

Double Targeted Universal Adversarial Perturbations

Philipp Benz*, Chaoning Zhang*, Tooba Imtiaz, and In So Kweon

Korea Advanced Institute of Science and Technology (KAIST)

pbenz@kaist.ac.kr, chaoningzhang1990@gmail.com

* Equal contribution

Abstract. Despite their impressive performance, deep neural networks (DNNs) are widely known to be vulnerable to adversarial attacks, which makes it challenging for them to be deployed in security-sensitive applications, such as autonomous driving. Image-dependent perturbations can fool a network for one specific image, while universal adversarial perturbations are capable of fooling a network for samples from all classes without selection. We introduce a double targeted universal adversarial perturbations (DT-UAPs) to bridge the gap between the instance-discriminative image-dependent perturbations and the generic universal perturbations. This universal perturbation attacks one targeted source class to sink class, while having a limited adversarial effect on other non-targeted source classes, for avoiding raising suspicions. Targeting the source and sink class simultaneously, we term it double targeted attack (DTA). This provides an attacker with the freedom to perform precise attacks on a DNN model while raising little suspicion. We show the effectiveness of the proposed DTA algorithm on a wide range of datasets and also demonstrate its potential as a physical attack.¹

1 Introduction

Despite the recent success of deep learning [1,2,3,4,5], deep neural networks (DNNs) remain vulnerable to adversarial attacks [6,7,8,9,10,11]. This poses a threat for deploying DNNs in security-sensitive applications, such as autonomous driving and robotics. Various attack methods [12] have been proposed in the past few years, which can be roughly divided into two main categories: image-dependent attacks [6,7,13,14,15] and universal attacks [16,17,18,19,20]. Image-dependent attacks construct perturbations tailored for a specific input image to be misclassified by the network; while universal attack methods aim to generate one single universal adversarial perturbation (UAP) that can fool the network for most samples of all classes.

Bridging the gap between the discriminative nature of image-dependent perturbations and the non-discriminative universal perturbation, we propose to attack a certain source class while limiting the influence of the attack on other,

¹ Code: <https://github.com/phibenz/double-targeted-uap.pytorch>

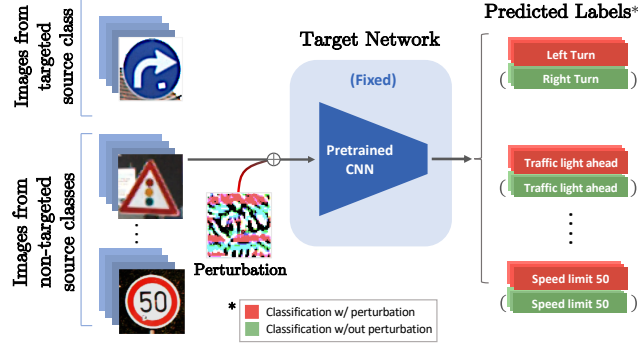


Fig. 1: Overview of the Double Targeted Attack (DTA). In this example, the perturbation causes the network to classify images of the targeted source class **right turn** as the sink class **left turn**. Image classifications from the non-targeted source classes remain unaltered. The DT-UAP is added to all image samples.

non-targeted classes. More specifically, we aim to fool the network with a single perturbation that can systematically shift a certain source class to a different sink class of choice. Since the proposed attack targets both the source and the sink class, we name it double targeted attack (DTA). To avoid confusion, while other works [15,20] use the term “target class”, we adopt “sink class” instead, since the proposed DTA also has target class(es) on the source side.

In this work, we focus on the exploration of universal perturbations due to their merit of being image-agnostic. This property eases the attack procedure for real-time applications such as autonomous driving or robotics, as the perturbation can be constructed in advance, and applying the prepared perturbation only requires one summation [16]. UAPs attack all classes, making it obvious to an observer that a system is under attack. For achieving a more covert universal attack, class-discriminative universal adversarial perturbation (CD-UAP) has been introduced in [8] to attack chosen class(es) on the source side. It would be more challenging yet meaningful to not only being class-discriminative on the source side, but also targets on the sink side. Compared with existing UAP attacks, DTA can be more dangerous in practice, since it allows precise attacks with flexible control over the targeted source class and the sink class. Applying double targeted universal adversarial perturbations (DT-UAP) can have fatal implications in practice. For instance, in the context of autonomous driving, an attacker can intentionally craft a perturbation to fool a network to misclassify traffic signs from “turn left” to “turn right” as shown in Fig. 1.

Technically, the proposed DTA does not strictly fall into the group of universal attacks, since it does not attack all classes. However, the DTA crafts one single perturbation that can be applied to the entire data distribution, which is similar to the existing UAPs [16]. It is a non-trivial task to craft the DT-UAP because there is an inherent conflict between two objectives. For the samples

from the targeted source class, the goal of the crafted perturbation is to shift their classification output to the sink class. This will inevitably have a similar influence on the non-targeted source classes, which conflicts with the goal of the attack being discriminative between the targeted source class and other non-targeted source classes. Inspired by [8], we have designed an algorithm that explicitly deals with the trade-off between them.

To demonstrate its effectiveness, we evaluate the proposed DTA on five classification datasets from different domains for various DNN architectures.

Our results establish the existence of DT-UAPs to attack data samples discriminatively. Though the designed DTA algorithm is mainly for perturbing samples to be misclassified from one targeted source class into one sink class, it can also be extended for shifting multiple targeted source classes to one sink class. We validate this specific attack scenario on the ImageNet dataset. Overall, our proposed algorithm has been validated to be effective to achieve discriminative targeted attacks with extensive experiments on different datasets and scenarios. Finally, we also demonstrate the potential of DTA being applied as a physical attack.

