we follow default settings in [5] for the first group attacks. For all iterative attack methods including ours, we set the iteration number to be 10. For blur-based baselines, we set the standard variation of Gaussian blur and the kernel size of Defocus blur to be 15.0, which is the same with our method for a fair comparison. For the image quality comparison in Sec. 3.2, we tune the hyper-parameter of all attacks to cover the Succ. Rate from low to high and show the relationship between image quality and Succ. Rate. This helps to see if our method could maintain high attack Succ. Rate and transferability while keeping image natural.

Setup of ABBA. In the experimental part of Sec. 3.2, we implement our methods by setting the hyper-parameters, *i.e.*, maximum translation ϵ_θ to 0.4 and maximum valid kernel size ϵ to 15.0 with the iteration 10. The step sizes are set as 0.04 and 1.5 for updating the kernels and translation parameters, respectively. Such setups are around the medium values among the range of our hyper-parameters, well balancing the attack success rate and visual quality. We discuss the effect of ϵ_θ and ϵ in Sec. 3.6. In addition, we have four variants of our method, *i.e.*, ABBA_{pixel}, ABBA_{obj}, ABBA_{bg}, and ABBA_{image}. The first one is the attack introduced in Sec. 3.6 while the rest three methods blur different regions of input images, *e.g.*, ABBA_{obj} only adds motion blur to the object region by fixing the kernels of background pixels to be $\{k_{pp} = 1, k_{pq} = 0 | S(p) = 0, q \in \mathcal{N}(p), q \neq p\}$ and ABBA_{image} adds motion blur to the whole image while forcing object and background to share the kernels and translations.

Metrics. We use the success rate (Succ. Rate), *i.e.*, the rate of adversarial examples that fool attacked DNNs, to evaluate the effectiveness of attacks. Regarding quality of generated adversarial examples, we use BRISQUE [33] instead of L_p norms, PSNR, or SSIM due to the following reasons: 1) our ABBA cannot be fairly evaluated by the L_1 , L_2 , L_∞ , PSNR, or SSIM metrics since the perturbation is not additive and no longer well aligned pixel-to-pixel. 2) BRISQUE [33] is a natural scene statistic-based distortion-generic blind/no-reference (NR) image quality assessment and widely used to evaluate losses of ‘naturalness’ in the image due to the presence of distortions including additive noise and blur. It should be note a more natural image often has a smaller BRISQUE value.

3.2 Comparison with Baselines on Transferability

Tab. 1 summarize the comparison results. We discuss them from two aspects: 1) the comparison with additive-perturbation-based attacks. 2) the advantages over blur-based methods. For the first aspect, compared with all additive-perturbation-based attacks, our methods, *i.e.*, ABBA and ABBA_{pixel}, achieve the highest success rate on all defence models, demonstrating the higher transferability of ABBA over baselines. Compared with two SOTA methods, *i.e.*, DIM and TIDIM, ABBA_{pixel} achieves slightly lower success rate when attacking Inc-v4 and IncRes-v2 while obtaining higher results than TIDIM when attacking Xception. In summary, our methods ABBA_{pixel} and ABBA have competitive transferability with SOTA additive-perturbation-based attacks while achieving significant advantages in attacking defence models. For the second aspect, ABBA achieves higher success rate than GaussBlur, DefocusBlur, and Interp_{blur} on all normal trained models and defence models. We also compare them for attacking object or background regions. Obviously, the success rates of all

⁴We use the released code in github.com/dongyp13/Translation-Invariant-Attacks to get the results of FGSM, MIFGSM, DIM, TIFGSM, TIMIFGSM, and TIDIM.

	Attacking Results (Inc-v3)				Defence Results (Inc-v3)			
	Inc-v3	Inc-v4	IncRes-v2	Xception	Inc-v3 _{em3}	Inc-v3 _{em4}	IncRes-v2 _{em}	HGD
GaussBlur	34.7	22.7	18.4	26.1	23.6	23.8	19.3	16.9
GaussBlur _{obj}	13.6	6.0	5.2	7.1	8.6	7.8	6.3	4.6
GaussBlur _{bg}	18.8	10.8	9.2	12.0	13.0	13.1	10.9	8.7
DefocBlur	30.0	16.8	11.1	18.8	17.5	18.3	15.0	12.9
DefocBlur _{obj}	10.0	3.0	2.9	3.6	5.2	4.6	3.8	2.7
DefocBlur _{bg}	16.9	9.2	7.0	10.5	10.1	10.3	9.2	7.8
Interpblur	34.7	3.6	0.5	3.4	7.1	7.1	4.3	1.4
ABBA _{obj}	21.0	4.9	4.2	7.0	10.1	10.5	8.3	4.9
ABBA _{bg}	30.9	11.6	10.1	12.9	1.2	0.8	1.2	0.5
ABBA _{image}	62.4	29.8	28.8	34.1	43.2	43.8	38.9	28.4
ABBA _{apixel}	89.2	65.5	65.8	71.2	69.8	72.5	68.0	63.1
ABBA	65.6	31.2	29.7	33.5	46.6	48.7	41.2	31.0
Interpnoise	95.8	20.5	18.6	22.9	16.8	16.1	9.4	3.3
FGSM	79.6	35.9	30.6	42.1	15.6	14.7	7.0	2.1
MIFGSM	97.8	47.1	46.4	47.7	20.5	17.4	9.5	6.9
DIM	98.3	73.8	67.8	71.6	24.2	24.3	13.0	9.7
TIFGSM	75.4	37.3	32.1	38.6	28.2	28.9	22.3	18.4
TIMIFGSM	97.9	52.4	47.9	44.6	35.8	35.1	25.8	25.7
TIDIM	98.5	75.2	69.2	61.3	46.9	47.1	37.4	38.3

Table 1: Adversarial comparison results on NeurIPS’17 adversarial competition dataset according to the success rate. The adversarial examples are generated from Inc-v3. There are two comparison groups. For the first one, we compare blur-based methods, *i.e.*, Interpretation-based blur (Interp_{blur}), GaussBlur, and DefocusBlur with our **ABBA** by considering the effects of attacking different regions, *i.e.*, object or background regions, of inputs. In addition to above methods, the second group comparison contains additive-perturbation-based attacks, *i.e.*, Interpretation-based noise (Interp_{noise}) [17], FGSM [1], MIFGSM [3], DIM [31], and TIFGSM, TIMIFGSM, and TIDIM [5]. We highlight the top three results with pink, yellow, and blue, respectively.

methods decrease significantly when we add blur to only object or background regions. Besides these transferability results, we also add an analysis in the supplementary material about the interpretable explanation of the high transferability of our method by comparing with FGSM and MIFGSM.

