

# Targeted Adversarial Perturbations for Monocular Depth Prediction

**Alex Wong**

Department of Computer Science  
University of California, Los Angeles  
alexw@cs.ucla.edu

**Safa Cicek**

Department of Electrical and Computer Engineering  
University of California, Los Angeles  
safacicek@ucla.edu

**Stefano Soatto**

Department of Computer Science  
University of California, Los Angeles  
soatto@cs.ucla.edu

## Abstract

We study the effect of adversarial perturbations on the task of monocular depth prediction. Specifically, we explore the ability of small, imperceptible additive perturbations to selectively alter the perceived geometry of the scene. We show that such perturbations can not only globally re-scale the predicted distances from the camera, but also alter the prediction to match a different target scene. We also show that, when given semantic or instance information, perturbations can fool the network to alter the depth of specific categories or instances in the scene, and even remove them while preserving the rest of the scene. To understand the effect of targeted perturbations, we conduct experiments on state-of-the-art monocular depth prediction methods. Our experiments reveal vulnerabilities in monocular depth prediction networks, and shed light on the biases and context learned by them.

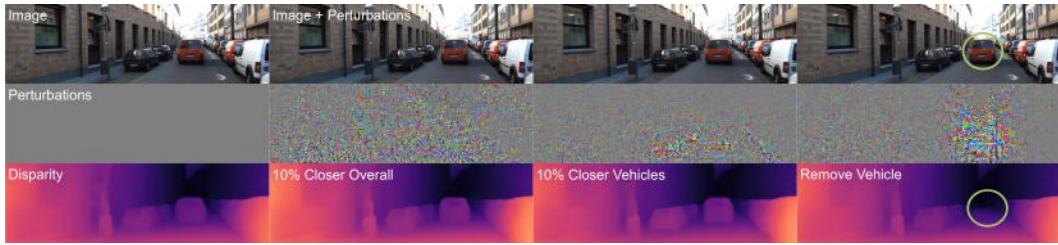


Figure 1: **Altering the predicted scene with adversarial perturbations.** Top to bottom: input image; adversarial perturbations with upper norm of  $2 \times 10^{-2}$ ; predicted scene visualized as disparity. Left to right: original image and predicted scene; overall scene altered to be 10% closer; all vehicles altered to be 10% closer; vehicle in the center of the road is removed by perturbations.

## 1 Introduction

Consider the image shown in the top-left of Fig. 1, captured from a moving car. The corresponding depth of the scene, inferred by a deep neural network and visualized as disparity, is shown underneath. Can adding a small perturbation cause the perceived vehicle in front of us disappear? Indeed, this is shown on the rightmost panel of the same figure: The perturbed image, shown on the top-right, is indistinguishable from the original. Yet, the perturbation, amplified and shown in the center row, causes the depth map to be altered in a way that makes the car in front of us disappear.

Adversarial perturbations are small signals that, when added to images, are imperceptible yet can cause the output of a deep neural network to change catastrophically [Szegedy et al., 2013]. We know that they can fool a network to mistake a tree for a peacock [Moosavi-Dezfooli et al., 2016]. But, as autonomous vehicles are increasingly employing learned perception modules, mistaking a stop sign for a speed limit [Eykholt et al., 2018] or causing obstacles to disappear is not just an interesting academic exercise. We explore the possibility that small perturbations can alter not just the class label associated to an image, but the inferred depth map, for instance to make the entire scene appear closer or farther, or portions of the scene, like specific objects, to become invisible or be perceived as being elsewhere in the scene.

When semantic segmentation is available, perturbations can target a specific category in the predicted scene. Some categories (e.g. traffic lights, humans) are harder to attack than others (e.g. roads, nature). When instance segmentation is available, perturbations can manipulate individual objects, for instance make a car disappear or move it to another location. We call these phenomena collectively as *stereopagnosia*, as the solid geometric analogue of prosopagnosia [Damasio et al., 1982].

Stereopagnosia sheds light on the role of context in the representation of geometry with deep networks. When attacking a specific category or instance, while most of the perturbations are localized, some are distributed throughout the scene, far from the object of interest. Even when the target effect is localized (e.g., make a car disappear), the perturbations are non-local, indicating that the network exploits non-local context, which represents a vulnerability. Could one perturb regions in the image, for instance displaying billboards, thus making cars seemingly disappear?

We note that, although the adversarial perturbations we consider are not universal, that is, they are tailored to a specific scene and its corresponding image, they are somewhat robust. Blurring the image after applying the perturbations reduces, but does not eliminate, stereopagnosia. To understand generalizability of adversarial perturbations, we examine the transferability of the perturbations between two monocular depth prediction models with different architectures and losses.

## 2 Related Work

Adversarial perturbations have been studied extensively for classification (Sec. 2.1). We focus on regression, where there exists some initial work. However, we study the *targeted* case where the entire scene, a particular object class, or even an instance is manipulated by the choice of perturbation.

### 2.1 Adversarial Perturbations

The early works [Szegedy et al., 2013, Goodfellow et al., 2014] show the existence of small, imperceptible additive noises that can alter the predictions of deep learning based classification networks. Since then many more advanced attacks [Moosavi-Dezfooli et al., 2016] have been proposed. [Moosavi-Dezfooli et al., 2017] showed the existence of universal perturbations i.e. constant additive perturbations to degrade the accuracy over the entire dataset.

More recently, [Naseer et al., 2019] studied transferability of attacks across datasets and models. [Peck et al., 2017] derived lower bounds on the magnitudes of perturbations. [Najafi et al., 2019] studied the attacks in semi-supervised learning setting. [Qin et al., 2019, Tramèr and Boneh, 2019] proposed methods to enhance robustness to adversarial attacks. [Laidlaw and Feizi, 2019] extended adversarial attacks beyond small additive perturbations. [Ilyas et al., 2019] showed that the existence of adversarial attacks makes deep networks more predictive.

Despite the exponentially growing literature on adversarial attacks for the classification task, there only have been a few works extending analysis of adversarial perturbations to dense-pixel prediction tasks. [Xie et al., 2017a] studied adversarial perturbations for detection and segmentation. [Hendrik Metzen et al., 2017] demonstrated targeted universal attacks for semantic segmentation. [Mopuri et al., 2018] examined universal perturbations in a data-free setting for segmentation and depth prediction to alter predictions in arbitrary directions. Unlike them, we study *targeted* attacks where network is fooled to predict a specific target.

Our goal is to analyze the robustness of the monocular depth prediction networks to different targeted attacks to explore possible explanations of what is learned by these models. With a similar motivation, [Hu et al., 2019] identified the smallest set of image pixels from which the network can estimate a depth map with small error. Unlike them, we analyze the monocular depth networks by studying their robustness against *targeted adversarial* attacks.

## 2.2 Monocular Depth Prediction

[Eigen et al., 2014, Eigen and Fergus, 2015, Liu et al., 2015, Liu et al., 2016, Laina et al., 2016] trained deep networks with ground-truth annotations to predict depth from a single image. However, high quality depth maps are often unavailable and, when available, are expensive to acquire. Hence, trends shifted to weaker supervision from crowd-sourced data [Chen et al., 2016], and ordinal relationships amongst depth measurements [Zoran et al., 2015, Fu et al., 2018].

Recently, supervisory trends shifted to unsupervised (self-supervised) learning, which relies on stereo-pairs or video sequences during training, and provides supervision in the form of image reconstruction. While depth from video-based methods is up to an *unknown* scale, stereo-based methods can predict depth in *metric* scale because the pose (baseline) between the cameras is known.

To learn depth from stereo-pairs, [Garg et al., 2016] predicted disparity by reconstructing one image from its stereo-counterpart. Monodepth [Godard et al., 2017] predicted both left and right disparities from a single image and laid the foundation for [Poggi et al., 2018, Pillai et al., 2019, Wong and Soatto, 2019]. To learn depth from videos, [Mahjourian et al., 2018, Zhou et al., 2017] also jointly learned pose between temporally adjacent frames to enable image reconstruction by reprojection. [Wang et al., 2018, Yang et al., 2018] leveraged visual odometry, [Fei et al., 2018] used gravity, [Casser et al., 2019, Luo et al., 2018] considered motion segmentation, and [Yin and Shi, 2018] jointly learned depth, pose and optical flow. Monodepth2 [Godard et al., 2019] explored both stereo and video-based methods and proposed a reprojection loss to discard potential occlusions.

To study the effect of adversarial perturbations, we examine the robustness of Monodepth2, the state-of-the-art, and its predecessor, Monodepth. While Monodepth2 proposed both stereo and video-based models, we choose their stereo model because the predicted depth is in *metric* scale, which enables us to study perturbations to alter the scale of the scene without changing its topology.

In Sec. 3, we discuss our method. We show perturbations for altering entire predictions to a target scene in Sec. 4 and localized attacks on specific categories and object instances in Sec. 5. We discuss the transferability of such perturbations in Sec. 6 and their robustness against defenses in Sec. B of Supp. Mat.

