# Adversarial Attacks on Monocular Depth Estimation

Ziqi Zhang‡, Xinge Zhu†, Yingwei Li§, Xiangqun Chen‡, Yao Guo‡

‡Peking University,　　†The Chinese University of Hong Kong,　　§Johns Hopkins University

*Abstract*— Recent advances of deep learning have brought exceptional performance on many computer vision tasks such as semantic segmentation and depth estimation. However, the vulnerability of deep neural networks towards adversarial examples have caused grave concerns for real-world deployment. In this paper, we present a systematic study of adversarial attacks on monocular depth estimation, an important task of 3D scene understanding in scenarios such as autonomous driving and robot navigation. In order to understand the impact of adversarial attacks on depth estimation, we first define a taxonomy of different attack scenarios for depth estimation, including *non-targeted attacks*, *targeted attacks* and *universal attacks*. We then adapt several state-of-the-art attack methods for classification on the field of depth estimation. Besides, multi-task attacks are introduced to further improve the attack performance for universal attacks. Experimental results show that it is possible to generate significant errors on depth estimation. In particular, we demonstrate that our methods can conduct targeted attacks on given objects (such as a car), resulting in depth estimation 3-4× away from the ground truth (e.g., from 20m to 80m).

## I. INTRODUCTION

Monocular depth prediction, *i.e.* predicting the per-pixel distances to the camera, is a key task for 3D scene understanding. Learning 3D scene geometries has many applications such as robot assisted surgery, robot navigation and autonomous driving [1], [2], [3], [4], [5]. With its widespread applications, more and more works have significantly promoted the monocular depth estimation performance using DCNN-based models [6], [7], [8], [9], [10], [11], [12].

Meanwhile, the advances in deep learning have largely pushed forward the state-of-the-art in various tasks in computer vision, such as image classification and semantic segmentation. However, several studies [13], [14] have demonstrated the vulnerability of deep learning based methods to deliberately generated adversarial examples, bringing reliability concerns to the applications in many safety-critical domains, such as autonomous driving and video surveillance. Hence, more and more work [15], [16] about adversarial attacks have been proposed to evaluate the robustness and reliability of neural network models before they are deployed in the real world.

Recent studies on adversarial attacks have served as the foundation of adversarial training, *i.e.* a common technique to enhance model robustness [17], [18], [19], [20]. Specifically, adversarial examples generated from the attack methods are combined into training data. Models trained with the mixture of clean data and adversarial data thus become more robust to adversarial attacks. Similar to attacks on classification, a systematic study of attack on depth estimation could also



(a) Original mean object depth is 8.2m



(b) Adversarial mean object depth is 3.7m (0.5×)



(c) Adversarial mean object depth is 35.5m (4.3×)



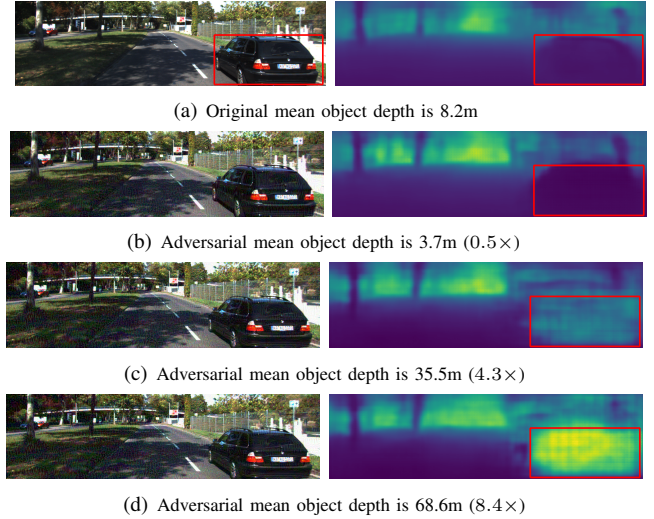(d) Adversarial mean object depth is 68.6m (8.4×)

Fig. 1: An illustration of the **targeted attack** on monocular depth estimation. The goal of this attack is to mislead the depth of specific objects (the black car in the red box in this case). The left column contains RGB images and the right column includes the corresponding depth prediction results. The first row shows the original images while three different adversarial examples are shown in the following rows. In the depth prediction image, bright yellow represents large depth value while dark blue represents small value. The depth estimation results are distorted up to 8× from the original prediction. Best viewed in color.

strengthen the robustness of depth models.

Although a number of attack methods [14], [21], [22], [23], [24], [13], [16] have been proposed, they mainly focus on classification [22], [14], [21] or segmentation [24], [25], [26]. Both tasks aim to pick a label among a fixed number of categories, so these techniques are not designed specifically to attack tasks such as monocular depth estimation, which aims to regress a precise value for each pixel.

In order to investigate the problem of adversarial attacks on monocular depth estimation, we first define a taxonomy of different attack scenarios for depth estimation, including *non-targeted attacks* (for a specific image), *targeted attacks* (for a specific object in an image) and *universal attacks* (for an arbitrary image). We then adapted several state-of-the-art attack methods to perform these attacks.

Moreover, we introduce a new multi-task attack strategy to improve the performance in the universal attack scenario. Multi-task strategies have been widely applied [27], [28], [29], [30] in supervised learning, where various supervi-

sion signals provide a more thorough understanding. In our setting, the segmentation task and depth estimation form the multi-task attack strategy (*i.e.* attacking both tasks simultaneously), where the high-level task (semantic segmentation) and low-level task (depth estimation) could offer complementary information to further boost the attacking performance for universal attacks.

We have performed extensive experiments on the popular KITTI dataset [31]. Results demonstrate that the average depth estimation errors (*i.e.* RMSE) can be distorted up to $10\times$ compared with the original prediction, which means that the adversarial attacks are real and practical. For targeted attacks, the distance estimation results can be distorted up to 3-4$\times$ (e.g., from 20m to 80m) for given objects. Fig. 1 displays an example. In addition, our proposed multi-task method achieves a better performance in the universal attacks, compared to single-task setting.