3.3 Comparison with Baselines on Image Quality

We conduct an analysis about the attack success rate and image quality (*i.e.*, measured by BRISQUE [33] that evaluates losses of ‘naturalness’ in the image). Note that, smaller BRISQUE corresponds to more natural images. We also consider two comparison groups: 1) the blur-based attacks, *i.e.*, GaussBlur and DefocBlur, and 2) SOTA additive-perturbation-based attacks, *i.e.*, FGSM, MIFGSM, DIM, TIFGSM, TIMIFGSM, and TIDIM. For each compared method, we tune their hyper-parameters⁵ to cover success rates from low to high on the NeurIPS’17 adversarial competition dataset. For each hyper-parameter, we calculate the average BRISQUE of the adversarial images that successfully fool DNNs. As a result, for each attacked model, we can draw a plot for an attack method and the points more near the top left corner are better (*i.e.*, looks more naturally while fooling DNNs). As shown in Fig. 4, in general, the quality of images generated by all baseline methods gradually gets worse as their success rate becomes larger. In contrast, the BRISQUE of our ABBA always stay at a small value even if the success rate increases, demonstrating that ABBA can produce visual natural adversarial examples with high attack success rate. When transferring the adversarial examples to other models, we observe that success rate of additive-perturbation-based attacks, especially the FGSM, MIFGSM, TIFGSM, and DIM, decrease sharply, while the blur-based attacks are not impacted so largely. Comparing with SOTA methods (*i.e.*, TIDIM and TIMIFGSM) on the transferability, we see that adversarial examples of ABBA and the two methods have similar BRISQUE scores when their success rate is small. When success rate of TIDIM and TIFGSM further increases, their BRISQUE values become smaller than ours.

3.4 Adversarial Blur Examples in the Simulation World and Real World

Results in AirSim environment. We validate our ABBA_{physical} on the AirSim simulator [34] that supports hardware-in-loop with cameras physically and visually realistic simulations.

AirSim also provides APIs to retrieve relevant data (*e.g.*, real-time depth map) and control cameras in a platform independent way, which meets the requirements of ABBA_{physical} introduced in Sec. 2.4. In particular, we choose the open-source Neighborhood environment [35], containing 70 cars with various styles. For each car, we select a observation view where a DNN (*e.g.*, Inc-v3) could classify it correctly. Then, we conduct the ABBA_{physical} experiment for each car with the following steps: 1) Set a camera to capture a original image of the car at the selected observation view. 2) Use our method introduced in Sec. 2.4 to calculate an adversarial blurred image (*i.e.*, ABBA’s result) and the camera translation with the depth map and camera intrinsic parameters. 3) Move the camera along a straight line to the translation destination while taking N pictures and averaging them to the final

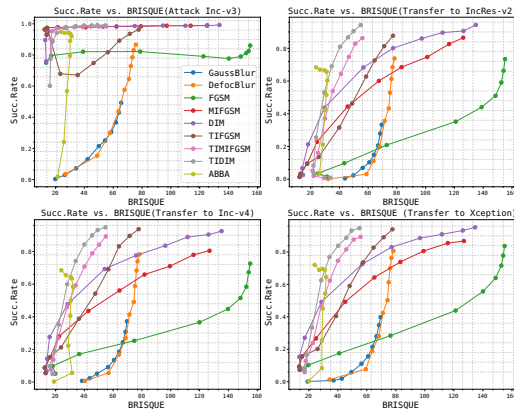


Figure 4: Succ. Rate vs. BRISQUE. Note that, smaller BRISQUE corresponds to more natural images.

⁵For FGSM, MIFGSM, DIM, TIFGSM, TIMIFGSM, and DIM, we tune their success rate by changing their maximum perturbation, *i.e.*, ϵ_a in Eq. (1). We tune the blur-based attacks by tuning the kernel size from 3 to 21. For ABBA, we change the ϵ and ϵ_θ from 5 to 50 and 0.1 to 1, respectively.

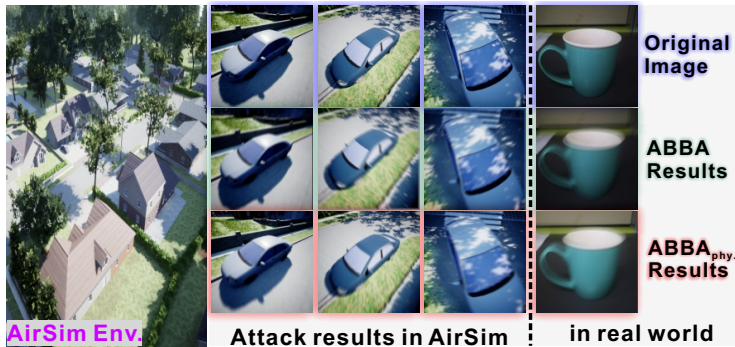


Figure 5: Attack results in the AirSim environment (medium) and real world (right). The left sub-figure shows a snapshot of the neighborhood environment. The results of $ABBA_{\text{physical}}$ are produced by physically moving the camera and mobile phone in the AirSim and real world, respectively. All the adversarial examples fool the Inc-v3 model.

blurred image (*i.e.*, $ABBA_{\text{physical}}$'s result), which is equivalent to the motion blur generation process [22]. 4) Test whether the real blurred image and adversarial blurred image could fool DNNs (*i.e.*, Inc-v3, Inc-v4, IncRes-v2, and Xception in Table 2) successfully. We summarize the attack success rates in Table 2 and observe that: 1) $ABBA_{\text{physical}}$ generating adversarial blur examples by physically moving cameras can achieve high success rate on Inc-v3, which demonstrates that there might exist real-world motion blur can fool DNNs easily. 2) Although our $ABBA$ has very high transferability across DNNs, the physical adversarial examples cannot fool other DNNs easily.

Results in the real world environment. We further perform a preliminary experiment to validate our $ABBA_{\text{physical}}$ in the real world through a mobile phone: 1) we capture a sharp image with a mobile phone, *i.e.*, the cup in the first row of Fig. 5. 2) we use $ABBA$ to generate an adversarial blur image with the Inc-v3 model, *i.e.*, the second row of Fig. 5, and obtain the image translation parameters that indicate the mobile phone's moving direction and distance. 3) we move the mobile phone along the direction indicated by the translation parameters and shoot a real-blurred image in the same scene with longer exposure time. We find both blurred images are misclassified by the Inc-v3 model. Here, we ignore the cup's depth since we shoot the sharp and real-blurred images at almost the same position and empirically tune the moving distance of the mobile phone. Such operation could be replaced by high-precision robot arms in the future.

Table 2: Success rate of $ABBA$ and $ABBA_{\text{physical}}$ for attacking DNNs in the AirSim environment.

	Adversarial examples from Inc-v3			
	Inc-v3	Inc-v4	IncRes-v2	Xception
Succ. Rate of $ABBA$	97.6	87.8	82.9	90.2
Succ. Rate of $ABBA_{\text{physical}}$	85.3	9.7	14.6	11.0

3.5 Effect of Deblurring Methods

Here, we discuss the effect of SOTA deblurring methods to our adversarial examples and 'normal' motion blurred images. The 'normal' motion blur is usually synthesized by averaging neighbouring video frames [22], which is equivalent to set the kernel weights as $\frac{1}{N}$ to average the translated examples in Eq. (4). We can regard such normal motion blur as an attack and add them to all images in the testing data. We use DeblurGAN [15] and DeblurGANv2 [16] to handle our adversarial examples and the normal motion blurred images and calculate the relative decrease of the success rate, *i.e.*, $r = \frac{s-s'}{s}$ where s and s' represent the success rate before and after deblurring. Smaller r means the attack is more resilient against deblurring methods. As shown in Fig. 6, in general, compared with the normal motion blurred images, our adversarial examples are harder to be handled by the state-of-the-art deblurring methods. For DeblurGANv2, the relative decrease r of normal motion blur is usually larger than 0.5 when kernel size is in [15, 35]. This means DeblurGANv2 can effectively remove the normal motion blur and improve the classification accuracy. In contrast, the r of our method is always smaller than the r of normal motion blur, and gradually decreases in the range of [15, 35], indicating that it is more difficult to use DeblurGANv2 to defend our attack as kernel size becomes larger. Similar results are obtained on DeblurGAN where the difference of r between our method and normal motion blur is much smaller than that on DeblurGANv2.