## 3 Finding Targeted Adversarial Perturbations

Given a pretrained depth prediction network,  $f_d : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}_+^{H \times W}$ ,  $f_d : x \mapsto d(x)$  our goal is to find a small additive perturbation  $v(x) \in \mathbb{R}^{H \times W \times 3}$ , as a function of the input image  $x$ , which can change its prediction to a target depth  $f_d(x + v(x)) = d^t(x) \neq d(x)$  with some norm constraint  $\|v(x)\|_\infty \leq \xi$  and high probability  $\mathbb{P}(f_d(x + v(x)) = d^t(x)) \geq 1 - \delta$ .

We begin by examining Dense Adversarial Generation (DAG) proposed by [Xie et al., 2017a] for finding adversarial perturbations for the semantic segmentation task. The perturbations from DAG can be formulated as the sum of a gradient ascent term (that pushes the predictions away from those of the original image) and a gradient descent term (that pulls predictions towards the target predictions). In the case of semantic segmentation, this formulation works well because the gradient ascent term suppresses the probability for the original predictions, which naturally increases the probability of the target predictions (zero-sum) driven by the gradient descent term. However, such is not the case for regression tasks, which requires the network to predict a real-valued scalar (as opposed to probability mass) for a targeted scene. Hence, the gradient ascent term maximizes the difference between the original and predicted depth, which results in DAG “overshooting” the target depth.

Instead, we use a simple objective function, similar to [Hendrik Metzen et al., 2017], but we modify it for the regression task by minimizing the normalized difference between predicted and target depth,

$$\ell(x, v(x), d^t(x), f_d) := \frac{\|f_d(x + v(x)) - d^t(x)\|_1}{d^t(x)}. \quad (1)$$

We minimize this objective function with respect to an image  $x$  by following an iterative optimization procedure as detailed in Alg. 1. The  $\text{CLIP}(v_n(x), -\xi, \xi)$  operation clamps any value of  $v(x)$  larger than  $\xi$  to  $\xi$  and any value smaller than  $-\xi$  to  $-\xi$ . For all the experiments,  $\xi \in \{2 \times 10^{-3}, 5 \times 10^{-3}, 1 \times 10^{-2}, 2 \times 10^{-2}\}$ .

### 3.1 Implementation Details

We evaluate adversarial targeted attacks on KITTI semantic split [Alhaija et al., 2018]. This is a dataset of 200 outdoor scenes, captured by car-mounted stereo cameras and a LIDAR sensor, with

---

**Algorithm 1** Proposed method to calculate targeted adversarial perturbations for a regression task.

---

**Parameters:** Learning rate  $\eta$ , noise upper norm  $\xi$ .  
**Inputs:** Image  $x$ , target depth map  $d^t(x)$ , pretrained depth network  $f_d$ .  
**Outputs:** Perturbation  $v_N(x)$ .  
**Init:**  $v_0(x) = 0$ .  
**for**  $n = 0 : N - 1$  **do**  
     $v_n(x) = \text{CLIP}(v_n(x), -\xi, \xi)$   
    Calculate  $\ell(x, v(x), d^t(x), f_d)$  as defined in Eqn. 1.  
     $v_{n+1}(x) = v_n(x) - \eta \nabla \ell(x, v(x), d^t(x), f_d)$   
**end for**

---

ground-truth semantic segmentation and instance labels. The semantic and instance labels in this split enables our experiments in Sec. 5 for targeting specific categories or instances in a scene.

The depth models (Monodepth, Monodepth2) are trained on the KITTI dataset [Geiger et al., 2012] using Eigen split [Eigen and Fergus, 2015]. The Eigen split contains 32 out of the total 61 scenes, and is comprised of 23,488 stereo pairs with an average size of  $1242 \times 375$ . Images are resized to  $640 \times 192$  as a preprocessing step and perturbations are computed with 500 steps of SGD. Entire optimization for each frame takes  $\approx 12$ s (Monodepth2 takes 22ms for each forward pass and 11s  $\approx 500 \times 22$ ms in total) using a GeForce GTX 1080. Details on hyper-parameters are provided in Sec. C of Supp. Mat.

For all the experiments, we use absolute relative error (ARE), computed with respect to the target depth  $d^t(x)$ , as our evaluation metric:

$$\text{ARE} = \frac{\|f_d(x + v(x)) - d^t(x)\|_1}{d^t(x)}. \quad (2)$$

## 4 Attacking the Entire Scene

Given a depth network  $f_d$ , our goal is to find adversarial perturbations to alter the predictions to a target scene  $d^t(x)$  for an image  $x$ . For this, we examine three settings (i) scaling the entire scene by a factor, (ii) symmetrically flipping the scene, and (iii) altering the scene to a preset scene.

### 4.1 Scaling the Scene

For autonomous navigation, misjudging an obstacle to be farther away than it is could prove disastrous. Hence, to alter the distances in the predicted scene without changing its topology or structure, we examine perturbations that will scale the scene (bringing the scene closer to or farther away from the camera) by a factor of  $1 + \alpha$ . The target scene is defined as:

$$d^t(x) = \text{scale}(f_d(x)) = (1 + \alpha)f_d(x) \quad (3)$$

for  $\alpha \in \{-0.10, -0.05, +0.05, +0.10\}$  or  $-10\%, -5\%$  (closer),  $+5\%, +10\%$  (farther), respectively. Column two of Fig. 1 shows the scene scaled 10% closer to the camera by applying visually imperceptible perturbations with  $\xi = 2 \times 10^{-2}$ . On average, scaling the scene by  $-10\%, -5\%, +5\%, +10\%$  with  $\xi = 2 \times 10^{-2}$  require an  $\|v(x)\|_1$  of 0.0160, 0.0124, 0.0126, and 0.0161, respectively. We note that scaling the scene by  $\pm 5\%$  requires less perturbations than  $\pm 10\%$  and the magnitude required for both directions is approximately symmetric. Also, perturbations are typically located along the object boundaries with concentrations on the road. For a side by side visualization of comparisons between different scaling factors, please see Sec. H.1 in Supp. Mat.

In Fig. 2-(a, b), we compare our approach with DAG, (Sec. 3) re-purposed for depth prediction task. While both are bounded by the same upper norm, DAG consistently produces results with higher error and generally with a higher standard deviation. As seen in Fig. 2-(a, b), even with  $\xi = 2 \times 10^{-3}$ , we are able to find perturbations that can scale the scene to be  $\approx 1\%$  from  $\pm 5\%$  and  $\approx 3\%$  from  $\pm 10\%$ . With  $\xi = 2 \times 10^{-2}$ , we are able to fully reach all four targets with less than  $\approx 0.5\%$  error.

### 4.2 Symmetrically Flipping the Scene

We now examine the problem setting where the target scene still retains the same structures given by the image, however, they are mirrored across the y- (horizontal flip) or x-axis (vertical flip).

For the *horizontal flip* scenario, we denote the target depth as  $d^t(x) = \text{fliph}(f_d(x))$  where the *fliph* operator horizontally flips the predicted depth map  $f_d(x)$  across the y-axis.

Fig. 3-(b) shows that the perturbations can fool the network into predicting horizontal flipped scenes. For scenes with different structures on either side,  $v(x)$  fools the network into *creating and removing*

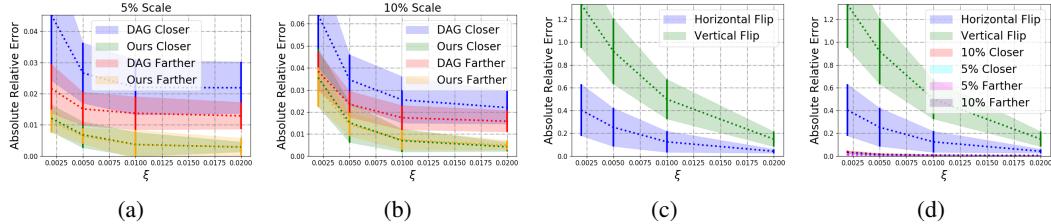


Figure 2: **ARE with various upper norm  $\xi$  for scaling and flipping Monodepth2 predictions.** (a) and (b) Comparisons between DAG and the proposed method for scaling the scene by  $\pm 5\%$  and  $\pm 10\%$ . (c) Results for horizontally and vertically flipping the predictions. (d) comparison between scaling and flipping tasks. Vertically flipping proves to be the most challenging.