2 Related work

2.1 Image-Dependent Attacks

Adversarial attacks, which craft one perturbation specifically for one input image to fool a network are called image-dependent attacks. Szegedy *et al.* optimized such perturbations by using box-constrained L-BFGS [6]. Goodfellow *et al.* then introduced the Fast Gradient Sign Method (FGSM), an efficient one-step attack to generate adversarial examples [7]. The iterative variant of FGSM (I-FGSM) updates the perturbation by only a fraction of the allowed upper bound in each iteration [13]. Integrating the momentum term into the iterative process of I-FGSM (MI-FGSM) further improved the success rate of adversarial attacks [21]. DeepFool [14] is also an iterative attack, manipulating the models' decision boundaries in the perturbation crafting process. Incorporating the minimization of the perturbation magnitude into the optimization function, Carlini and Wagner (C&W) introduced another three variants of image-dependent attacks [15]. Another effective multi-step attack variant was introduced by Madry *et al.* using projected gradient descent (PGD) to craft adversaries [22]. The proposed DTA differentiates itself by attacking an entire class instead of only a single image.

2.2 Universal Attacks

A universal adversarial perturbation (UAP) is a single perturbation, which enables fooling a network for most input samples. Accumulating image-dependent perturbations by iteratively applying DeepFool [14], Moosavi *et al.* crafted the first UAPs [16]. In another variant, UAPs are crafted by leveraging the Jacobian matrices of the networks' hidden layers [17]. Assuming no access to the

original training data, Fast Feature Fool proposed to generate data-free UAPs by optimizing the feature change caused by the applied UAP [18]. Generative Adversarial Perturbations (GAP) were proposed by Poursaeed *et al.* [20], using generative models to craft image-dependent and universal perturbations. Data-free targeted UAP has been introduced in [9], showing UAP have dominant features over images. The almost absent computational overhead (single summation) in the deployment of UAPs, makes them a favorable choice for the attack of real-world applications. Despite being universal, our proposed DTA differentiates itself from the existing universal attacks in its class-discriminative nature, i.e. by having a different influence on a sample depending on whether or not it belongs to the targeted source class. CD-UAP has been introduced in [8], our DT-UAP also targets on the sink side and thus constitutes a more challenging task. Moreover, we show that our DTA can also be used in physical attack [23,24].

2.3 Attack on autonomous driving and robotics

Deep learning has achieved the maturity to be deployed in safety and security-critical applications, such as autonomous driving [25] and robotics [26]. The threat of adversarial attacks in these applications has also been widely explored. For example, Melis *et al.* [27] demonstrated the vulnerability of robots to the adversarially manipulated input images with the techniques in [6], and argue that secure robotics need to adopt strategies to enforce DNNs to learn more robust representations. Attack on the learning policy of robotics has been explored in [28]. Considering adversarial attacks in the context of autonomous driving, [29] generates UAPs to attack road sign classifiers. Another work [30] performs an attack in autonomous driving with traffic signs. Besides the classical classification dataset to evaluate the adversarial attack method, we also evaluate the proposed method on a traffic sign dataset and another robotics-related dataset.

3 Double Targeted Attack

3.1 Problem Formulation

The purpose of the proposed attack is to craft a single perturbation to shift one targeted source class to a different sink class. The source class to be attacked as well as the sink class are determined by the attacker to realize a flexible and precise attack. We term it double targeted attack (DTA).

Let $x \sim X$ denote a single sample from a distribution in \mathbb{R}^d , and $\hat{F}(x) = p$ being a classification function, mapping input $x \in \mathbb{R}^d$ to a predicted class $p \in [1, C]$ for a classification problem of C classes. Here the classification function is represented through a DNN parameterized by the weights θ . For most samples from the targeted source class $x_t \sim X_t$, we seek a perturbation δ that satisfies the constraint

$$\hat{F}(x_t + \delta) = y_{\text{sink}} \quad \text{subject to} \quad \|\delta\|_p \leq \epsilon, \quad (1)$$

where the sink class satisfies $y_{\text{sink}} \neq F(x_t)$, and ϵ indicates the magnitude limit for the l_p norm of the crafted perturbation δ . Note that limiting X_t in Eq. 1 to a single image results in an image-dependent targeted attack. Meanwhile, it is equivalent to a non-discriminative targeted universal attack if the targeted samples X_t comprise the entire dataset X .

Empirically, we find that a perturbation crafted under the constraint of Eq. 1 also shifts samples from the non-targeted source classes into the sink class with a high targeted fooling ratio. To incorporate covertness within the proposed attack, this effect of non-targeted samples $x_{nt} \sim X_{nt}$ shifting to the sink class should be minimized. The crafted perturbation should ideally shift instances from the chosen source class to a different sink class while having limited influence on the samples from the non-targeted source classes. More specifically, the proposed DTA has two objectives: (1) to increase the targeted fooling ratio for the samples from the chosen source class to the chosen sink class; (2) to decrease the targeted fooling ratio for samples from the non-targeted source class(es) into the sink class, where the targeted fooling ratio is defined as the ratio of samples fooled into the sink class. These two objectives contradict each other, leading to an inevitable trade-off. In the following subsection, we state the loss function for DTA and design the algorithm for explicitly handling this trade-off between the two objectives.

3.2 DTA Loss design

To achieve selectivity among the targeted source class and non-targeted source classes, we explicitly design different loss functions for the two. For the targeted class and the non-targeted classes, the loss is indicated by \mathcal{L}_t and \mathcal{L}_{nt} , respectively. The final loss \mathcal{L} can then be calculated as:

$$\mathcal{L} = \mathcal{L}_t + \alpha \mathcal{L}_{nt}, \quad (2)$$

where α is a hyper-parameter for weighting the trade-off between \mathcal{L}_t and \mathcal{L}_{nt} . In practice, this hyper-parameter can be fine-tuned by the attacker for a specific task. For simplicity, we set α to 1 in all of our experiments. We empirically found that this setting works well when the same number of samples are sampled from X_t and X_{nt} in every iteration update.