The main contributions of this work are as follows.

- This paper presents a systematic study of adversarial attacks on the task of monocular depth estimation. We define and benchmark various experimental settings for depth attacks, including non-targeted attacks, targeted attacks and universal attacks.
- We have adapted several state-of-the-art adversarial attack methods for classification to perform attacks on depth estimation. A new attack strategy, multi-task attack, is introduced to enforce attacking performance on single task with the support of auxiliary task by enriching supervision signals.
- Experimental results demonstrate that adversarial attacks on depth estimation is a real threat as we are able to generate significant errors in the estimation results. In particular, we are able to attack specific objects such as pedestrians and cars, bringing out much larger results than the clean estimation. In addition, our proposed multi-task attack strategy achieves better performance than existing methods.

## II. RELATED WORK

### A. Monocular Depth Estimation

Monocular depth estimation plays an important role for understanding spatial structures from 2D images. Traditional methods use triangulation to compute spatial position of each point corresponding in stereo images. Saxena *et al* [32] proposed a supervised learning approach leveraging local- and global-features. After that, various of methods based on hand-crafted features have been proposed [33], [34], [35], [36].

**DCNN based techniques**. Eigen *et al* first predicted dense depth estimation with deep neural networks [37]. Supervised deep learning-based approaches thus advanced the state of the art with different techniques. Xie *et al* utilized the skip-connection strategy to fuse low-resolution feature maps and high-resolution feature maps [38]. Laina *et al* proposed a fully convolutional architecture with residual learning and an up-sampling module [9]. Fu *et al* proposed an ordinal

regression loss to recast depth network learning as a classification problem [39]. In some work, CRF was integrated into deep architectures as well [40], [12], [41] .

**Multi-task strategy**. Recent work further demonstrated that depth estimation can be learned in a multi-task setting. Xu *et al* proposed to simultaneously solve the problem of depth estimation and scene parsing in a joint CNN [42]. Zhang *et al* proposed a task-recursive-learning framework for semantic segmentation and depth estimation [43].

### B. Adversarial Attacks

**White-box and black-box attacks**. White-box attack means that attackers have the access of the target model. Szegedy *et al* suggested that adding human imperceptible perturbations leads deep neural networks to wrong predictions [14]. Later, utilizing the linear property of neural network, a fast attack method was developed and named as Fast Gradient Sign Method [13]. Later, an iteration based method [16] was built to generate stronger adversarial examples. On the contrary, in black-box attack, attackers cannot access the target model. People developed a variety of methods to attack a black-box model, including query-based [44], [45], decision-based [46], and transfer-based methods [47], [22], [48], [49], [50], [51]. Several work studied the adversarial attacks on semantic segmentation [25], [52], [24].

**Per-image and universal attacks**. Per-image attacks mean to generate an adversarial example for each given image specifically. Oppositely, universal attacks generates a universal perturbation that can be directly added to any test image to fool the classifier. Moosavi-Dezfooli *et al* suggested the existence of universal adversarial perturbations and generated them by iteratively optimizing the per-instance adversarial loss [53]. Shafahi *et al* developed a much faster universal adversarial attack method [54]. Mopuri *et al* generated data-independent universal perturbation by maximizing spurious activations at each layer [55]. Metzen generated universal adversarial perturbations against semantic segmentation [24].

To the best of our knowledge, no existing work have studied adversarial attacks on tasks such as monocular depth estimation.

### C. Monocular Depth Estimation and Adversarial Attacks

Some work relates to both monocular depth estimation and adversarial attacks. Van *et al* [56] crafted some special cases to study the internal mechanism of how networks see depth from an image. The difference between this work and [56] is that this paper studies human-undetectable adversarial perturbation while Van *et al* constructs obvious fake images that can be recognized from one sight. Mopuri *et al* [55] proposed a data-free method to craft universal adversarial examples for a given CNN. Although their method is effective on depth estimation and semantic segmentation respectively, it aims at a single network and can not generate universal examples that can attack both tasks from a different structure

at the same time. On the other hand, the technique described in Section III-C generates universal examples that attack two tasks from two different networks. Hu *et al* [57] studied the white-box adversarial attacks (I-FGSM) on depth estimation in an indoor situation and proposed to defense by saliency map. By comparison, this paper studies a more dangerous situation: the depth of a certain object can be manipulated in autonomous driving. Besides, this paper utilizes multi-task (depth estimation and semantic segmentation) to attack better while Hu *et al* employs multi-task (depth estimation and saliency map) to defense.

## III. METHODOLOGY

In this section, we first formulate the problem of adversarial attacks on depth estimation in Section III-A. Then Section III-B presents three state-of-the-art methods we adapted from classification. Finally, multi-task strategy is illustrated in Section III-C .

### A. Problem Formulation

Let $\mathbf{x} \in \mathbb{R}^{3 \times \mathbf{h} \times \mathbf{w}}$ be an input image and $\mathbf{y}^{\text{true}} \in \mathbb{R}^{\mathbf{h} \times \mathbf{w}}$ be the ground truth depth image. $f$ represents deep neural network parameterized by $\theta$ and $L(f(\mathbf{x}; \theta), \mathbf{y}^{\text{true}})$ is the loss function. We denote $\xi$ as the adversarial perturbation and let $\mathbf{x}^{\text{adv}} = \mathbf{x} + \xi$ be the corresponding adversarial example. To make the adversarial example imperceptible, we restrict the perturbation $\|\xi\|_\infty < \epsilon$, where $\epsilon$ is a given perturbation constraint. More specifically, this paper considers three different types of depth attacks: *non-targeted attacks*, *targeted attacks* and *universal attacks*.

