

# A UNIFIED APPROACH TO INTERPRETING AND BOOSTING ADVERSARIAL TRANSFERABILITY

Xin Wang<sup>a\*</sup>, Jie Ren<sup>a\*</sup>, Shuyun Lin<sup>a</sup>, Xiangming Zhu<sup>a</sup>, Yisen Wang<sup>b</sup>, Quanshi Zhang<sup>a†</sup>

<sup>a</sup>Shanghai Jiao Tong University

<sup>b</sup>Key Lab. of Machine Perception (MoE), School of EECS, Peking University, Beijing, China

## ABSTRACT

In this paper, we use the interaction inside adversarial perturbations to explain and boost the adversarial transferability. We discover and prove the negative correlation between the adversarial transferability and the interaction inside adversarial perturbations. The negative correlation is further verified through different DNNs with various inputs. Moreover, this negative correlation can be regarded as a unified perspective to understand current transferability-boosting methods. To this end, we prove that some classic methods of enhancing the transferability essentially decrease interactions inside adversarial perturbations. Based on this, we propose to directly penalize interactions during the attacking process, which significantly improves the adversarial transferability. Our code is available online<sup>1</sup>.

## 1 INTRODUCTION

Adversarial examples of deep neural networks (DNNs) have attracted increasing attention in recent years (Ma et al., 2018; Madry et al., 2018; Wang et al., 2019; Ilyas et al., 2019; Duan et al., 2020; Wu et al., 2020b; Ma et al., 2021). Goodfellow et al. (2014) found the transferability of adversarial perturbations, and used perturbations generated on a source DNN to attack other target DNNs. Although many methods have been proposed to enhance the transferability of adversarial perturbations (Dong et al., 2018; Wu et al., 2018; 2020a), the essence of the improvement of the transferability is still unclear.

This paper considers the interaction inside adversarial perturbations as a new perspective to interpret adversarial transferability. Interactions inside adversarial perturbations are defined using the Shapley interaction index proposed in game theory (Michel & Marc, 1999; Shapley, 1953). Given an input sample  $x \in \mathbb{R}^n$ , the adversarial attack aims to fool the DNN by adding an imperceptible perturbation  $\delta \in \mathbb{R}^n$  on  $x$ . Each unit in the perturbation map is termed a *perturbation unit*. Let  $\phi_i$  denote the importance of the  $i$ -th perturbation unit  $\delta_i$  to attacking.  $\phi_i$  is implemented as the Shapley value, which will be explained later. The interaction between perturbation units  $\delta_i, \delta_j$  is defined as the change of the  $i$ -th unit’s importance  $\phi_i$  when the  $j$ -th unit is perturbed *w.r.t* the case when the  $j$ -th unit is not perturbed. If the perturbation  $\delta_j$  on the  $j$ -th unit increases the importance  $\phi_i$  of the  $i$ -th unit, then there is a positive interaction between  $\delta_i$  and  $\delta_j$ . If the perturbation  $\delta_j$  decreases the importance  $\phi_i$ , it indicates a negative interaction.

In this paper, we discover and partially prove a clear negative correlation between the transferability and the interaction between adversarial perturbation units, *i.e.* adversarial perturbations with lower transferability tend to exhibit larger interactions between perturbation units. We verify such a correlation based on both the theoretical proof and comparative studies. Furthermore, based on the correlation, we propose to penalize interactions during attacking to improve the transferability.

\*Equal contribution

†Correspondence. This study is conducted under the supervision of Dr. Quanshi Zhang. zqs1022@sjtu.edu.cn. Quanshi Zhang is with the John Hopcroft Center and the MoE Key Lab of Artificial Intelligence, AI Institute, at the Shanghai Jiao Tong University, China.

<sup>1</sup><https://github.com/xherdan76/A-Unified-Approach-to-Interpreting-and-Boosting-Adversarial-Transferability>

In fact, our research group led by Dr. Quanshi Zhang has proposed game-theoretic interactions, including interactions of different orders (Zhang et al., 2020) and multivariate interactions (Zhang et al., 2021c). As a basic metric, the interaction can be used to explain signal processing in trained DNNs from different perspectives. For example, we have build up a tree structure to explain the hierarchical interactions between words in NLP models (Zhang et al., 2021a). We have also used interactions to explain the generalization power of DNNs (Zhang et al., 2021b). The interaction can also explain the utility of adversarial training (Ren et al., 2021). As an extension of the system of game-theoretic interactions, in this study, we explain the adversarial transferability based on interactions.

In this paper, the background for us to investigate the correlation between adversarial transferability and the interaction is as follows. First, we prove that multi-step attacking usually generates perturbations with larger interactions than single-step attacking. Second, according to (Xie et al., 2019), multi-step attacking tends to generate more over-fitted adversarial perturbations with lower transferability than single-step attacking. We consider that the more dedicated interaction reflects more over-fitting towards the source DNN, which hurts adversarial transferability. In this way, we propose the hypothesis that the transferability and the interaction are negatively correlated.

- Comparative studies are conducted to verify this negative correlation through different DNNs.
- Unified explanation. Such a negative correlation provides a unified view to understand current transferability-boosting methods. We theoretically prove that some classic transferability-boosting methods (Dong et al., 2018; Wu et al., 2018; 2020a) essentially decrease interactions between perturbation units, which also verifies the hypothesis of the negative correlation.
- Boosting adversarial transferability. Based on above findings, we propose a loss to decrease interactions between perturbation units during attacking, namely the interaction loss, in order to enhance the adversarial transferability. The effectiveness of the interaction loss further proves the negative correlation between the adversarial transferability and the interaction inside adversarial perturbations. Furthermore, we also try to only use the interaction loss to generate perturbations without the loss for the classification task. We find that such perturbations still exhibit moderate adversarial transferability for attacking. Such perturbations may decrease interactions encoded by the DNN, thereby damaging the inference patterns of the input.

Our contributions are summarized as follows. (1) We reveal the negative correlation between the transferability and the interaction inside adversarial perturbations. (2) We provide a unified view to understand current transferability-boosting methods. (3) We propose a new loss to penalize interactions inside adversarial perturbations and enhance the adversarial transferability.

## 2 RELATED WORK

**Adversarial transferability.** Attacking methods can be roughly divided into two categories, white-box attacks (Szegedy et al., 2013; Goodfellow et al., 2014; Papernot et al., 2016; Carlini & Wagner, 2017; Kurakin et al., 2017; Su et al., 2017; Madry et al., 2018) and black-box attacks (Liu et al., 2016; Papernot et al., 2017; Chen et al., 2017a; Bhagoji et al., 2018; Ilyas et al., 2018; Bai et al., 2020). A specific type of the black-box attack is based on the adversarial transferability (Dong et al., 2018; Wu et al., 2018; Xie et al., 2019; Wu et al., 2020a), which transfers adversarial perturbations on a surrogate/source DNN to a target DNN.

Thus, some previous studies focused on the transferability of adversarial attacking. Liu et al. (2016) demonstrated that non-targeted attacks were easy to transfer, while the targeted attacks were difficult to transfer. Wu et al. (2018) and Demontis et al. (2019) explored factors in affecting the transferability, such as network architectures, model capacity, and gradient alignment. Several methods have been proposed to enhance the transferability of adversarial perturbations. The momentum iterative attack (MI Attack) (Dong et al., 2018) incorporated the momentum of gradients to boost the transferability. The variance-reduced attack (VR Attack) (Wu et al., 2018) used the smoothed gradients to craft perturbations with high transferability. The diversity input attack (DI Attack) (Xie et al., 2019) applied the adversarial attacking to randomly transformed input images, which included random resizing and padding with a certain probability. The skip gradient method (SGM Attack) (Wu et al., 2020a) used the gradients of the skip connection to improve the transferability. Dong et al. (2019) proposed the translation-invariant attack (TI Attack) to evade robustly trained DNNs. Li et al.

(2020) used the dropout erosion and the skip connection erosion to improve the transferability. In comparison, we explain the transferability based on game theory, and discover the negative correlation between the transferability and interactions as a unified explanation for some above methods.

Interaction. The interaction between input variables has been widely investigated. Michel & Marc (1999) proposed the Shapley interaction index based on the Shapley value (Shapley, 1953) in game theory. Daria Sorokina (2008) defined the interaction of input variables of additive models. Scott Lundberg (2017) quantified interactions between each pair of input variables for tree-ensemble models. Some studies mainly focused on interactions to analyze DNNs. Tsang et al. (2018) measured statistical interactions based on DNN weights. Murdoch et al. (2018) proposed to extract interactions in LSTMs by disambiguating information of different gates, and Singh et al. (2019) extended this method to CNNs. Jin et al. (2020) quantified the contextual independence of words to hierarchically explain the LSTMs. Janizek et al. (2020) extended the method of Integrated Gradients (Sundararajan et al., 2017) to quantify pairwise interactions of input features based on the Hessian matrix, which required the DNN to use the SoftPlus operation replace the ReLU operation. Chen et al. (2020) extended the attribution in (Chen & Ji, 2020) to use the Shapley interaction index to generate hierarchical explanations of NLP tasks. In comparison, in this study, we use the Shapley interaction index to explain and improve the transferability of adversarial perturbations.

### 3 THE RELATIONSHIP BETWEEN TRANSFERABILITY AND INTERACTIONS

Preliminaries: the Shapley value. The Shapley value was first proposed in game theory (Shapley, 1953). Considering multiple players in a game, each player aims to win a high reward. The Shapley value is considered as a unique and unbiased approach to fairly allocating the total reward gained by all players to each player (Weber, 1988). The Shapley value satisfies four desirable properties, the linearity, dummy symmetry and efficiency (please see the Appendix A.1 for details). Let  $f = \{1, 2, \dots, n\}$  denote the set of all players, and let  $v(S)$  denote the reward function.  $v(S)$  represents the reward obtained by a set of players  $S$ . The Shapley value  $(\phi_i)$  unbiasedly measures the contribution of the  $i$ -th player to the total reward gained by all players in as follows.

$$\phi_i = v(\{i\}) - v(\emptyset); \quad (\phi_{ij}) = \sum_{S: i \in S, j \notin S} \frac{|S|!(n - |S| - 1)!}{n!} (v(S \cup \{j\}) - v(S)) \quad (1)$$

Adversarial attack. Given an input sample  $x \in [0, 1]^n$  with the true label  $y \in \{1, 2, \dots, C\}$ , we use  $h(x) \in \mathbb{R}^C$  to denote the output of the DNN before the softmax layer. To simplify the story, in this study, we mainly focus on the untargeted adversarial attack. The goal of the untargeted adversarial attack is to add a human-imperceptible perturbation  $\delta \in \mathbb{R}^n$  on the sample  $x$ , and make the DNN classify the perturbed sample  $x + \delta$  into an incorrect category, i.e.  $\arg \max_{y' \in \mathbb{Y}} h_{y'}(x + \delta) \neq y$ . The objective of adversarial attacking is usually formulated as follows.

$$\text{maximize } \ell(h(x + \delta); y) \quad \text{s.t. } \|\delta\|_p \leq \epsilon; \quad x + \delta \in [0, 1]^n; \quad (2)$$

where  $\ell(h(x + \delta); y)$  is referred to the classification loss, and  $\epsilon$  is a constant of the norm constraint. Please see Appendix C for technical details of solving Equation (2).

#### 3.1 THEORETICAL UNDERSTANDING OF THE ADVERSARIAL ATTACK IN GAME THEORY

In adversarial attacking, given the perturbation  $\delta \in \mathbb{R}^n$ , we use  $S = \{1, 2, \dots, n\}$  to denote all units/dimensions in the perturbation. We use the Shapley value in Equation (1) to measure the contribution of each perturbation unit  $i \in S$  to the attack. To this end, it requires us to define the utility of a subset of perturbation units  $S' \subseteq S$  for attacking, which can be formulated as  $u(S') = \max_{y' \in \mathbb{Y}} h_{y'}(x + \delta^{(S')}) - h_y(x + \delta^{(S')})$ , according to Equation (2).  $u(S)$  is the value of the  $S$ -th element of  $u(\cdot) \in \mathbb{R}^C$ .  $\delta^{(S)} \in \mathbb{R}^n$  is the perturbation which only contains perturbation units  $S$ , i.e.  $\delta_i \in S, \delta_i = \delta_i; \delta_i \notin S, \delta_i = 0$ . In this way,  $u(\cdot) = \max_{y' \in \mathbb{Y}} h_{y'}(x + \delta^{(\cdot)}) - h_y(x + \delta^{(\cdot)})$  denotes the utility of all perturbation units, and  $u(\emptyset) = \max_{y' \in \mathbb{Y}} h_{y'}(x) - h_y(x)$  denotes the baseline score without perturbations. Thus, the overall contribution of perturbation units can be measured as  $u(\cdot) - u(\emptyset)$ . We apply the Shapley value in Equation (1) to assign the overall contribution to each perturbation unit as  $\phi_i = u(\{i\}) - u(\emptyset)$ , where  $(\phi_{ij})$  denotes the contribution of the  $i$ -th perturbation unit.

Interactions. Perturbation units do not contribute to the adversarial utility independently. For example, perturbation units may form a certain pattern, an edge in the image. Thus, perturbations units in the edge must appear together. The absence of a few units in the pattern may invalidate this pattern. Let us consider two perturbation units. According to (Michel & Marc, 1999), the Shapley interaction index between units is defined as the additional contribution as follows.

$$I_{ij}(\cdot) = (S_{ij}) - [(i) + (j)]; \quad (3)$$

where  $(i)$  and  $(j)$  represent the individual contributions of unit  $i$  and  $j$ , respectively, when the perturbation units work individually. Note that  $(i)$  is computed in the scenario of considering the unit always absent.  $(i) = v(\emptyset) - v(\{i\})$ , due to the absence of perturbation unit  $i$ .  $(S_{ij})$  denotes the joint contribution of  $i, j$ , when perturbation units  $i, j$  are regarded as a singleton unit  $S_{ij} = \{i, j\}$ . In this case, units  $i, j$  are supposed to be always perturbed or not perturbed simultaneously, and we can consider that there are only two players in the game. Thus, the set of all perturbation units is considered as  $\{i, j\}$ . The joint contribution of  $S_{ij}$  is denoted by  $(S_{ij})$ , s.t.  $(S_{ij}) = v(\{i, j\}) - v(\emptyset)$ .

The interaction defined in Equation (3) is equivalent to the change of the unit's importance  $i$  when the unit exists w.r.t the case when the unit is absent. Please see Appendix D for details.

If  $I_{ij}(\cdot) > 0$ , it means  $i$  and  $j$  cooperate with each other, i.e. the interaction is positive; if  $I_{ij}(\cdot) < 0$ , it means  $i$  and  $j$  conflict with each other, i.e. the interaction is negative. The absolute value of  $|I_{ij}(\cdot)|$  indicates the interaction strength. The interaction is symmetric that  $I_{ij}(\cdot) = I_{ji}(\cdot)$ .

We are given an input sample  $x \in \mathbb{R}^n$  and a DNN  $h(\cdot)$  trained for classification. With the definition of interactions, in adversarial attacking, we have the following propositions:

**Proposition 1.** (Proof in Appendix E) The adversarial perturbation generated by the multi-step attack via gradient descent is given as  $\delta_{\text{multi}}^m = \sum_{t=0}^{m-1} r^t \nabla_x (h(x + \delta_{\text{multi}}^t); y)$ , where  $\delta_{\text{multi}}^t$  denotes the perturbation after the  $t$ -th step of updating, and  $m$  is referred to as the total number of steps. The adversarial perturbation generated by the single-step attack is given as  $\delta_{\text{single}} = r \nabla_x (h(x); y)$ . Then, the expectation of interactions between perturbation units  $i, j$ ,  $E_{a,b}[I_{ab}(\delta_{\text{multi}}^m)]$ , is larger than  $E_{a,b}[I_{ab}(\delta_{\text{single}})]$ , i.e.  $E_{a,b}[I_{ab}(\delta_{\text{multi}}^m)] > E_{a,b}[I_{ab}(\delta_{\text{single}})]$ .

Note that when we compare interactions inside different perturbations, magnitudes of these perturbations should be similar. It is because the comparison of interactions between adversarial perturbations of different magnitudes is not fair. Therefore, we use the stepsize in the single-step attack to roughly (not accurately) balance the magnitude of perturbations. The fairness is further discussed in Appendix E.1.