For the targeted class, the loss \mathcal{L}_t should shape the perturbation to fool the network by shifting the prediction from the source class into the sink class. This can be realized through (1) decreasing the logit value for the originally predicted class \hat{L}_p with $p = \arg \max(\hat{L}(x_t))$ to not being the highest logit anymore, while (2) increasing the logit for the sink class \hat{L}_{sink} , to be the dominant logit, where $\hat{L}(\cdot)$ indicates the function mapping to the logit values and \hat{L}_i is the specific logit value of class i . Thus, \mathcal{L}_t can be decomposed into two parts as follows:

$$\mathcal{L}_t = \mathcal{L}_{t1} + \mathcal{L}_{t2}, \text{ with} \quad (3)$$

$$\mathcal{L}_{t1} = \max(\hat{L}_p(x_t + \delta) - \max_{i \neq p}(\hat{L}_i(x_t + \delta)), 0) \quad (4)$$

$$\mathcal{L}_{t2} = \max(\max_{i \neq y_{\text{sink}}}(\hat{L}_i(x_t + \delta) - \hat{L}_{\text{sink}}(x_t + \delta)), -D) \quad (5)$$

Algorithm 1: Double Targeted Attack Algorithm

Input: Data distribution X , Classifier \hat{F} , Loss function \mathcal{L} , Mini-batch size m ,
Number of iterations I , Perturbation magnitude ϵ

Output: Perturbation vector δ

$X_t \subseteq X$ ▷ Subset

$X_{nt} \subseteq X$ ▷ Subset

$\delta \leftarrow 0$ ▷ Initialize

for $iteration = 1, \dots, I$ **do**

$B_t \sim X_t: |B_t| = \frac{m}{2}$ ▷ Randomly sample

$B_{nt} \sim X_{nt}: |B_{nt}| = \frac{m}{2}$ ▷ Randomly sample

$B \leftarrow B_t \cup B_{nt}$ ▷ Concatenate

$g_\delta \leftarrow \mathbb{E}_B[\nabla_\delta \mathcal{L}]$ ▷ Calculate gradient

$\delta \leftarrow \text{Optim}(g_\delta)$ ▷ Update perturbation

$\delta \leftarrow \frac{\delta}{\|\delta\|_p} \epsilon$ ▷ Projection

end

where the hyper-parameter D constitutes an intensity value of the dominance of the targeted logit value. A higher D implies a higher chance that the sample will be classified as the sink class. For the non-targeted source classes, we adopt the widely used cross-entropy function as:

$$\mathcal{L}_{nt} = \mathcal{X}(\hat{L}(x_{nt} + \delta), \mathbb{1}(\hat{F}(x_{nt}))) \quad (6)$$

with $\mathbb{1}(\cdot)$ indicating a one-hot encoded vector of C classes. In practice, an attacker can change the hyper-parameters according to the requirements. For instance, the attacker can increase the parameter α in Eq. 2 to increase the covertness of the proposed attack accompanied by a relatively low targeted fooling ratio for the targeted class, or increase the parameter D in order to achieve stronger classifications into the sink class.

To balance the two contradicting objectives, clamping of the logit values was adopted in \mathcal{L}_t . Without this clamping operation, the loss part of the targeted classes \mathcal{L}_t can prevail by shifting the samples from the targeted source class to the sink class, while disregarding the other objective of limiting the influence on samples from the non-targeted classes. Since this loss clamping is applied to every targeted source class sample in the batch, it can also facilitate avoiding any sample dominating over other samples for contributing to the gradient of the universal perturbation. A similar clamping technique has been applied in [15] but with the objective to achieve a minimum-magnitude (image-dependent) perturbation that can attack a specific sample.

3.3 DTA Algorithm

With the loss functions defined above, the procedure to craft DT-UAPs with DTA is shown in Algorithm 1. For each perturbation update iteration, we include samples from both the targeted source class and the non-targeted source

classes. More specifically, we randomly select the same number (half of the mini-batch size) of samples from the targeted source class and the non-targeted source classes to form B_t and B_{nt} , which can be concatenated to one batch B . We then calculate the loss parts \mathcal{L}_t and \mathcal{L}_{nt} referring to Eq. 3 and Eq. 6, respectively. The total loss \mathcal{L} can then be calculated referring to Eq. 2. This procedure illustrates how the loss \mathcal{L} in Algorithm 1 is calculated. The perturbation can then be updated with the loss gradient calculated with respect to the perturbation. Note that the gradient thus computed is the expected gradient, i.e. the average of the gradients in this mini-batch. For the update of the perturbation, we can adopt any existing optimizer, but we empirically found that the ADAM [31] optimizer converges the fastest for our method. In the final step, the perturbation is projected to the l_p -ball with radius ϵ in order to satisfy the magnitude constraint. This process is repeated for I iterations. Mini-batch training and balancing the sample amount from the two data distributions result in a simple yet effective algorithm. Our algorithm is mainly inspired by [8,9]. Their algorithm has been shown to outperform UAP [16] and GAP [20] by a large margin, achieving SOTA performance for universal attack. Here, we tailor it to suite our purpose of being double targeted.

4 Results and Analysis

4.1 Experimental Setup

We apply the proposed DTA to various deep convolutional neural network architectures and construct perturbations on various datasets: CIFAR-10 [32], GTSRB [33], EuroSAT [34], YCB [35] and large-scale ImageNet [36]. CIFAR-10 and ImageNet are two commonly used benchmark datasets for image classification tasks. The GTSRB dataset consists of 43 classes of different German traffic signs and is a commonly used dataset for autonomous driving applications. The EuroSAT dataset is used for land cover classification tasks via satellite images categorized into 10 classes. The YCB dataset is a benchmark dataset for robotic manipulation and consists of a total of 98 classes of daily life objects.