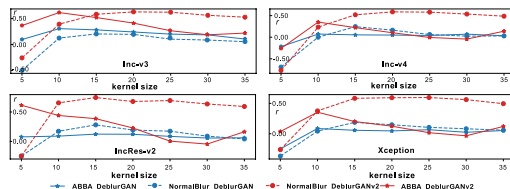


Figure 6: The relative decrease of attack success rate before and after deblurring. Two state-of-the-art deblurring methods, *i.e.*, DeblurGAN [15] and DeblurGANv2 [16], are used to deblur the our adversarial blurred examples and normal motion blurred images.

3.6 Hyper-parameter Analysis and Ablation Study

Effect of ϵ and ϵ_θ : We calculate the success rate of our method with different ϵ and ϵ_θ in Eq. (5) and find that the success rates become gradually higher with the increase of ϵ and ϵ_θ . More results and discussions can be found in the supplementary material. **Effect of motion directions:** we fix ϵ_θ and ϵ and tune the motion direction of object and background by setting different x-axis and y-axis translations to see the success rate variation. We find that the success rate reaches the highest value around 45° moving direction. We visualize and discuss the results in the supplementary material. **Effect of blurred regions and importance of adaptive translations.** We conduct an ablation study by adding motion blur to different regions of an image. The success rate results are reported in Tab. 1. Compared with $ABBA_{\text{pixel}}$, $ABBA$ makes good balance between the attack success rate and visual effects. Compared with other variants, $ABBA$ that jointly tunes the object and background translations can obtain much better transferability across normal trained and defense-based models. We show and discuss their visualization results in the supplementary material.

4 Conclusions

In this paper, we initiated the first step to comprehensively investigate the potential hazards of motion blur for DNNs. We propose the kernel-prediction-based attack that can fool DNNs with a high success rate, which is further regularized by visual saliency to make the motion-blurred image visually more natural. Besides, we also validate that our adversarial blur might indeed exist in the real world. Comprehensive evaluation demonstrates the usefulness of our methods. Our results call for the attention of future research that takes motion blur effects into consideration during real-time image perception DNN designs. We also hope our work facilitates more general solutions of robust DNN towards addressing common camera motion blur effects. Moreover, we have demonstrated that the proposed kernel-prediction-based adversarial attack can be extended to other kinds of attacking methods, *e.g.*, adversarial attacks based on denoising [36], raining [37], camera exposure [38], and bias field in the medical imaging [39]. We will also discuss the effects of motion blur to visual object tracking [40, 41, 42], through the proposed $ABBA$ and recent attacking method [43] against tracking.

Acknowledgments and Disclosure of Funding

We appreciate the anonymous reviewers for their valuable feedback. This research was supported by Singapore National Cybersecurity R&D Program No. NRF2018NCR-NCR005-0001, National Satellite of Excellence in Trustworthy Software System No. NRF2018NCR-NSOE003-0001, NRF Investigatorship No. NRFI06-2020-0022. It was also supported by JSPS KAKENHI Grant No.20H04168, 19K24348, 19H04086, and JST-Mirai Program Grant No. JPMJMI18BB, Japan, and the National Natural Science Foundation of China (NSFC) under Grant 61671325, Grant 61572354, Grant 61672376, and Grant U1803264. We gratefully acknowledge the support of NVIDIA AI Tech Center (NVAITC) to our research.

5 Broader Impact

In this work, we make an early attempt to investigate the motion-blur effects to DNNs, which is a common phenomenon in the real-world image capturing process of a camera. We present the very *first* attack based on manipulating the motion blur of the images. Through comprehensive experiments, we have demonstrated that very successful attacks can be well disguised in naturally-looking motion blur patterns of the image, unveiling the system vulnerabilities of any image capture stage that deals with camera or object motions.

Considering that image capturing and sensing is an integral and essential part of almost every computer vision application that interacts with the real world, the message we are trying to convey here is an important one, *i.e.*, attackers can intentionally make use of motion blur, either by tampering with the camera sensing hardware or the image processing software to embed such an attack. Even unintentionally, the motion blur effects still commonly exist in the real-world application, posing threats to the DNNs behind the camera. This work is the first attempt to identify and showcase that such as attack based off image motion blur is not only feasible, but also leads to high attack success rate while simultaneously maintaining high realism in the image motion blur patterns. In a larger sense, this work can and will provide new thinking into how to better design the image capturing pipeline in order to mitigate potential risk caused by the vulnerabilities discussed herein, especially for mission- and safety-critical applications that are involved with moving objects or

moving sensors such as autonomous driving scenarios, mobile face authentication with a hand-held device, computer-aided diagnostics in medical imaging, robotics, *etc.*

Bad actors can potentially make use of this newly proposed attack mode as a wheel to pose risks on existing imaging systems that are not yet prepared for this new type of attack and effect based on image motion blur. We, as researchers, believe that our proposed method can accelerate the research and development of the DNN resilient mechanism against such motion blur effects. Therefore, our work can serve as an asset and a stepping stone for future-generation trustworthy design of computer vision DNNs and systems.

In addition to the societal impact discussed above, the proposed method can also influence various research directions. For example, our proposed ABBA method:

- hints new data augmentation technique for training powerful DNN-based deblurring methods.
- hints new DNN design, detection/defense techniques to be resilience against motion blur effects.
- hints new direction of analyzing the effect of motion blur to video analysis tasks, *e.g.*, real-time visual object detection, tracking, segmentation, and action recognition.

A Supplementary Material

In the submitted manuscript, we have reported the attack results of Inc-v3 on four normal trained models and four defense models, and compared with 14 attack instances on the transferability and the image quality. In this supplementary material, (1) we present the evaluation details of our method regarding the transferability on five more defense models, comparing the visualization results with SOTA attacks, and discussing the attack results of another three DNNs. (2) We also conducted an in-depth hyper-parameter analysis and ablation study of our method, and posted an interpretable explanation about the difference between our method and baselines on the transferability. (3) Finally, we validated the generalization of our method by attacking an STN-based CNN.

Overall, the results of this supplementary material further demonstrated that the proposed adversarial blur attack can fool DNNs effectively while generating visually natural blurred images. All experimental results and discussions infer that motion blur as a common effect in the real world has a high risk of fooling SOTA DNNs and our attack methods initiate the first step to study the potential hazards of motion blur for DNNs.

A.1 Attack Results on Eight Defense Models

Besides the results on the four defense models reported in the submission, we also compared our method with baselines on another five defense models including R&P [44], NeurIPS-r3⁶, and three models from the state-of-the-art feature denoise-based (FD) defense method [45] (*i.e.*, ResNetXt101 with all denoising (FD_{R101}), ResNet152 with four denoising blocks (FD_{R152}), and adversarial trained baseline model ResNet152 (FD_{R152B})). The R&P method transforms input images through random resizing and padding, which ranked the second in the NeurIPS 2017 defense competition. NeurIPS-r3 is the third rank submission of NeurIPS 2017 defense competition and combines adversarial trained VGG16, Inc-v3, IncRes-v2, and ResNet152v2 models in an ensemble way. Besides, NeurIPS-r3 also performs transformations, *i.e.*, shear, shifting, zoom, rotation, JPEG compression, and noise corruption, on input images. The FD method ranked the first in Competition on Adversarial Attacks and Defenses (CAAD)-2018.