**Non-targeted attacks.** Its goal is to maximize the prediction error for a given image, such that the original model would predict incorrectly (output incorrect depth values), which is the most typical attack type. It maximizes a loss function as described by Eq. 1:

$$\max_{\mathbf{x}^{\text{adv}} - \mathbf{x}} L(f(\mathbf{x}^{\text{adv}}; \theta), \mathbf{y}^{\text{true}}) \ \ s.t. \ \|\mathbf{x} - \mathbf{x}^{\text{adv}}\|_\infty < \epsilon. \quad (1)$$

**Targeted attacks.** Taking into consideration the characteristics of depth estimation, this attack aims to mislead the model to produce incorrect depth estimation for specific (masked) objects towards a predefined depth value, whose objective is formulated in Eq. 2:

$$\min_{\mathbf{x}^{\text{adv}} - \mathbf{x}} L(f(\mathbf{x}^{\text{adv}}; \theta), C \cdot M + \mathbf{y}^{\text{true}} \odot (1 - M))$$
$$s.t. \ \|\mathbf{x} - \mathbf{x}^{\text{adv}}\|_\infty < \epsilon, \quad (2)$$

where $M$ is a binary object mask and $C$ is a predefined depth value. The goal is to cause mis-estimation of certain objects while preserving the other parts. For example, a targeted attack could induce an autonomous driving car to predict a rider ahead to be farther away.

**Universal attacks.** This third attack type aims to train a universal adversarial perturbation (UAP) that can be added to a broad class of images. The training goal is described

in Eq. 3, where $N$ is the training image number. It empowers attackers who cannot generate per-instance adversarial examples on the go with an image-agnostic perturbation.

$$\max_{\mathbf{x}^{\text{adv}} - \mathbf{x}} \frac{1}{N} \sum_{i=0}^{N} L(f(\mathbf{x}_\mathbf{i}^{\text{adv}}; \theta), \mathbf{y}_\mathbf{i}^{\text{true}})$$
$$s.t. \ \|\mathbf{x}_\mathbf{i} - \mathbf{x}_\mathbf{i}^{\text{adv}}\|_\infty < \epsilon \quad (3)$$

### B. Adapting Existing Attacking Methods

In order to reveal the impact of adversarial attacks on depth estimation, we first adapted three state-of-the-art attack methods to attack depth estimation. These tasks include Fast Gradient Sign Method (FGSM) [13], Iterative FGSM (I-FGSM) [16], and Momentum I-FGSM (MI-FGSM) [22]. With the loss functions defined earlier, we implement each of the methods as follow.

**FGSM.** It computes the direction of the loss gradient, and then adds the max possible perturbations on the original image, by

$$\mathbf{x}^{\text{adv}} = \mathbf{x} + \epsilon \cdot \text{sign} \left( \nabla_\mathbf{x} L(f(\mathbf{x}; \theta), \mathbf{y}^{\text{true}}) \right) \quad (4)$$

where $\text{sign}(\cdot)$ denotes the sign function.

**I-FGSM.** It initializes an adversarial example $\mathbf{x}_0^{\text{adv}} = \mathbf{x}$ and then iteratively updates it by

$$\mathbf{x}_{\text{t+1}}^{\text{adv}} = \text{Clip}_\mathbf{x}^\epsilon \{\mathbf{x}_\text{t}^{\text{adv}} + \alpha \cdot \text{sign} \left( \nabla_\mathbf{x} L(f(\mathbf{x}_\text{t}^{\text{adv}}; \theta), \mathbf{y}^{\text{true}}) \right)\}, \ (5)$$

The clip function $\text{Clip}_\mathbf{x}^\epsilon$ ensures the generated adversarial example is within the $\epsilon$-ball of the original image $x$ with ground-truth $\mathbf{y}^{\text{true}}$.

**MI-FSGM.** This method integrates the momentum term into the attack to stabilize update directions and escape from poor local maxima. At the $t^{\text{th}}$ iteration, the accumulated gradient is

$$\mathbf{g}_{\text{t+1}} = \mu \cdot \mathbf{g}_\text{t} + \frac{\nabla_\mathbf{x} L(f(\mathbf{x}; \theta), \mathbf{y}^{\text{true}})}{||\nabla_\mathbf{x} L(f(\mathbf{x}; \theta), \mathbf{y}^{\text{true}})||_1}, \quad (6)$$

where $\mu$ is the momentum decay factor. The sign of $g_{n+1}$ is then used to generate the adversarial example, by

$$\mathbf{x}_{\text{t+1}}^{\text{adv}} = \text{Clip}_\mathbf{x}^\epsilon \{\mathbf{x}_\text{t}^{\text{adv}} + \alpha \cdot \text{sign}(\mathbf{g}_{\text{t+1}})\}. \quad (7)$$

### C. Multi-Task Strategy for Universal Attacks

Typically, universal attacks are not as effective as non-targeted or targeted attacks in term of the increased error, as they often attack a set of images simultaneously. To further improve the performance of universal attacks, we introduce the multi-task attack strategy. Specifically, In our setting, both segmentation and depth estimation tasks are available. The proposed strategy aims to generate the universal adversarial perturbation to attack both tasks simultaneously. The conjugation of low-level information provided by depth estimation and semantic information from segmentation could offer more complementary signals than the single task, thus boosting the performance of depth attacking.