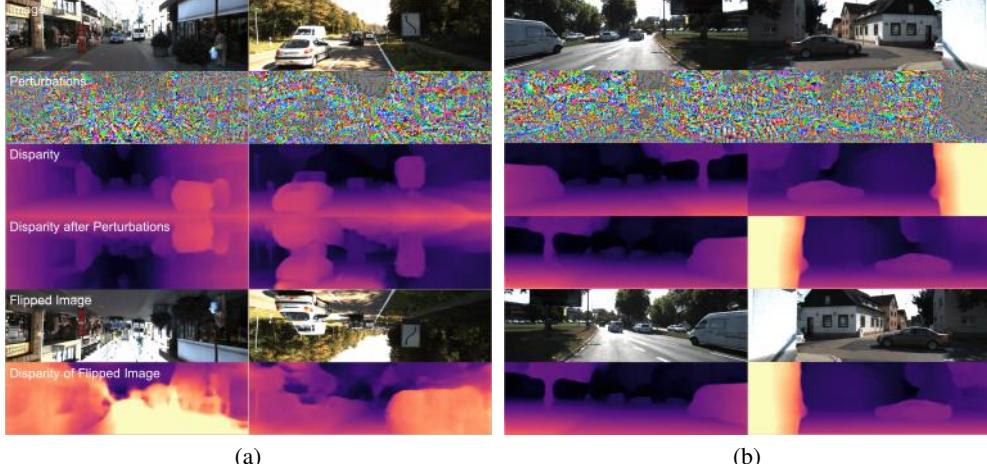


Figure 3: (a) **Examples of success (left) and failure (right) cases for vertical flip.** For the failure case, the car and road still remain on the bottom of the predictions. This is likely because the network is biased to predict closer structures on the bottom half of the image and farther ones on the top half (last two rows). (b) **Examples of horizontal flip.** Here, we observe the noise required to create and remove surfaces. Surprisingly, removing the white wall (right) requires very little perturbations.

surfaces. We note that the amount of noise required to horizontally flip the scene is much more than that of scaling the scene (i.e. for  $\xi = 2 \times 10^{-2}$ ,  $\|v(x)\|_1 = 0.161$  for scaling  $+10\%$  and  $0.326$  for flipping horizontally), which illustrates the difficulty in altering the scene structures. Interestingly, the amount of noise required to remove the white wall (Fig. 3-(b)) is significantly less than the rest.

For the *vertical flip* scenario, we denote the target depth as  $d^t(x) = \text{flipv}(f_d(x))$  where *flipv* operator vertically flips the predicted depth map  $f_d(x)$  across the x-axis.

As seen in Fig. 3-(a), perturbations cannot fully flip the predictions vertically. Even on successful attempts (left), there are still artifacts in the output. For failure cases (right), portions of the cars still remain on the bottom half of the predictions. This experiment reveals the potential *biases* learned by the network. To verify this, we feed vertically flipped images to the network. As seen in the last two rows of Fig. 3-(a), the network still assigns closer depth values to the bottom half of the image (now sky) and farther depth values to the top half (now road and cars).

In Fig. 2-(c, d), we plot the ARE achieved by the proposed method for different target depth maps: horizontal flip, vertical flip and different scales. Both flipping tasks are much harder than the scaling tasks. Particularly, fooling the network to produce vertical flipped predictions is the most challenging task as the error is  $\approx 15\%$ , even with  $\xi = 2 \times 10^{-2}$ .

### 4.3 Altering Predictions to Fit a Preset Scene

We now examine perturbations for altering the predicted scene  $d(x_1) = f_d(x_1)$  to an entirely different pre-selected one  $d^t(x_1) = f_d(x_2)$  obtained from images sampled from the same training distribution  $x_1, x_2 \sim \mathbb{P}(x)$ :  $d(x_1) = f_d(x_1) \neq d^t(x_1) = f_d(x_2)$ .

Fig. 4 shows that cars can be removed and road signs can be replaced with trees (leftmost), walls (column two) and vegetation (column three) can be added to the scene, and an urban street with vehi-

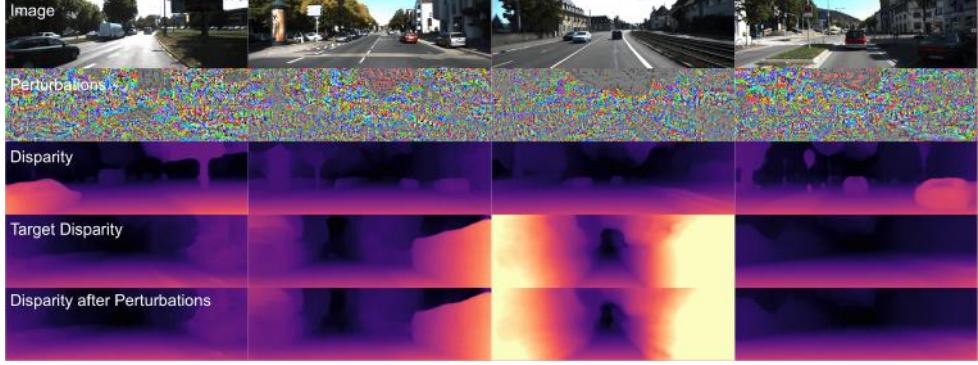


Figure 4: **Altering the predicted scene to a preset scene.** Adversarial perturbations can remove a car and replace a road sign with trees (leftmost), add walls (column two), and vegetation (column three) to open streets, and transform an urban environment with vehicles to an open road (rightmost).

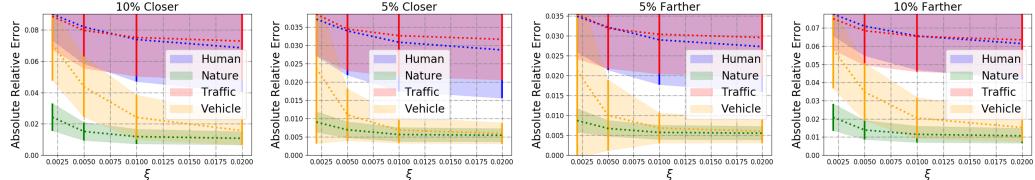


Figure 5: **ARE for scaling different categories closer and farther.** It is easier to fool the network to predict vehicle and nature categories closer and farther than is to fool human and traffic categories.

cles can be transformed to an open road (rightmost). While perturbations are visually imperceptible, we note that  $\|v(x)\|_1 = 0.362$ , which is  $\approx 2 \times$  the amount required for horizontal flip. However, the existence of such perturbations demonstrates just how vulnerable depth prediction networks can be.

Additionally, this experiment also confirms the biases learned by network discussed in Sec. 4.2. While perturbations can alter the scene to a preset one with *structures not present* in the image, we have difficulties finding perturbations that can vertically flip the predicted scene.

## 5 Localized Attacks on the Scene

Given semantic and instance segmentation [Alhaija et al., 2018], we now examine adversarial perturbations to target localized regions in the scene. Our goal is to fool the network into (i) predicting depths that are closer or farther by a factor of  $1 + \alpha$  for all objects belonging to a semantic category, (ii) removing specific instances from the scene, and (iii) moving specific instances to different regions of the scene, all the while keeping the rest of the scene unchanged.

### 5.1 Category Conditioned Scaling

Unlike Sec. 4.1, we want to alter a subset of the scene, partitioned by semantic segmentation, such that predictions belonging to an object category (e.g. vehicle, nature, human) are brought closer to or farther from the camera by a factor of  $1 + \alpha$  for  $\alpha \in \{-0.10, -0.05, +0.05, +0.10\}$ .

We assume a binary category mask  $M \in \{0, 1\}^{H \times W}$  derived from a semantic segmentation where all pixels belonging to a category are marked with 1 and 0 otherwise. We denote the target depth as

$$d^t(x) = (\mathbf{1} - M) \circ f_d(x) + (1 + \alpha)M \circ f_d(x) \quad (4)$$

where  $\mathbf{1}$  is an  $H \times W$  matrix of 1s. Column three of Fig. 1 illustrates this problem setting where the perturbations fool Monodepth2 into predicting all vehicles to be 10% closer to the camera. Fig. 5 shows a comparative study between different categories. Unlike Sec. 4.1, it is more difficult to alter a specific portion of the scene without affecting the rest. Moreover, each category exhibits a different level of robustness to adversarial noise. Some categories are harder to attack than others, e.g. traffic signs and human categories ( $\approx 3\%$  error for  $\alpha = \pm 5\%$  and  $\approx 6\%$  for  $\alpha = \pm 10\%$ ) are harder to alter than vehicle and nature ( $\approx 0.5\%$  error for  $\alpha = \pm 5\%$  and  $\approx 1\%$  for  $\alpha = \pm 10\%$ ). For interested readers, please see Sec. H.3 in Supp. Mat. for visualizations, additional experiments, and performance comparisons amongst all categories.