In the Table 1 of our submission, we report the defense results of baselines, *i.e.*, FGSM, MIFGSM, DIM, TIFGSM, TIMIFGSM, and TIDIM that use L_2 norm bound and show low transferability. Here, we report the results of above baselines with L_∞ norm bound and show the comparison results in Table I. Note that, all other results of above baselines in our submission and supplementary material are based on L_∞ norm bound.

As reported in Table I, our method, *i.e.*, ABBA_{pixel}, achieves the highest transferability across all defense models and ABBA has competitive results with the state-of-the-art baseline TIDIM. Such results reveal a potential big shortcoming of existing studies of defense methods, *i.e.*, only considering the adversarial noise while ignoring other potential factors in physical environment. Note that,

⁶<https://github.com/anlthms/nips-2017>

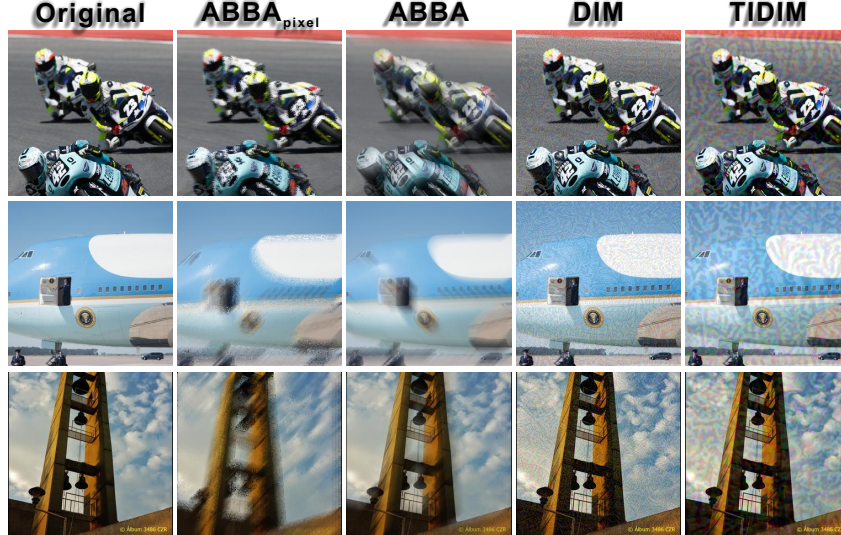


Figure 1: Three visualization results of $ABBA_{\text{pixel}}$, $ABBA$, DIM, and TIDIM. All adversarial examples mislead the Inc-v3 model.

compared with adversarial noise, motion blur frequently happens in our daily life and widely exists among various computer vision-based applications, thus its influence to DNNs should be carefully studied and addressed.

Table 1: Adversarial comparison results on NeurIPS’17 adversarial competition dataset according to the success rate. We use nine defense models to evaluate all attacks. The adversarial examples are generated from Inc-v3. There are two comparison groups. For the first one, we compare blur-based methods, *i.e.*, Interpretation-based blur ($Interp_{\text{blur}}$), GaussBlur, and DefocusBlur with our $ABBA$ by considering the effects of attacking different regions, *i.e.*, object or background regions, of inputs. In addition to above methods, the second group comparison contains additive-perturbation-based attacks, *i.e.*, Interpretation-based noise ($Interp_{\text{noise}}$) [17], FGSM [1], MIFGSM [3], DIM [31], and TIFGSM, TIMIFGSM, and TIDIM [5]. We highlight the top three results with pink, yellow, and blue, respectively.

	Defence Results (Adv. Examples from Inc-v3)								
	Inc-v3 _{env3}	Inc-v3 _{env4}	IncRes-v2 _{ens}	HGD	R&P	NeurPIS-r3	FD _{R101}	FD _{R152}	FD _{R152B}
GaussBlur	23.6	23.8	19.3	16.9	17.2	17.6	35.6	35.8	35.9
GaussBlur _{obj}	8.6	7.8	6.3	4.6	4.8	5.1	13.9	13.9	14.6
GaussBlur _{bg}	13.0	13.1	10.9	8.7	10.0	9.3	19.5	19.2	20.1
DefocBlur	17.5	18.3	15.0	12.9	14.6	14.2	31.1	30.9	31.1
DefocBlur _{obj}	5.2	4.6	3.8	2.7	3.3	2.9	10.8	10.5	11.1
DefocBlur _{bg}	10.1	10.3	9.2	7.8	9.0	8.1	19.5	17.6	18.5
$Interp_{\text{blur}}$	7.1	7.1	4.3	1.4	2.9	2.9	25.5	25.8	28.6
$ABBA_{\text{obj}}$	10.1	10.5	8.3	4.9	6.2	7.1	18.7	18.4	19.1
$ABBA_{\text{bg}}$	1.2	0.8	1.2	0.5	0.6	0.7	43.5	44.1	45.5
$ABBA_{\text{image}}$	43.2	43.8	38.9	28.4	34.1	35.0	61.1	61.9	62.4
$ABBA_{\text{pixel}}$	69.8	72.5	68.0	63.1	65.0	65.7	79.6	81.0	82.1
$ABBA$	46.6	48.7	41.2	31.0	36.7	38.5	64.2	64.6	65.6
$Interp_{\text{noise}}$	16.8	16.1	9.4	3.3	4.1	4.4	39.6	41.4	46.8
FGSM	15.6	14.7	7.0	2.1	6.5	9.8	39.2	41.4	45.3
MIFGSM	20.5	17.4	9.5	6.9	8.7	12.9	39.0	40.2	44.6
DIM	24.2	24.3	13.0	9.7	13.3	18.0	39.1	40.3	45.1
TIFGSM	28.2	28.9	22.3	18.4	19.8	24.5	39.7	41.8	45.4
TIMIFGSM	35.8	35.1	25.8	25.7	23.9	26.7	39.3	41.2	45.8
TIDIM	46.9	47.1	37.4	38.3	36.8	41.4	40.0	42.2	45.8

A.2 Visualization Comparison with Baselines

We show several adversarial examples of $ABBA_{\text{pixel}}$, $ABBA$, DIM, and TIDIM in Fig. 1, Fig. VIII and Fig. IX. All examples can mislead the Inc-v3 model.

Obviously, our method $ABBA$ can generate visually natural motion-blurred examples on various objects and these examples are very similar to real images captured by real-world cameras where the motion blur is caused by object or camera moving. In contrast, the adversarial examples of DIM and TIDIM have obvious unreal patterns. The noise-like pattern of DIM is drastically different from the natural noise usually caused by the camera sensor, *e.g.*, Gaussian noise. The perturbation pattern of TIDIM is more perceptible than that of DIM, although TIDIM achieves much higher transferability than DIM. Compared with $ABBA$, our another method, *i.e.*, $ABBA_{\text{pixel}}$, breaks local pattern of the original input. However, $ABBA_{\text{pixel}}$ ’s examples look more imperceptible than TIDIM’s results. More comparison results are shown in Fig. VIII and Fig. IX.



Figure II: Comparison between adversarial-blurred images and blurred images for training deblurring models.