Proposition 1 shows that, in general, adversarial perturbations generated by the multi-step attack tend to exhibit larger interactions than those generated by the single-step attack. In addition, Appendix E.4 shows that the multi-step attack usually generates perturbations with larger interactions than noisy perturbations of the same magnitude. Besides, Xie et al. (2019) demonstrated that the multi-step attack tends to over-fit the source DNN, which led to low transferability. Intuitively, large interactions mean a strong cooperative relationship between perturbation units, which indicates the significant over-fitting towards adversarial perturbations oriented to the source DNN. In this way, we propose the hypothesis that the adversarial transferability and the interactions inside adversarial perturbations are negatively correlated.

### 3.2 EMPIRICAL VERIFICATION OF THE NEGATIVE CORRELATION

To verify the negative correlation between the transferability and interactions, we conduct experiments to examine whether adversarial perturbations with low transferability tend to exhibit larger interactions than those perturbations with high transferability. Given a source DNN and an input sample  $x$ , we generate the adversarial example  $x + \delta$ . Then, given a target DNN  $h^t(\cdot)$ , we measure the transfer utility of  $\delta$  as  $\text{Transfer Utility} = [\max_{y \in \mathcal{Y}} h_{y^0}^{(t)}(x + \delta) - h_{y^0}^{(t)}(x)] / [\max_{y \in \mathcal{Y}} h_{y^0}^{(t)}(x) - h_{y^0}^{(t)}(x)]$  as mentioned in Section 3.1. The interaction is given as  $\text{Interaction} = E_{i,j}[I_{ij}(\cdot)]$ , which is computed on the source DNN. Note that the computational cost of  $I_{ij}(\cdot)$  is NP-hard. However, we prove that we can simplify the computation of the average interaction over all pairs of units as

Figure 1: The negative correlation between the transfer utility and the interaction. The correlation is computed as the Pearson correlation. The blue shade in each sub figure represents the 95% confidence interval of the linear regression.

follows, which significantly reduces the computational cost. Please see Appendix F for the proof.

$$E_{i;j} [I_{ij}(\cdot)] = \frac{1}{n-1} E_i [v(\cdot) - v(\text{avg}) - v(\text{fig}) + v(\cdot)]: \quad (4)$$

Using 50 images randomly sampled from the validation set of the ImageNet dataset (Russakovsky et al., 2015), we generate adversarial perturbations on four types of DNNs, including ResNet-34/152 (RN-34/152) (He et al., 2016) and DenseNet-121/201 (DN-121/201) (Huang et al., 2017). We transfer adversarial perturbations generated on each ResNet to DenseNets. Similarly, we also transfer adversarial perturbations generated on each DenseNet to ResNets. Figure 1 shows the negative correlation between the transfer utility and the interaction. Each sub figure corresponds to a specific pair of source DNN and target DNN. In each sub figure, each point represents the average transfer utility and the average interaction of adversarial perturbations through all testing images. Different points represent the average interaction and the average transfer utility computed using different hyper-parameters. Given an input image, adversarial perturbations are generated by solving the relaxed form of Equation (2) via the gradient descent,  $\min_x \ell(h(x + \delta); y) + c \|\delta\|_p^2$  s.t.  $x + \delta \in [0, 1]^n$ , where  $c \in \mathbb{R}$  is a scalar. In this way, we gradually change the value of  $c$  and set different values of  $p$  as different hyper-parameters to generate different adversarial perturbations, thereby drawing different points in each sub figure. Fair comparisons require adversarial perturbations generated with different hyper-parameters to be comparable with each other. Thus, we select a constant  $c_2 = 1$  as the stopping criteria of all adversarial attacks. Please see Appendix G for more details.

#### 4 UNIFIED UNDERSTANDING OF TRANSFERABILITY BOOSTING ATTACKS

In this section, we prove that some classical methods of improving the adversarial transferability essentially decrease interactions between perturbation units, although these methods are not originally designed to decrease the interaction. Without loss of generality, let us be given an input sample  $x \in \mathbb{R}^n$  and a DNN  $h(\cdot)$  trained for classification.

- VR Attack (Wu et al., 2018) smooths the classification loss with the Gaussian noise during attacking. In the VR Attack, the gradient of the input sample is computed as follows:  $E_{N(0, \sigma^2)} [\nabla_x \ell(h(x + \delta); y)]$ . The following proposition proves that the VR Attack is prone to decrease interactions inside perturbation units.

**Proposition 2.** (Proof in Appendix H) The adversarial perturbation generated by the multi-step attack is given as  $\delta_{\text{multi}}^m = \sum_{t=0}^{m-1} r^t \delta^t(x; y)$ . The adversarial perturbation generated by the VR Attack is computed as  $\delta_{\text{vr}}^m = \sum_{t=0}^{m-1} r^t \delta_{\text{vr}}^t(x; y)$ , where  $\delta_{\text{vr}}^t(x; y) = E_{N(0, \sigma^2)} [\nabla_x \ell(h(x + \delta_{\text{vr}}^t); y)]$ . Perturbation units of  $\delta_{\text{vr}}^m$  tend to exhibit smaller interactions than  $\delta_{\text{multi}}^m$ , i.e.  $E_x E_{a;b} [I_{ab}(\delta_{\text{vr}}^m)] \leq E_x E_{a;b} [I_{ab}(\delta_{\text{multi}}^m)]$ .

Besides the theoretical proof, we also conduct experiments to compare interactions of perturbation units generated by the baseline multi-step attack (implemented as (Madry et al., 2018)) with those

<sup>2</sup>We set  $p = 2$  as the setting 1, and  $p = 5$  as the setting 2. To this end, the performance of adversarial perturbations is not the key issue in the experiment. Instead, we just randomly select to examine the trustworthiness of the negative correlation under various attacking conditions (even in extreme attacking conditions).

of perturbation units generated by the VR Attack. Table 5 shows that the VR Attack exhibits lower interactions between perturbation units than the baseline multi-step attack.

- **MI Attack** (Dong et al., 2018) incorporates the momentum of gradients when updating the adversarial perturbation. In the MI Attack, the gradient used in step  $t$  is computed as follows.

$$g^t = g^{t-1} + r_x \cdot \nabla_x h(x + g^{t-1}; y); y = r_x \cdot \nabla_x h(x + g^{t-1}; y_{-1}).$$

Note that the original MI Attack and the multi-step attack cannot be directly compared, since that magnitudes of the generated perturbations cannot be fairly controlled. The values of interactions are sensitive to the magnitude of perturbations. Comparing perturbations with different magnitudes is not fair. Thus, we slightly revise the MI Attack as  $\alpha > 0$ ;  $g_{mi}^t = g_{mi}^{t-1} + (1 - \alpha) r_x \cdot \nabla_x (h(x + g_{mi}^{t-1}); y)$ ;  $g_{mi}^0 = 0$ , where  $\alpha = (t-1)/t$ . We investigate the interaction of adversarial perturbations generated by the original multi-step attack and the MI Attack. We prove the following proposition, which shows that the MI Attack decreases the interaction between perturbation units in most cases.

**Proposition 3.** (Proof in Appendix I) The adversarial perturbation generated by the multi-step attack is given as  $p_{multi}^m = \sum_{t=0}^{m-1} r_x \cdot \nabla_x (h(x + p_{multi}^t); y)$ . The adversarial perturbation generated by the multi-step attack incorporating the momentum is computed as  $p_{mi}^m = \sum_{t=0}^{m-1} g_{mi}^t$ . Perturbation units of  $p_{mi}^m$  exhibit smaller interactions than  $p_{multi}^m$ , i.e.  $E_{a,b}[I_{ab}(p_{mi}^m)] \leq E_{a,b}[I_{ab}(p_{multi}^m)]$ .

- **SGM Attack** (Wu et al., 2020a) exploits the gradient information of the skip connection in ResNets to improve the transferability of adversarial perturbations. The SGM Attack revises the gradient in the backpropagation, which can be considered as to add a specific dropout operation in the backpropagation. We notice that Zhang et al. (2021b) has proved that the dropout operation can decrease the significance of interactions, so as to decrease the significance of the overfitting of DNNs. Thus, this also proves that the SGM Attack decreases interactions between perturbation units.

Besides the theoretical proof, we also conduct experiments to compare interactions of perturbation units generated by the baseline multi-step attack (implemented as Madry et al. (2018)) with those of perturbation units generated by the SGM Attack. Table 5 shows that the SGM Attack exhibits lower interactions than the baseline multi-step attack.

## 5 THE INTERACTION LOSS FOR TRANSFERABILITY ENHANCEMENT

**Interaction loss.** Based on findings in previous sections, we propose a loss to directly penalize interactions during attacking, in order to improve the transferability of adversarial perturbations. Based on Equation (2), we jointly optimize the classification loss and the interaction loss to generate adversarial perturbations. This method is termed the interaction-reduced attack (IR Attack).

$$\max_x [L(h(x + p); y) - \lambda \cdot \text{interaction}]; \quad \text{interaction} = E_{i,j} [I_{ij}(p)] \quad \text{s.t. } k_p \leq k; x + p \in [0, 1]^n; \quad (5)$$

where  $\lambda \cdot \text{interaction}$  is the interaction loss, and  $\lambda$  is a constant weight for the interaction loss. Although the computation of the interaction loss can be simplified according to Equation (4), the computational cost of the interaction loss is intolerable, when the dimension of images is high. Therefore, as a trade-off between the accuracy and the computational cost, we divide the input image into  $16 \times 16$  grids. We measure and penalize interactions at the grid level, instead of the pixel level. Moreover, we apply an efficient sampling method to approximate the expectation operation during the computation of interactions in Equation (4). Figure 2 visualizes interactions between adjacent perturbation units at the grid level generated with and without the interaction loss.

**Experiments.** For implementation, we generated adversarial perturbations on six different source DNNs, including Alexnet (Krizhevsky et al., 2012), VGG-16 (Simonyan & Zisserman, 2015), ResNet-34/152 (RN-34/152) (He et al., 2016) and DenseNet-121/201 (DN-121/201) (Huang et al., 2017). For each source DNN, we tested the transferability of the generated perturbations on seven target DNNs, including VGG-16, ResNet-152 (RN-152), DenseNet-201 (DN-201), SENet-154 (SE-154) (Hu et al., 2018), InceptionV3 (IncV3) (Szegedy et al., 2016), InceptionV4 (IncV4) (Szegedy et al., 2017), and Inception-ResNetV2 (IncResV2) (Szegedy et al., 2017). In addition, three state-of-the-art DNNs, including the Dual-Path-Network (DPN-68) (Chen et al., 2017b), the NASNet-LARGE (NASN-L) (Zoph et al., 2018), and the Progressive NASNet (PNASN) (Liu et al., 2018), were used as target DNNs to evaluate the ensemble source model (will be introduced in the next paragraph). Besides unsecured target DNNs mentioned above, we also used three secured target

Figure 2: Visualization of interactions between neighboring perturbation units generated with and without the interaction loss. The color in the visualization is computed as  $\frac{w_{ij}}{E_{j \in N_i} |w_{ij}|}$ , where  $N_i$  denotes the set of adjacent perturbation units of the perturbation unit  $i$ . Here, we ignore interactions between non-adjacent units to simplify the visualization. It is because adjacent units usually encode much more significant interactions than other units. The interaction loss forces the perturbation to encode more negative interactions.

Table 1: The success rates of  $L_1$  and  $L_2$  black-box attacks crafted on six source models, including AlexNet, VGG16, RN-34/152, DN-121/201, against seven target models. Transferability of adversarial perturbations can be enhanced by penalizing interactions.

Source	Method	VGG-16	RN152	DN-201	SE-154	IncV3	IncV4	IncResV2
AlexNet	PGDL <sub>1</sub>	67.0 1.6	27.8 1.1	32.3 0.4	28.2 0.7	29.1 1.5	23.0 0.4	18.6 1.5
	PGDL <sub>1</sub> +IR	78.7 1.0	42.0 1.5	50.3 0.4	41.2 0.6	43.7 0.5	36.4 1.5	29.0 1.0
VGG-16	PGDL <sub>1</sub>	–	43.0 1.8	48.3 2.0	52.9 2.7	39.3 0.7	49.3 1.1	29.7 2.0
	PGDL <sub>1</sub> +IR	–	63.1 1.6	70.0 1.1	71.2 1.5	57.6 1.0	68.6 3.2	49.2 1.2
RN-34	PGDL <sub>1</sub>	65.4 2.9	59.2 2.7	63.5 3.3	33.1 2.9	27.4 3.6	23.9 1.7	21.1 1.1
	PGDL <sub>1</sub> +IR	84.0 0.5	84.7 2.3	88.5 0.9	64.4 1.6	56.9 3.1	59.3 4.3	49.2 1.1
RN-152	PGDL <sub>1</sub>	51.6 3.2	–	61.5 2.4	33.9 1.5	28.1 0.9	25.0 1.2	22.4 1.0
	PGDL <sub>1</sub> +IR	72.3 1.2	–	82.1 1.3	61.1 0.9	53.6 0.8	50.6 3.5	46.0 2.3
DN-121	PGDL <sub>1</sub>	68.6 1.1	63.6 3.2	86.9 1.5	46.1 1.5	37.3 1.6	37.1 2.1	28.9 2.8
	PGDL <sub>1</sub> +IR	85.0 0.3	84.8 0.4	95.1 0.2	70.3 1.7	61.1 2.5	62.1 2.0	53.5 0.3
DN-201	PGDL <sub>1</sub>	64.4 1.4	67.8 0.2	–	50.9 0.8	39.5 3.3	36.5 0.9	34.2 0.4
	PGDL <sub>1</sub> +IR	78.6 2.5	85.0 1.1	–	73.9 0.5	61.6 1.8	63.7 0.6	56.4 2.1
AlexNet	PGDL <sub>2</sub>	85.1 1.5	58.9 1.0	60.2 2.1	55.1 1.5	56.0 3.7	49.6 3.4	44.6 3.3
	PGDL <sub>2</sub> +IR	91.6 1.1	72.0 1.6	76.8 1.0	69.0 1.0	73.0 0.8	63.1 2.1	59.4 1.9
VGG-16	PGDL <sub>2</sub>	–	76.7 0.9	82.3 2.9	83.5 1.9	77.5 3.6	82.1 2.2	69.4 2.1
	PGDL <sub>2</sub> +IR	–	86.5 0.9	88.9 1.5	89.6 1.2	85.2 1.1	88.3 1.4	80.4 0.4
RN-34	PGDL <sub>2</sub>	88.2 1.4	86.2 0.4	89.6 1.3	66.9 1.1	64.2 2.9	60.0 1.9	55.2 1.8
	PGDL <sub>2</sub> +IR	95.2 0.2	95.4 0.1	96.7 0.6	86.7 1.2	84.3 0.6	81.8 1.9	80.4 1.9
DN-121	PGDL <sub>2</sub>	89.4 1.1	86.8 1.0	97.6 1.0	75.6 1.7	70.1 2.9	70.4 4.4	66.5 4.7
	PGDL <sub>2</sub> +IR	94.2 0.1	93.3 0.8	97.7 0.3	87.8 0.7	84.5 0.7	84.2 0.1	82.4 0.1

Table 2: The success rates of black-box attacks crafted on the ensemble model (RN-34+RN-152+DN-121) against nine target models.

Source	Method	VGG-16	RN-152	DN-201	SE-154	IncV3	IncV4	IncResV2	DPN-68	NASN-L	PNASN
Ensemble	PGDL <sub>1</sub>	86.6 1.2	99.9 0.1	97.0 0.7	70.7 1.6	64.2 0.3	57.7 2.4	53.1 0.7	61.6 0.5	59.6 0.4	72.3 0.3
	PGDL <sub>1</sub> +IR	91.5 0.1	92.4 1.6	92.1 1.7	86.1 0.3	81.6 0.9	79.9 1.7	78.4 1.3	82.5 1.0	82.3 1.6	85.6 0.5

models for testing, which were learned via ensemble adversarial training: IncV3<sub>ens3</sub> (ensemble of three IncV3 networks), IncV4<sub>ens4</sub> (ensemble of four IncV3 networks), and IncResV2<sub>ens3</sub> (ensemble of three IncResV2 networks), which were released by Truett et al. (2017).

Ensemble source models. Besides above adversarial transferring from a single-source model, we also conducted the proposed IR Attack in the scenario of the ensemble-based attacking (Liu et al., 2016), in order to generate adversarial perturbations on the ensemble of RN-34, RN-152, and DN-121.