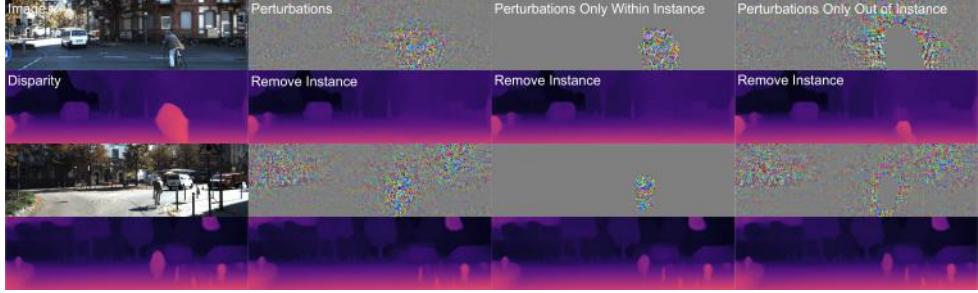


Figure 6: **Selectively removing instances of human (bikers and pedestrians).** Removing a localized target requires attacking non-local contextual information. Moreover, one can attack an instance without perturbing it at all. We demonstrate this by constraining the perturbation to be either completely within the instance mask or completely out of the mask.

## 5.2 Instance Conditioned Removing

We now consider the case where instance labels are available for removing a specific instance (e.g. car, pedestrian) from the scene. By examining this scenario, we hope to shed light on the possibility that a depth prediction network can “miss” a human or car, which may cause incorrect rendering in augmented reality or an accident in the autonomous navigation scenario.

Similar to Sec. 4.1, we assume a binary mask  $M$ , but in this case, of specific instance(s) in the scene, e.g. all pixels belonging to a specific pedestrian are marked with 1 and 0 otherwise. To obtain  $d^t(x)$ , we first remove the depth values in  $f_d(x)$  belonging to  $M$  by multiplying  $f_d(x)$  by  $\mathbf{1} - M$ . Then, we use the depth values  $f_d(x)$  on the contour of  $M$  to linearly interpolate the depth in the missing region:

$$d^t(x) = (\mathbf{1} - M) \circ f_d(x) + M \circ d_M^t(x) \quad (5)$$

where  $d_M^t(x) := \text{interp}(\text{contour}(f_d(x), M))$ . Fig. 6 shows examples of pedestrian and biker removal in the driving scenario where perturbations completely remove the targeted instance. With this attack, the road ahead becomes clear, which makes the agent susceptible to causing an accident.

Even though perturbations are concentrated on the targeted instance region, non-zero perturbations can be observed in the surrounding regions. While the target effect is localized (e.g., make a pedestrian disappear), the perturbation is non-local, implying that the network exploits non-local context, which presents a vulnerability to attacks against a target instance by perturbing other parts of the image.

## 5.3 Instance Conditioned Removing with Spatial Constraints on the Perturbations

Motivated by our results in Sec. 5.2, we extend the instance conditioned removal task to more constrained scenarios where the perturbations have to either exist (i) completely within the target instance mask  $M$  or (ii) completely outside of it.

For perturbations within the targeted instance mask  $M$ , we constrain  $v(x)$  to satisfy  $\|M \circ v(x)\|_\infty \leq \xi$  and  $\|(\mathbf{1} - M) \circ v(x)\|_\infty = 0$ . When constrained within  $M$ , perturbations can only remove *some instances* successfully (e.g. biker is completely removed in row two, column three of Fig. 6). In other cases, the perturbations can only remove the outer part of the instance, leaving parts of the instance in the scene (row four). This shows that depth prediction networks leverage global context; without attacking the contextual information located outside of  $M$  (e.g. without perturbing the entire image as in column two of Fig. 6), it is not always possible to completely remove the target instance.

Second, we want to answer the question posed in Sec. 1. Can perturbations remove an object by attacking anywhere (e.g. a billboard), *but the object* (e.g. a car)? In this more challenging case, the perturbations are constrained to be outside of the instance mask:  $\|(\mathbf{1} - M) \circ v(x)\|_\infty \leq \xi$  and  $\|M \circ v(x)\|_\infty = 0$ . Column four of Fig. 6 shows that even though there are no “direct attacks on the object” (perturbations in the masked region), the perturbations can *still remove parts of the target instance*. While some of the target instance still remains, our experiment demonstrates that depth prediction networks are indeed susceptible to *attacks against a target instance that does not require perturbing the instance at all*.

## 5.4 Instance Conditioned Translation

In this case study, we examine perturbations for moving an instance (e.g. vehicle, pedestrian) horizontally or vertically in the image space. As Sec. 5.2 and 5.3 have demonstrated the ability to remove localized objects from the scene, we now show that it is possible for perturbations to move

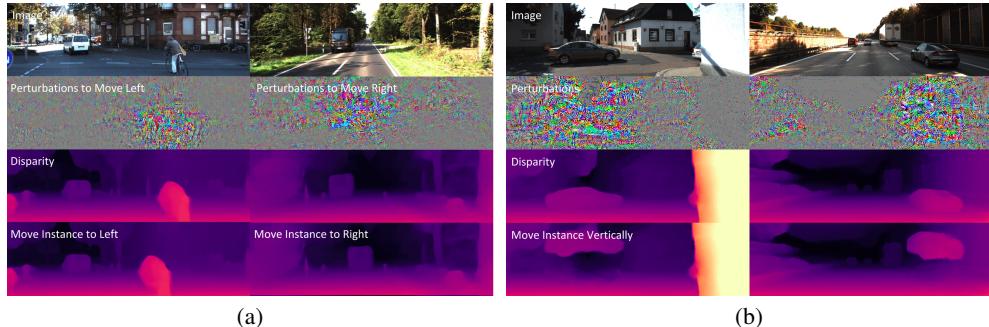


Figure 7: (a) **Moving horizontally.** Selected instances is moved by  $\approx 8\%$  in the left and right directions while rest of the scene is preserved. (b) **Flying cars.** A vehicle instance is moved  $\approx 42\%$  upward while rest of the scene is preserved. Noise is concentrated around the targeted instance.

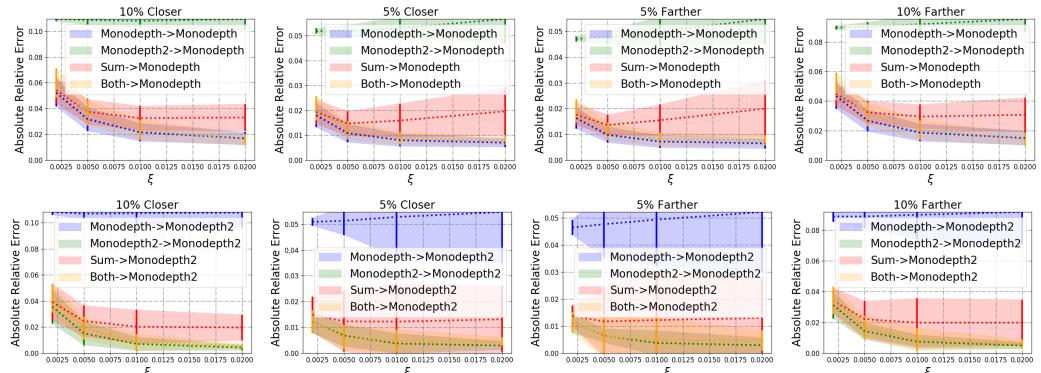


Figure 8: **Transferability across models.** Perturbations are (i) optimized for Monodepth and Monodepth2 separately, (ii) optimized for both together and (iii) summed over perturbations calculated for Monodepth and Monodepth2 separately. Each is tested on Monodepth and Monodepth2.

such objects to different locations in the scene (removing the instance and creating it elsewhere) while keeping the rest of the scene unchanged.

Fig. 7-(a) shows that perturbations can fool a network to move the target instance by  $\approx 8\%$  across the image in the left and right directions. When moved left, the biker (left column) is now in front our vehicle. When moved right, the truck (right column) is in the wrong lane and looks to be on-coming traffic. Moreover, Fig. 7-(b) shows that perturbations can move select instances by  $\approx 42\%$  in the upward direction, creating the illusion that there are “flying cars” in the scene.

## 6 Transferability Across Different Models

Transferability is important for black-box scenarios, a practical setting where the attacker does not have access to the target model or its training data. To examine transferability, we test our perturbations crafted for Monodepth2 [Godard et al., 2019] to fool its predecessor Monodepth [Godard et al., 2017] (different architecture and loss function) in Monodepth2 $\rightarrow$ Monodepth, and vice versa in Monodepth $\rightarrow$ Monodepth2, for the scene scaling task. (Sec. 4.1).

To this end, we also optimized perturbations for Monodepth to scale the entire scene. Overall, the perturbations optimized for one model does not transfer to another and, interestingly, transferability decays with increasing norm (Fig. 8), which may be due to perturbations overfitting to the model. We summed the perturbations for Monodepth and Monodepth2 (“Sum” in Fig. 8) and found that their summation can affect both models with reduced effects as the upper norm increases. For  $\xi = 2 \times 10^{-3}$ , the potency is nearly unaffected, meaning, for small norms, their summation can attack both models equally well. Lastly, by optimizing for both models (“Both” in Fig. 8), the same perturbation can fool both as if it was optimized for the models individually, with performance indistinguishable from Monodepth $\rightarrow$ Monodepth and Monodepth2 $\rightarrow$ Monodepth2 across all norms. This shows that both models share a space that is vulnerable to adversarial attacks. Hence, crafting perturbations for an array of potential models may be an avenue towards achieving absolute transferability across models.