For the different datasets, we evaluate DTA with at least two different networks. Overall, we explore various DNN architectures, including VGG-16 [37], ResNet-20/50 [38], Inception-V3 [39] and MobileNet-V2 [40]. To evaluate our approach, we use the metric of the targeted fooling ratio κ , which is defined as the ratio of samples fooled into the sink class. We apply the targeted fooling ratio to the targeted source class and non-targeted source classes, indicated by κ_t and κ_{nt} , respectively. Consequently, the higher (lower) κ_t (κ_{nt}), the better. For the following experiments, we set the number of iterations to $I = 500$, adopt the l_∞ norm and cap the perturbation magnitude at $\epsilon = 15$ for images in the range $[0, 255]$. All our experiments are performed using the PyTorch (v.0.4.1) [41] framework on a single GPU TITAN X (Pascal). Note that for crafting the perturbation, we only use the correctly classified images from the training dataset and report the results on all samples from the validation dataset.

Table 1: Experimental results for the Double Targeted Attack (DTA) for the datasets CIFAR-10, GTSRB, EuroSAT, YCB and ImageNet under 10 scenarios S_0 to S_9 . For each scenario, the targeted fooling ratios for the targeted source samples (κ_t) and the non-targeted source samples (κ_{nt}) are reported. All numbers are reported in %.

Dataset	Model	S_0		S_1		S_2		S_3		S_4		S_5		S_6		S_7		S_8		S_9		Avg	
		κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
CIFAR-10	VGG-16	77.5	20.5	83.5	22.0	78.2	14.7	81.4	21.5	73.0	18.6	79.1	14.2	75.1	15.1	76.7	24.6	75.0	20.3	86.2	16.6	78.6	18.8
	ResNet-20	78.8	26.1	84.6	28.0	84.0	24.3	84.2	26.9	77.1	22.0	82.1	21.3	83.8	14.7	72.9	33.2	80.0	27.8	89.8	22.3	81.7	24.7
GTSRB	VGG-16	89.0	0.2	100	1.1	87.1	1.2	72.2	0.6	91.0	1.3	83.6	2.4	88.3	1.1	80.0	0.7	95.0	1.9	81.1	1.7	86.7	1.2
	ResNet-20	84.3	0.5	100	1.6	53.1	0.2	77.8	1.8	87.6	2.9	77.1	4.4	70.0	2.7	88.3	1.2	80.0	0.3	64.4	0.7	78.3	1.6
EuroSAT	ResNet-50	96.2	33.0	98.8	18.0	95.2	31.1	96.6	22.1	99.2	28.7	95.0	24.0	94.4	44.3	96.3	17.6	96.3	24.5	91.2	22.7	95.9	26.6
	Inception-V3	94.3	28.7	95.2	18.9	93.8	41.4	99.2	56.3	93.0	29.4	93.0	24.2	91.6	34.6	96.0	21.8	96.8	31.6	89.2	18.8	94.2	30.6
YCB	ResNet-50	100	14.5	100	24.2	100	32.4	96.7	38.0	100	33.5	99.2	38.3	100	44.4	99.2	41.7	100	19.0	100	33.1	99.5	31.9
	Inception-V3	100	16.6	100	30.0	100	38.7	99.2	31.2	100	12.9	98.3	20.0	100	32.2	100	36.6	100	17.3	100	39.2	99.8	27.5
ImageNet	VGG-16	72.0	10.3	96.0	19.5	90.0	19.5	82.0	28.3	74.0	15.9	82.0	13.0	66.0	8.9	64.0	12.9	66.0	21.5	70.0	26.1	76.2	17.6
	ResNet-50	74.0	13.9	94.0	21.4	82.0	15.2	72.0	20.9	62.0	13.6	84.0	15.5	72.0	9.8	66.0	21.4	66.0	17.3	62.0	18.1	73.4	16.7
	Inception-V3	78.0	10.0	86.0	15.7	86.0	12.2	78.0	15.6	58.0	9.5	76.0	12.9	70.0	8.9	72.0	15.7	62.0	18.9	66.0	17.8	73.2	13.7
	MobileNet-V2	74.0	11.3	94.0	17.0	88.0	20.4	70.0	15.3	72.0	16.0	84.0	15.0	74.0	14.5	74.0	21.7	72.0	18.8	70.0	21.9	77.2	17.2

Table 2: Targeted source class to sink class mapping for the datasets CIFAR-10, GTSRB, YCB, EuroSAT, and ImageNet.

S	CIFAR-10	GTSRB	YCB	EuroSAT	ImageNet
S_0	bird → airplane	turn right ahead → turn left ahead	large clamp → strawberry	Herb. Vegetation → Annual Crop	wig → lab coat
S_1	deer → frog	end prev. limitation → end no passing	flat screwdriver → mini soccer ball	Industrial → Permanent Crop	photocopier → castle
S_2	frog → cat	no passing → no Lkw permitted	cups type f → larger marker	Permanent Crop → Highway	flagpole → sewing machine
S_3	ship → cat	wild animals possible → bicycle lane	hammer → lego duplo type i	River → Highway	jersey → rain barrel
S_4	truck → horse	no vehicles permitted → speed limit 70	cups type c → toy airplane part i	Sea Lake → Residential	theater curtain → brass
S_5	airplane → deer	no passing → speed limit 60	tuna fish can → plastic nut	Residential → Pasture	drilling platf. → pomegranate
S_6	horse → dog	slippery road → uneven surfaces	tomato soup can → cups type g	Permanent Crop → River	fireboat → aircraft carrier
S_7	dog → frog	pedestrian crossing → double curves	cups type h → chain	Pasture → Permanent Crop	torch → golfcart
S_8	dog → deer	speed limit 20 → Speed Limit 120	marbles type 3 → key	Pasture → Industrial	candle → howler monkey
S_9	airplane → automobile road narrows right	children crossing	cups type j → toy airplane part k	Annual Crop → Forest	ruddy turnstone → kuvasz

4.2 Quantitative Results

We evaluate the effectiveness of the proposed DTA by randomly selecting 10 source-to-sink shift scenarios indicated by S_0 to S_9 for each dataset. The results are summarized in Table 1, where we report the targeted fooling ratio for both the targeted class κ_t and the non-targeted classes κ_{nt} for each scenario. The exact mapping of the targeted source class to the sink class can be found in Table 2.