Due to the lack of ground truth (no depth ground truth and large-scale segmentation labels coexist simultaneously

in one dataset), we conduct the non-targeted attack to depth estimation and use the least-likely method (LLM) [16] to perform the segmentation task. LLM takes the least likely label $\mathbf{y}^{\text{LL}} = \arg\min_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$ as the attack target and increases the prediction probability of that label as described in Eq. 8.

$$\min_{\mathbf{x}^{\text{adv}} - \mathbf{x}} L(f(\mathbf{x}^{\text{adv}};\theta), \mathbf{y}^{\text{LL}}) \ \ s.t. \ \|\mathbf{x} - \mathbf{x}^{\text{adv}}\|_{\infty} < \epsilon \quad (8)$$

We utilize MI-FGSM to illustrate the generation of universal adversarial perturbation with the multi-task strategy. Note MI-FGSM can be replaced by any method in section III-B easily. To compute the adversarial perturbation for a minibatch, we use $\overline{L} = \mathbb{E}_{\mathbf{x} \in B} L(f(\mathbf{x}, \theta), \mathbf{y})$ as the target loss function, where B represents a minibatch. In specific, we define $\overline{L}_t^{\text{depth}} = \mathbb{E}_{\mathbf{x} \in B_t} L^{\text{depth}}(f(\mathbf{x};\theta), \mathbf{y}^{\text{true}})$ and $\overline{L}_t^{\text{semantic}} = \mathbb{E}_{\mathbf{x} \in B_t} L^{\text{semantic}}(f(\mathbf{x}, \theta), \mathbf{y}^{\text{LL}})$. To combine information from both tasks, we compute the averaged loss $L_t$ as Eq. 9 and averaged gradient $\overline{\mathbf{g}}_t$ as Eq. 10, where $w_{\text{depth}}$ and $w_{\text{semantic}}$ are predefined weights for each task.

$$L_t = w_{\text{depth}} \cdot |\overline{L}_t^{\text{depth}}| + w_{\text{semantic}} \cdot |\overline{L}_t^{\text{semantic}}| \quad (9)$$

$$\overline{\mathbf{g}}_t = w_{\text{depth}} \cdot \frac{\nabla_{\mathbf{x}} \overline{L}_t^{\text{depth}}}{\|\nabla_{\mathbf{x}} \overline{L}_t^{\text{depth}}\|_1} + w_{\text{semantic}} \cdot \frac{\nabla_{\mathbf{x}} \overline{L}_t^{\text{semantic}}}{\|\nabla_{\mathbf{x}} \overline{L}_t^{\text{semantic}}\|_1} \quad (10)$$

Our multi-task universal attack algorithm is summarized in Algorithm 1.

In our implementation, two tasks are involved in the multi-task attack and note that it is easy to extend this strategy to more tasks to take advantage of more supervision signals.

## IV. EXPERIMENTS

### A. Experimental Setup

**Dataset.** We use the KITTI depth estimation dataset [31], which consists of $85,898$ training images and $6,852$ validation images, created by aggregating LiDAR scans from consecutive frames before projecting into one image. It contains semi-dense depth ground truth of about $30\%$ annotated pixels. We center crop every image to the size of $321 \times 929$. In the per-image attack, all images in the validation set are used. In the universal attack, we randomly select $17,913$ images as the training set and test on the validation set. We train universal adversarial perturbation for 2 epochs.

**Evaluation metrics.** We adopt root mean squared error (RMSE) [58] and masked mean depth (MMD) to evaluate the performance on **non-targeted attack** and **targeted attack**, respectively. Specifically, given the ground truth depth $Y^* = \{y*\}$ and predicted depth $Y = \{y\}$, RMSE is computed by $\sqrt{\frac{1}{|Y^*|} \sum_{y^* \in Y^*} \|y^* - y\|^2}$. The **larger** RMSE is, the farther predicted depth is from the ground truth depth, which means a better attack performance. Denoting $M = \{m\}$ as the predefined mask, MMD is computed by $\frac{1}{|M|} \sum_{m \in M} y \cdot m$, which is the mean depth value in the masked area (served as the attack targets in the image). It represents the average

---

**Algorithm 1** Multi-Task Universal Perturbation Generation

---

**Require:** Training samples $\mathbf{X}$, perturbation bound $\epsilon$, learning rate $\gamma$, momentum $\mu$, epoch number $N_{\text{ep}}$ and iteration number T

**Require:** Loss function of the two tasks $L^{\text{depth}}$ and $L^{\text{semantic}}$ and weights of the two tasks $w_{\text{depth}}$ and $w_{\text{semantic}}$.

1: Randomly initialize $\delta$
2: **for** epoch $= 1 \ldots N_{\text{ep}}$ **do**
3:   **for** minibatch $B \subset X$ **do**
4:     Initialize $B_0 = B$, $\mathbf{g_0} = \mathbf{0}$
5:     **for** iteration $t = 0 \ldots T - 1$ **do**
6:       Compute single task loss $\overline{L}_t^{\text{depth}}$ and $\overline{L}_t^{\text{semantic}}$
7:       Compute single task gradient $\nabla_{\mathbf{x}} \overline{L}_t^{\text{depth}}$ and $\nabla_{\mathbf{x}} \overline{L}_t^{\text{semantic}}$
8:       Compute average loss according to Eq. 9
9:       Compute average gradient according to Eq. 10
10:      Update gradient $\mathbf{g_{t+1}} = \mu \cdot \mathbf{g_t} + \overline{\mathbf{g}}_t$
11:      Compute $\mathbf{x_{t+1}}$ by Eq. 7 and update minibatch $B_{t+1} = \{\mathbf{x_{t+1}}\}$
12:     **end for**
13:     Compute perturbation $\delta_B = \mathbb{E}_{\mathbf{x_T} \in B_T}(\mathbf{x_T} - \mathbf{x_0})$
14:     Update universal perturbation $\delta \leftarrow \delta + \gamma \cdot \delta_B$
15:   **end for**
16: **end for**
17: **return** Universal perturbation $\delta$

---

distance of objects to the camera. Note that the predefined $C$ in Eq. 2 is set to 100 meters in our experiments, thus the **larger** MMD is, the better performance the attack method shows (*i.e.* the results are closer to $C = 100$). The predefined masks is the masks of target objects. How these objects are selected will be explained in the following paragraph.