## 7 Conclusion

Depth prediction networks are indeed vulnerable to adversarial perturbations. Not only can such perturbations alter the perception of the scene, but can also affect specific instances, making the network behave unpredictably, which can be catastrophic in applications that involve interaction with physical space. These perturbations also shed light on the network’s dependency on non-local context for local predictions, making non-local targeted attacks possible. While we have exposed vulnerabilities, we hope that our findings on the network’s biases, the effect of context, and robustness of the perturbations can help design more secure and interpretable models that are not susceptible to such attacks.

## References

- [Alhaija et al., 2018] Alhaija, H., Mustikovela, S., Mescheder, L., Geiger, A., and Rother, C. (2018). Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision (IJCV)*.
- [Casser et al., 2019] Casser, V., Pirk, S., Mahjourian, R., and Angelova, A. (2019). Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8001–8008.
- [Chen et al., 2016] Chen, W., Fu, Z., Yang, D., and Deng, J. (2016). Single-image depth perception in the wild. In *Advances in Neural Information Processing Systems*, pages 730–738.
- [Cordts et al., 2016] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., and Schiele, B. (2016). The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [Damasio et al., 1982] Damasio, A. R., Damasio, H., and Van Hoesen, G. W. (1982). Prosopagnosia: anatomic basis and behavioral mechanisms. *Neurology*, 32(4):331–331.
- [Deng et al., 2009] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.
- [Eigen and Fergus, 2015] Eigen, D. and Fergus, R. (2015). Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658.
- [Eigen et al., 2014] Eigen, D., Puhrsch, C., and Fergus, R. (2014). Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374.
- [Eykholt et al., 2018] Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., and Song, D. (2018). Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1625–1634.
- [Fei et al., 2018] Fei, X., Wong, A., and Soatto, S. (2018). Geo-supervised visual depth prediction. *arXiv preprint arXiv:1807.11130*.
- [Fu et al., 2018] Fu, H., Gong, M., Wang, C., Batmanghelich, K., and Tao, D. (2018). Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011.
- [Garg et al., 2016] Garg, R., BG, V. K., Carneiro, G., and Reid, I. (2016). Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer.
- [Geiger et al., 2012] Geiger, A., Lenz, P., and Urtasun, R. (2012). Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE.
- [Godard et al., 2017] Godard, C., Mac Aodha, O., and Brostow, G. J. (2017). Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 270–279.
- [Godard et al., 2019] Godard, C., Mac Aodha, O., Firman, M., and Brostow, G. J. (2019). Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3828–3838.
- [Goodfellow et al., 2014] Goodfellow, I. J., Shlens, J., and Szegedy, C. (2014). Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- [Hendrik Metzen et al., 2017] Hendrik Metzen, J., Chaithanya Kumar, M., Brox, T., and Fischer, V. (2017). Universal adversarial perturbations against semantic image segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2755–2764.

- [Hu et al., 2019] Hu, J., Zhang, Y., and Okatani, T. (2019). Visualization of convolutional neural networks for monocular depth estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3869–3878.
- [Ilyas et al., 2019] Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., and Madry, A. (2019). Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, pages 125–136.
- [Kingma and Ba, 2014] Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Laidlaw and Feizi, 2019] Laidlaw, C. and Feizi, S. (2019). Functional adversarial attacks. In *Advances in Neural Information Processing Systems*, pages 10408–10418.
- [Laina et al., 2016] Laina, I., Rupprecht, C., Belagiannis, V., Tombari, F., and Navab, N. (2016). Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE.
- [Liu et al., 2015] Liu, F., Shen, C., and Lin, G. (2015). Deep convolutional neural fields for depth estimation from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5162–5170.
- [Liu et al., 2016] Liu, F., Shen, C., Lin, G., and Reid, I. (2016). Learning depth from single monocular images using deep convolutional neural fields. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2024–2039.
- [Luo et al., 2018] Luo, C., Yang, Z., Wang, P., Wang, Y., Xu, W., Nevatia, R., and Yuille, A. (2018). Every pixel counts++: Joint learning of geometry and motion with 3d holistic understanding. *arXiv preprint arXiv:1810.06125*.
- [Mahjourian et al., 2018] Mahjourian, R., Wicke, M., and Angelova, A. (2018). Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5667–5675.
- [Moosavi-Dezfooli et al., 2017] Moosavi-Dezfooli, S.-M., Fawzi, A., Fawzi, O., and Frossard, P. (2017). Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1765–1773.
- [Moosavi-Dezfooli et al., 2016] Moosavi-Dezfooli, S.-M., Fawzi, A., and Frossard, P. (2016). Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2574–2582.
- [Mopuri et al., 2018] Mopuri, K. R., Ganeshan, A., and Babu, R. V. (2018). Generalizable data-free objective for crafting universal adversarial perturbations. *IEEE transactions on pattern analysis and machine intelligence*, 41(10):2452–2465.
- [Najafi et al., 2019] Najafi, A., Maeda, S.-i., Koyama, M., and Miyato, T. (2019). Robustness to adversarial perturbations in learning from incomplete data. In *Advances in Neural Information Processing Systems*, pages 5542–5552.
- [Naseer et al., 2019] Naseer, M. M., Khan, S. H., Khan, M. H., Khan, F. S., and Porikli, F. (2019). Cross-domain transferability of adversarial perturbations. In *Advances in Neural Information Processing Systems*, pages 12885–12895.
- [Peck et al., 2017] Peck, J., Roels, J., Goossens, B., and Saeys, Y. (2017). Lower bounds on the robustness to adversarial perturbations. In *Advances in Neural Information Processing Systems*, pages 804–813.
- [Pillai et al., 2019] Pillai, S., Ambrus, R., and Gaidon, A. (2019). Superdepth: Self-supervised, super-resolved monocular depth estimation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9250–9256. IEEE.
- [Poggi et al., 2018] Poggi, M., Tosi, F., and Mattoccia, S. (2018). Learning monocular depth estimation with unsupervised binocular assumptions. In *2018 International Conference on 3D Vision (3DV)*, pages 324–333. IEEE.
- [Qin et al., 2019] Qin, C., Martens, J., Gowal, S., Krishnan, D., Dvijotham, K., Fawzi, A., De, S., Stanforth, R., and Kohli, P. (2019). Adversarial robustness through local linearization. In *Advances in Neural Information Processing Systems*, pages 13824–13833.
- [Silberman et al., 2012] Silberman, N., Hoiem, D., Kohli, P., and Fergus, R. (2012). Indoor segmentation and support inference from rgbd images. In *European conference on computer vision*, pages 746–760. Springer.
- [Szegedy et al., 2013] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2013). Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- [Tramèr and Boneh, 2019] Tramèr, F. and Boneh, D. (2019). Adversarial training and robustness for multiple perturbations. In *Advances in Neural Information Processing Systems*, pages 5858–5868.

- [Wang et al., 2018] Wang, C., Miguel Buenaposada, J., Zhu, R., and Lucey, S. (2018). Learning depth from monocular videos using direct methods. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2022–2030.
- [Wong and Soatto, 2019] Wong, A. and Soatto, S. (2019). Bilateral cyclic constraint and adaptive regularization for unsupervised monocular depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5644–5653.
- [Xie et al., 2017a] Xie, C., Wang, J., Zhang, Z., Zhou, Y., Xie, L., and Yuille, A. (2017a). Adversarial examples for semantic segmentation and object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1369–1378.
- [Xie et al., 2017b] Xie, S., Girshick, R., Dollár, P., Tu, Z., and He, K. (2017b). Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500.
- [Yang et al., 2018] Yang, N., Wang, R., Stückler, J., and Cremers, D. (2018). Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *European Conference on Computer Vision*, pages 835–852. Springer.
- [Yin et al., 2019] Yin, W., Liu, Y., Shen, C., and Yan, Y. (2019). Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5684–5693.
- [Yin and Shi, 2018] Yin, Z. and Shi, J. (2018). Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1983–1992.
- [Zhou et al., 2017] Zhou, T., Brown, M., Snavely, N., and Lowe, D. G. (2017). Unsupervised learning of depth and ego-motion from video. In *CVPR*, volume 2, page 7.
- [Zoran et al., 2015] Zoran, D., Isola, P., Krishnan, D., and Freeman, W. T. (2015). Learning ordinal relationships for mid-level vision. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 388–396.