Table II: Adversarial comparison results on NeurIPS’17 adversarial competition dataset. There is no available Xception model based on the author’s implementations [5] of baselines, *i.e.*, FGSM, MIFGSM, DIM, TIFGSM, TIFMIFGSM, and TIDIM. Hence, we leave these baselines’ results empty for the Xception model. We highlight the top three results with pink, yellow, and blue, respectively.

	Attacking Results (Inc-v3)				Attacking Results (Inc-v4)				Attacking Results (IncRes-v2)				Attacking Results (Xception)			
	Inc-v3	Inc-v4	IncRes-v2	Xception	Inc-v3	Inc-v4	IncRes-v2	Xception	Inc-v3	Inc-v4	IncRes-v2	Xception	Inc-v3	Inc-v4	IncRes-v2	Xception
GaussBlur	34.7	22.7	18.4	26.1	14.2	26.7	10.9	17.2	12.1	11.8	20.1	13.8	16.1	15.7	11.9	32.5
GaussBlur _{obj}	13.6	6.0	5.2	7.1	3.5	9.5	2.2	3.9	3.2	2.8	6.4	2.7	3.7	3.4	2.6	10.9
GaussBlur _{bg}	18.8	10.8	9.2	12.0	6.7	13.4	5.5	7.3	6.7	6.5	11.8	6.8	7.6	7.1	6.3	16.3
DefocBlur	30.0	16.8	11.1	18.8	18.7	36.2	13.2	22.3	15.8	14.7	23.4	17.4	18.5	18.7	12.9	36.8
DefocBlur _{obj}	10.0	3.0	2.9	3.6	3.9	10.3	3.1	4.4	3.8	3.2	7.5	3.2	4.4	4.4	2.6	11.8
DefocBlur _{bg}	16.9	9.2	7.0	10.5	10.4	20.1	8.3	12.8	8.9	9.4	15.3	10.1	10.2	11.4	8.4	21.6
InterPblur	34.7	3.6	0.5	3.4	2.7	26.7	0.8	3.1	3.1	3.1	20.1	3.4	3.0	3.1	0.8	32.5
ABA _{obj}	21.0	4.9	4.2	7.0	11.6	28.9	9.7	11.5	11.2	11.9	29.0	12.7	9.1	9.6	7.7	30.2
ABA _{bg}	30.9	11.6	10.1	12.9	14.0	31.7	13.3	15.7	14.0	14.0	25.8	13.2	12.5	14.3	11.4	33.3
ABA _{image}	62.4	29.8	28.8	34.1	32.0	66.7	28.8	36.2	33.0	30.7	63.4	37.0	28.9	28.4	26.1	66.7
ABA _{pixel}	89.2	65.5	65.8	71.2	77.7	88.1	71.3	76.0	81.8	78.3	92.0	80.6	74.0	67.5	66.8	86.2
ABA	65.6	31.2	29.7	33.5	39.5	74.9	37.3	43.2	38.4	38.6	71.6	44.2	32.3	35.2	35.9	73.1
InterPnoise	95.8	20.5	15.6	22.9	5.2	92.6	1.6	6.0	6.7	6.0	91.8	8.3	3.5	2.3	0.4	93.4
FGSM	79.6	35.9	30.6	42.1	43.1	72.6	32.5	45.2	44.3	36.1	64.3	45.4				
MIFGSM	97.8	47.1	46.4	47.7	67.1	98.8	54.3	58.5	74.8	64.8	100.0	61.7				
DIM	98.3	73.8	67.8	71.6	81.8	98.2	74.2	79.1	86.1	83.5	99.1	80.8				
TIFGSM	75.4	37.3	32.1	38.6	45.3	68.1	33.7	39.4	49.7	41.5	63.7	44.0				
TIFMIFGSM	97.9	52.4	47.9	44.6	68.6	98.8	55.3	50.8	76.1	69.5	100.0	59.9				
TIDIM	98.5	75.2	69.2	61.3	80.7	98.7	73.2	65.5	86.4	85.5	98.8	71.0				

Besides above visualization results, we further conduct an experiment to compare our adversarial blur images with the blur images for training deblurring models [15]. Specifically, given a sharp image, *e.g.*, the left sub-figures in Fig. II, we use ABBA to generate corresponding adversarial blur images and compare them with the blur images for training. Obviously, both blur looks realistic, which demonstrates the capability of ABBA to generate visually natural blur images.

A.3 Attack Results of Inc-v3, Inc-v4, IncRes-v2, and Xception

Besides the attack results of Inc-v3 reported in our submission, we further show the results of Inc-v4, IncRes-v2, and Xception in Table II. Note that, there is no available Xception model based on the authors’ implementations [5] of FGSM, MIFGSM, DIM, TIFGSM, TIFMIFGSM, and TIDIM. Hence, we leave these baselines’ results empty for the Xception model. Similar to the results of Inc-v3, for the transferability results, our method, *i.e.*, ABBA_{pixel}, achieves slightly lower success rate than the state-of-the-art additive-perturbation-based attacks, *i.e.*, DIM and TIDIM, when attacking Inc-v3, Inc-v4, and IncRes-v2, and obtains higher success rate than TIDIM when attacking the Xception model. For the whitebox attacks, TIMIFGSM and MIFGSM usually achieve the highest success rate.

A.4 Hyper-parameter Analysis and Ablation Study

Effect of ϵ and ϵ_θ . We calculate the success rate of our method with different ϵ and ϵ_θ in the Eq. (5) of our submission, respectively. Specifically, we try ϵ with the range $[5, N]$ where $N = 51$ and ϵ_θ in $[0, 1]$. As shown in Fig. III (a), the success rates become gradually higher with the increase of ϵ and ϵ_θ . The highest success rates are 94.8%, 68.5%, 68.4%, and 72.1% on Inc-v3, Inc-v4, IncRes-v2, and

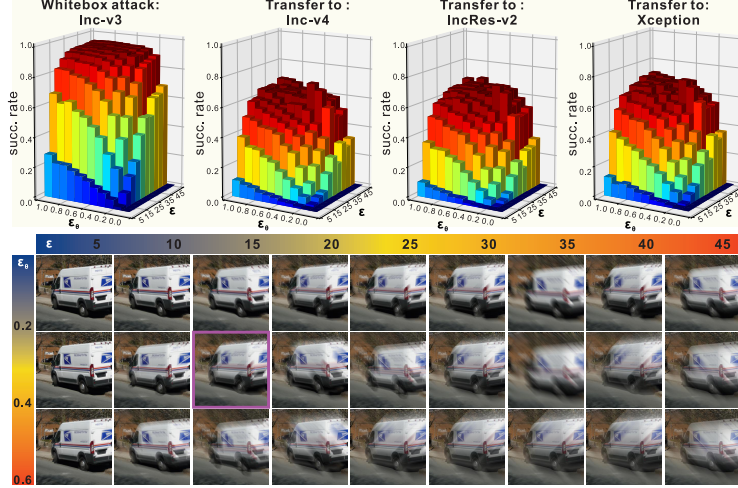


Figure III: Up: shows the success rate of ABBA w.r.t. the variation of both ϵ and ϵ_θ in Eq. (5) of our submission where ϵ is within $[5, 50]$ with step size 5 and ϵ_θ is in $[0, 1]$ with step size 0.1. Down: shows an example of .