**Target Objects.** We select cars, people, and riders with small depth values as the attack target. Firstly, we use Mask-RCNN[59] to predict the instance segmentation for each image and select the instance mask of cars, people and riders. Then we use the sparse depth ground truth of KITTI dataset to estimate the average depth of each instance. Finally, we select instances whose average depth is smaller than 50m.

**Baseline models.** We follow the encoder-decoder framework in [9] to set up the models in our experiments. Four different backbones are used as the encoding part, including ResNet-18 [60], VGG-16 [61], ResNet-50 and ResNet-101. For the decoding parts, they are composed of 4 UpConv layers followed by a bilinear upsampling layer in all models [58]. Note that we refer each model by the backbone architecture. The total training epoch is 20 for all models and we show the baseline performance in Table II.

**Semantic segmentation model used in multi-task attack.** We utilize the state-of-the-art semantic segmentation model, PSPNet [62], to provide additional supervision signals. We use Cityscapes [63] to train the network. The total epoch is 100 and random crop ($321 \times 929$ in our experiment to suit the image size in KITTI) is used during training. When optimizing the universal attack with multi-task strategy, two

| | Metric | RMSE ↑ (non-targeted attack) | | | | MMD ↑ (targeted attack) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Method | ResNet-18 | VGG-16 | ResNet-50 | ResNet-101 | ResNet-18 | VGG-16 | ResNet-50 | ResNet-101 |
| ResNet-18 | FGSM | 11.54 (2.7×) | 10.30 (2.4×) | 11.02 (2.7×) | **10.18** (2.2×) | 23.51 (1.1×) | 18.45 (0.9×) | 19.34 (0.9×) | 19.66 (0.9×) |
| | I-FGSM | **35.05** (8.3×) | 10.30 (2.4×) | 16.02 (3.9×) | 7.95 (1.7×) | 72.33 (3.5×) | 26.95 (1.3×) | 34.42 (1.7×) | 25.46 (1.2×) |
| | MI-FGSM | 32.37 (7.7×) | **11.60** (2.6×) | **16.70** (4.0×) | 9.63 (2.0×) | **72.35** (3.5×) | 27.38 (1.3×) | **35.54** (1.7×) | **25.57** (1.2×) |
| VGG-16 | FGSM | 8.69 (2.1×) | 11.5 (2.6×) | 9.45 (2.3×) | **9.30** (2.0×) | 22.90 (1.1×) | 24.63 (1.2×) | 22.98 (1.1×) | 19.77 (0.9×) |
| | I-FGSM | 7.13 (1.7×) | **37.40** (8.5×) | 10.48 (2.5×) | 6.34 (1.3×) | 23.39 (1.1×) | 74.82 (3.6×) | 25.73 (1.2×) | **23.09** (1.1×) |
| | MI-FGSM | **9.00** (2.1×) | 36.23 8.2×) | **12.52** (3.0×) | 8.26 (1.8×) | **24.08** (1.1×) | **76.47** (3.7×) | **27.07** (1.3×) | 22.76 (1.1×) |
| ResNet-50 | FGSM | 10.27 (2.4×) | 9.98 (2.3×) | 11.58 (2.8×) | **9.59** (2.0×) | 20.48 (1.0×) | 19.77 (1.0×) | 25.27 (1.2×) | 20.42 (1.0×) |
| | I-FGSM | 8.26 (1.9×) | 8.22 (1.9×) | **42.27** (10.1×) | 6.42 (1.4×) | 24.77 (1.2×) | 24.58 (1.2×) | 78.72 (3.8×) | 23.92 (1.1×) |
| | MI-FGSM | **10.51** (2.5×) | **10.40** (2.3×) | 38.37 (9.2×) | 8.54 (1.8×) | **25.34** (1.2×) | **25.21** (1.2×) | **79.78** (3.8×) | **23.96** (1.1×) |
| ResNet-101 | FGSM | **9.26** (2.2×) | 8.69 (2.0×) | 10.22 (2.4×) | 12.54 (2.7×) | 23.88 (1.1×) | 20.45 (1.0×) | 23.61 (1.1×) | 31.14 (1.5×) |
| | I-FGSM | 7.39 (1.8×) | 7.13 (1.6×) | 9.49 (2.3×) | **25.25** (5.4×) | 23.87 (1.1×) | 23.59 (1.1×) | 26.16 (1.3×) | 56.35 (2.6×) |
| | MI-FGSM | 9.08 (2.2×) | **9.00** (2.0×) | **11.66** (2.8×) | 24.52 (5.2×) | **24.34** (1.2×) | **23.65** (1.1×) | **26.89** (1.3×) | **56.65** (2.7×) |

TABLE I: The overall evaluation (RMSE represents the average depth estimation error for the whole image, while MMD indicates the mean estimated depth value of the masked target) results for adversarial attacks in white-box (shaded cells) and black-box (others) settings. The number in the parentheses are relative to the baseline results in Table II. A higher value indicates better attack effect.

| Model | ResNet18 | VGG16 | ResNet50 | ResNet101 |
|---|---|---|---|---|
| RMSE(m) | 4.22 | 4.42 | 4.17 | 4.70 |
| MMD(m) | 20.76 | 20.57 | 20.82 | 21.34 |

TABLE II: Baseline performance of different models.

different networks are used for the depth estimation and semantic segmentation.

**Universal multi-task attack setting.** To implement universal multi-task attack, root mean squared error(RMSE) loss is selected as $L^{depth}$ in Equ. 9 and $L^{semantic}$ is the widely used softmax loss [62].

**Attack settings.** We use $L_2$ distance as the loss function $L$. If not specified intentionally, we set the perturbation constraint $\epsilon = 16$. For iterative attack methods, we set iteration number to $\min(\epsilon+4, 1.25\epsilon)$ following [16] and step size $\alpha = 1$. For MI-FGSM, we set $\mu = 1$ as [22] suggests.