# Supplementary Materials

## A Summary of Contents

In Sec. B, the robustness of perturbations against defenses is discussed. Additional implementation details that we could not fit into main text due to space constraints are given in Sec. C. More experimental results on changing the scale of the scene are provided in Sec. D. In Sec. E, existence of the successful adversarial attacks for indoor scenes (NYU-V2) is shown for state-of-the-art indoor monocular depth prediction model. In Sec. F, we examine how predictions behave when linear operations are applied to perturbations (sum of two perturbations and linear scaling of a perturbation). Limitations and failure cases for the perturbations are analyzed in Sec. G. Finally, in Sec. H, more qualitative and quantitative results are provided for the experiments whose compressed versions are presented in the main text.

## B Robustness of the Targeted Attacks Against Defense Mechanisms

In the main text, we have shown that depth prediction networks are prone to adversarial attacks. In this section, we will examine the robustness of the perturbations against common defense mechanisms: (i) Gaussian blurring and (ii) adversarial training.

### B.1 Defense through Gaussian Blurring

In Fig. 9, we show the effect of Gaussian blurring as a simple defense mechanism on our targeted attacks by blurring the image with additive perturbations. Even though Gaussian blur does reduce the effectiveness of the perturbations (increased ARE over all scales), the resulting scene is still only  $\approx 3\%$  away from a target depth that is 10% closer or farther than the original predictions for  $\xi = 2 \times 10^{-3}$ . This is the performance that the method achieves for the case of  $\xi = 2 \times 10^{-3}$  without blurring. In other words, the effect of the blurring can be suppressed simply by increasing the upper norm of the noise by  $10\times$  for scaling the scene by  $\pm 10\%$ .

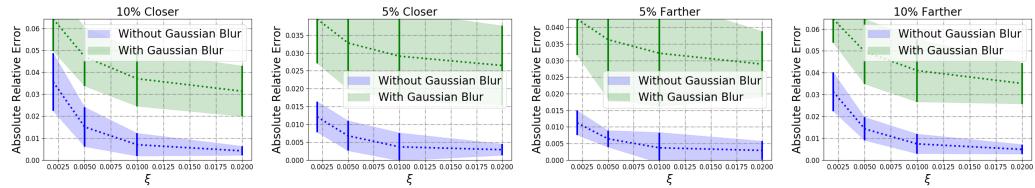


Figure 9: **Gaussian blur.** Absolute relative error (ARE) achieved by adversarial perturbations of different norms ( $\xi$ ) for different scales. For each scale, we plot the ARE with and without Gaussian blur. Even though absolute relative error increases with the Gaussian blur, the proposed method can still find a small norm noise to alter the scene.

### B.2 Defense through Adversarial Training

To examine the robustness of adversarial perturbations to adversarial training, we crafted adversarial perturbations for scaling the scene by a factor of  $1 + \alpha$  where  $\alpha \in \{-0.10, -0.05, +0.05, +0.10\}$  for the KITTI Eigen split [Eigen and Fergus, 2015] (consisting of 22600 stereo pairs). We trained Monodepth2 [Godard et al., 2019] by minimizing the normalized discrepancy between the predicted depth of a perturbed image ( $f_d(x + v(x))$ ) and its prediction for the original image ( $f_d(x)$ ).

$$\ell(x, v(x), f_d) = \frac{\|f_d(x) - f_d(x + v(x))\|_1}{f_d(x)} \quad (6)$$

Fig. 10 shows the effect of the perturbations on Monodepth2 after adversarial training. While training does reduce the influence of adversarial perturbations on the scene scaling task, it does not make the network invariant to the adversarial perturbations. With perturbations of  $\xi = 2 \times 10^{-3}$ , the predict scene is  $\approx 7\%$  from the target

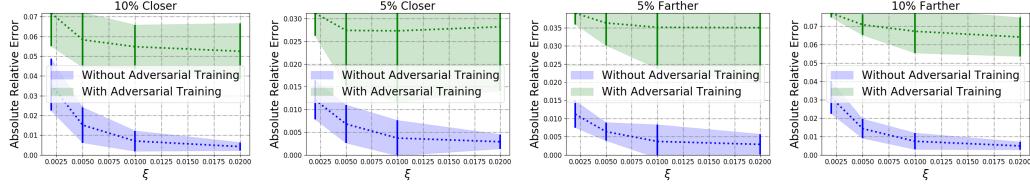


Figure 10: **Adversarial training.** Absolute relative error (ARE) achieved by adversarial perturbations of different norms ( $\xi$ ) for different scales. For each scale, we plot the ARE with and without adversarial training. Even though absolute relative error increases with the adversarial training, the perturbations can still affect the predicted scene with small norm noise.

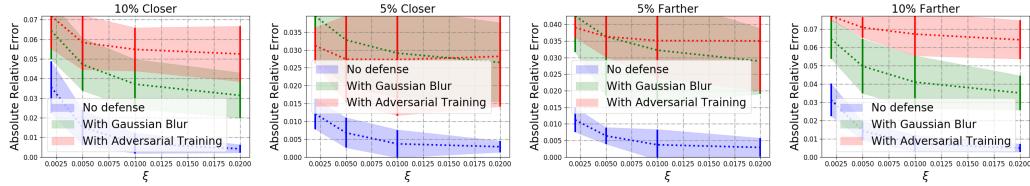


Figure 11: **Adversarial training vs Gaussian blur.** Absolute relative error (ARE) achieved by adversarial perturbations of different norms ( $\xi$ ) for different scales. For each scale, we plot the ARE without any defense, with Gaussian blur and with adversarial training. Both Gaussian blur and adversarial training makes the depth prediction network more robust to perturbations. Performances of the two defense mechanisms are comparable for small norms ( $\xi$ ), but adversarial training is more effective on larger norms.

scene that is 10% closer than or farther from the original and  $\approx 3\%$  from the target scene that is 5% closer or farther. For  $\xi = 2 \times 10^{-2}$ , perturbations can still fool the network to predict a target scene scaled by  $\pm 10\%$  with  $\approx 5\%$  absolute relative error and  $\approx 2\%$  error for fooling the network to predict a target scene that is scaled by  $\pm 10\%$ .

To compare the two defense mechanisms, we refer to Fig. 11. For smaller norms, e.g.  $\xi = 2 \times 10^{-3}$ , we observe a similar performance in using Gaussian blur (Sec. B.1) and adversarial training as defenses against adversarial perturbations; whereas, adversarial training is clearly better for larger  $\xi$ . This may be due to Gaussian blur's ability to destroy the perturbation for small norms and, hence, able to mitigate the effect of the perturbations. However, for larger norms, the blurring does not corrupt the perturbations enough and therefore does not reduce the effect of perturbations by as much.

## C Additional Implementation Details for Outdoor Scenario

In this section, we provide the additional implementation details for crafting adversarial perturbations for Monodepth [Godard et al., 2017] and Monodepth2 [Godard et al., 2019] on the KITTI dataset (outdoor driving scenario) as discussed in the main text.

### C.1 Hyper-parameters

Upper Norm	$\xi = 2 \times 10^{-3}$	$\xi = 5 \times 10^{-3}$	$\xi = 1 \times 10^{-2}$	$\xi = 2 \times 10^{-2}$
Monodepth	$\eta = 1.0$	$\eta = 2.0$	$\eta = 3.0$	$\eta = 4.0$
Monodepth2	$\eta = 0.1$	$\eta = 1.0$	$\eta = 3.0$	$\eta = 5.0$

Table 1: **Learning rates.** We achieve the best performances with the given learning rates.

Regarding hyper-parameters for crafting *adversarial perturbations*: We search the learning rate for each noise norm from the set  $\{0.1, 1.0, 2.0, 3.0, 4.0, 5.0, 10.0\}$ . We report the best performing ones in Table 1. Regarding our choice for the number of steps to run, we experimented with 200, 400, 500, 800, and 1000 steps and found little difference in performance measured by ARE between 500 steps, and 800 and 1000 steps. While an increase number of steps will obtain slight performance improvements, conscious of the time complexity, we chose 500 for our experiments.

Regarding hyper-parameters for *adversarial training*: As a defense against adversarial perturbations, we optimized Eqn. 6 for Monodepth2 using Adam [Kingma and Ba, 2014] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . We used a batch size of 4 and starting learning rate of  $1 \times 10^{-5}$ . We decreased the learning rate to  $5 \times 10^{-6}$  after 10 epochs and to  $2.5 \times 10^{-6}$  after 20 epochs and  $1 \times 10^{-6}$  after 30 epochs for a total of 40 epochs. Training takes approximately 4 hours using an Nvidia GeForce GTX 1080 GPU.

## C.2 Monodepth and Monodepth2

We study the effects of adversarial perturbations on the state-of-the-art monocular depth prediction method, Monodepth2 [Godard et al., 2019] and its predecessor Monodepth [Godard et al., 2017]. The two models utilize different network architectures and trained with different loss functions. In this section, we provide details on the two methods.