Overall, the results in Table 1 indicate that DTA achieves reasonable performance for different mapping scenarios on a wide range of datasets. This conclusion stems from two major observations. First, the targeted fooling ratio for the targeted classes (κ_t) is quite high. Second, there is a significant gap between κ_t and κ_{nt} , which indicates that the crafted perturbation is discriminative between targeted class and non-targeted classes. We further analyze the performance of each dataset.

CIFAR-10 With an average κ_t of around 80%, DTA performs reasonably well on CIFAR-10, fooling most of the targeted source class into the sink class. The gap between κ_t and κ_{nt} is about 58%, indicating sufficient selectivity.

GTSRB For the task of road sign classification, our proposed DTA can even achieve a 100% targeted fooling ratio for scenario S_1 while maintaining a very low targeted fooling ratio of 1.1% and 1.6% for VGG-16 and ResNet-20, respectively, on the non-targeted source samples. Overall, DTA exhibits high κ_t values, while maintaining the lowest κ_{nt} values among all examined datasets. Therefore, DTA achieves the highest gap between κ_t and κ_{nt} for the GTSRB dataset. The low κ_{nt} indicates that the perturbations for attacking GTSRB are especially covert. We speculate that the reason behind the high performance on the GTSRB dataset is that the in-class variation is very small, making the discriminative attack a relatively easy task.

EuroSAT and YCB The results of DTA on the EuroSAT and YCB datasets exhibit similar behavior, with very high values for κ_t , above 94%, while having κ_{nt} values of around 30%. With a gap of more than 60%, DTA poses a strong, covert threat for applications deploying satellite images and classification tasks for robotic manipulation.

ImageNet The results show that DTA is able to fool a network for a single class out of the 1000 into a sink class for all 4 investigated DNNs, namely VGG-16, ResNet-50, Inception-V3, and MobileNet-V2. For specific scenarios such as S_3 or S_4 , there can be a relatively large performance gap among different DNN architectures. Overall, with an average κ_t of around 75% and an average κ_{nt} of 16%, different DNNs have comparable performance.

4.3 Qualitative Results

In this subsection, we illustrate perturbations and perturbed samples generated by the proposed DTA. Fig. 2 shows the original targeted source image, along with the amplified universal perturbation and the resulting adversarial image. It can be observed that the DTA produces patterns with different characteristics for each dataset. The adversarial image is still identifiable as a source class instance to a human observer, however, the DNN classifies the manipulated image (from the targeted source class) with high confidence into the sink class.

4.4 Universal Multi2One Targeted Perturbation

Finally, we extend the DT-UAPs to a more challenging scenario to demonstrate an extension of the DTA. To this end, we alter the objective from one targeted source class to instead support multiple source classes (MS) while still leading the samples from these classes to one sink class. Due to this property of classifying multiple source classes to one sink class, we term the resulting perturbation a universal Multi2One targeted perturbation. Crafting such perturbations is more challenging since multiple source classes add complexity which has to be compensated by the universal perturbation. We evaluate this attack for 4 scenarios, which are detailed in Table 4 under the same settings as before. The results in

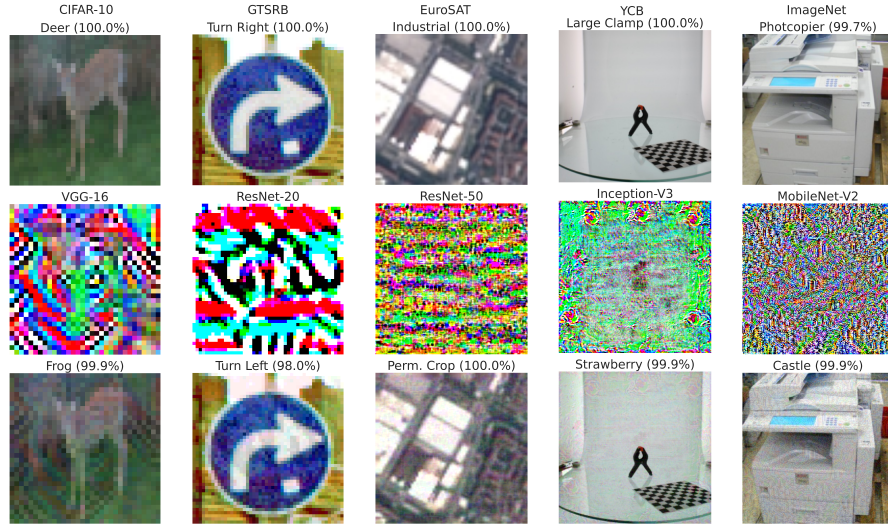


Fig. 2: Examples of adversarial perturbations for various datasets and networks. The figure shows the original images (top), an amplified version of the corresponding perturbations (middle) and the resulting adversarial examples (bottom). The confidence values of the network and the predicted labels are stated above the images. The target network is indicated above the amplified perturbation.

Table 3 show that our proposed DTA also achieves reasonable performance in the case of shifting multiple targeted source classes into the sink class.

4.5 Ablation Analysis

In the following, we perform ablation studies for the proposed DTA algorithm. All ablation experiments are performed on ResNet-20 for CIFAR-10 and ResNet-50 for the ImageNet dataset.