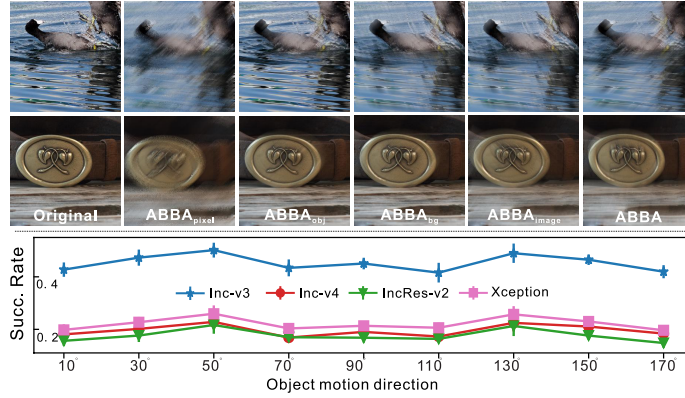


Figure IV: Up: two examples of $ABBA_{\text{pixel}}$, $ABBA_{\text{obj}}$, $ABBA_{\text{bg}}$, $ABBA_{\text{image}}$, and $ABBA$. Bottom: Success rates of our method with respect to the object motion directions.

Xception, respectively. We also visualize adversarial examples of an image that has been successfully attacked on all $\epsilon > 0$ and $\epsilon_\theta > 0$. Obviously, as ϵ and ϵ_θ increase, the visual effects of adversarial examples gradually become worse and the perturbations are more easily perceived. According to numerous attacking on different images, we choose the $\epsilon = 15.0$ and $\epsilon_\theta = 0.4$ to balance the success rate and visual effects when comparing with baselines on transferability in Sec. 3.2 in our submission.

Effect of motion directions. we fix $\epsilon_\theta = 0.4$ and $\epsilon = 15.0$ and tune the motion direction of object and background by setting different x-axis and y-axis translations. For each object motion direction, we calculate the mean and standard variation of the success rates on different background moving directions. As shown in Fig. IV (B), the success rate increases as the object motion direction becomes larger in $[10^\circ, 50^\circ]$ while decreasing as the direction is smaller in $[50^\circ, 70^\circ]$. The success rate variation has symmetrical trend in the range of $[90^\circ, 170^\circ]$. Such results are mostly caused by the L_∞ used for constraining the translation. The motion direction is directly related to the translation and the success rate reaches the highest value around 45° .

Effect of blurred regions and importance of adaptive translations. As reported in Tab. 1 in our submission and cases shown in Fig. IV (U), $ABBA_{\text{pixel}}$ achieves the highest attack success rate and transferability among all variants, which, however, changes the original image obviously and looks unnatural. $ABBA_{\text{obj}}$ and $ABBA_{\text{bg}}$ have the worst success rate on all models although they tend to generate visually natural motion blur. $ABBA_{\text{image}}$ and $ABBA$ make good balance between the attack success rate and visual effects. In particular, $ABBA$ that jointly tunes the object and background translations can obtain much better transferability across normal trained and defense-based models. Note that, when compared with the results using fixed motion directions in Fig. IV (B), $ABBA$ obtains the highest success rate among all motion direction, further demonstrating usefulness of adaptive translations.

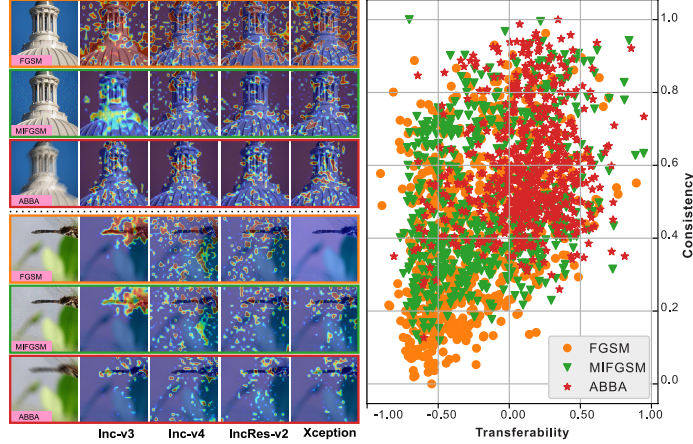


Figure V: Left: the interpretable maps of six adversarial examples generated by FGSM, MIFGSM, and ABBA, respectively, with four models. Right: the transferability & consistency distributions of adversarial examples generated by the three attacks.

A.5 Interpretable Explanation of the Transferability

In the following, we explore the difference between ABBA, FGSM, and MIFGSM on the transferability. Note that, we implement FGSM and MIFGSM on the same platform (*i.e.*, pytorch with foolbox 2.3.0) with ABBA for fair comparison. We modify the method in [17] that generates an interpretable map for a classification model $f(\cdot)$ with a given perturbation. Then, we observe that the transferability of an adversarial example generated by an attack correlates with the consistency of interpretable maps of different models. Specifically, given an adversarial example \mathbf{X}^{adv} generated by an attack and the original image \mathbf{X}^{real} , we can calculate an interpretable map \mathbf{M}^f for $f(\cdot)$ by optimizing:

$$\arg \min_{\mathbf{M}^f} f_y(\mathbf{M}^f \odot \mathbf{X}^{\text{adv}} + (1 - \mathbf{M}^f) \odot \mathbf{X}^{\text{real}}) + \lambda_1 \|\mathbf{M}^f\|_1 + \lambda_2 \text{TV}(\mathbf{M}^f) \quad (6)$$

where $f_y(\cdot)$ denotes the score at label y that is the ground truth label of \mathbf{X}^{real} and $\text{TV}(\cdot)$ is the total-variation norm. Intuitively, optimizing Eq. (6) is to find the region that causes misclassification. We optimize Eq. (6) via gradient decent in 150 iterations and fix $\lambda_1 = 0.05$ and $\lambda_2 = 0.2$. We can calculate four interpretable maps for each adversarial example based on four models, *i.e.*, Inc-v3, Inc-v4, IncRes-v2, and Xception, as shown in Fig. V(L).

We observe that the interpretable maps of our method have similar distributions across the four models while the maps of FGSM and MIFGSM do not exhibit this phenomenon. To further validate this observation, we calculate the standard variation across the four maps at each pixel and get a value by mean pooling. We normalize the value and regard it as the consistency measure for the four maps. As shown in Fig. V(R), the consistency of adversarial examples of our method is generally higher than that of FGSM and MIFGSM. We further study the transferability of an adversarial example across models. Given an adversarial example from Inc-v3 and a model $f(\cdot)$, we calculate a score to measure the transferability under this model: $t^f = f_c(\mathbf{X}^{\text{adv}}) - f_y(\mathbf{X}^{\text{adv}})$ where $c \neq y$ is the label having maximum score among non-ground-truth labels. If $t^f > 0$ means the adversarial example fool $f(\cdot)$ successfully, and vice versa. As shown in Fig. V(R), the transferability of adversarial examples of our method is generally higher than that of FGSM and MIFGSM.

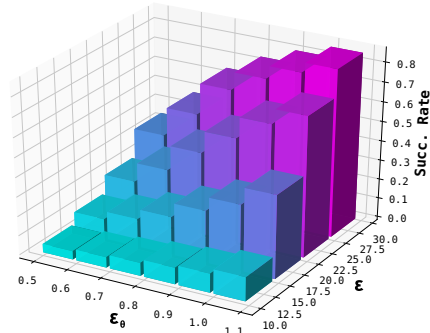


Figure VI: Success rate of ABBA w.r.t. different ϵ_θ and ϵ that are maximum translations and maximum number of valid kernel elements.