Regarding *Monodepth*: Monodepth uses a ResNet50 encoder architecture as its backbone and a standard decoder with skip connections. Monodepth predicts both left and right disparities from a single image (assuming it is the left image of a stereo-pair) and uses image reconstruction as supervision. Additionally, it is trained with a standard local smoothness term weighted by image gradients and a left-right disparity consistency term as its regularizers.

Regarding *Monodepth2*: Monodepth2, unlike Monodepth, uses ResNet18 encoder (pretrained on ImageNet) as its backbone network architecture. Rather than simply minimizing an image reconstruction loss, Monodepth2 leverages a heuristic to discount occluded pixels and also uses a criterion to discount static frames. Similar to Monodepth, Monodepth2 also minimizes a local smoothness regularizer weighted by image gradients.

## D Scaling with Larger Factors

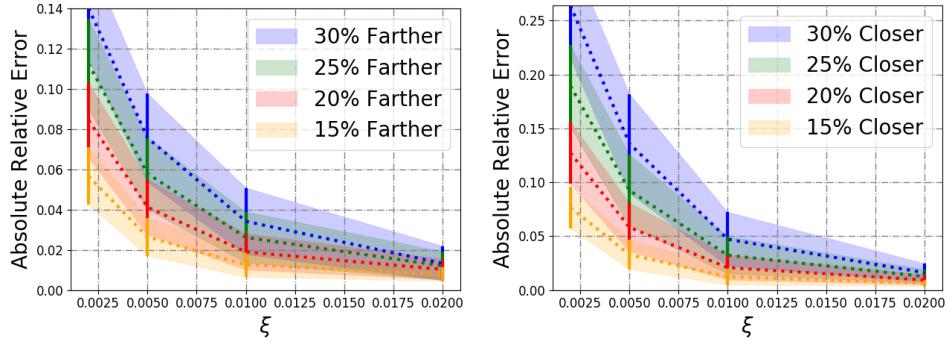


Figure 12: **ARE with various upper norm  $\xi$  for scaling Monodepth2 predictions by larger factors.** We increased the scaling to 15%, 20%, 25% and 30% closer and farther. We can see that for  $\xi = 2 \times 10^{-2}$ , perturbations are still able to scale the entire scene by  $\approx 1\%$  error.

In Sec. 4.1 and Fig. 2-(a, b) of the main text, we showed that it is possible, even for small norms such as  $\xi = 2 \times 10^{-3}$  to scale the scene to be 5% or 10% closer or farther with small error. In this section, we demonstrate that it is possible to scale the scene by larger amounts (up to 30% closer or farther). Fig. 12 shows that with  $\xi = 2 \times 10^{-3}$ , it is only able to scale the up to 15% with reasonable error; whereas perturbations with  $\xi = 5 \times 10^{-3}$  can achieve this up to 20% closer or farther. However, using larger norms ( $1 \times 10^{-2}$  and  $\xi = 2 \times 10^{-2}$ ), one can scale the scene up to 30% with small errors (less than 5% ARE for  $\xi = 1 \times 10^{-2}$  and  $\approx 1\%$  ARE for  $\xi = 2 \times 10^{-2}$ ).

To see how far we can push for each upper norm, Fig. 17 shows various scales that each upper norm is capable of achieving. We note that  $\xi = 2 \times 10^{-2}$ , is still able to obtain less than 2% ARE when scaling the scene by 45%; however, standard deviation starts to grow larger as the scaling increases.

## E Adversarial Attacks for Indoor Scenes

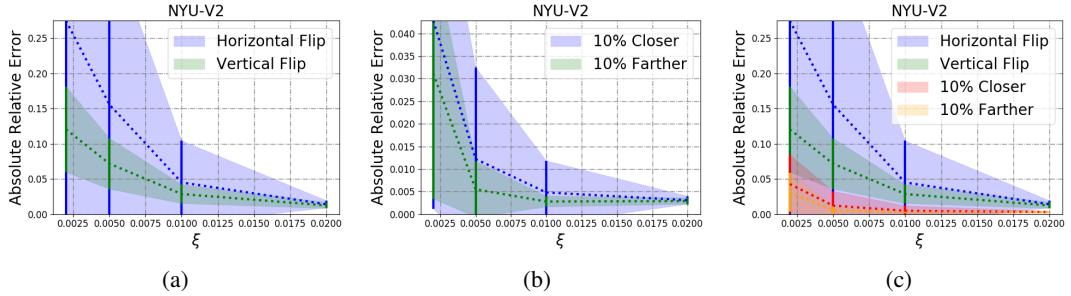


Figure 13: **Indoor quantitative results.** ARE with various upper norm  $\xi$  for scaling and flipping VNL predictions. (a) Results for horizontally and vertically flipping the predictions. (b) Results for scaling the scene by  $\pm 10\%$ . (c) comparison between scaling and flipping tasks.

To show the applicability of the adversarial method on indoor scenes, we examine the adversarial perturbations for Kinect Dataset NYU Depth V2 (NYU-V2) [Silberman et al., 2012]. We tested the effectiveness of adversarial perturbations on Virtual Normal Loss (VNL) [Yin et al., 2019] which is the state-of-the-art monocular depth prediction method for NYU-V2, trained in the supervised setting.

### E.1 Implementation Details

NYU-V2 consists of 1449 RGBD images gathered from a wide range of buildings, comprising 464 different indoor scenes across 26 scene classes. The images were hand-selected from 435,103 video frames, to ensure diversity. 1449 labeled samples are split into 795 training and 654 test images.

In the method proposed by [Yin et al., 2019], a 3D point cloud is reconstructed from the estimated depth map. Then, three non-collinear points are randomly sampled with large distances to form a virtual plane. The deviation between ground truth and prediction for the direction of the normal vector corresponding to the plane is penalized. The pre-trained ResNeXt-101 [Xie et al., 2017b] model on ImageNet [Deng et al., 2009] is used as the backbone architecture. During training, images are cropped to the size  $384 \times 384$  for NYU-V2. We use the same image resolution for our experiments. The training set is randomly sampled from 29K images of the raw unlabeled training set.

The time it takes to forward an image with this method is  $\approx 0.15$  seconds ( $\approx 7$  times more than Monodepth2). Due to computational limitations, we choose the first 20 images out of 654 images of the test split for our experiments. We run SGD for 1000 steps. The learning rate is kept at 10.0 for the entire optimization.

### E.2 Scaling and Symmetrically Flipping the Scene

In Fig. 13, we compare the performance for different target depth maps: scaling the scene by  $\pm 10\%$ , horizontal and vertical flipping. Unlike outdoor case (KITTI), in the indoor (NYU-V2) horizontal flipping is a harder task than vertical flipping. Achieving a vertically flipped scene was expected to be simpler as layouts of indoor scenes are more diverse. Hence, the depth network does not overfit to a particular layout type e.g. the one in which there are large depth values only at the top of the image. Horizontal flipping being relatively harder for indoor scenes can be explained by the large divergence between the depth distributions of the original predictions and the target depths. The reason is that for the indoor scenes, most scenes are not symmetric in the horizontal direction, unlike the outdoor driving scenario where the left and right parts of the scene from the ego-view are usually symmetric.

Since [Yin et al., 2019] normalizes images with the deviation of the dataset which approximately scales the image by 5, it also effectively scales the noise with the same deviation. But, since the relative norm is still the same, we use the same norm values  $\xi \in \{2 \times 10^{-3}, 5 \times 10^{-3}, 1 \times 10^{-2}, 2 \times 10^{-2}\}$  when plotting the ARE in Fig. 13.

In Fig. 14, we present qualitative results for NYU-V2 for  $\xi = 2 \times 10^{-2}$ . Small, white borders around RGB images exist in the raw dataset. For all the tasks, including vertical flip, adversarial perturbations manage to fool the model to predict the target depth with small errors. For horizontally flipped target depth (a), predictions have more artifacts than vertical flipped depth (b) and scaled depths (c,d).