*B. Single Model Attack*

In this experiment, we follow [22] to report the white-box and black-box attacks in the single-model setting, where all four networks attack each other in both **non-targeted** and **targeted** attacks. Table I reports the overall results, where models in each row generate adversarial perturbations that are evaluated on models in each column.

For non-targeted attacks, we are showing the mean error (RMSE) values to indicate the extent of distortion by these attacks. We can see that for white-box attacks, the estimation errors are up to $10\times$ of the original error in the baseline models (Table II). The black-box attacks are not so effective, but still achieve an average error of up to $4\times$ compared to the baselines.

For targeted attacks, we show the average distance in the estimation results in Table I. We can see that the estimated distance is up to $4\times$ (from around 20m to 80m) farther than the original prediction (Table II). As an illustrative example, Fig. 1 presents a case to attack a car nearby. This adversary phenomenon would be really dangerous if a pedestrian or car is estimated to be 68 meters away when it is actually only 8 meters away, in scenarios such as automonous driving.

In addition, Table I also demonstrates the robustness of different backbones. As shown in the shadow cells, ResNet-50 is the most fragile one in both non-targeted and targeted attacks. Considering its high performance on clean images, this observation is consistent with the previous claim [64] that models with higher performance are more vulnerable.

*C. Move Objects to an Arbitrary Distance*

This section shows that the depth value of an object can be manipulated to any number, *i.e.* move any object to an arbitrary distance. Fig. 2 displays a sample attack result. The example image is 2011_09_26/0036/image_02/0000000050.png of the KITTI depth dataset. The attack target is the black car on the right side of the image (red box in Fig 2(a)) and the original average depth of the target vehicle is 11.4m. After attack, the average depth of the target object can achieve a minimum of 4.4m and a maximum of 80.8m. This means the attacker can not only pull the target closer (Fig. 2(b), from 11.4m to 4.4m), but also push it further to nearly arbitrary distance (Fig. 2(d) to Fig. 2(g), from 11.4m to a maximum of 80.8m).

Note that this attack is concealed very well. On the one hand, the operation on the RGB image is very subtle (the left column of Fig. 2). On the other hand, the depth of other parts of the image is barely changed. In Fig. 2, the depth detail of the pole is preserved well (orange boxes in Fig 2) .

| Model | Method | ResNet-18 | VGG-16 | ResNet-50 |
|---|---|---|---|---|
| Single-task | FGSM | 5.967 | 6.784 | 6.18 |
| | I-FGSM | 6.108 | 6.877 | 6.464 |
| | MI-FGSM | **7.173** | **7.503** | **7.257** |
| Multi-task | FGSM | 6.693 | 8.098 | 8.197 |
| | I-FGSM | 5.566 | 6.021 | 5.863 |
| | MI-FGSM | **9.378** | **8.582** | **9.757** |

TABLE III: The average depth errors (RMSE) of universal adversarial attacks with respect to different attack methods and settings. Single-task represents utilizing only non-targeted depth information and multi-task means adopting information of both depth and segmentation.

(a) Original Image , Mean Object Depth=11.4m

(b) Adversarial Example, Mean Object Depth=4.4m

(c) Adversarial Example, Mean Object Depth=20.2m

(d) Adversarial Example, Mean Object Depth=31.0m

(e) Adversarial Example, Mean Object Depth=44.7m

(f) Adversarial Example, Mean Object Depth=54.8m

(g) Adversarial Example, Mean Object Depth=68.5m

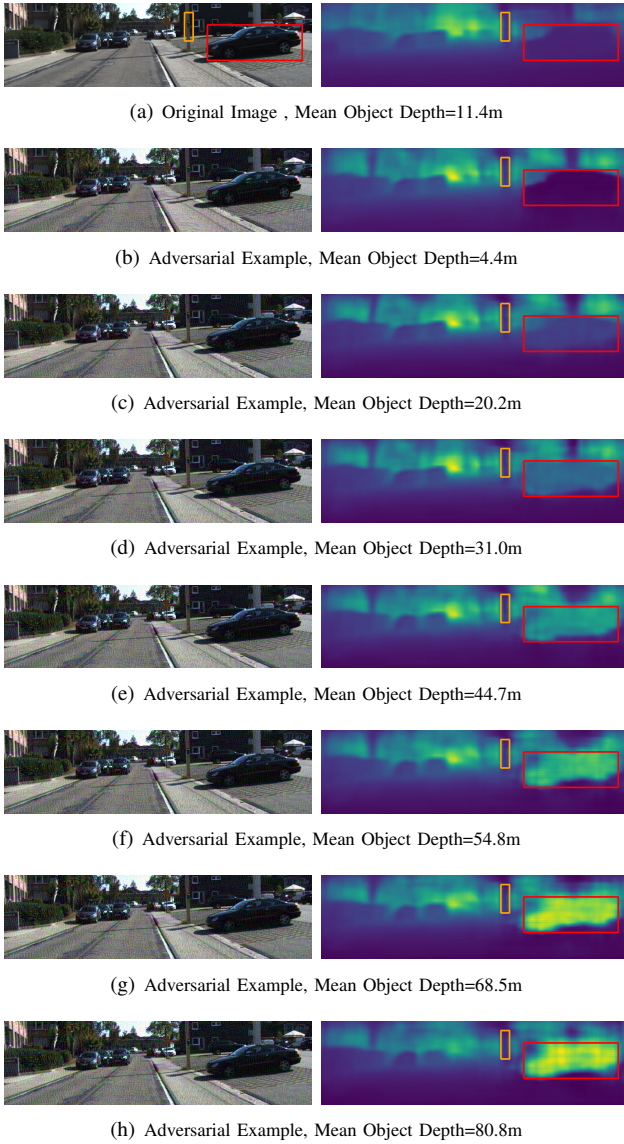(h) Adversarial Example, Mean Object Depth=80.8m