Baselines. The first baseline method, the PGD Attack (Madry et al., 2018), directly solved the Equation (2), which was widely used for adversarial attacks. Besides this baseline attack, the other four baselines were the MI Attack (Dong et al., 2018), the VR Attack (Wu et al., 2018), the SGM Attack (Wu et al., 2020a), and the TI Attack (Dong et al., 2019). Our method was implemented according to Equation (5), namely the IR Attack. Because the SGM Attack was one of the top-ranked methods of boosting the adversarial transferability, we further added the interaction loss to the SGM Attack as another implementation of our method (namely the SGM+IR Attack). We also used the interaction loss to boost the performance of the MI Attack and the VR Attack (namely MI+IR and VR+IR, respectively). Please see Appendix M.1 for details. Moreover, as Section 4 states, the MI Attack, VR Attack, and SGM Attack also decrease interactions during attacking. Thus, we combined the IR Attack with all these interaction-reducing techniques together as a new implementation of our method, namely the HybridIR Attack. All attacks were conducted with 100

Table 3: Transferability against the secured models: the success rates of black-box attacks crafted on RN-34 and DN-121 source models against three secured models.

Source	Method	IncV3 <sub>ens 3</sub>	IncV3 <sub>ens 4</sub>	IncRes <sub>ens 3</sub>	Source	Method	IncV3 <sub>ens 3</sub>	IncV3 <sub>ens 4</sub>	IncRes <sub>ens 3</sub>						
RN-34	PGDL <sub>1</sub>	9.8	0.1	10.0	0.5	5.7	0.3	DN-121	PGDL <sub>1</sub>	12.8	0.1	11.2	1.7	6.9	1.0
	PGDL <sub>1</sub> +IR	26.5	2.9	22.1	1.3	14.3	0.4		PGDL <sub>1</sub> +IR	28.0	1.8	26.5	2.1	17.4	1.3
	TI <sup>4</sup>	21.4	0.8	20.9	0.9	14.9	1.4		TI <sup>4</sup>	26.8	1.3	26.1	1.5	19.4	1.6
	TI <sup>4</sup> + IR	33.6	0.4	33.2	0.3	24.0	0.5		TI <sup>4</sup> + IR	38.0	2.5	42.2	7.7	29.0	1.4

Table 4: The success rates of black-box attacks crafted by different methods on four source models (RN-34/152, DN-121/201) against seven target models. Transferability of adversarial perturbations can be enhanced by penalizing interactions.

Source	Method	VGG-16	RN152	DN-201	SE-154	IncV3	IncV4	IncResV2							
RN-34	MI	80.1	0.5	73.0	2.3	77.7	0.5	48.9	0.8	46.2	1.2	39.9	0.5	34.8	2.5
	VR	88.8	0.2	86.4	1.6	87.9	2.4	62.1	1.5	58.4	3.0	56.3	2.3	49.7	0.9
	SGM	91.8	0.6	89.0	0.9	90.0	0.4	68.0	1.4	63.9	0.3	58.2	1.1	54.6	1.2
	SGM+IR	94.7	0.6	91.7	0.6	93.4	0.8	72.7	0.4	68.9	0.9	64.1	1.3	61.3	1.0
	HybridIR	96.5	0.1	94.9	0.3	95.6	0.6	79.7	1.0	77.1	0.8	73.8	0.1	70.2	0.5
RN-152	MI	70.3	0.6	-	74.8	1.4	51.7	0.8	47.1	0.9	40.5	1.6	36.8	2.7	
	VR	83.9	3.4	-	91.1	0.9	70.0	3.7	63.1	0.9	58.8	0.1	56.2	1.3	
	SGM	88.2	0.5	-	90.2	0.3	72.7	1.4	63.2	0.7	59.1	1.5	58.1	1.2	
	SGM+IR	92.0	1.0	-	92.5	0.4	79.3	0.1	69.6	0.8	66.2	1.0	63.6	0.9	
	HybridIR	95.3	0.4	-	96.9	0.2	84.7	0.7	80.0	1.2	77.5	0.8	75.6	0.6	
DN-121	MI	83.0	4.9	72.0	0.7	91.5	0.2	58.4	2.6	54.6	1.6	49.2	2.4	43.9	1.5
	VR	91.5	0.5	88.7	0.5	98.8	0.2	75.1	1.3	74.3	1.7	75.6	3.0	69.8	1.3
	SGM	88.7	0.9	88.1	1.0	98.0	0.4	78.0	0.9	64.7	2.5	65.4	2.3	59.7	1.7
	SGM+IR	91.7	0.2	90.4	0.4	94.3	0.1	87.0	0.4	78.8	1.3	79.5	0.2	75.8	2.7
	HybridIR	96.9	0.4	96.8	0.4	99.1	0.4	90.9	0.5	88.4	0.8	87.8	0.8	87.1	0.4
DN-201	MI	77.3	0.8	74.8	1.4	-	64.6	1.0	56.5	2.5	51.1	2.1	47.8	1.9	
	VR	87.3	1.1	90.4	1.2	-	78.0	1.5	75.8	2.1	75.8	1.3	71.3	1.2	
	SGM	87.3	0.3	92.4	1.0	-	82.9	0.2	72.3	0.3	71.3	0.6	68.8	0.5	
	SGM+IR	89.5	0.9	91.8	0.7	-	87.3	1.2	82.5	0.8	80.3	0.3	81.5	0.5	
	HybridIR	94.4	0.1	96.9	0.5	-	91.7	0.2	89.6	0.6	88.3	0.3	87.3	0.7	

steps on randomly selected 1000 images of the validation set in the ImageNet dataset. We set  $\epsilon = 16=255$  for the  $L_1$  attack, and set  $\epsilon = 16=255$  following the setting in (Dong et al., 2018) for the  $L_2$  attack. The step size was set to  $\epsilon/255$  for all attacks. Considering the efficiency of signal processing in DNNs with different depths, we set  $\eta = 1$  for the IR Attack, when the source DNN was ResNet. We set  $\eta = 2$ , for other source DNNs. To enable fair comparisons, the transferability of each baseline was computed based on the best adversarial perturbation during the 100 steps via the leave-one-out (LOO) validation. Please see the Appendix K for the motivation and the evidence of the LOO evaluation of transferability. All attacks were conducted with three different random samplings of grids or different initial perturbations.

Table 1 reports the success rates of the baseline attack (PGD (Madry et al., 2018)) and the IR Attack, namely PGDL<sub>1</sub> +IR of  $L_1$  attacks and PGDL<sub>2</sub> +IR of  $L_2$  attacks. Compared with the baseline attack, the transferability was significantly improved by the interaction loss on various source models against different target models. Let us focus on the attack. For most source models and target models, the transferability enhancement brought by the interaction loss was more than 10%. In particular, when the source DNN and the target DNN were DN-201 and IncV4, respectively, the baseline attack achieved the transferability of 36.5%. With the interaction loss, the transferability was improved to 63.7% (27% gain). As Table 2 shows, in most cases, the IR Attack on the ensemble model generated more transferable perturbations than the PGD Attack. Besides, as Table 3 shows, our interaction loss also improved the transferability against the secured target DNNs. Such improvement further verified the negative correlation between transferability and interactions. Note that we did not use the LOO in Table 3, in order to make experimental settings in this table consistent with the evaluation used by Turner et al. (2017). Table 4 shows the improvement of the transferability obtained by the interaction loss on other attacking methods. The interaction loss could further boost the transferability of state-of-the-art transfer attacks. Without the interaction loss, the highest transferability made by the SGM Attack against the IncResV2 was 68.8% (when the source is DN-201). When the interaction loss was added, the transferability was improved to 81.5% (

<sup>3</sup>Previous studies usually set the number of steps to 10 or 20. Here, we set the number of steps to 100 together with the leave-one-out validation for fair comparisons of different attacks.

<sup>4</sup>The TI Attack was designed oriented to the secured DNNs which were robustly trained via adversarial training. Thus, we applied the TI Attack to the secured models in Table 3.



Figure 3: (a) The success rates of black-box attacks with the IR Attack using different values of  $\alpha$ . The success rates increased, when the value of  $\alpha$  increased. (b) The transferability of adversarial perturbations generated by only using the interaction loss (without the classification loss). Such adversarial perturbations still exhibited moderate adversarial transferability. Points localized at the last epoch represent the transferability of noise perturbations as the baseline.

12% gain). Moreover, the HybridIR Attack, which combined all methods of reducing interactions together, improved success rates from the range of 54.9% to the range of 70.2%.

We can understand behaviors of the proposed interaction loss as follows. Different methods generate adversarial perturbations in different manifolds, thereby exhibiting different transferability. Based on the current perturbation, the interaction loss can point out the optimization direction towards further decrease of interactions in a local manner due to its optimization power. Thus, the interaction loss further boosts the transferability.

To further demonstrate the broad applicability of the interaction loss, besides untargeted attacks on the ImageNet dataset, we also conducted targeted attacks on the CIFAR-10 dataset (Krizhevsky & Hinton, 2009). Experimental results consistently showed that the adversarial transferability can be enhanced by reducing interactions in targeted attacks. Please see Appendix. M.2 for details.

**Effects of the interaction loss** We tested the transferability of perturbations generated by the IR Attack with different weights of the interaction loss. In particular, the baseline attack (PGD) can be considered as the IR Attack when  $\alpha = 0$ . We conducted attacks on two source DNNs (RN-34, DN-121), and transferred adversarial perturbations to seven target DNNs (VGG16, RN-152, DN-201, SE-154, IncV3, IncV4, IncResV2). The attacks were conducted with 100 steps on 100 validation images in ImageNet. Figure 3 (a) shows the black-box success rates with different values of  $\alpha$ . The transferability of the IR Attack increased along with the increase of the weight

**Attack only with the interaction loss** To further understand the effects of the interaction loss, we generated perturbations by exclusively using the interaction loss (without the classification loss). We used the RN-34 and DN-121 as source DNNs and tested the transferability on seven target DNNs. The attacks were conducted with 100 steps on 100 ImageNet validation images. Figure 3 (b) shows the curve of the transferability in different epochs. We compared such adversarial perturbations with noise perturbations generated as  $\text{sign}(\text{noise})$ , where  $\text{noise} \sim \mathcal{N}(0, \sigma^2)$ , and  $\sigma = 16 = 255$ , which was the same as the value used in the attack. We found that perturbations generated by only using the interaction loss still exhibited moderate adversarial transferability. This phenomenon may be explained as that such perturbations decrease most interactions in the DNN, thereby damaging the inference patterns in the input image.

## 6 CONCLUSION

In this paper, we have analyzed the transferability of adversarial perturbations from the perspective of interactions based on game theory. We have proved that the multi-step attack tends to generate adversarial perturbations with large interactions. We have discovered and partially proved the negative correlation between the transferability and interactions inside adversarial perturbations. I.e. adversarial perturbations with higher transferability usually exhibit more negative interactions. We have proved that some classical methods of enhancing the transferability essentially decrease interactions between perturbation units, which provides a unified view to understand the enhancement of transferability. Moreover, we have proposed a new loss to directly penalize interactions between perturbation units during attacking, which significantly improves the transferability of previous methods. Furthermore, we have found that adversarial perturbations generated only using the interaction loss without the classification loss still exhibited moderate transferability, which provides a new perspective to understand the transferability of adversarial perturbations.

## ACKNOWLEDGMENTS

All members in Shanghai Jiao Tong university, including Xin Wang, Jie Ren, Shuyun Lin, Xiangming Zhu, and Dr. Quanshi Zhang are supported by National Natural Science Foundation of China (61906120 and U19B2043) and Huawei Technologies. Dr. Yisen Wang is partially supported by the National Natural Science Foundation of China under Grant 62006153, and CCF-Baidu Open Fund (OF2020002). Xin Wang is supported by Wu Wen Jun Honorary Doctoral Scholarship, AI Institute, Shanghai Jiao Tong University.

## REFERENCES

- Marco Ancona, Cengiz Oztireli, and Markus Gross. Explaining deep neural networks with a polynomial time algorithm for shapley value approximation. *International Conference on Machine Learning* pp. 272–281, 2019.
- Yang Bai, Yuyuan Zeng, Yong Jiang, Yisen Wang, Shu-Tao Xia, and Weiwei Guo. Improving query efficiency of black-box adversarial attack. *ECCV*, 2020.
- Arjun Nitin Bhagoji, Warren He, Bo Li, and Dawn Song. Practical black-box attacks on deep neural networks using efficient query mechanisms. *European Conference on Computer Vision* pp. 158–174. Springer, 2018.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (S&P)*. 39–57. IEEE, 2017.
- Hanjie Chen and Yangfeng Ji. Learning variational word masks to improve the interpretability of neural text classifiers. In *EMNLP*, 2020.
- Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. Generating hierarchical explanations on text classification via feature interaction detection. *ACL*, 2020.
- Jianbo Chen, Le Song, Martin J. Wainwright, and Michael I. Jordan. L-shapley and c-shapley: Efficient model interpretation for structured data. *arXiv:1808.02610*2018a.
- 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. *arXiv:1708.03999*2017a.
- Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: Elastic-net attacks to deep neural networks via adversarial examples. *ICAAI*, 2018b.
- Yunpeng Chen, Jianan Li, Huaxin Xiao, Xiaojie Jin, Shuicheng Yan, and Jiashi Feng. Dual path networks. In *Advances in neural information processing systems* pp. 4467–4475, 2017b.
- Mirek Riedewald Daria Sorokina, Rich Caruana. Detecting statistical interactions with additive groves of trees. In *ICML*, 2008.
- Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, Cristina Nita-Rotaru, and Fabio Roli. Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. *28th USENIX Security Symposium USENIX Security 19*, pp. 321–338, 2019.
- Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Xiaolin Hu, Jianguo Li, , and Jun Zhu. Boosting adversarial attacks with momentum. *CVPR*, 2018.
- Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* pp. 4312–4321, 2019.
- Ranjie Duan, Xingjun Ma, Yisen Wang, James Bailey, A Kai Qin, and Yun Yang. Adversarial camouflage: Hiding physical-world attacks with natural styles. *CVPR*, 2020.

- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572, 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation network. Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7132–7141, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700–4708, 2017.
- Qian Huang, Isay Katsman, Horace He, Zeqi Gu, Serge Belongie, and Ser-Nam Lim. Enhancing adversarial example transferability with an intermediate level attack. Proceedings of the IEEE International Conference on Computer Vision, pp. 4733–4742, 2019.
- Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. ICML, 2018.
- Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. NeurIPS, 2019.
- Nathan Inkawhich, Wei Wen, Hai Helen Li, and Yiran Chen. Feature space perturbations yield more transferable adversarial examples. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7066–7074, 2019.
- Nathan Inkawhich, Kevin Liang, Lawrence Carin, and Yiran Chen. Transferable perturbations of deep feature distributions. International Conference on Learning Representations, 2020.
- Joseph D Janizek, Pascal Sturmfels, and Su-In Lee. Explaining explanations: Axiomatic feature interactions for deep networks. arXiv preprint arXiv:2002.04138, 2020.
- Xisen Jin, Zhongyu Wei, Junyi Du, Xiangyang Xue, and Xiang Ren. Towards hierarchical importance attribution: Explaining compositional semantics for neural sequence models. ICLR, 2020.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. pp. 1097–1105, 2012.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In arXiv:1607.02533, 2017.
- Yann LeCun, léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 6(11):2278–2324, 1998.
- Yingwei Li, Song Bai, Yuyin Zhou, Cihang Xie, Zhishuai Zhang, and Alan Yuille. Learning transferable adversarial examples via ghost networks. Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, 2020.
- Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. Proceedings of the European Conference on Computer Vision (ECCV), pp. 19–34, 2018.
- Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. ICLR, 2016.

- Xingjun Ma, Bo Li, Yisen Wang, Sarah M Erfani, Sudanthi Wijewickrema, Grant Schoenebeck, Dawn Song, Michael E Houle, and James Bailey. Characterizing adversarial subspaces using local intrinsic dimensionality. In *ICLR*, 2018.
- Xingjun Ma, Yuhao Niu, Lin Gu, Yisen Wang, Yitian Zhao, James Bailey, and Feng Lu. Understanding adversarial attacks on deep learning based medical image analysis systems. *Pattern Recognition*, 110:107332, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- Grabisch Michel and Roubens Marc. An axiomatic approach to the concept of interaction among players in cooperative games. *International Journal of Game Theory*, 1999.
- W. James Murdoch, Peter J. Liu, and Bin Yu. Beyond word importance: Contextual decomposition to extract interactions from lstm. In *ICLR*, 2018.
- Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. *IEEE European Symposium on Security & Privacy*, 2016.
- Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. *arXiv:1602.02697*, 2017.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Aubert, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32, pp. 8026–8037. Curran Associates, Inc., 2019.
- Jie Ren, Die Zhang, Yisen Wang, Lu Chen, Zhanpeng Zhou, Xu Cheng, Xin Wang, Yiting Chen, Jie Shi, and Quanshi Zhang. Game-theoretic understanding of adversarially learned feature. *arXiv preprint arXiv:2103.07364*, 2021.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- Su-In Lee Scott Lundberg. Consistent feature attribution for tree ensemble models. *ICML WHI Workshop*, 2017.
- Lloyd S Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2(28): 307–317, 1953.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- Chandan Singh, W. James Murdoch, and Bin Yu. Hierarchical interpretations for neural network predictions. In *ICLR*, 2019.
- Jiawei Su, Danilo Vasconcellos Vargas, and Sakurai Kouichi. One pixel attack for fooling deep neural networks. In *arXiv:1710.08864*, 2017.
- Mukund Sundararajan and Amir Najmi. The many shapley values for model explanation. *arXiv preprint arXiv:1908.08474*, 2019.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pp. 3319–3328, 2017.
- 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.

- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *Thirty-first AAAI conference on artificial intelligence* 2017.
- Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defense. *arXiv preprint arXiv:1705.07204* 2017.
- Michael Tsang, Dehua Cheng, and Yan Liu. Detecting statistical interactions from neural network weights. In *ICLR*, 2018.
- Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanquan Gu. On the convergence and robustness of adversarial training. *ICML*, 2019.
- Robert J Weber. Probabilistic values for games. *The Shapley Value. Essays in Honor of Lloyd S. Shapley*, pp. 101–119, 1988.
- Dongxian Wu, Yisen Wang, Shu-Tao Xia, James Bailey, and Xingjun Ma. Skip connections matter: On the transferability of adversarial examples generated with resonance. *International Conference on Learning Representations* 2020a.
- Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. In *NeurIPS* 2020b.
- Lei Wu, Zhanxing Zhu, and Cheng Tai. Understanding and enhancing the transferability of adversarial examples. *arXiv preprint arXiv:1802.09707* 2018.
- Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L Yuille. Improving transferability of adversarial examples with input diversity. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2730–2739, 2019.
- Die Zhang, Huilin Zhou, Hao Zhang, Xiaoyi Bao, Da Huo, Ruizhao Chen, Xu Cheng, Mengyue Wu, and Quanshi Zhang. Building interpretable interaction trees for deep nlp models. *AAAI*, 2021a.
- Hao Zhang, Xu Cheng, Yiting Chen, and Quanshi Zhang. Game-theoretic interactions of different orders. *arXiv preprint arXiv:2010.14978* 2020.
- Hao Zhang, Sen Li, Yinchao Ma, Mingjie Li, Yichen Xie, and Quanshi Zhang. Interpreting and boosting dropout from a game-theoretic view. *ICLR*, 2021b.
- Hao Zhang, Yichen Xie, Longjie Zheng, Die Zhang, and Quanshi Zhang. Interpreting multivariate interactions in dnns. *IAAAI*, 2021c.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.

## A MOTIVATIONS FOR USING THE SHAPLEY INTERACTION INDEX

In this section, we discuss the motivations of using the Shapley interaction index to define the interaction.

### A.1 FOUR PROPERTIES OF SHAPLEY VALUES

Let  $N = \{1, 2, \dots, n\}$  denote the set of all players, and the reward function  $v$ . Without ambiguity, we use  $(ij)$  to denote the Shapley value of the player  $i$  in the game with all players and reward function  $v$ , which is given as follows.

$$(ij) = \sum_{S \subseteq N \setminus \{i, j\}} \frac{|S|!(n - |S| - 1)!}{n!} (v(S \cup \{i, j\}) - v(S)) \quad (6)$$

The Shapley value satisfies the following four properties (Weber, 1988):

- **Linearity property:** If there are two games and the corresponding reward functions are  $v(S)$  and  $w(S)$ , i.e.  $v(S)$  and  $w(S)$  measure the reward obtained by players in these two games. Let  $(ij)$  and  $(ij)_w$  denote the Shapley value of the player  $i$  in the game  $v$  and game  $w$ , respectively. If these two games are combined into a new game, and the reward function becomes  $v_w(S) = v(S) + w(S)$ , then the Shapley value comes to be  $(ij)_{v_w} = (ij)_v + (ij)_w$  for each player in  $N$ .
- **Dummy property:** A player  $i \in N$  is referred to as a dummy player if  $v(S \cup \{i\}) = v(S) + v(\{i\})$ . In this way,  $(ij) = (ij)_{v(\{i\})}$ , which means that player  $i$  plays the game independently.
- **Symmetry property:** If  $v(S \cup \{i, j\}) = v(S \cup \{j, i\})$ , then Shapley values of player  $i$  and  $j$  are equal, i.e.  $(ij) = (ji)$ .
- **Efficiency property:** The sum of each individual's Shapley value is equal to the reward won by the coalition  $N$ , i.e.  $\sum_i (ij) = v(N) - v(\emptyset)$ . This property guarantees the overall reward can be allocated to each player in the game.

### A.2 MOTIVATIONS

**Theoretical rigor.** We use the Shapley interaction index defined based on the Shapley value, because the Shapley value has a solid theoretical foundation in the game theory, which is the attribution satisfying the above four desirable axioms.

**Whether the metric depends on network architectures.** Because adversarial transferability is a general property for the attack, a convincing metric for adversarial transferability is supposed not to be directly related to the network architecture. To this end, the computation of the interaction defined on the Shapley value does not depend on the network architecture. In comparison, previous definitions of the interaction are usually oriented to model architectures. For example, the interaction proposed by Tsang et al. (2018) requires the DNN to be fully-connected. The two interaction metrics proposed by Murdoch et al. (2018) and Jin et al. (2020) are designed for LSTMs. The Hessian-based interaction (Janizek et al., 2020) requires the DNN to use the softPlus operation to replace the ReLU operation.

**Computational cost.** The computational cost of the Shapley-based interaction-reduction loss is relatively low. Because of the efficiency axiom of the Shapley value, we prove that the time cost of computing the interaction loss  $\mathcal{L}_{\text{interaction}} = \frac{1}{n-1} E_i [v(\emptyset) - v(N \setminus \{i\}) - v(\{i\}) + v(N)]$  is linear, i.e.  $O(n)$ , where  $n$  is the dimension of features. The linear complexity makes it possible to apply the interaction to high-dimensional data and deep neural networks. In contrast, the complexity of computing all possible pairwise interactions defined in (Daria Sorokina, 2008) is  $O(n^2)$ .

## B COMPARISONS BETWEEN INTERACTIONS INSIDE PERTURBATIONS OF DIFFERENT ATTACKS

We have theoretically proved that some classical attacking methods of boosting the adversarial transferability essentially decrease interactions inside perturbations. Besides the theoretical proof in Ap-

Table 5: The average interaction inside adversarial perturbations generated by different attacks.

Method	RN-34	RN-152	DN-121	DN-201
Baseline (PGD Attack)	0.422	0.926	0.909	0.784
SGM Attack	-0.012	0.037	0.395	0.308
VR Attack	0.097	0.270	0.242	0.137

pendix I and Appendix H, we also conduct experiments to compare interactions of perturbation units when we generate adversarial perturbations with and without these attacking methods. Such experiments further verify that these methods of boosting the transferability essentially decrease interactions. We conduct attacks with the validation set in the ImageNet dataset on four DNNs, and measure the average interaction inside perturbation units. As Table 5 shows, the SGM Attack and the VR Attack decrease interactions inside perturbations.

## C ADVERSARIAL ATTACK

In general, the objective of adversarial attacking can be formulated as the following optimization problem.

$$\text{maximize } \ell(h(x + \delta); y) \quad \text{s.t. } \|\delta\|_{k_p} \leq \epsilon; \quad x + \delta \in [0; 1]^n; \quad (7)$$

where  $\ell(h(x + \delta); y)$  is the classification loss. There are many ways to solve the above optimization problem under different norm constraints  $\|\delta\|_{k_p}$  (Goodfellow et al., 2014; Carlini & Wagner, 2017; Kurakin et al., 2017; Madry et al., 2018; Chen et al., 2018b; Wang et al., 2019).

Optimization-based approach. One approach to approximately solving Equation (7) is to solve the following relaxed form:

$$\text{minimize } f = \ell(h(x + \delta); y) + c \|\delta\|_{k_p} \quad \text{s.t. } \delta \in [0; 1]^n; \quad (8)$$

where  $c > 0$  is a scalar constant to balance the classification loss and the norm constraint. Szegedy et al. (2013); Carlini & Wagner (2017) have demonstrated the effectiveness of this method.

Projected gradient descent (PGD) (Madry et al., 2018). The PGD Attack is usually considered as one of the simplest and the most widely used baseline for adversarial attacking. In this paper, this method is called the Baseline. The PGD Attack directly optimizes the classification loss in Equation (7). Considering the norm constraint, after each step of updating, the PGD Attack projects the adversarial perturbation back to the  $\epsilon$ -ball, if the perturbation goes beyond the ball.

PGD updates adversarial perturbations in each step with the following equation:

$$\delta^{t+1} = \begin{cases} \text{proj}_{\|\cdot\|_{k_p} \leq \epsilon}(\delta^t + \text{sign}(r \cdot \nabla \ell(h(x + \delta^t); y))) & ; \quad p = 1 \\ \text{proj}_{\|\cdot\|_{k_p} \leq \epsilon}(\delta^t + \frac{r \cdot \nabla \ell(h(x + \delta^t); y)}{\|\nabla \ell(h(x + \delta^t); y)\|_{k_2}}) & ; \quad p = 2; \end{cases} \quad (9)$$

where  $\delta^t$  denotes the perturbation of the  $t$ th step.  $\text{proj}_{\|\cdot\|_{k_p} \leq \epsilon}(\cdot)$  and  $\text{proj}_{\|\cdot\|_{k_2} \leq \epsilon}(\cdot)$  are projection operations, which project the perturbation back to the  $\epsilon$ -ball, if the perturbation goes beyond the ball in the step size. Given  $\delta \in \mathbb{R}^n$ , we have:

$$\text{proj}_{\|\cdot\|_{k_p} \leq \epsilon}(\delta) = \begin{cases} \delta & ; \quad \|\delta\|_{k_p} \leq \epsilon \\ \frac{\epsilon}{\|\delta\|_{k_p}} \delta & ; \quad \|\delta\|_{k_p} > \epsilon \end{cases}; \quad \text{proj}_{\|\cdot\|_{k_2} \leq \epsilon}(\delta) = \begin{cases} \delta & ; \quad \|\delta\|_{k_2} \leq \epsilon \\ \frac{\epsilon}{\|\delta\|_{k_2}} \delta & ; \quad \|\delta\|_{k_2} > \epsilon \end{cases}; \quad (10)$$

## D EQUIVALENT FORMS OF THE INTERACTION

In Section 3.1, the interaction between units  $i$  and  $j$  is defined as the additional contribution as follows.

$$I_{ij}(\delta) = (S_{ij} \cdot \delta) \cdot [(i, j) \cdot \delta + (j, i) \cdot \delta]; \quad (11)$$

where  $(S_{ij} \cdot \delta)$  denotes the joint contribution of  $i$  and  $j$ , when perturbation units  $i$  and  $j$  are regarded as a singleton unit  $S_{ij} = \{i, j\}$ , as follows.

$$(S_{ij} \cdot \delta) = \sum_{S \ni i, j} \frac{|S|! (n - |S| - 2)!}{(n - 1)!} (v(S[f; i, j]) - v(S));$$

where  $S_{ij} = f; j; g$  represents the coalition of perturbation units. In this game, because perturbation units  $i; j$  are regarded as a singleton player, we can consider there are only 1 players in the game, and consequently the set of players changes to  $\{f; j; g\} \cup S_{ij}$ .

$(ij \setminus n \setminus f; j; g)$  and  $(j \setminus n \setminus f; g)$  represent the individual contributions of unit  $i$  and  $j$ , respectively, when the perturbation units  $j$  work individually. The individual contribution of perturbation unit  $i$ , when perturbation unit  $j$  is absent, is given as follows.

$$(ij \setminus n \setminus f; j; g) = \sum_{S \setminus n \setminus f; j; g} \frac{j! S_j! (n - j - S_j - 2)!}{(n - 1)!} (v(S \setminus f; i; g) - v(S));$$

In this game, because the perturbation unit  $i$  is always absent, we can consider there are only 1 players in the game. Consequently the set of players changes to  $\{g\}$ .

Similarly, the individual contribution of perturbation unit  $j$  when perturbation unit  $i$  is absent, is given as follows.

$$(j \setminus n \setminus f; i; g) = \sum_{S \setminus n \setminus f; i; g} \frac{j! S_j! (n - j - S_j - 2)!}{(n - 1)!} (v(S \setminus f; j; g) - v(S));$$

In Section 1, the interaction between perturbation units  $i$  and  $j$  is defined as the change of the importance  $\rho_i$  of the  $i$ -th unit when the  $j$ -th unit  $j$  is perturbed w.r.t the case when the  $j$ -th unit  $j$  is not perturbed. If the perturbation  $j$  on the  $j$ -th unit increases the importance of the  $i$ -th unit, then there is a positive interaction between  $i$  and  $j$ . If the perturbation  $j$  decreases the importance of  $i$ , it indicates a negative interaction. Mathematically, this definition can be written as follows.

$$I_{ij}^0(\cdot) = \rho_{i; w=j} - \rho_{i; w=0j}; \quad (12)$$

where  $\rho_{i; w=j}$  represents the importance of  $i$ , when  $j$  is always present;  $\rho_{i; w=0j}$  represents the importance of  $i$ , when  $j$  is always absent. When perturbation unit  $j$  is always present, the contribution of perturbation unit  $i$  is given as follows.

$$\rho_{i; w=j} = \sum_{S \setminus n \setminus f; i; j; g} \frac{j! S_j! (n - j - S_j - 2)!}{(n - 1)!} (v(S \setminus f; i; j; g) - v(S \setminus f; j; g));$$

In this game, because the perturbation unit  $j$  is always present, we can consider there are only 1 players.

When perturbation unit  $i$  is always absent, the contribution of perturbation unit  $j$  is given as follows.

$$\rho_{i; w=0j} = \sum_{S \setminus n \setminus f; i; g} \frac{j! S_j! (n - j - S_j - 2)!}{(n - 1)!} (v(S \setminus f; i; g) - v(S));$$

In this game, because the perturbation unit  $i$  is always absent, we can consider there are only 1 players.

The interaction in Equation (11) is equal to the interaction in Equation (12),

$$I_{ij}(\cdot) = I_{ij}^0(\cdot)$$

## E PROOF OF PROPOSITION 1

To simplify the problem setting, we do not consider some tricks in adversarial attacking, such as gradient normalization and the clip operation. In multi-step attacking, the final perturbation generated after  $t$  steps is given as follows.

$$t_{\text{multi}} \stackrel{\text{def}}{=} \sum_{t^0=0}^{t-1} \alpha^t \cdot \nabla_x \ell(h(x + t_{\text{multi}}^{t^0}); y);$$

where  $\alpha$  represents the step size, and  $\ell(x; y)$  is referred as the classification loss.



To simplify the notation, we use  $g(x)$  to denote  $x^T(h(x); y)$ , i.e.  $g(x) \stackrel{\text{def}}{=} x^T(h(x); y)$ . Furthermore, we define the update of the perturbation with the multi-step attack at each step as follows.

$$x_{\text{multi}}^t \stackrel{\text{def}}{=} g(x + x_{\text{multi}}^{t-1}); \quad (13)$$

In this way, the perturbation can be written as follows.

$$x_{\text{multi}}^t = x_{\text{multi}}^1 + x_{\text{multi}}^2 + \dots + x_{\text{multi}}^{t-1}. \quad (14)$$

Lemma 1. Given the sample  $x \in \mathbb{R}^n$  and the adversarial perturbation  $\delta \in \mathbb{R}^n$ , we use  $S = \{1, 2, \dots, n\}$  to denote the set of all perturbation units. The score function is denoted by  $v(S) = L(x + \delta^{(S)})$ , where  $\delta^{(S)}$  satisfies  $\delta_i \in S; \delta_i = 0; \delta_i \notin S; \delta_i = 0$ . The Shapley interaction between perturbation units  $a, b$  can be written as  $I_{ab} = g_a H_{ab}(x)_b + R_2(\delta)$ , where  $H_{ab}(x) = \frac{\partial^2 L(x)}{\partial x_a \partial x_b}$  represents the element of the Hessian matrix, and  $R_2(\delta)$  denotes terms with elements of higher than the second order.