Loss Function We perform an ablation study for the loss function design. In Table 5 the performance of DTA for different loss function configurations is shown. We observe that our chosen loss design $\mathcal{L}_t + \mathcal{L}_{nt}$ achieves the best performance. In particular, we observe that excluding the non-targeted loss part \mathcal{L}_{nt} results in a very high κ_t close to 100% for both, the CIFAR-10 and ImageNet dataset. However, the κ_{nt} also increases drastically compared to the result obtained using $\mathcal{L}_t + \mathcal{L}_{nt}$. The average κ_{nt} for ImageNet is 74.5%, and that for CIFAR-10 is even higher with a value of 98.7%. This clearly shows that under the absence of \mathcal{L}_{nt} , DTA fails to achieve the objective of being discriminative between samples from the targeted source class and non-targeted source classes. With the existence of \mathcal{L}_{nt} , the absence of either \mathcal{L}_{t1} or \mathcal{L}_{t2} also leads to inferior

Table 3: Experimental results for the universal Multi2One targeted perturbation on ImageNet under 4 scenarios MS_0 to MS_3 . For each scenario, κ_t κ_{nt} are reported. All numbers are reported in %.

Model	MS_0		MS_1		MS_2		MS_3		Avg	
	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
VGG16	63.3	24.9	69.3	33.7	76.0	25.8	69.3	26.8	69.5	27.8
ResNet-50	64.0	30.1	63.3	32.3	78.7	29.2	62.7	23.2	67.2	28.7
Inception-V3	58.0	19.4	56.7	23.8	66.7	19.0	66.7	20.8	62.0	20.8
MobileNet-V2	68.0	27.2	66.0	28.0	74.0	25.6	66.0	24.4	68.5	26.3

Table 4: Targeted source classes to sink class mapping for the Multi2One attack on ImageNet.

MS	ImageNet
MS_0	affenpinscher, black grouse, alp \rightarrow mosque
MS_1	necklace, four-poster, jersey \rightarrow llama
MS_2	wig, photocopier, flagpole \rightarrow castle
MS_3	granny smith, dragonfly, drilling platform \rightarrow brass

performance. Moreover, with the existence of \mathcal{L}_{nt} , we further explore another variant of \mathcal{L}_t adopting the cross-entropy (CE) loss indicated as $\mathcal{L}_t^{\text{CE}}$. Similar to \mathcal{L}_t , $\mathcal{L}_t^{\text{CE}}$ is decomposed into two parts $\mathcal{L}_{t1}^{\text{CE}}$ and $\mathcal{L}_{t2}^{\text{CE}}$. $\mathcal{L}_{t1}^{\text{CE}}$ aims to reduce the logit value of the source class logit by calculating the negative cross-entropy between the network output and the one hot encoded source class label and $\mathcal{L}_{t2}^{\text{CE}}$ aims to increase the sink class logit by calculating the cross-entropy between the network output and the one hot encoded sink class label. We observe that this setup also achieves inferior performance compared to $\mathcal{L}_t + \mathcal{L}_{nt}$. The reason for this inferior performance can be attributed to the nature of the CE loss manipulating all logits, and not clamping the loss values.

Dominance Value D Further, we investigate the influence of the dominance value D for clamping the loss part \mathcal{L}_{t2} . Fig. 3 (left) shows the targeted fooling rates κ_t and κ_{nt} plotted over various dominance values. We observe that the value of D has a significant influence on the behavior of the proposed DTA. Increasing D increases both κ_t and κ_{nt} . More specifically, κ_t increases and saturates with further increasing D , while κ_{nt} increases almost linearly with the increase of D . The results show that it is beneficial to choose an appropriate D for achieving high κ_t with relatively low κ_{nt} . However, here we only aim to show the influence of the hyper-parameter D on the behavior of the proposed DTA and do not intend to find the optimal value which is dependent on the choice of models and dataset.

Perturbation Magnitude ϵ One constraint of adversarial perturbations is to be bound to a certain magnitude range. Here we investigate the influence of

Table 5: Analysis of the influence of different loss function configurations. For each scenario of the 4 different scenarios, κ_t and κ_{nt} are reported. All numbers are reported in %.

\mathcal{L}	Dataset	S_0		S_1		S_2		S_3		Avg	
		κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
$\mathcal{L}_t + \mathcal{L}_{nt}$	CIFAR-10	78.8	26.1	84.6	28.0	84.0	24.3	84.2	26.9	82.9	26.3
	ImageNet	74.0	13.9	94.0	21.4	82.0	15.2	72.0	20.9	80.5	17.9
$\mathcal{L}_t^{\text{CE}} + \mathcal{L}_{nt}$	CIFAR-10	77.4	46.1	89.4	49.6	88.8	44.6	88.1	57.0	85.9	49.3
	ImageNet	66.0	10.4	98.0	29.5	92.0	29.2	78.0	27.0	83.5	24.0
\mathcal{L}_t	CIFAR-10	99.2	97.9	99.6	98.4	99.7	99.3	100	99.3	99.6	98.7
	ImageNet	100	68.5	100	76.7	94.0	67.3	98.0	85.4	98.0	74.5
$\mathcal{L}_{t1} + \mathcal{L}_{nt}$	CIFAR-10	18.1	3.4	17.1	2.7	23.0	5.8	2.8	4.9	15.3	4.2
	ImageNet	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.1	0.0	0.1
$\mathcal{L}_{t2} + \mathcal{L}_{nt}$	CIFAR-10	81.9	32.4	88.9	38.9	89.4	34.4	90.0	35.3	87.6	35.3
	ImageNet	78.0	23.7	96.0	31.4	90.0	22.7	76.0	30.1	85.0	27.0

Table 6: Influence of α in Eq. 2 on the targeted fooling ratios κ_t and κ_{nt} .

Dataset	0.1		0.5		1		2		10	
	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
CIFAR-10	98.0	68.4	90.1	37.8	84.6	28.0	70.8	15.3	29.4	5.1
ImageNet	98.0	60.6	94.0	31.5	94.0	21.4	90.0	9.0	0.0	0.1

the perturbation magnitude ϵ and report the results in Fig 3 (right). A sharp increase of κ_t can be observed for ϵ values between 2.5 and 10 saturating around a targeted fooling ratio of 90% for further increased ϵ values, while κ_{nt} increases more steadily with increasing ϵ values.