Fig. 2: Visualization of targeted attack on 2011_09_26/0036/image_02/0000000050.png of KITTI depth dataset to change the depth result of the target object (the black car in the red boxes) to arbitrary value without changing other part. Besides, the modification to the rgb image is subtle and the depth detail of other parts is preserved well (the pole in orange boxes)

*D. Universal Attack*

In this experiment, we perform the universal adversarial attack in both **multi-task** and single-task settings. We set $w_{\mathrm{depth}} = 0.5$, $w_{\mathrm{semantic}} = 0.5$ for multi-task setting and $w_{\mathrm{depth}} = 1$, $w_{\mathrm{semantic}} = 0$ for single-task in Eq. 10. Table III reports the white-box attack result, where single-task denotes to only use depth information and multi-task represents to leverage multiple supervision signals. As shown in Table III, with the help of the multi-task strategy, the results of FGSM and MI-FGSM are improved significantly compared to single-task. This demonstrates that more comprehensive



(a) Image

(b) Image + UAP

(c) Depth Prediction

(d) Adversarial Depth Prediction

(e) Semantic Segmentation

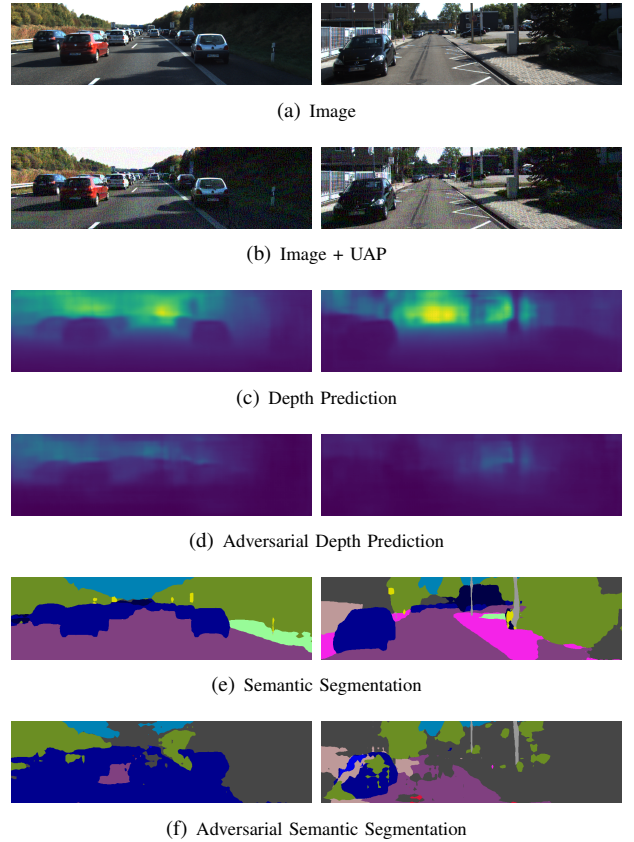(f) Adversarial Semantic Segmentation

Fig. 3: Visualization of the effect of the universal perturbation. The left and right represents two different images.

supervision signals do benefit the universal adversarial perturbation generation. The reason may be that the high-level information of semantic segmentation enriches supervision signal of low-level depth information, thus reinforces attack effect on universal adversarial perturbation. Fig. 3 displays an example of the universal attack, where the results of both depth estimation and semantic segmentation are distorted dramatically.

## V. CONCLUSION

In this paper, we present to the best of our knowledge the first systematic investigation of adversarial attacks on monocular depth estimation. We propose three types of attack, *i.e.* non-targeted attacks, targeted attacks and universal attacks and adapt three state-of-the-art methods to perform these attacks. Muti-task setting enriches supervision signals by adopting high-level information for universal attacks. Extensive experiments demonstrate the vulnerability of depth estimation models and the effectiveness of multi-task strategy compared to different methods. In particular, we demonstrate the targeted attack on certain objects is able to twist depth estimation up to an average of $4\times$ from the ground truth. Our work reveals the severe impact of adversarial attacks on depth estimation, as these attacks may result in grave security concerns if applied in cases such as autonomous driving. We hope this work can provide a baseline and guidance for adversarial attacks research on monocular depth estimation.

REFERENCES

[1] E. Coupeté, F. Moutarde, and S. Manitsaris, "Gesture recognition using a depth camera for human robot collaboration on assembly line," *Procedia Manufacturing*, vol. 3, pp. 518–525, 2015.

[2] G. N. DeSouza and A. C. Kak, "Vision for mobile robot navigation: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 24, no. 2, pp. 237–267, 2002.

[3] D. Stoyanov, M. V. Scarzanella, P. Pratt, and G.-Z. Yang, "Real-time stereo reconstruction in robotically assisted minimally invasive surgery," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2010, pp. 275–282.

[4] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, "Sparsity invariant cnns," in *International Conference on 3D Vision (3DV)*, 2017.

[5] Y. Ma, X. Zhu, S. Zhang, R. Yang, W. Wang, and D. Manocha, "Trafficpredict: Trajectory prediction for heterogeneous traffic-agents," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 6120–6127.

[6] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2650–2658.

[7] S. Kim, K. Park, K. Sohn, and S. Lin, "Unified depth prediction and intrinsic image decomposition from a single image via joint convolutional neural fields," in *European conference on computer vision*. Springer, 2016, pp. 143–159.

[8] Y. Kuznietsov, J. Stuckler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6647–6655.

[9] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 2016, pp. 239–248.

[10] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 2024–2039, 2016.

[11] A. Roy and S. Todorovic, "Monocular depth estimation using neural regression forest," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5506–5514.

[12] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille, "Towards unified depth and semantic prediction from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2800–2809.

[13] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *ICLR*, 2015.

[14] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *ICLR*, 2014.

[15] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE S&P*, 2017.