Proof. The Shapley interaction between perturbation units is

$$I_{ab}(\delta) = \sum_{S \ni a, b} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} [v(S \cup \{a, b\}) - v(S \cup \{b\}) - v(S \cup \{a\}) + v(S)];$$

where  $v(S) = L(x + \delta^{(S)})$ . Here, the classification loss can be approximated as  $v(S) = L(x) + g^T(x) + \frac{1}{2} \delta^T H(x) \delta + R_2(\delta)$  using Taylor series. Thus  $v(S^0) =$

$$v(S^0) = L(x) + \sum_{a \in S^0} g_a(x) \delta_a + \frac{1}{2} \sum_{a, b \in S^0} g_{ab} H_{ab}(x) \delta_a \delta_b + R_2^0(\delta);$$

where  $R_2^0(\delta)$  denotes terms with elements of higher than the second order.

In this way, the Shapley interaction  $I_{ab}$  is given as

$$\begin{aligned} I_{ab}(\delta) &= \sum_{S \ni a, b} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} [v(S \cup \{a, b\}) - v(S \cup \{b\}) - v(S \cup \{a\}) + v(S)] \\ &= \sum_{S \ni a, b} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} \{ [L(x) + \sum_{a^0 \in S \cup \{a, b\}} g_{a^0}(x) \delta_{a^0} + \frac{1}{2} \sum_{a^0, b^0 \in S \cup \{a, b\}} g_{a^0 b^0} H_{a^0 b^0}(x) \delta_{a^0} \delta_{b^0} + R_2^{(S \cup \{a, b\})}(\delta)] \\ &\quad - [L(x) + \sum_{a^0 \in S \cup \{b\}} g_{a^0}(x) \delta_{a^0} + \frac{1}{2} \sum_{a^0, b^0 \in S \cup \{b\}} g_{a^0 b^0} H_{a^0 b^0}(x) \delta_{a^0} \delta_{b^0} + R_2^{(S \cup \{b\})}(\delta)] \\ &\quad - [L(x) + \sum_{a^0 \in S \cup \{a\}} g_{a^0}(x) \delta_{a^0} + \frac{1}{2} \sum_{a^0, b^0 \in S \cup \{a\}} g_{a^0 b^0} H_{a^0 b^0}(x) \delta_{a^0} \delta_{b^0} + R_2^{(S \cup \{a\})}(\delta)] \\ &\quad + L(x) + \sum_{a^0 \in S} g_{a^0}(x) \delta_{a^0} + \frac{1}{2} \sum_{a^0, b^0 \in S} g_{a^0 b^0} H_{a^0 b^0}(x) \delta_{a^0} \delta_{b^0} + R_2^{(S)}(\delta) \} \\ &= \sum_{S \ni a, b} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} g_a H_{ab}(x)_b \delta_a \delta_b \\ &\quad + \sum_{S \ni a, b} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} [R_2^{(S \cup \{a, b\})}(\delta) - R_2^{(S \cup \{a\})}(\delta) - R_2^{(S \cup \{b\})}(\delta) + R_2^{(S)}(\delta)] \\ &= \sum_{s=0}^n \sum_{S \ni a, b; |S|=s} \frac{s!(n - s - 2)!}{(n - 1)!} [g_a H_{ab}(x)_b \delta_a \delta_b + R_2(\delta)] \end{aligned}$$

$$\begin{aligned}
&= \sum_{s=0}^{\infty} \frac{X^2 (n-2)!}{s!(n-s-2)!} \frac{s!(n-s-2)!}{(n-1)!} [{}_a H_{ab}(x)_b] + \hat{R}_2(\cdot) \\
&= {}_a H_{ab}(x)_b + \hat{R}_2(\cdot);
\end{aligned}$$

where  $\hat{R}_2(\cdot)$  denotes terms with elements of higher than the second order.  $\square$

Lemma 2. The update of the perturbation with the multi-step attack at step  $t$  defined in Equation (13) can be written as  $x_{\text{multi}}^t = [I + H(x)]^{t-1} g(x) + \hat{R}_1^t$ , where  $g(x) \stackrel{\text{def}}{=} r_x(h(x); y)$  represents the gradient, and  $H(x) \stackrel{\text{def}}{=} r_x^2(h(x); y)$  represents the Hessian matrix.  $\hat{R}_1^t$  denotes terms with elements in  $x_{\text{multi}}^t$  of higher than the first order.

Proof. If  $t = 1$ ,  $x_{\text{multi}}^1 = g(x)$ .

Let  $0 < t < \infty$ ;  $x_{\text{multi}}^{t-1} = [I + H(x)]^{t-2} g(x) + \hat{R}_1^{t-1}$ , then we have

$$\begin{aligned}
x_{\text{multi}}^t &= g(x + x_{\text{multi}}^{t-1}) \quad \text{According to Equation (13)} \\
&= g(x + x_{\text{multi}}^{t-1} + x_{\text{multi}}^{t-2} + \dots + x_{\text{multi}}^1) \quad \text{According to Equation (14)} \\
&= g(x + [I + H(x)] + [I + H(x)]^2 + \dots + [I + H(x)]^{t-2}) g(x) + \hat{R}_1^t;
\end{aligned}$$

where  $\hat{R}_1^t$  denotes terms of elements in  $x_{\text{multi}}^t$  of higher than the first order.

Using the Taylor series, we get

$$x_{\text{multi}}^t = g(x) + H(x)T(x) + \frac{1}{2} H^2(x)T^2(x) + \dots + \frac{1}{(t-1)!} H^{t-1}(x)T^{t-1}(x) + \hat{R}_1^t; \quad (15)$$

where  $T(x)$  in Equation (15) is given as follows.

$$T(x) = [I + H(x)] + [I + H(x)]^2 + \dots + [I + H(x)]^{t-2} g(x); \quad (16)$$

Multiply  $[I + H(x)]$  on both sides of Equation (16), and we get

$$(I + H(x))T(x) = [I + H(x)]^2 + [I + H(x)]^3 + \dots + [I + H(x)]^{t-1} g(x); \quad (17)$$

Then, according to Equation (17) and Equation (16), we get

$$H(x)T(x) = [I + H(x)]^{t-1} g(x); \quad (18)$$

Substituting Equation (18) back to Equation (15), we have

$$x_{\text{multi}}^t = [I + H(x)]^{t-1} g(x) + \hat{R}_1^t;$$

In this way, we have proved that  $x_{\text{multi}}^t = [I + H(x)]^{t-1} g(x) + \hat{R}_1^t$ .  $\square$

Proposition 1. The adversarial perturbation generated by the multi-step attack via gradient descent is given as  $x_{\text{multi}}^m = \prod_{t=0}^{m-1} r_x^{-1}(h(x + x_{\text{multi}}^t); y)$ , where  $x_{\text{multi}}^t$  denotes the perturbation after the  $t$ -th step of updating, and  $m$  is referred to as the total number of steps. The adversarial perturbation generated by the single-step attack is given as  $x_{\text{single}} = r_x^{-1}(h(x); y)$ . The expectation of interactions between perturbation units in  $x_{\text{multi}}^m$ ,  $E_{a,b}[I_{ab}(x_{\text{multi}}^m)]$ , is larger than  $E_{a,b}[I_{ab}(x_{\text{single}})]$ , i.e.  $E_{a,b}[I_{ab}(x_{\text{multi}}^m)] > E_{a,b}[I_{ab}(x_{\text{single}})]$ .

### E.1 FAIRNESS OF COMPARISONS OF INTERACTIONS INSIDE DIFFERENT PERTURBATIONS

Proposition 1 is valid for different loss functions of generating of adversarial perturbations. In this section, we discuss the fairness of comparisons of interactions inside different perturbations.

When we compare interactions inside different perturbations, magnitudes of these perturbations should be similar, because the comparison of interactions between adversarial perturbations of different magnitudes is not fair. For fair comparisons, in Section 3.1, this paper controls the magnitude of the single-step attack by setting the step size of the single-step attack, where  $\alpha$  denotes the step size and the total number of steps of the multi-step attack, respectively. The equivalent step size  $\alpha/m$  makes the magnitude of perturbations generated by the single-step attack to be similar to that of perturbations generated by the multi-step attack, when we use the target score before the softmax layer to generate adversarial perturbations,  $\text{sub}(h(x); y) = \max_{y \in \mathcal{Y}} h(x) - h_y(x)$ . In this case, the magnitude of the gradient  $\nabla_x \text{sub}(h(x); y)$  is relatively stable. In particular, this type of loss has been widely used. For example, one of the most widely used attacking (Carlini & Wagner, 2017), uses the score before the softmax layer for targeted attacking.

### E.2 PROOF OF PROPOSITION 1

Proof. According to Lemma 2, the update of the perturbation with the multi-step attack at the step  $t$  is given as follows.

$$x_{\text{multi}}^t = [I + H(x)]^{t-1} g(x) + \hat{R}_1^t; \quad (19)$$

where  $\hat{R}_1^t$  denotes terms with elements in  $\mathcal{R}_{\text{multi}}^1$  of higher than the first order, and represents the step size.

To simplify the notation without causing ambiguity, we write  $g(x)$  and  $H(x)$  as  $g$  and  $H$ , respectively. In this way, according to Equation (14) and Equation (19),  $x_{\text{multi}}^m$  can be written as follows.

$$\begin{aligned} x_{\text{multi}}^m &= \prod_{t=1}^m [I + [I + H] + [I + H]^2 + \dots + [I + H]^{m-1}] g + \sum_{t=1}^m \hat{R}_1^t \\ &= mI + \frac{m(m-1)}{2}H + \dots + g + \sum_{t=1}^m \hat{R}_1^t; \end{aligned} \quad (20)$$

where  $m$  represents the total number of steps. According to Lemma 1, the Shapley interaction between perturbation units  $a$  and  $b$  in  $x_{\text{multi}}^m$  is given as follows.

$$I_{ab}(x_{\text{multi}}^m) = \frac{m}{m_{\text{multi}};a} H_{ab} \frac{m}{m_{\text{multi}};b} + \hat{R}_2(x_{\text{multi}}^m); \quad (21)$$

where  $\hat{R}_2(x_{\text{multi}}^m)$  denotes terms with elements in  $\mathcal{R}_{\text{multi}}^2$  of higher than the second order.

According to Equation (20) and Equation (21), we have

$$\begin{aligned}
I_{ab}(\frac{m}{\text{multi}}) &= H_{ab}^h mg_a + \frac{2m(m-1)}{2} \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) + \sum_{t=1} X^n \left\{ \frac{\alpha(\frac{t}{\text{multi}}; a)}{Z} \right\} \quad \text{in} \\
&\quad \text{terms of } \frac{t}{\text{multi}}; a \\
&\quad \text{of higher than the first order,} \\
&\quad \text{which corresponds to the term of} \\
&\quad \hat{R}_1^t \text{ in Equation (20)} \\
mg_b + \frac{2m(m-1)}{2} \sum_{a^0=1} X^n (H_{a^0b} g_{a^0}) + \sum_{t=1} X^n \left\{ \frac{\alpha(\frac{t}{\text{multi}}; b)}{Z} \right\} \quad \text{in} \\
&\quad \text{terms of } \frac{t}{\text{multi}}; b \\
&\quad \text{of higher than the first order,} \\
&\quad \text{which corresponds to the term of} \\
&\quad \hat{R}_1^t \text{ in Equation (20)} \\
&= \frac{2m^2 g_a g_b H_{ab}}{Z} \quad \text{first-order terms w.r.t. elements in } H \\
&+ \frac{3(m-1)m^2}{2} \sum_{b^0=1} X^n g_b (H_{ab^0} g_{b^0}) + \frac{3(m-1)m^2}{2} \sum_{a^0=1} X^n g_a (H_{a^0b} g_{a^0}) H_{ab} \quad \# \\
&\quad \text{second-order terms w.r.t. elements in } H \\
&+ \frac{4(m-1)^2 m^2}{4} \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) (H_{a^0b} g_{a^0}) H_{ab} + \dots \quad \# \\
&\quad \text{R}_2^{\text{multi}}(H) \\
&+ \left[ \sum_{t=1} X^n \left\{ \frac{\alpha(\frac{t}{\text{multi}}; a)}{Z} \right\} \right] H_{ab} \frac{m}{\text{multi}}; b + \left[ \sum_{t=1} X^n \left\{ \frac{\alpha(\frac{t}{\text{multi}}; b)}{Z} \right\} \right] H_{ab} \frac{m}{\text{multi}}; a + \hat{R}_2(\frac{m}{\text{multi}}) \\
&\quad \text{R}_2^0(\frac{m}{\text{multi}}) \\
&= 2m^2 g_a g_b H_{ab} + \frac{3(m-1)m^2}{2} \sum_{a^0=1} X^n g_a H_{ab} (H_{a^0b} g_{a^0}) \\
&+ \frac{3(m-1)m^2}{2} \sum_{b^0=1} X^n g_b H_{ab} (H_{ab^0} g_{b^0}) + \hat{R}_2^0(\frac{m}{\text{multi}}) + R_2^{\text{multi}}(H); \tag{22}
\end{aligned}$$

where  $R_2^{\text{multi}}(H)$  represents terms with elements in  $H$  of higher than the second order, and  $\hat{R}_2(\frac{m}{\text{multi}})$  represents terms with elements in  $\frac{m}{\text{multi}}$  of higher than the second order.

Let us consider the single-step attack. When we compare interactions inside different perturbations, magnitudes of these perturbations should be similar, because the comparison of interactions between adversarial perturbations of different magnitudes is not fair. For fair comparisons, in Section 3.1, this paper controls the magnitude of the single-step attack, as follows. The single-step attack only uses the gradient information on the original input which generates adversarial perturbations as:

$$\text{single} = mg:$$

Therefore, according to Lemma 1, the interaction between perturbations  $\text{single}$  is given as follows.

$$\begin{aligned}
I_{ab}(\text{single}) &= \text{single}_a H_{ab} \text{single}_b + \hat{R}_2(\text{single}) \\
&= m^2 2g_a g_b H_{ab} + \hat{R}_2(\text{single}); \tag{23}
\end{aligned}$$

where  $\hat{R}_2(\text{single})$  denotes terms with elements in  $\text{single}$  of higher than the second order. In this way, according to Equation (22) and Equation (23), the expectation of the difference between  $I_{ab}(\frac{m}{\text{multi}})$  and  $I_{ab}(\text{single})$  is given as follows.

Figure 4: (a) Histograms of the value of the Hessian element  $H_{ab}(x)$  w.r.t. different values of  $a; b$ ; (b) Histograms of the value of  $\frac{g_b H_{ab}}{g_a}$  w.r.t. different values of  $b$ . Because the Hessian of the DNN with the ReLU activation is not well defined, we replace the ReLU activation with the Softplus activation  $f(x) = \frac{1}{2} \log(1 + e^x)$ . We train VGG-16, ResNet-32, and DensetNet-121 on the CIFAR-10 dataset (Krizhevsky et al., 2009), and use the cross-entropy loss as the classification loss.

$$\begin{aligned}
 & E_{a,b} [I_{ab}(\frac{m}{\text{multi}}) - I_{ab}(\text{single})] \\
 &= E_{a,b} \left[ \frac{3(m-1)m^2}{2} g_a H_{ab} \prod_{a^0=1}^n (H_{a^0 b} g_{a^0}) + \frac{3(m-1)m^2}{2} g_b H_{ab} \prod_{b^0=1}^n (H_{ab^0} g_{b^0}) \right] \\
 &+ R_2^0(\frac{m}{\text{multi}}) + R_2^{\text{multi}}(H) - R_2(\text{single}) \\
 &= \frac{3(m-1)m^2}{2} E_{a,b} \left[ \underbrace{g_a H_{ab} \prod_{a^0=1}^n (H_{a^0 b} g_{a^0})}_{U_{ab}} + \underbrace{g_b H_{ab} \prod_{b^0=1}^n (H_{ab^0} g_{b^0})}_{U_{ba}} \right] + E_{a,b} [R_{ab}];
 \end{aligned}$$

where

$$R_{ab} = R_2^0(\frac{m}{\text{multi}}) + R_2^{\text{multi}}(H) - R_2(\text{single});$$

Assumption 1 Magnitudes of elements in the Hessian matrix  $H(x)$  is small that  $|H_{ab}(x)| \leq 1$ , where  $1 \leq a, b \leq n$ . Therefore  $H^k(x) = 0$ , if  $k > 2$ .

We verify the assumption by directly measuring the value of  $H_{ab}(x)$ . As Figure 4 (a) shows, the value of  $H_{ab}(x)$  is very small that  $|H_{ab}(x)| \leq 1$ .