Weighting Factor α One way an attacker can control the behavior of DTA is by manipulating the weighting factor α in Eq. 2. In Table 6 we evaluate the influence of α on κ_t and κ_{nt} . Higher values of α lead to lower values of κ_{nt} , since α weights the contribution of \mathcal{L}_{nt} to the final loss value. Even though this behavior is desired, κ_t decreases simultaneously. For an effective attack, an attacker might consider a large gap between κ_t and κ_{nt} , where neither a too large nor too small α is beneficial.

Number of Available Training Samples Finally, we investigate the influence of the available number of training samples on the attack behavior. In Table 7 we report the influence of the number of available training samples per class on the attack performance. With the same number of training iterations, we find that a smaller number of training samples per class lead to lower κ_t and κ_{nt} and the gap between κ_t and κ_{nt} decreases accordingly. However, with as small

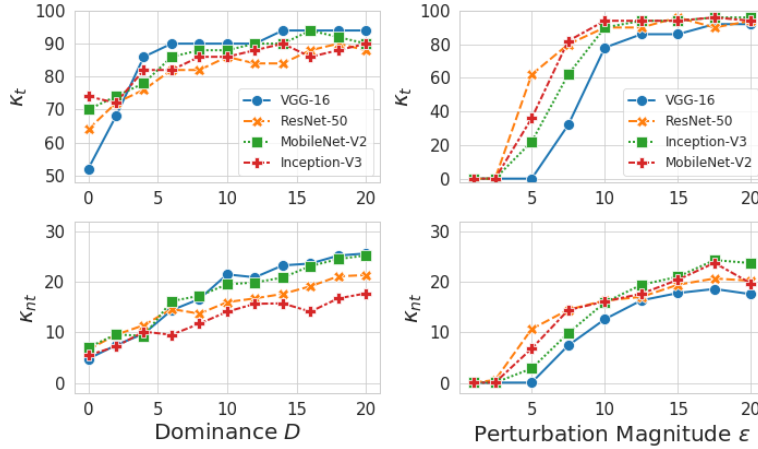


Fig. 3: Analysis of the influence of the dominance value D (left) and perturbation magnitude ϵ (right) on the targeted fooling ratios κ_t (top) and κ_{nt} (bottom) for the ImageNet dataset.

Table 7: Influence of number of training samples per class on the targeted fooling ratios κ_t and κ_{nt} .

Dataset	50		100		250		500		1000	
	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
CIFAR-10	46.7	18.4	60.2	20.3	73.9	22.7	80.9	27.6	83.2	27.8
ImageNet	50.0	2.9	64.0	4.3	86.0	12.5	94.0	19.0	96.0	17.7

as 50 samples per class, the algorithm still works reasonably well. For example, for ImageNet κ_t is 50% while κ_{nt} is as low as 2.9%.

Table 8: Quantitative results for the generated DT-Patch on ImageNet

Hammer \rightarrow Hummingbird		Screwdriver \rightarrow Go-Kart		Coffee Mug \rightarrow Chocolate Sauce	
κ_t	κ_{nt}	κ_t	κ_{nt}	κ_t	κ_{nt}
80.0	42.7	92.0	44.9	96.0	41.6

5 Double Targeted Patch

We extend DT-UAP to a physical-world attack [23,24] by generating a physical patch. We apply the concept of the DTA to attack one source class to a sink class.

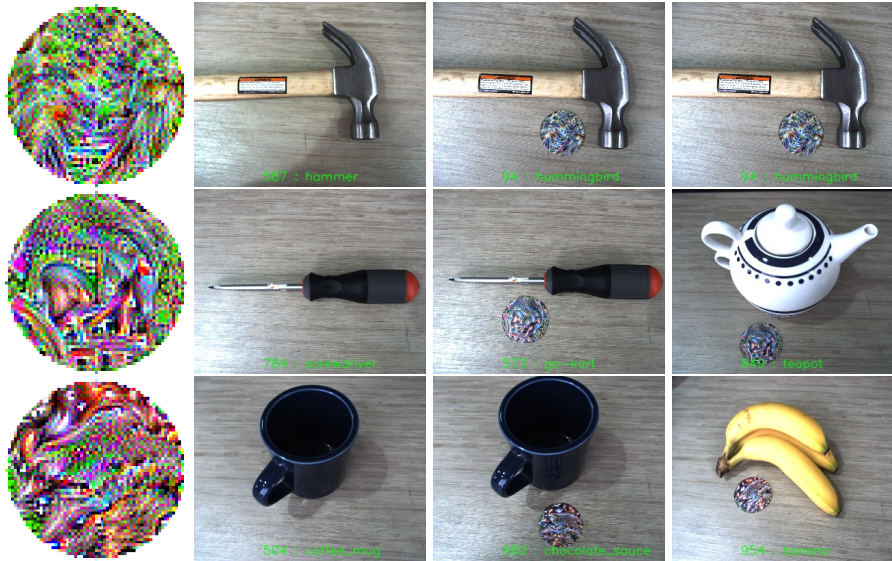


Fig. 4: Real-world examples of the DT-Patch for three different scenarios (see Table 8).

We choose VGG-16 trained on ImageNet as the target network. For generating a physical patch, we restrict the perturbation to a circular area, as well as its magnitude to lie in image range, i.e. $x + \delta \in [0, 1]$. We show three cases by choosing the source-sink class pairs as indicated in Table 8. Despite being a more challenging scenario than the original adversarial patch, we observe from Table 8 that κ_t is larger than κ_{nt} by a non-trivial margin. This indicates that the patch fulfills the objective. The qualitative results in Figure 4 show the applied patch fooling the source into the sink class, while having no influence on a sample from a non-targeted class.