A.6 Attack results of STN-based model

As introduced in Sec. 2.3 in the submission, we employ the spatial transformer network (STN) to tune the translation parameters of the object and background, and it may post a question if our method could also be useful in attacking STN-based CNN models. Towards more comprehensive



Figure VII: Visualization examples of ABBA for attacking STN-based CNN with two group hyper-parameters, *i.e.*, $\epsilon_\theta = 0.5$, $\epsilon = 15.0$ and $\epsilon_\theta = 1.0$, $\epsilon = 15.0$.

evaluation and answering this question, we further conduct an experiment to attack the STN-based CNN implemented by [46] on the MNIST dataset. *First*, we use our ABBA to attack the STN-based CNN with different hyper-parameters, *i.e.*, the maximum translations (ϵ_θ) and the maximum number of valid kernel elements (ϵ). As shown in Fig. VI, the success rate of our method gradually increases as the ϵ_θ and ϵ become larger, which demonstrates the effectiveness of our method. *Second*, we show several examples of our attack results in Fig. VII with two groups of hyper-parameters, *i.e.*, $\epsilon_\theta = 0.5$, $\epsilon = 15.0$ and $\epsilon_\theta = 1.0$, $\epsilon = 15.0$. As shown in Fig. VII, with the small $\epsilon_\theta = 0.5$, our method can generate slightly blurred handwritten digits that look naturally but fool the STN-based CNN. When using larger $\epsilon_\theta = 1.0$, the adversarial images become more blurred but still natural.

References

- [1] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [2] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv:1607.02533*, 2017.
- [3] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li. Boosting adversarial attacks with momentum. In *CVPR*, pages 9185–9193, 2018.
- [4] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017.
- [5] Y. Dong, T. Pang, H. Su, and J. Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *CVPR*, pages 4307–4316, 2019.
- [6] Shaokai Ye, Sia Huat Tan, Kaidi Xu, Chenglong Bao Yanzhi Wang, and Kaisheng Ma. Brain-inspired reverse adversarial examples. In *arXiv:1905.12171*, 2019.
- [7] Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Song. Spatially transformed adversarial examples. In *ICLR*, 2018.
- [8] Rima Alaifari, Giovanni S. Alberti, and Tandri Gauksson. ADef: an iterative algorithm to construct adversarial deformations. In *ICLR*, 2019.
- [9] Run Wang, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, Yihao Huang, and Yang Liu. Amora: Black-box adversarial morphing attack. *arXiv preprint arXiv:1912.03829*, 2019.
- [10] Anand Bhattad, Min Jin Chong, Kaizhao Liang, Bo Li, and David Forsyth. Unrestricted adversarial examples via semantic manipulation. In *ICLR*, 2020.
- [11] B. Mildenhall, J. T. Barron, J. Chen, D. Sharlet, R. Ng, and R. Carroll. Burst denoising with kernel prediction networks. In *CVPR*, pages 2502–2510, 2018.
- [12] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive convolution. In *CVPR*, pages 2270–2279, 2017.
- [13] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive separable convolution. In *ICCV*, pages 261–270, 2017.

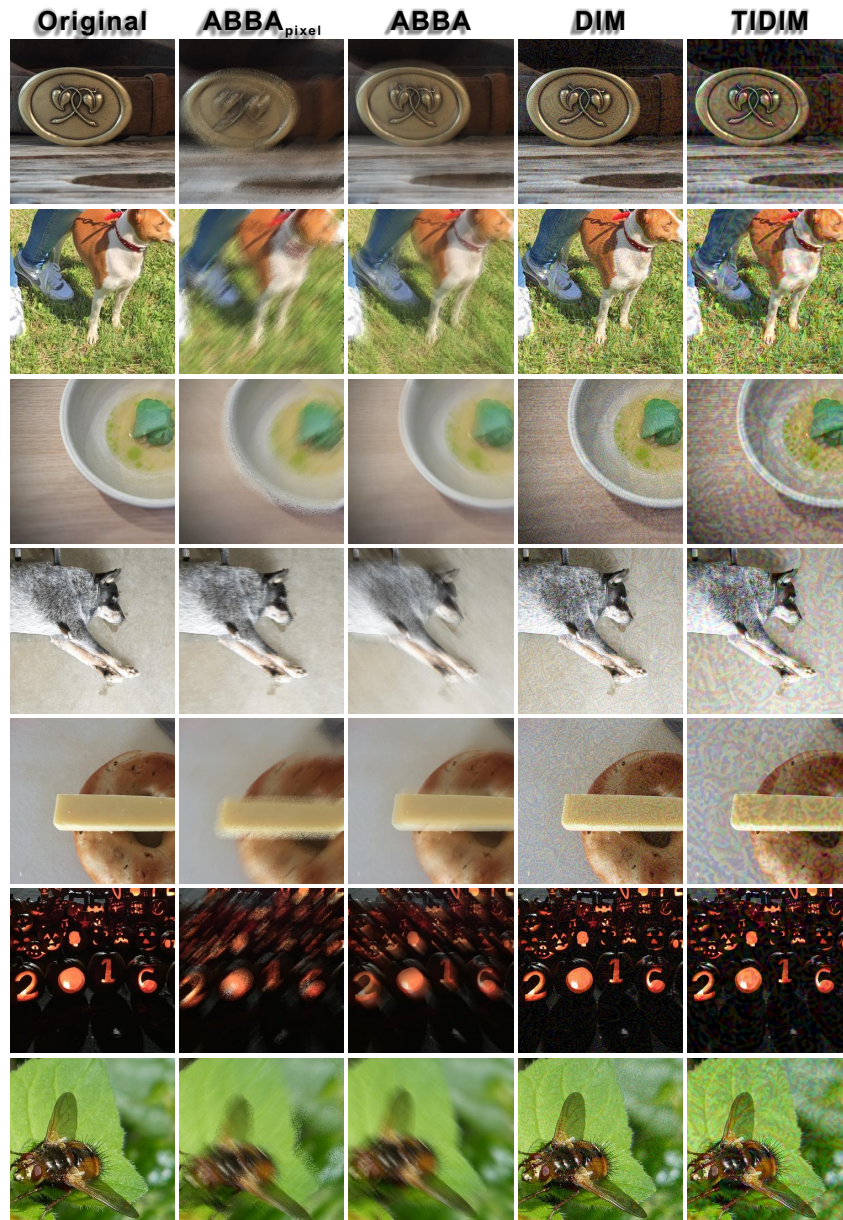


Figure VIII: Seven visualization results of ABBA_{pixel}, ABBA, DIM, and TIDIM. All adversarial examples fool the Inc-v3 model.

- [14] T. Brooks and J. T. Barron. Learning to synthesize motion blur. In *CVPR*, pages 6833–6841, 2019.
- [15] O. Kupyn, V. Budzan, M. Mykhailych, D. Mishkin, and J. Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *CVPR*, 2018.
- [16] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *ICCV*, 2019.
- [17] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *ICCV*, pages 3449–3457, 2017.
- [18] K. He, J. Sun, and X. Tang. Guided image filtering. *TPAMI*, 35(6):1397–1409, 2013.
- [19] Steve Bako, Thijs Vogels, Brian McWilliams, Mark Meyer, Jan Novák, Alex Harvill, Pradeep Sen, Tony DeRose, and Fabrice Rousselle. Kernel-predicting convolutional networks for denoising monte carlo renderings. *ACM Trans. Graph.*, 36(4):97:1–97:14, 2017.
- [20] Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Jiashi Feng, and Jianmin Jiang. A simple pooling-based design for real-time salient object detection. In *CVPR*, 2019.
- [21] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and koray kavukcuoglu. Spatial transformer networks. In *NIPS*, 2015.