According to Assumption 1, we have  $R_2^{\text{multi}}(H) = 0$ . Note that the magnitude of  $\frac{m}{\text{multi}}$  and the magnitude of  $\text{single}$  are small, then  $R_2^0(\frac{m}{\text{multi}}) = 0$ , and  $R_2(\text{single}) = 0$ . In this way, we have  $E_{a,b} [R_{ab}] = E_{a,b} [R_2^0(\frac{m}{\text{multi}}) + R_2(\text{single}) + R_2^{\text{multi}}(H)] = 0$ .

Moreover, for the expectation of  $U_{ab}$ , we have

$$E_{a,b} [U_{ab}] = \frac{1}{n(n-1)} \sum_{b=1}^n \sum_{a \neq b}^n g_a H_{ab} \prod_{a^0=1}^n (g_{a^0} H_{a^0 b})$$

Figure 5: Histograms of the value  $\mathbb{E}_b[I_{ab}]$  w.r.t. different values of  $a$

$$= \frac{1}{n(n-1)} \sum_{b=1}^{n-1} \sum_{a=1}^{b-1} g_a H_{ab} \sum_{a^0=1}^{b-1} \{Z_{A^0}\} + \sum_{b=1}^{n-1} \sum_{a^0=1}^{b-1} g_{a^0} H_{a^0 b} \sum_{a=1}^{b-1} \{Z_A\}$$

Let us focus on terms  $A$  and  $B$ . Note that  $A$  is the sum of  $n$  terms ( $n$  is large). In comparison,  $B$  is just a single term in  $A$ . Therefore, the sign of  $A - B$  is usually dominated by the term  $A$ . In this way, we get  $\text{Prob}[\text{sign}(A - B) = \text{sign}(A)] \approx 1$ . Therefore,  $\text{Prob}[(A - B)A \leq 0] \approx 1$ . We verify this assumption by measuring the value of  $\frac{g_b H_{bb}}{\sum_{a=1}^{b-1} g_a H_{ab}}$ . If  $\text{Prob} \left[ \frac{g_b H_{bb}}{\sum_{a=1}^{b-1} g_a H_{ab}} \geq 1 \right] \approx 1$ , then we have  $\text{Prob}[\text{sign}(A - B) = \text{sign}(A)] \approx 1$ . As Figure 4 (b) shows, the value of  $\frac{g_b H_{bb}}{\sum_{a=1}^{b-1} g_a H_{ab}}$  is very small that  $\frac{g_b H_{bb}}{\sum_{a=1}^{b-1} g_a H_{ab}} \ll 1$ . To this end, we have  $(A - B)B \leq 0$ , and we get

$$\mathbb{E}_{a,b}[U_{ab}] \leq 0 \tag{24}$$

Due to the symmetry of  $a$  and  $b$ , we have  $\mathbb{E}_{a,b}[U_{ba}] = \mathbb{E}_{a,b}[U_{ab}]$ . Therefore,

$$\begin{aligned} & \mathbb{E}_{a,b}[I_{ab}(\text{multi}) - I_{ab}(\text{single})] \\ &= \frac{3(m-1)m^2}{2} \mathbb{E}_{a,b}[U_{ab} + U_{ba}] + \mathbb{E}_{a,b}[R_{ab}] \\ &= \frac{3(m-1)m^2}{2} \mathbb{E}_{a,b}[U_{ab}] + 0 \end{aligned}$$

0:

□

### E.3 VERIFICATION OF PROPOSITION 1

We verify that perturbations generated by the multi-step attack tend to exhibit larger interaction than those generated by the single-step attack by measuring the value of  $\mathbb{E}_b[I_{ab}]$ . As shown in Appendix F, we prove that  $\mathbb{E}_b[I_{ab}] = v(\cdot) - v(\cdot - n f_{ag}) - v(f_{ag}) + v(\cdot; \cdot)$ . Because the image data is high-dimensional, the cost of computing  $\mathbb{E}_b[I_{ab}]$  is high. As Appendix J.1 demonstrates, given the input image, we can measure the interaction at the grid level, instead of the pixel level, to reduce the computational cost. Therefore, we divide the input image into 166 ( $L = 16$ ) grids, and use Equation (39) to compute the interaction  $\mathbb{E}_{(p,q)} I_{(p,q);(p^0,q^0)}(\cdot) = v(\cdot) - v(\cdot - n f_{pqg}) - v(f_{pqg}) + v(\cdot; \cdot)$ , where  $(p; q)$  denotes the coordinate of a grid. The experiments were conducted with ImageNet validation images on ResNet-32 and DenseNet-121.

For fair comparisons, the magnitude of perturbations generated by the single-step attack is controlled to be same as that generated by the multi-step attack. As Figure 5 (left) shows, perturbations generated by the multi-step attack tend to exhibit larger interaction than those generated by the single-step attack.

#### E.4 PERTURBATIONS GENERATED BY THE MULTISTEP ATTACK TEND TO EXHIBIT LARGER INTERACTION THAN GAUSSIAN NOISE

Moreover, we compare the interaction inside perturbation units generated by the multi-step attack with the Gaussian noise perturbation. Similarly, for fair comparisons, the magnitude of the Gaussian noise is controlled to be similar to that generated by the multi-step attack. As Figure 5 (right) shows, perturbations generated by the multi-step attack tend to exhibit larger interaction than Gaussian noise.

#### F EXPECTATION OF THE SHAPLEY INTERACTION

In Equation (3), the Shapley interaction between two perturbation units is given as follows.

$$I_{ij}(\mathbf{x}) = (S_{ij} - \sum_{j \in S} v_{ij}(\mathbf{x})) - (v_{ij}(\mathbf{x}) + \sum_{j \in S} v_{ij}(\mathbf{x}));$$

where  $(S_{ij} - \sum_{j \in S} v_{ij}(\mathbf{x}))$  is the Shapley value of the singleton unit  $S_{ij} = \{i, j\}$ , when perturbation units  $i, j$  form a coalition.  $v_{ij}(\mathbf{x})$  and  $\sum_{j \in S} v_{ij}(\mathbf{x})$  are Shapley values of perturbation units  $i, j$ , when these two perturbation units work individually. In this way, we can write the Shapley interaction in a closed form as follows.

$$I_{ij}(\mathbf{x}) = \sum_{S \ni i, j} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} [v(S \setminus \{i, j\}) - v(S \setminus \{j\}) - v(S \setminus \{i\}) + v(S)]; \quad (25)$$

where  $v(S) = \max_{y \in \mathbb{R}^n} h_{y^0}^{(S)}(\mathbf{x} + \mathbf{e}^{(S)}) - h_y^{(S)}(\mathbf{x} + \mathbf{e}^{(S)})$ . The expectation of interaction is given as follows.

$$E_{i,j} [I_{ij}(\mathbf{x})] = \frac{1}{n-1} E_i [v(\mathbf{x}) - v(\mathbf{x} \setminus \{i\}) - v(\mathbf{x} \setminus \{j\}) + v(\mathbf{x})]; \quad (26)$$

which is proved as follows.

Proof. As proved in Appendix D,  $I_{ij}(\mathbf{x}) = I_{ij}^0(\mathbf{x})$ : Therefore, the interaction between players  $i$  and  $j$  is given as follows.

$$\begin{aligned} I_{ij}(\mathbf{x}) &= \sum_{S \ni i, j} \frac{|S|!(n - |S| - 2)!}{(n - 1)!} [v(S \setminus \{i, j\}) - v(S \setminus \{j\}) - v(S \setminus \{i\}) + v(S)] \\ &= \sum_{j \in S} v_{ij}(\mathbf{x}) - \sum_{j \in S} v_{ij}(\mathbf{x}); \end{aligned}$$

The expectation of the interaction can be written as follows.

$$E_{i,j} [I_{ij}(\mathbf{x})] = \frac{1}{(n-1)} E_i \left[ \sum_{j \in S} v_{ij}(\mathbf{x}) - \sum_{j \in S} v_{ij}(\mathbf{x}) \right];$$

According to the efficiency property of Shapley values (please refer to Appendix A.1 for details):

$$\begin{aligned} \sum_{j \in S} v_{ij}(\mathbf{x}) &= v(\mathbf{x}) - v(\mathbf{x} \setminus \{i\}) \\ \sum_{j \in S} v_{ij}(\mathbf{x}) &= v(\mathbf{x} \setminus \{i\}) - v(\mathbf{x}); \end{aligned}$$

In this way,

$$E_{i,j} [I_{ij}(\mathbf{x})] = \frac{1}{n-1} E_i [v(\mathbf{x}) - v(\mathbf{x} \setminus \{i\}) - v(\mathbf{x} \setminus \{i\}) + v(\mathbf{x})];$$

□

## G DETAILS OF OBSERVING THE NEGATIVE CORRELATION BETWEEN THE TRANSFERABILITY AND THE INTERACTION

In Section 3.2, we directly measure the transfer utility and interactions of different adversarial perturbations. Here, we give more details of the experiments. We measure the transfer utility as  $\text{Transfer Utility} = [\max_{y \in \mathcal{Y}} h_{y^0}^{(t)}(x + \delta) - h_y^{(t)}(x + \delta)] / [\max_{y \in \mathcal{Y}} h_{y^0}^{(t)}(x) - h_y^{(t)}(x)]$ . We measure the interaction as  $\mathbb{E}_{i,j} [I_{ij}(\delta)] = \frac{1}{n-1} \mathbb{E}_i [v(\delta) - v(\text{fig}) - v(i) + v(\text{fig} + i)]$ . As Appendix J.1 demonstrates, to reduce the computational cost, given the input image, we can measure the interaction at the grid level, instead of the pixel level. Therefore, we divide the input image into 16 ( $L = 16$ ) grids, and use Equation (39) to compute the interaction as  $\mathbb{E}_{(p,q)} [I_{(p,q)}(\delta)] = \frac{1}{L^2-1} \mathbb{E}_{(p,q)} [v(\delta) - v(\text{fig} + \delta_{(p,q)}) - v(\delta_{(p,q)}) + v(\text{fig} + \delta_{(p,q)})]$ , where  $(p, q)$  denotes the coordinate of a grid.

Using the validation set of the ImageNet dataset (Russakovsky et al., 2015), we generate adversarial perturbations on four types of DNNs, including ResNet-34/152 (RN-34/152) (He et al., 2016) and DenseNet-121/201 (DN-121/201) (Huang et al., 2017). We transfer adversarial perturbations generated on each ResNet to DenseNets. Similarly, we also transfer adversarial perturbations generated on each DenseNet to ResNets. Given an input image, adversarial perturbations are generated using Equation (8), i.e.  $\min_{\delta} \ell(h(x + \delta); y) + c \|\delta\|_p$  s.t.  $\delta \in [0, 1]^n$ , where  $c \in \mathbb{R}$  is a scalar constant. In this way, we gradually change the value of different hyper-parameters to generate different adversarial perturbations.  $c_k = k + c_0$ , where  $c_0 \in \mathbb{R}$  is a constant. Moreover, to ensure adversarial perturbations generated with different values of  $c$  change smoothly, we use the perturbation generated with  $c_1$  to initialize the perturbation for  $c_k$ , i.e.  $\delta_{\text{init}}^{(c_k)} = \delta^{(c_{k-1})}$ , where  $c_1 \in \mathbb{R}$  is a constant. In our experiments, we set  $c_0 = 0.6$ . For fair comparisons, we need to ensure adversarial perturbations generated with different hyper-parameters be comparable with each other. Thus, we select a constant and let  $k_2 = \frac{c_1}{c_0}$  as the stopping criteria of all adversarial attacks. We set the number of steps as 1000. The threshold is set to ensure that attacks with different hyper-parameters are almost converged when the norm of the perturbation  $\|\delta\|_2$  reaches  $k_2$ . Note that different attacking methods may successfully attack different sets of testing samples, so we select testing samples that can be successfully attacked by all attacking methods with different  $c_k$  values (i.e. those having reached the stopping criteria under all attacks). The interaction and the transfer utility reported in Figure 1 are measured on the selected samples for fair comparisons.

## H PROOF OF PROPOSITION 2

To simplify the problem setting, we do not consider some tricks in adversarial attacking, such as gradient normalization and the clip operation. In VR attack (Wu et al., 2018), the final perturbation generated after steps is given as follows.

$$\delta_{\text{vr}}^t \stackrel{\text{def}}{=} \sum_{t^0=0}^{t-1} r_x \nabla_x \hat{h}(x + \delta_{\text{vr}}^{t^0}; y);$$

where

$$\hat{h}(x; y) = \mathbb{E}_{N(0, \sigma^2 I)} [\ell(h(x + \delta); y)]; \quad (27)$$

According to Equation (27), the gradient and the Hessian matrix  $\nabla_x \hat{h}(x; y)$  is given as follows.

$$\begin{aligned} \hat{g}(x) &= r_x \nabla_x \hat{h}(x; y) \\ &= \mathbb{E}_{N(0, \sigma^2 I)} [r_x \nabla_x \ell(h(x + \delta); y)]; \\ \hat{H}(x) &= r_x^2 \nabla_x^2 \hat{h}(x; y) \\ &= \mathbb{E}_{N(0, \sigma^2 I)} [r_x^2 \nabla_x^2 \ell(h(x + \delta); y)]; \end{aligned} \quad (28)$$

where  $\sigma$  represents the step size.

Lemma 3. Given the Gaussian smoothed  $\hat{h}(x) = \mathbb{E}_{N(0, \sigma^2 I)} [\ell(h(x); y)]$ , where  $\ell(h(x); y)$  is the original classification loss,  $\delta_a \in [b, c]$ ,  $\delta_b \in [a, c]$ , we have



$$E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b = 0, \quad \text{and } E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbb{E}_{\mathbf{g}_c \sim \mathcal{N}(\mathbf{0}; \sigma_c^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{H}_{cb} \mathbf{g}_c = 0.$$

Proof. According to Equation (28), we have

$$\mathbf{g}_a(x) = E_{\mathbf{N}(\mathbf{0}; \sigma_a^2)} [\mathbf{g}_a(x + \mathbf{z})] = E_{x \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} [\mathbf{g}_a(x^0)];$$

$$\mathbf{H}_{ab}(x) = E_{\mathbf{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \frac{\partial \mathbf{g}(x + \mathbf{z})}{\partial \mathbf{x}} = E_{x \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \frac{\partial \mathbf{g}(x^0)}{\partial \mathbf{x}} = E_{x \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} [\mathbf{H}_{ab}(x^0)];$$

This indicates that the gradient and the Hessian matrix in the VR attack are both smoothed by the Gaussian noise. Because the Lipschitz constants of  $\mathbf{g}_a(x)$  and  $\mathbf{H}_{ab}(x)$  are usually limited to a certain range, we can ignore the tiny probability of large gradients and large elements in the Hessian matrix, and roughly assume that  $\mathbf{g}_a(x) \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)$ , and  $\mathbf{H}_{ab}(x) \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)$ , where  $\sigma_a, \sigma_{H_{ab}} \in \mathbb{R}$  are two constants denoting the standard deviation. Thus,  $\mathbf{g}_a(x)$  and  $\mathbf{H}_{ab}(x)$  can be written as follows.

$$\mathbf{g}_a(x) = \hat{\mathbf{g}}_a(x) + \mathbf{g}_a; \quad \mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2);$$

$$\mathbf{H}_{ab}(x) = \hat{\mathbf{H}}_{ab}(x) + \mathbf{H}_{ab}; \quad \mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2); \quad (29)$$