[16] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," *arXiv preprint arXiv:1607.02533*, 2016.

[17] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," *arXiv preprint arXiv:1802.00420*, 2018.

[18] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *ICLR*, 2018.

[19] U. Shaham, Y. Yamada, and S. Negahban, "Understanding adversarial training: Increasing local stability of supervised models through robust optimization," *Neurocomputing*, vol. 307, pp. 195–204, 2018.

[20] A. Sinha, H. Namkoong, and J. Duchi, "Certifying some distributional robustness with principled adversarial training," *arXiv preprint arXiv:1710.10571*, 2017.

[21] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in *ICLR*, 2017.

[22] Y. Dong, F. Liao, T. Pang, H. Su, J. Hu, J. Li, and J. Zhu, "Boosting adversarial attacks with momentum," in *CVPR*, 2018.

[23] N. Carlini and D. Wagner, "Adversarial examples are not easily detected: Bypassing ten detection methods," in *AISec*, 2017.

[24] J. H. Metzen, M. C. Kumar, T. Brox, and V. Fischer, "Universal adversarial perturbations against semantic image segmentation," in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 2774–2783.

[25] A. Arnab, O. Miksik, and P. H. Torr, "On the robustness of semantic segmentation models to adversarial attacks," 2018.

[26] R. Sun, X. Zhu, C. Wu, C. Huang, J. Shi, and L. Ma, "Not all areas are equal: Transfer learning for semantic segmentation via hierarchical region selection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4360–4369.

[27] X. Zhu, H. Zhou, C. Yang, J. Shi, and D. Lin, "Penalizing top performers: Conservative loss for semantic segmentation adaptation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 568–583.

[28] Z. Chen, V. Badrinarayanan, C.-Y. Lee, and A. Rabinovich, "Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks," *arXiv preprint arXiv:1711.02257*, 2017.

[29] I. Kokkinos, "Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6129–6138.

[30] X. Zhu, J. Pang, C. Yang, J. Shi, and D. Lin, "Adapting object detectors via selective cross-domain alignment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 687–696.

[31] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[32] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 5, pp. 824–840, 2009.

[33] C. Hane, L. Ladicky, and M. Pollefeys, "Direction matters: Depth estimation with a surface normal classifier," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 381–389.

[34] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.

[35] R. Ranftl, V. Vineet, Q. Chen, and V. Koltun, "Dense monocular depth estimation in complex dynamic scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4058–4066.

[36] Y. Xu, X. Zhu, J. Shi, G. Zhang, H. Bao, and H. Li, "Depth completion from sparse lidar data with depth-normal constraints," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2811–2820.

[37] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Advances in neural information processing systems*, 2014, pp. 2366–2374.

[38] J. Xie, R. Girshick, and A. Farhadi, "Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 842–857.

[39] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2002–2011.

[40] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5162–5170.

[41] D. Xu, W. Wang, H. Tang, H. Liu, N. Sebe, and E. Ricci, "Structured attention guided convolutional neural fields for monocular depth estimation," in *CVPR*, 2018.

[42] D. Xu, W. Ouyang, X. Wang, and N. Sebe, "Pad-net: Multi-tasks guided prediciton-and-distillation network for simultaneous depth estimation and scene parsing," in *CVPR*, 2018.

[43] Z. Zhang, Z. Cui, C. Xu, Z. Jie, X. Li, and J. Yang, "Joint task-recursive learning for semantic segmentation and depth estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 235–251.

[44] A. N. Bhagoji, W. He, B. Li, and D. Song, "Practical black-box attacks on deep neural networks using efficient query mechanisms," in *ECCV*, 2018.

[45] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, 2017.

[46] C. Guo, J. S. Frank, and K. Q. Weinberger, "Low frequency adversarial perturbation," *arXiv preprint arXiv:1809.08758*, 2018.

[47] S. Baluja and I. Fischer, "Learning to attack: Adversarial transformation networks," in *AAAI*, 2018.

[48] O. Poursaeed, I. Katsman, B. Gao, and S. Belongie, "Generative adversarial perturbations," in *CVPR*, 2017.

[49] C. Xiao, B. Li, J.-Y. Zhu, W. He, M. Liu, and D. Song, "Generating adversarial examples with adversarial networks," in *IJCAI*, 2018.

[50] C. Xie, Z. Zhang, J. Wang, Y. Zhou, Z. Ren, and A. Yuille, "Improving transferability of adversarial examples with input diversity," *arXiv preprint arXiv:1803.06978*, 2018.

[51] W. Zhou, X. Hou, Y. Chen, M. Tang, X. Huang, X. Gan, and Y. Yang, "Transferable adversarial perturbations," in *ECCV*, 2018.

[52] O. Poursaeed, I. Katsman, B. Gao, and S. Belongie, "Generative adversarial perturbations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4422–4431.

[53] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *CVPR*, 2017.

[54] A. Shafahi, M. Najibi, Z. Xu, J. Dickerson, L. S. Davis, and T. Goldstein, "Universal adversarial training," *arXiv preprint arXiv:1811.11304*, 2018.

[55] K. R. Mopuri, U. Garg, and R. V. Babu, "Fast feature fool: A data independent approach to universal adversarial perturbations," in *BMVC*, 2017.

[56] T. van Dijk and G. C. de Croon, "How do neural networks see depth in single images?" *arXiv preprint arXiv:1905.07005*, 2019.

[57] J. Hu and T. Okatani, "Analysis of deep networks for monocular depth estimation through adversarial attacks with proposal of a defense method," *arXiv preprint arXiv:1911.08790*, 2019.

[58] F. Mal and S. Karaman, "Sparse-to-dense: Depth prediction from sparse depth samples and a single image," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–8.

[59] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.

[60] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[61] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[62] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

[63] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.

[64] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness may be at odds with accuracy," *stat*, vol. 1050, p. 11, 2018.