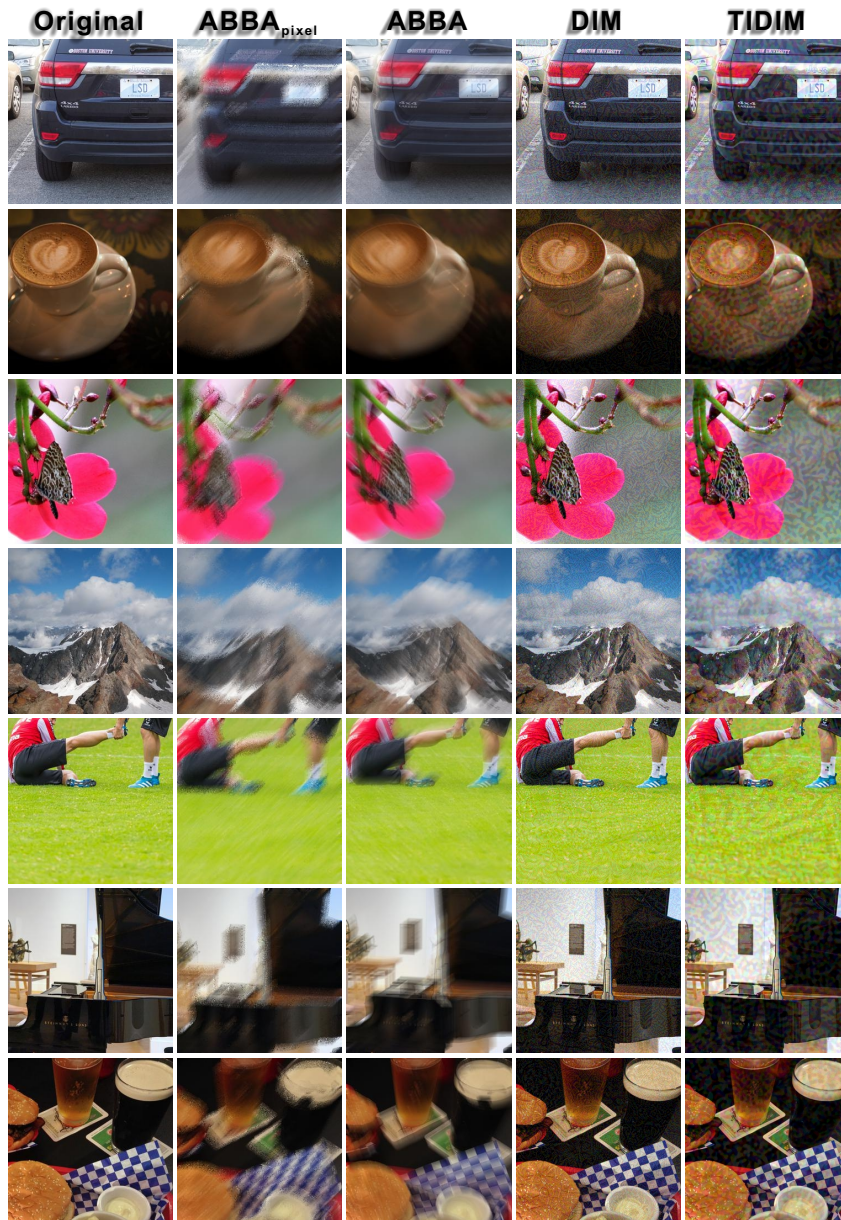


Figure IX: Seven visualization results of $ABBA_{pixel}$, $ABBA$, DIM, and TIDIM. All adversarial examples fool the Inc-v3 model.

- [22] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, pages 3883–3891, 2017.
- [23] M. Noroozi, P. Chandramouli, and P. Favaro. Motion deblurring in the wild. In *Pattern Recognition*, pages 65–77, 2017.
- [24] Shital Shah, Debadepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics*, 2017.
- [25] 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. *CoRR*, 2018.
- [26] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- [27] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, 2017.
- [28] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2017.
- [29] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In *ICLR*, 2018.

- [30] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *CVPR*, 2018.
- [31] 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 *CVPR*, 2019.
- [32] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models, 2017.
- [33] A. Mittal, A. K. Moorthy, and A. C. Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on Image Processing*, 21(12):4695–4708, 2012.
- [34] Microsoft. Airsim simulator. website. <https://microsoft.github.io/AirSim/>.
- [35] Microsoft AirSim. Released airsim environments. website. <https://github.com/microsoft/AirSim/releases>.
- [36] Yupeng Cheng, Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Shang-Wei Lin, Weisi Lin, Wei Feng, and Yang Liu. Pasadena: Perceptually aware and stealthy adversarial denoise attack. *arXiv preprint arXiv:2007.07097*, 2020.
- [37] Liming Zhai, Felix Juefei-Xu, Qing Guo, Xiaofei Xie, Lei Ma, Wei Feng, Shengchao Qin, and Yang Liu. It’s Raining Cats or Dogs? Adversarial Rain Attack on DNN Perception. *arXiv preprint arXiv:2009.09205*, 2020.
- [38] Yupeng Cheng, Felix Juefei-Xu, Qing Guo, Huazhu Fu, Xiaofei Xie, Shang-Wei Lin, Weisi Lin, and Yang Liu. Adversarial Exposure Attack on Diabetic Retinopathy Imagery. *arXiv preprint arXiv:2009.09231*, 2020.
- [39] Binyu Tian, Qing Guo, Felix Juefei-Xu, Wen Le Chan, Yupeng Cheng, Xiaohong Li, Xiaofei Xie, and Shengchao Qin. Bias Field Poses a Threat to DNN-based X-Ray Recognition. *arXiv preprint arXiv:2009.09247*, 2020.
- [40] Qing Guo, Wei Feng, Ce Zhou, Rui Huang, Liang Wan, and Song Wang. Learning dynamic siamese network for visual object tracking. In *Proceedings of the IEEE international conference on computer vision*, pages 1763–1771, 2017.
- [41] Qing Guo, Wei Feng, Ce Zhou, Chi-Man Pun, and Bin Wu. Structure-regularized compressive tracking with online data-driven sampling. *IEEE Transactions on Image Processing*, 26(12):5692–5705, 2017.
- [42] Qing Guo, Ruize Han, Wei Feng, Zhihao Chen, and Liang Wan. Selective spatial regularization by reinforcement learned decision making for object tracking. *IEEE Transactions on Image Processing*, 29:2999–3013, 2020.
- [43] Qing Guo, Xiaofei Xie, Felix Juefei-Xu, Lei Ma, Zhongguo Li, Wanli Xue, Wei Feng, and Yang Liu. Spark: Spatial-aware online incremental attack against visual tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [44] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. In *Proceedings of the International Conference on Learning Representations*, 2018.
- [45] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [46] Chen-Hsuan Lin. Spatial transformer networks. website. <https://github.com/chenhsuanlin/inverse-compositional-STN>.