To simplify the notation without causing ambiguity, we write  $\mathbf{g}(x)$  and  $\mathbf{H}(x)$  as  $\hat{\mathbf{g}}$  and  $\hat{\mathbf{H}}$ , respectively. Moreover, we write  $\mathbf{g}(x)$  and  $\mathbf{H}(x)$  as  $\mathbf{g}$  and  $\mathbf{H}$ , respectively. In this way, we have

$$E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} (\mathbf{g}_a + \mathbf{g}_b)^T (\hat{\mathbf{H}}_{ab} + \mathbf{H}_{ab})$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b + \mathbf{g}_a^T \hat{\mathbf{H}}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \hat{\mathbf{H}}_{ab} \mathbf{g}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b + \mathbf{g}_a^T \hat{\mathbf{H}}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \hat{\mathbf{H}}_{ab} \mathbf{g}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b$$

$$= E_x \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} = \mathbf{0};$$

According to Equation (29), we have  $\mathbf{g}_a = \hat{\mathbf{g}}_a + \mathbf{g}_a, \mathbf{g}_b = \hat{\mathbf{g}}_b + \mathbf{g}_b$ . Thus, we have

$$E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_b$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} (\hat{\mathbf{g}}_a + \mathbf{g}_a) (\hat{\mathbf{g}}_b + \mathbf{g}_b)^T (\hat{\mathbf{H}}_{ab} + \mathbf{H}_{ab})$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \hat{\mathbf{g}}_a \hat{\mathbf{g}}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \hat{\mathbf{g}}_a \mathbf{g}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a \hat{\mathbf{g}}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a \mathbf{g}_b$$

$$+ \mathbf{H}_{ab} \hat{\mathbf{g}}_a \hat{\mathbf{g}}_b + \mathbf{H}_{ab} \hat{\mathbf{g}}_a \mathbf{g}_b + \mathbf{H}_{ab} \mathbf{g}_a \hat{\mathbf{g}}_b + \mathbf{H}_{ab} \mathbf{g}_a \mathbf{g}_b$$

$$= E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \hat{\mathbf{g}}_a \hat{\mathbf{g}}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \hat{\mathbf{g}}_a \mathbf{g}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a \hat{\mathbf{g}}_b + \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{g}_a \mathbf{g}_b$$

$$+ \mathbf{H}_{ab} \hat{\mathbf{g}}_a \hat{\mathbf{g}}_b + \mathbf{H}_{ab} \hat{\mathbf{g}}_a \mathbf{g}_b + \mathbf{H}_{ab} \mathbf{g}_a \hat{\mathbf{g}}_b + \mathbf{H}_{ab} \mathbf{g}_a \mathbf{g}_b$$

$$= E_x \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} + \mathbf{0} = \mathbf{0};$$

Moreover, according to Equation (29), we have

$$E_x \mathbb{E}_{\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}; \sigma_a^2)} \mathbb{E}_{\mathbf{H}_{ab} \sim \mathcal{N}(\mathbf{0}; \sigma_{H_{ab}}^2)} \mathbb{E}_{\mathbf{g}_c \sim \mathcal{N}(\mathbf{0}; \sigma_c^2)} \mathbf{g}_a^T \mathbf{H}_{ab} \mathbf{H}_{cb} \mathbf{g}_c$$



where

$$R_{ab}^{vr} = \frac{3(m-1)m^2}{2} \sum_{a^{02f} 1;2;\dots;n \text{ gnf ag}} \sum_{h=1}^n \sum_{i=1}^n \underbrace{(g_a g_a^0 H_{ab} H_{a^0b} - g_a g_a^0 H_{ab} H_{a^0b})}_{V_{ab}} + \sum_{b^{02f} 1;2;\dots;n \text{ gnf bg}} \sum_{h=1}^n \sum_{i=1}^n \underbrace{(g_b g_b^0 H_{ab} H_{ab^0} - g_b g_b^0 H_{ab} H_{ab^0})}_{V_{ba}} + 2m^2 \sum_{h=1}^n \sum_{i=1}^n (g_a g_b H_{ab} - g_a g_b H_{ab}) + R_2^{vr} \binom{m}{vr} R_2^0 \binom{m}{multi} + R_2^{vr}(H) R_2^{multi}(H)$$

The expectation of  $R_{ab}^{vr}$  is give as follows.

$$E_x E_{a,b} [R_{ab}^{vr}] = \frac{3(m-1)m^2}{2} E_{a,b} [E_x [V_{ab}] + E_x [V_{ba}]] + \frac{2}{(m-1)} E_x \sum_{h=1}^n \sum_{i=1}^n (g_a g_b H_{ab} - g_a g_b H_{ab}) + E_x E_{a,b} [R_2^{vr} \binom{m}{vr} R_2^0 \binom{m}{multi} + R_2^{vr}(H) R_2^{multi}(H)] = 0$$

According to Assumption 1, we have  $R_2^{vr}(H) = 0$ , and  $R_2^{multi}(H) = 0$ . Note that the magnitude of  $\binom{m}{mi}$  and the magnitude of  $\binom{m}{multi}$  are small, then  $R_2^0 \binom{m}{vr} = 0$ , and  $R_2 \binom{m}{multi} = 0$ . According to Lemma 3, we have  $E_x \sum_{h=1}^n \sum_{i=1}^n (g_a g_b H_{ab} - g_a g_b H_{ab}) = 0$ ,  $E_x \sum_{h=1}^n \sum_{i=1}^n (g_a g_a^0 H_{ab} H_{a^0b} - g_a g_a^0 H_{ab} H_{a^0b}) = 0$ . Therefore, we get  $E_x [V_{ab}] = 0$ . In this way,  $E_x E_{a,b} [R_{ab}^{vr}] = 0$ .

Furthermore, according to Lemma 3, we have  $E_x \sum_{h=1}^n \sum_{i=1}^n g_a^2 H_{ab}^2 = 0$ .

Therefore,

$$E_x E_{a,b} [I_{ab}(\text{vr}) - I_{ab}(\text{multi})] = \frac{3(m-1)m^2}{2} E_{a,b} \sum_{h=1}^n \sum_{i=1}^n (g_a^2 H_{ab}^2 - g_a^2 H_{ab}^2) + E_x \sum_{h=1}^n \sum_{i=1}^n (g_b^2 H_{ab}^2 - g_b^2 H_{ab}^2) + E_x E_{a,b} [R_{ab}^{vr}] = \frac{3(m-1)m^2}{2} E_{a,b} \sum_{h=1}^n \sum_{i=1}^n (g_a^2 H_{ab}^2 - g_a^2 H_{ab}^2) + E_x \sum_{h=1}^n \sum_{i=1}^n (g_b^2 H_{ab}^2 - g_b^2 H_{ab}^2) + 0 = 0$$

□

## I PROOF OF PROPOSITION 3

To simplify the problem setting, we do not consider some tricks in adversarial attacking, such as gradient normalization and the clip operation. Note that the original MI Attack and the multi-step attack cannot be directly compared, since that magnitudes of the generated perturbations cannot be fairly controlled. The value of interactions is sensitive to the magnitude of perturbations. Comparing perturbations with different magnitudes is not fair. Thus, we slightly revise the MI Attack as

$$g_{mi}^t \stackrel{\text{def}}{=} g_{mi}^{t-1} + (1 - \alpha) r_x \cdot (h(x + g_{mi}^{t-1}); y); \quad (31)$$

where  $t$  denotes the step and  $\alpha = (t-1)/t$ .  $\cdot (h(x); y)$  is referred as the classification loss. To simplify the notation, we use  $g(x)$  to denote  $r_x \cdot (h(x); y)$ , i.e.  $g(x) \stackrel{\text{def}}{=} r_x \cdot (h(x); y)$ .

In MI attack, the final perturbation generated after  $t$  steps is given as follows.

$$g_{mi}^t \stackrel{\text{def}}{=} \sum_{i=0}^{t-1} g_{mi}^i$$

where  $\delta$  represents the step size.

Furthermore, we define the update of perturbation with the MI attack at each step as follows.

$$x_{mi}^t \stackrel{\text{def}}{=} g_{mi}^t. \quad (32)$$

In this way, the perturbation can be written as follows.

$$x_{mi}^t = x_{mi}^1 + x_{mi}^2 + \dots + x_{mi}^{t-1}. \quad (33)$$

Lemma 4. The update of the perturbation with the MI attack at step  $t$  defined in Equation (32) can be written as  $x_{mi}^t = I + \frac{t-1}{2}H(x) + R_1^t(H(x)) g(x) + R_1^t$ , where  $R_1^t(H(x))$  denotes terms of elements in  $H(x)$  higher than the first order, and  $R_1^t$  denotes terms with elements in  $g$  higher than the first order.

Proof. If  $t = 1$ ,  $x_{mi}^1 = g(x)$ .

Let  $8\delta t^0 < t$ ,  $x_{mi}^{t_0} = I + \frac{t_0-1}{2}H(x) + R_1^{t_0}(H(x)) g(x) + R_1^{t_0}$ .

According to Equation (31) and Equation (32), we have

$$x_{mi}^t = \frac{t-1}{t}g_{mi}^{t-1} + \frac{1}{t}g(x + x_{mi}^{t-1}) :$$

Applying the Taylor series to the term  $g(x + x_{mi}^{t-1})$ , we get

$$x_{mi}^t = \frac{t-1}{t}g_{mi}^{t-1} + \frac{1}{t}g(x) + H(x) x_{mi}^{t-1} + r_1^{t-1} ; \quad (34)$$

where  $r_1^{t-1}$  denotes terms of elements in  $g$  higher than the first order.

According to Equation (33) and Equation (34), we get

$$x_{mi}^t = \frac{t-1}{t}g_{mi}^{t-1} + \frac{1}{t}g(x) + H(x) x_{mi}^1 + x_{mi}^2 + \dots + x_{mi}^{t-1} + r_1^{t-1} :$$

Because  $8\delta t^0 < t$ ,  $x_{mi}^{t_0} = I + \frac{t_0-1}{2}H(x) + R_1^{t_0}(H(x)) g(x) + R_1^{t_0}$ , we have  $x_{mi}^{t-1} = I + \frac{t-2}{2}H(x) + R_1^{t-1}(H(x)) g(x) + R_2^{t-1}$ . According to Equation (32), we get  $x_{mi}^1 = I + \frac{t-2}{2}H(x) + R_1^{t-1}(H(x)) g(x) + R_2^{t-1}$ .

In this way, we get

$$\begin{aligned}
 x_{mi}^t &= \left( \frac{t-1}{t} I + \frac{t-2}{2} H(x) + R_1^{t-1}(H(x)) \right) g(x) \\
 &\quad + \frac{1}{t} \left( I + H(x) \right) \left( t-1 \right) + \frac{(t-2)(t-1)}{4} H(x) + \frac{X^1}{t^{0=1}} R_1^{t^0}(H(x)) g(x) + \frac{X^1}{t^{0=1}} R_1^{t^0} \\
 &\quad + r_1^{t-1} \\
 &= \frac{t-1}{2} R_1^{t-1}(H(x)) + \frac{(t-2)(t-1)}{4t} H^2(x) + \frac{1}{t} H(x) \frac{X^1}{t^{0=1}} R_1^{t^0}(H(x)) \\
 &\quad + I + \frac{t-1}{2} H(x) + \frac{Z}{5} g(x) + \frac{1}{t} \frac{X^1}{t^{0=1}} \left\{ \frac{Z}{R_1^t} \right\} + r_1^{t-1} \\
 &= I + \frac{t-1}{2} H(x) + R_1^t(H(x)) g(x) + R_1^t:
 \end{aligned}$$

where  $R_1^t(H(x))$  denotes terms of elements in  $H(x)$  higher than the first order, and  $R_1^t$  denotes terms with elements in  $t_{mi}^{t-1}$  of higher than the first order. In this way, we have proved that 1;  $x_{mi}^t = I + \frac{t-1}{2} H(x) + R_1^t(H(x)) g(x) + R_1^t$ .  $\square$

**Proposition 3.** The adversarial perturbation generated by multi-step attack is denoted by  $\sum_{t=0}^{m-1} r_{x^t}(h(x + t_{multi}); y)$ . The adversarial perturbation generated by multi-step attack incorporating the momentum is computed as  $\sum_{t=0}^{m-1} g_{mi}^t$ . Perturbation units of  $m_{mi}$  exhibit smaller interactions than  $m_{multi}$ , i.e.  $E_{ij} [I_{ij}(\frac{m}{mi})] < E_{ij} [I_{ij}(\frac{m}{multi})]$ .

**Proof.** According to Lemma 4, the update of the perturbation with the MI attack at the step  $t$  is given as follows.

$$x_{mi}^t = I + \frac{t-1}{2} H(x) + R_1^t(H(x)) g(x) + R_1^t: \quad (35)$$

where  $R_1^t(H(x))$  denotes terms of elements in  $H(x)$  of higher than the first order, and  $R_1^t$  denotes terms with elements in  $t_{mi}^{t-1}$  of higher than the first order.

To simplify the notation without causing ambiguity, we write  $g(x)$  and  $H(x)$  as  $g$  and  $H$ , respectively. In this way, according to Equation (33) and Equation (35), (35) can be written as follows.

$$x_{mi}^m = I + \frac{m(m-1)}{4} H + \sum_{t=1}^m R_1^t(H) g + \sum_{t=1}^m R_1^t: \quad (36)$$

where  $m$  represents the total number of steps. According to Lemma 1, the Shapley interaction between perturbation units in  $m_{mi}$  is given as follows.

$$I_{ab}(\frac{m}{mi}) = m_{mi;a} H_{ab} m_{mi;b} + R_2(\frac{m}{mi}); \quad (37)$$

where  $R_2(\frac{m}{mi})$  denotes terms with elements in  $m_{mi}$  of higher than the second order.

According to Equation (36) and Equation (37), we get

$$\begin{aligned}
I_{ab}(\frac{m}{m_i}) &= H_{ab} [ mg_a + \frac{2m(m-1)}{4} \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) + \sum_{t=1} X^n \underbrace{O(\frac{t}{m_i;a})}_{\substack{\text{terms of } t_{mi;a} \\ \text{of higher than the } \text{rst order,} \\ \text{which corresponds to the term of} \\ \mathbb{R}_1^t \text{ in Equation (36)}}} ] \\
&+ mg_b + \frac{2m(m-1)}{4} \sum_{a^0=1} X^n (H_{a^0b} g_{a^0}) + \sum_{t=1} X^n \underbrace{O(\frac{t}{m_i;b})}_{\substack{\text{terms of } t_{mi;b} \\ \text{of higher than the } \text{rst order,} \\ \text{which corresponds to the term of} \\ \mathbb{R}_1^t \text{ in Equation (36)}}} ] + \mathbb{R}_2(\frac{m}{m_i}) \\
&= \underbrace{\frac{2m^2 g_a g_b H_{ab}}{4}}_{\substack{\text{first-order terms w.r.t. elements in } H \\ \#}} + \underbrace{\frac{3(m-1)m^2}{4} g_b \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) + \frac{3(m-1)m^2}{4} g_a \sum_{a^0=1} X^n (H_{a^0b} g_{a^0})}_{\substack{\text{second-order terms w.r.t. elements in } H \\ \#}} H_{ab} \\
&+ \underbrace{\frac{4(m-1)^2 m^2}{16} \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) \sum_{a^0=1} X^n (H_{a^0b} g_{a^0})}_{\mathbb{R}_2^{mi}(H)} + \dots \\
&+ \underbrace{[O(\frac{t}{m_i;a})] H_{ab} \frac{m}{m_i;b} + [O(\frac{t}{m_i;b})] H_{ab} \frac{m}{m_i;a}}_{\mathbb{R}_2^0(\frac{m}{m_i})} + \mathbb{R}_2(\frac{m}{m_i}) \\
&= \frac{2m^2 g_a g_b H_{ab}}{4} + \frac{3(m-1)m^2}{4} g_a \sum_{a^0=1} X^n (H_{a^0b} g_{a^0}) \\
&+ \frac{3(m-1)m^2}{4} g_b \sum_{b^0=1} X^n (H_{ab^0} g_{b^0}) + \mathbb{R}_2^0(\frac{m}{m_i}) + \mathbb{R}_2^{mi}(H); \tag{38}
\end{aligned}$$

where  $\mathbb{R}_2^{mi}(H)$  denotes terms of elements in higher than the second order, and  $\mathbb{R}_2^0(\frac{m}{m_i})$  denotes terms of elements in  $\frac{m}{m_i}$  higher than the second order

According to Equation (22) and Equation (38), the expectation of the difference between  $I_{ab}(\frac{m}{m_i})$  and  $I_{ab}(\frac{m}{m_{\text{multi}}})$  is given as follows.

$$\begin{aligned}
&E_{a,b} [I_{ab}(\frac{m}{m_i}) - I_{ab}(\frac{m}{m_{\text{multi}}})] \\
&= \frac{3(m-1)m^2}{4} E_{a,b} \left[ \underbrace{g_a H_{ab} \sum_{a^0=1} X^n (H_{a^0b} g_{a^0})}_{U_{ab}} + \underbrace{g_b H_{ab} \sum_{b^0=1} X^n (H_{ab^0} g_{b^0})}_{U_{ba}} \right] + E_{a,b} \mathbb{R}_{ab}^{mi};
\end{aligned}$$

where

$$\mathbb{R}_{ab}^{mi} = \mathbb{R}_2^0(\frac{m}{m_i}) - \mathbb{R}_2^0(\frac{m}{m_{\text{multi}}}) + \mathbb{R}_2^{mi}(H) - \mathbb{R}_2^{mi}(H);$$

According to Assumption 1, we have  $\mathbb{R}_2^{mi}(H) \rightarrow 0$ , and  $\mathbb{R}_2^{multi}(H) \rightarrow 0$ . Note that the magnitude of  $\frac{m}{m_i}$  and the magnitude of  $\frac{m}{m_{\text{multi}}}$  are small, then  $\mathbb{R}_2^0(\frac{m}{m_i}) \rightarrow 0$ , and  $\mathbb{R}_2^0(\frac{m}{m_{\text{multi}}}) \rightarrow 0$ . Therefore,  $E_{a,b} \mathbb{R}_{ab}^{mi} = E_{a,b} [\mathbb{R}_2^0(\frac{m}{m_i}) - \mathbb{R}_2^0(\frac{m}{m_{\text{multi}}}) + \mathbb{R}_2^{mi}(H) - \mathbb{R}_2^{multi}(H)] \rightarrow 0$ . Moreover, similar to Equation (24) in the proof of Proposition 1, we have  $E_{a,b} [U_{ab}] = E_{a,b} [U_{ba}] \rightarrow 0$ .