6 Conclusion and Future Work

We proposed DTA to extend the existing UAP and CD-UAP for a more flexible attack control. The generated DT-UAP shifts one predefined source class into one predefined sink class, simultaneously attempts to minimize the targeted fooling ratio for samples from the non-targeted source classes. The effectiveness of DTA is demonstrated with extensive experiments on multiple datasets for different network architectures. We further presented an extension of DTA to the Multi2One scenario, driving multiple source classes into one sink class. With some preliminary results we found it also worked for a very challenging Multi2Multi scenario with limited success, and leave further explorations for future work.

References

1. Sutskever, I., Hinton, G.E., Krizhevsky, A.: Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* (2012) 1097–1105
2. Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Kingsbury, B., et al.: Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine* **29** (2012)
3. Collobert, R., Weston, J.: A unified architecture for natural language processing: Deep neural networks with multitask learning. In: *Proceedings of the 25th international conference on Machine learning*, ACM (2008) 160–167
4. Zhang, C., Rameau, F., Kim, J., Argaw, D.M., Bazin, J.C., Kweon, I.S.: Deepptz: Deep self-calibration for ptz cameras. In: *Winter Conference on Applications of Computer Vision (WACV)*. (2020)
5. Zhang, C., Rameau, F., Lee, S., Kim, J., Benz, P., Argaw, D.M., Bazin, J.C., Kweon, I.S.: Revisiting residual networks with nonlinear shortcuts. In: *British Machine Vision Conference (BMVC)*. (2019)
6. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* (2013)
7. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014)
8. Zhang, C., Benz, P., Imtiaz, T., Kweon, I.S.: Cd-uap: Class discriminative universal adversarial perturbation. In: *AAAI Conference on Artificial Intelligence (AAAI)*. (2020)
9. Zhang, C., Benz, P., Imtiaz, T., Kweon, I.S.: Understanding adversarial examples from the mutual influence of images and perturbations. In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. (2020)
10. Liu, H., Ji, R., Li, J., Zhang, B., Gao, Y., Wu, Y., Huang, F.: Universal adversarial perturbation via prior driven uncertainty approximation. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2019) 2941–2949
11. Benz, P., Zhang, C., Imtiaz, T., Kweon, I.S.: Data from model: Extracting data from non-robust and robust models. *CVPR Workshop on Adversarial Machine Learning in Computer Vision* (2020)
12. Akhtar, N., Mian, A.: Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access* **6** (2018) 14410–14430
13. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236* (2016)
14. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. (2016) 2574–2582
15. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: *2017 IEEE Symposium on Security and Privacy (SP)*, IEEE (2017) 39–57
16. Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P.: Universal adversarial perturbations. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. (2017) 1765–1773
17. Khrulkov, V., Oseledets, I.: Art of singular vectors and universal adversarial perturbations. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018) 8562–8570

18. Mopuri, K.R., Garg, U., Babu, R.V.: Fast feature fool: A data independent approach to universal adversarial perturbations. In: 2017 British Conference on Machine Vision (BMVC), IEEE (2017)
19. Metzen, J.H., Kumar, M.C., Brox, T., Fischer, V.: Universal adversarial perturbations against semantic image segmentation. In: 2017 IEEE International Conference on Computer Vision (ICCV), IEEE (2017) 2774–2783
20. Poursaeed, O., Katsman, I., Gao, B., Belongie, S.: Generative adversarial perturbations. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 4422–4431
21. Dong, Y., Liao, F., Pang, T., Su, H., Zhu, J., Hu, X., Li, J.: Boosting adversarial attacks with momentum. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2018) 9185–9193
22. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
23. Brown, T.B., Mané, D., Roy, A., Abadi, M., Gilmer, J.: Adversarial patch (2017)
24. Liu, A., Wang, J., Liu, X., Cao, B., Zhang, C., Yu, H.: Bias-based universal adversarial patch attack for automatic check-out. (2020)
25. Sallab, A.E., Abdou, M., Perot, E., Yogamani, S.: Deep reinforcement learning framework for autonomous driving. *Electronic Imaging* **2017** (2017) 70–76
26. Sünderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., Upcroft, B., Abbeel, P., Burgard, W., Milford, M., et al.: The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research* **37** (2018) 405–420
27. Melis, M., Demontis, A., Biggio, B., Brown, G., Fumera, G., Roli, F.: Is deep learning safe for robot vision. Adversarial examples against the iCub humanoid. CoRR, abs/1708.06939 (2017)
28. Clark, G., Doran, M., Glisson, W.: A malicious attack on the machine learning policy of a robotic system. In: 2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE), IEEE (2018) 516–521
29. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., Song, D.: Robust physical-world attacks on deep learning visual classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 1625–1634
30. Morgulis, N., Kreines, A., Mendelowitz, S., Weisglass, Y.: Fooling a real car with adversarial traffic signs. arXiv preprint arXiv:1907.00374 (2019)
31. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization (2014)
32. Krizhevsky, A.: Learning multiple layers of features from tiny images. Technical report (2009)
33. Stallkamp, J., Schlipsing, M., Salmen, J., Igel, C.: Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks* **32** (2012) 323–332
34. Helber, P., Bischke, B., Dengel, A., Borth, D.: Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* (2019)
35. Calli, B., Singh, A., Walsman, A., Srinivasa, S., Abbeel, P., Dollar, A.M.: The ycb object and model set: Towards common benchmarks for manipulation research. In: 2015 international conference on advanced robotics (ICAR), IEEE (2015) 510–517

36. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition, Ieee (2009) 248–255
37. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
38. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: European conference on computer vision, Springer (2016) 630–645
39. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 2818–2826
40. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 4510–4520
41. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. In: Advances in Neural Information Processing Systems. (2019) 8024–8035