Therefore,

$$\begin{aligned} & E_{a;b} [I_{ab}(m_i) - I_{ab}(m_{\text{multi}})] \\ &= \frac{3(m-1)m^2}{4} E_{a;b} [U_{ab} + U_{ba}] + E_{a;b} R_{ab}^{mi} \\ & \quad - \frac{3(m-1)m^2}{2} E_{a;b} [U_{ab}] + 0 = 0. \end{aligned}$$

□

Note that Proposition 3 just shows the revised MI Attack usually decreases the interaction between perturbation units. The proof towards all types of MI Attacks is still a challenge.

## J IMPLEMENTATION OF THE INTERACTION-REDUCED ATTACK (IR ATTACK)

### J.1 GRID-LEVEL INTERACTIONS FOR IMAGE DATA

Although the computation of  $E_{i;j} [I_{ij}(\cdot)]$  can be simplified using Equation (26), the computational cost of  $E_{i;j} [I_{ij}(\cdot)]$  is still high. Therefore, as Figure 6 shows, using the local property of images (Chen et al., 2018a), we can divide the entire image into  $L$  grids, and compute interactions at the grid level, instead of the pixel level. Let  $\mathcal{G} = \{g_1, g_2, \dots, g_L\}$  denote the set of grids. We use  $(p; q)$  to denote the coordinate of a grid. In this way, the expectation of interactions between perturbation grids is given as follows.

$$E_{(p;q);(p^0;q^0)} [I_{(p;q);(p^0;q^0)}(\cdot)] = \frac{1}{L^2 - 1} E_{(p;q)} [v(\cdot) - v(\cdot - \mathbf{f}_{pqg}) - v(\cdot + \mathbf{f}_{pqg}) + v(\cdot)]; \quad (39)$$

Figure 6: For the input image, we can divide the image into grids, and compute interactions at the grid level.

### J.2 SCALABILITY OF THE INTERACTION LOSS

In this section, we discuss about two kinds of scalability of the interaction loss.

- Is the computational cost of the interaction loss affordable when the number of players is large? To this end, we have proved in Equation (4) that the computational complexity of the expectation of the interaction is linear, which is scalable. In fact, we do not directly compute interaction using Equation (3). Instead, we compute the expectation of interactions with Equation (4). The computational cost of the IR Attack can be further reduced by calculating the grid-level interactions of images.

We further conducted experiments to measure the time cost of generating perturbations using the IR Attack. We conducted the IR attack for 100 steps on the ImageNet dataset. The time cost was measured using PyTorch 1.6 (Paszke et al., 2019) on Ubuntu 18.04, with the Intel(R) Core(TM) i7-9800X CPU @ 3.80GHz and a Titan RTX GPU. Table 6 shows the average computational cost of generating adversarial perturbations on an input image with size 224 by the IR Attack for 100 steps. It shows that the IR Attack is computationally applicable to high-dimensional data and deep neural networks.

Table 6: Average computational cost of generating adversarial perturbations over an input image by the IR Attack for 100 steps on different source DNNs.

	RN-34	RN-152	DN-121	DN-201
Time (seconds)	12.882	48.774	27.519	44.812

• Is the computation cost of the interaction loss affordable when we consider the continuous space of adversarial perturbations? It has been widely discussed (Ancona et al., 2019; Sundararajan & Najmi, 2019) that when applying the Shapley value, the feature space is regarded as binary. It is because as (Sundararajan & Najmi, 2019) shows that although there exist the Shapley-value-like attribution in a continuous space, only the Shapley value in the binary space is the unique attribution that satisfies the linearity axiom, the dummy axiom, the symmetry axiom, and the efficiency axiom that only in the binary space. Thus, when we compute the interaction, the perturbation can be regarded in the binary space, i.e., whether the perturbation unit is added to the input or not, which enables scalability.

## K EVALUATION OF THE TRANSFERABILITY VIA LEAVE-ONE-OUT VALIDATION

As Figure 7 shows, the highest transferability of the MI Attack is achieved in an intermediate step, rather than in the last step. This phenomenon presents a challenge for fair comparisons of the transferability between different attacking methods.

To this end, in order to enable fair comparisons of transferability between different methods, we estimate the adversarial perturbations with the highest transferability for each input image via the leave-one-out (LOO) validation as follows. Given a set of clean examples  $(x_i, y_i)_{i=1}^N$ , where  $y_i \in \{1, 2, \dots, C\}$ , we use  $x_i^t$  to denote the adversarial example at step  $t$  w.r.t. the clean example  $x_i$ , where  $t \in \{1, 2, \dots, T\}$ , and  $T$  is the number of total step. Given a target DNN  $h(\cdot)$  and an input example  $x$ , where  $h(\cdot)$  denotes the output before the softmax layer, we use  $C(x) = \arg \max_k h_k(x)$ ;  $k \in \{1, 2, \dots, C\}$  to denote the prediction of the example  $x$ .

$$x_i^t \stackrel{\text{def}}{=} x_i^{t_i}; \text{ s.t. } t_i = \arg \max_t E_{i \in \{1, 2, \dots, N\}} \mathbb{1}[C(x_i^t) \neq y_i];$$

where  $\mathbb{1}[\cdot]$  is the indicator function. Then the average transferability is given as follows.

$$\text{Transferability} = E_i [\mathbb{1}[C(x_i^t) \neq y_i]];$$

Figure 7: The curve of transferability in different steps.

## L ADDITIONAL RELATED WORK

Some studies paid attention to intermediate features to improve transferability. Activation Attack (Inkawhich et al., 2019) forced the intermediate features of the input image to be similar with the intermediate features of a target image, in order to generate highly transferable targeted example. Distribution Attack (Inkawhich et al., 2020) explicitly modeled the feature distribution of each class, and improve the targeted transferability by driving the feature of perturbed input image into the distribution of a specific target class. Intermediate Level Attack (Huang et al., 2019) improved the



transferability of an adversarial example by maximizing the feature perturbation of a pre-specified layer. In comparison, we explain and improve the transferability based on game theory. Moreover, we discover the negative correlation between the transferability and interactions.

## M ADDITIONAL EXPERIMENTS ON INTERACTION-REDUCED LOSS

### M.1 INTERACTION REDUCTION ON OTHER ATTACKS

To further demonstrate the effectiveness of the interaction loss, we have applied the interaction loss on other attacks besides the PGD Attack, including the MI Attack, the SGM Attack, and the VR Attack. More specifically, we added the interaction loss on the MI Attack (namely the MI+IR Attack), the SGM Attack (namely the SGM+IR Attack), and the VR Attack (namely the VR+IR Attack), respectively.

For the MI Attack and the SGM Attack, we directly applied Equation (7) to these attacks, because these attacks were compatible with the interaction loss. Besides, for the VR Attack, its objective function is given as follows.

$$\text{maximize } E_{N(0; \Sigma)} [\langle h(x + \delta); y \rangle] \quad \text{s.t. } \|\delta\|_p \leq \epsilon; x + \delta \in [0, 1]^n; \quad (40)$$

Therefore, the VR+IR Attack was implemented via sampling as follows.

$$\begin{aligned} \text{maximize } & \frac{1}{K} \sum_{k=1}^K [\langle h(x + \delta_k); y \rangle - E_{ij} [\langle I_{ij}(\delta) \rangle]]; \quad \|\delta_k\|_p \leq \epsilon; \delta_k \in [0, 1]^n; \\ \text{s.t. } & \|\delta_k\|_p \leq \epsilon; x + \delta_k \in [0, 1]^n; \end{aligned} \quad (41)$$

where the interaction loss was computed by considering the input image as  $x$ , rather than  $x + \delta$  in Equation (26). The VR Attack reported in Table 4 followed the original paper (Wu et al., 2018) to set  $K = 20$ . However, a crucial issue for applying the interaction loss to the VR attack was its extremely high computational cost. Therefore, for the implementation of the VR+IR Attack, we set  $K = 5$  and reduce the number of steps from 100 to 50. Just like experiments in Table 1, we also used the LOO strategy for evaluation.

Table 7, Table 8, and Table 9 compare the success rates of attacks with and without the interaction loss. The results demonstrated that the performance of the MI Attack, the SGM Attack, and the VR Attack can be further enhanced by directly adding the interaction loss to reduce interactions inside perturbations.

Table 7: The success rates of black-box attacks crafted by MI and MI+IR on four source models (RN-34/152, DN-121/201) against seven target models. The interaction loss can boost the transferability of MI.

Source	Method	VGG-16	RN152	DN-201	SE-154	IncV3	IncV4	IncResV2							
RN-34	MI	80.1	0.5	73.0	2.3	77.7	0.5	48.9	0.8	46.2	1.2	39.9	0.5	34.8	2.5
	MI+IR	90.0	0.5	85.7	0.3	88.5	0.6	67.0	0.1	66.9	1.8	60.2	0.7	53.9	2.3
RN-152	MI	70.3	0.6	–	–	74.8	1.4	51.7	0.8	47.1	0.9	40.5	1.6	36.8	2.7
	MI+IR	78.9	1.4	–	–	82.2	2.0	68.3	0.3	63.6	1.2	59.0	0.4	56.3	1.0
DN-121	MI	83.0	4.9	72.0	0.7	91.5	0.2	58.4	2.6	54.6	1.6	49.2	2.4	43.9	1.5
	MI+IR	89.0	0.8	83.2	1.5	93.4	0.6	74.2	0.7	69.6	0.9	64.7	0.5	58.2	2.3
DN-201	MI	77.3	0.8	74.8	1.4	–	–	64.6	1.0	56.5	2.5	51.1	2.1	47.8	1.9
	MI+IR	87.3	0.3	81.6	2.0	–	–	75.4	0.6	66.6	3.3	60.0	1.0	62.1	0.7

### M.2 ATTACKS ON THE CIFAR-10 DATASET

In Table 1, we have shown that reduction of interactions could improve the adversarial transferability on the ImageNet dataset (Russakovsky et al., 2015). To further demonstrate the broad applicability of such a negative correlation, we also conducted the targeted attack on the CIFAR-10 dataset (Krizhevsky & Hinton, 2009) to test the transferability of perturbations generated with the interaction loss.

Table 8: The success rates of black-box attacks crafted by SGM and SGM+IR on four source models (RN-34/152, DN-121/201) against seven target models. The interaction loss can boost the transferability of SGM.

Source	Method	VGG-16	RN152	DN-201	SE-154	IncV3	IncV4	IncResV2
RN-34	SGM	91.8 0.6	89.0 0.9	90.0 0.4	68.0 1.4	63.9 0.3	58.2 1.1	54.6 1.2
	SGM+IR	94.7 0.6	91.7 0.6	93.4 0.8	72.7 0.4	68.9 0.9	64.1 1.3	61.3 1.0
RN-152	SGM	88.2 0.5	–	90.2 0.3	72.7 1.4	63.2 0.7	59.1 1.5	58.1 1.2
	SGM+IR	92.0 1.0	–	92.5 0.4	79.3 0.1	69.6 0.8	66.2 1.0	63.6 0.9
DN-121	SGM	88.7 0.9	88.1 1.0	98.0 0.4	78.0 0.9	64.7 2.5	65.4 2.3	59.7 1.7
	SGM+IR	91.7 0.2	90.4 0.4	94.3 0.1	87.0 0.4	78.8 1.3	79.5 0.2	75.8 2.7
DN-201	SGM	87.3 0.3	92.4 1.0	–	82.9 0.2	72.3 0.3	71.3 0.6	68.8 0.5
	SGM+IR	89.5 0.9	91.8 0.7	–	87.3 1.2	82.5 0.8	80.3 0.3	81.5 0.5

Table 9: The success rates of black-box attacks crafted by VR and VR+IR on four source models (RN-34/152, DN-121/201) against seven target models. The interaction loss can boost the transferability of VR.

Source	Method	VGG-16	RN152	DN-201	SE-154	IncV3	IncV4	IncResV2
RN-34	VR	85.1	85.3	87.0	55.7	54.3	50.7	43.7
	VR+IR	90.8	92.2	93.3	75.4	75.4	67.5	66.1
DN-121	VR	88.8	88.4	98.2	72.9	73.5	72.5	63.6
	VR+IR	93.0	93.5	96.2	83.7	82.8	84.0	79.8

Following Wu et al. (2018), we chose three DNNs as the source DNN or the target DNN, which included: LeNet (LeCun et al., 1998), RN-20 (He et al., 2016), and DN-121 (Huang et al., 2017). We conducted the targeted attack under the norm constraint, and chose plane class as the target category. The norm constraint was set to 16/255, and the step size was set to 2/255. The transferability was computed based on the best adversarial perturbation during 50 steps via the leave-one-out (LOO) validation., which has been introduced in Appendix K.

As Table 10 shows, the transferability could be enhanced by reducing interactions on the targeted attack on the CIFAR-10 dataset. Particularly, when the source DNN is RN-20 and the target DNN is DN-121, the transferability improvement was about 30%, which was a considerable gain.

Table 10: The success rates of targeted black-box attacks on three source models, including LeNet, RN-20, DN-121, against three target models.

Source	Method	LeNet	RN-20	DN-121
LeNet	PGD	–	34.1 0.1	19.6 0.4
	PGD+IR	–	44.2 0.3	29.7 1.1
RN-20	PGD	10.8 0.8	–	41.9 1.3
	PGD+IR	19.7 0.3	–	71.8 1.0
DN-121	PGD	10.0 0.7	44.2 0.3	–
	PGD+IR	18.9 0.7	58.5 1.0	–

Table 11: The average interaction inside adversarial perturbations generated by PGD, DI and TI.

Method	RN-34	DN-121
Baseline (PGD Attack)	0.422	0.926
DI Attack	0.241	0.499
TI Attack	0.379	0.618

## N EMPIRICAL VERIFICATION OF OTHER TRANSFERABILITY-BOOSTING ATTACKS

We have theoretically analyzed the MI Attack, the VR Attack, and the SGM Attack. However, for other methods of improving adversarial transferability, such as Diversity Input (DI) (Xie et al.,

2019), which uses random data augmentation during attacking, it is difficult to mathematically prove that they essentially reduce interactions. Nevertheless, as Table 11 shows, we empirically demonstrated that two widely-used transferability-boosting attacks, DI and TI (Dong et al., 2019), also reduced interactions.

## O ADDITIONAL EXPERIMENTS ON EFFECTS OF THE INTERACTION LOSS

We conducted additional experiments to test the effects of the interaction loss. We conducted attacks on two source DNNs (RN-34, DN-121), and transferred adversarial perturbations to seven target DNNs (VGG16, RN-152, DN-201, SE-154, IncV3, IncV4, IncResV2).

We used the following two experimental settings to compare the transferability of adversarial perturbations generated with different values.

First, we re-drew the curves in Figure 3(a) by extending them from the range of  $[0; 1:2]$  to the range of  $[0; 2:0]$ , in order to show the performance of different values. We simply changed the value in the objective function (i.e. Equation (5)) without any other revisions. This was the most direct way to test the effects of  $\epsilon$ . Experimental results are shown in Figure 8.

Besides above experimental settings, we also compared adversarial perturbations generated with different  $\epsilon$  values, when we controlled each perturbation to have the same attacking utility. The attacking utility was defined as follows.

$$\text{Attacking Utility} = \max_{y \neq y'} h_{y'}(x + \delta) - h_y(x + \delta)$$

where  $y$  denote the label of the input image. This setting also ensured the fairness of comparisons from a new perspective. Please see Figure 9 for experimental results.

In sum, under both experimental settings, we found that the large  $\epsilon$  value usually yielded a high adversarial transferability in our experiments.

Figure 8: The success rates of black-box attacks with the IR Attack using different values under the first experimental setting.

Figure 9: The success rates of black-box attacks with the IR Attack using different values under the second experimental setting.