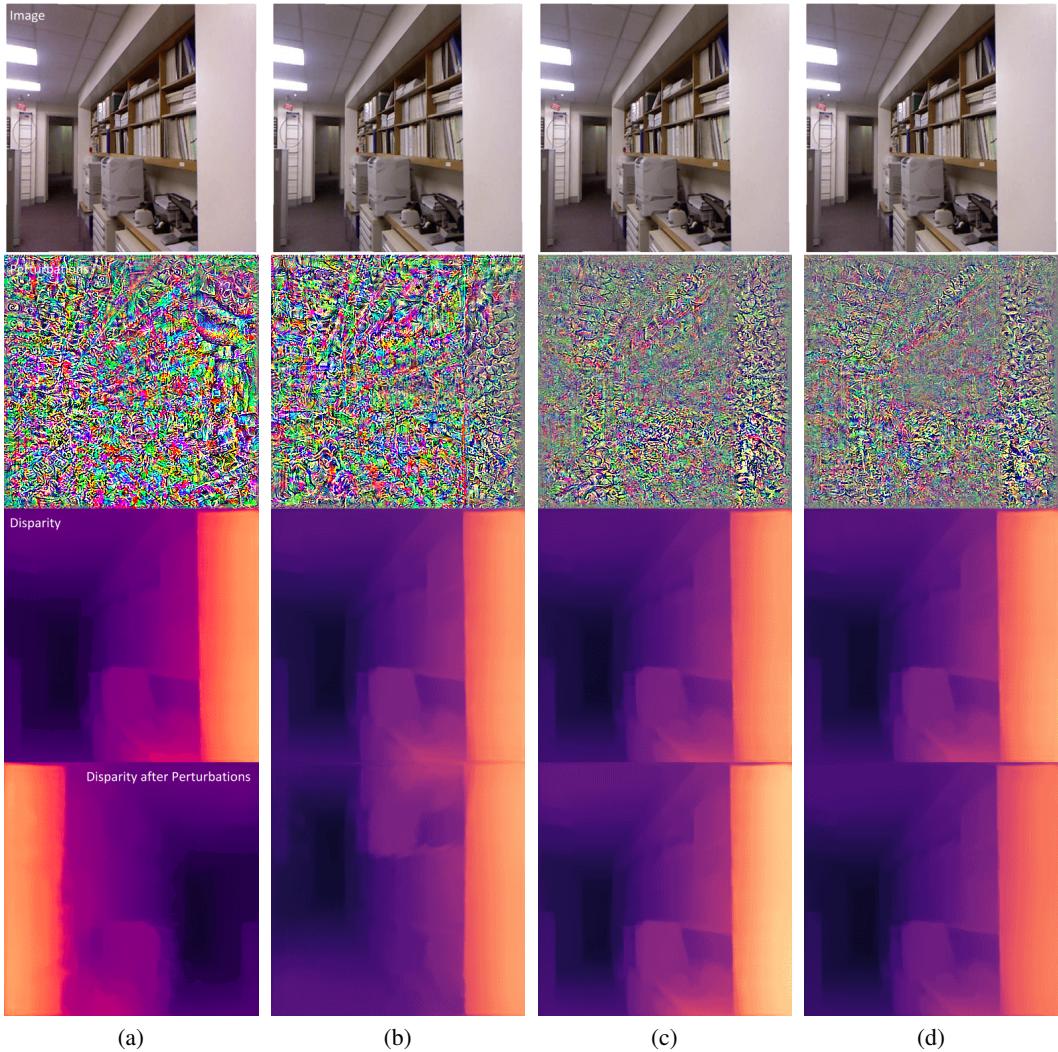


Figure 14: **Indoor qualitative results.** From left to right: (a) horizontal flip, (b) vertical flip, (c) scale 10% closer and (d) scale 10% farther. From top to bottom: RGB, noise, original disparity, disparity prediction for the perturbed image.

## F Linear Operations

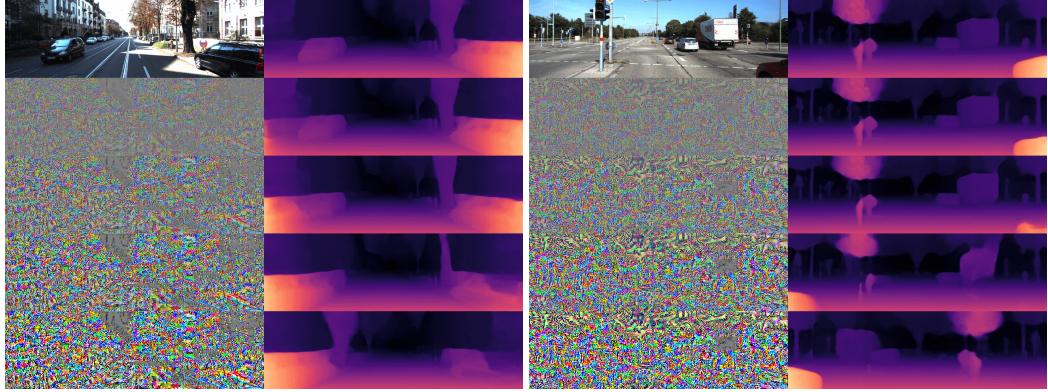


Figure 15: Disparity for  $x + \gamma v(x)$  where  $\gamma$  is 0.0, 0.25, 0.5, 0.75, 1.0 from top to bottom.  $v(x)$  is calculated for  $d^t = \text{fliph}(f_d(x))$ . So, the top is the original disparity map while bottom most is the flipped one. In between, portions of the scene are flipped smoothly.

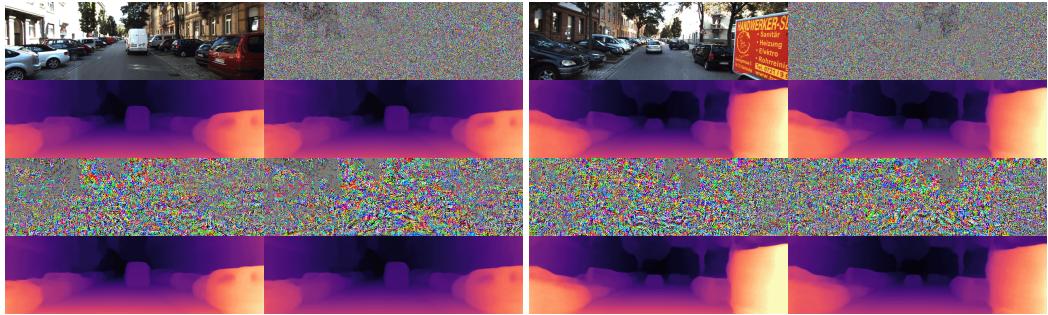


Figure 16: (1st row) left to right: RGB, sum of noises. (2nd row) left to right: original disparity, disparity when two noises  $v_1(x) + v_2(x)$  are added to the image. (3rd row) left to right: noise for 10% closer, noise for 10% farther. (4th row): Disparity predictions for the images perturbed with the noises in the 3rd row. When added, two noises cancel each other’s effect: the scale for  $f_d(x + v_1(x) + v_2(x))$  is close to the original one  $f_d(x)$ .

To better understand how predictions of the depth network changes within a ball of small radius, we examine the effect of linear operations on perturbations. Specifically, we visualize the predictions for the scaled perturbations and for the perturbations which we get after summing two perturbations calculated for two different target depth maps.

In Fig. 15, we take the perturbation  $v(x)$ , which we calculated to horizontally flip the prediction for the given image, and we visualize the prediction of the network for  $x + \gamma v(x)$  where  $\gamma \in \{0.0, 0.25, 0.5, 0.75, 1.0\}$ . As can be seen, between  $\gamma = 0$  and  $\gamma = 1$ , the scene is smoothly flipped. This implies that the adversarial perturbations can be used to control the depth prediction in a disentangled way. In other words, one causal factor (e.g. horizontal orientation) of the prediction can be independently controlled by tweaking  $\gamma$  only, keeping everything else the same.

As observed before, noise is small for the white regions. See the third column, where there is a gray rectangle in the noise corresponding to the white region of the truck. We speculate the reason behind this phenomenon as the white color being on the border of the support of RGB images. But, the noise is still large for black regions which are at the other extreme of the support (see perturbations corresponding to black vehicles). So, we left further understanding of this phenomenon as future work.

In Fig. 16, we take  $v_1(x)$  and  $v_2(x)$  which are optimized to scale the scene to 10% closer and 10% farther. Then, visualize the summed perturbation,  $v(x) = v_1(x) + v_2(x)$  and the prediction for  $x + v(x)$ . As can be seen, two noises cancel each other:  $\|v_1(x)\| \approx \|v_2(x)\| \gg \|v_1(x) + v_2(x)\|$ . Furthermore, the prediction for the image perturbed with the summed noise is close to the original prediction:  $f_d(x) \approx f_d(x + v_1(x) + v_2(x))$ . This shows that two perturbations with inverse functionalities can neutralize their effects when applied together.

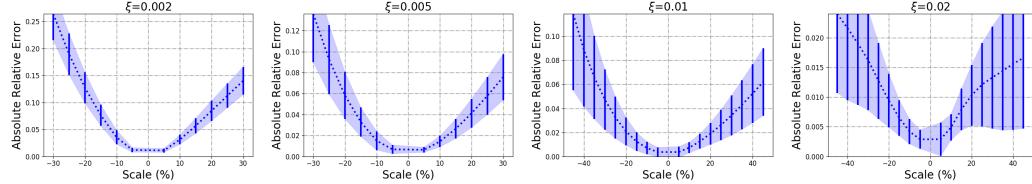


Figure 17: **ARE with various upper norm  $\xi$  for scaling Monodepth2 predictions.** This time error is plotted for large scale ratios 15%, 20%, 25% and 30% (up to 45% for larger  $\xi$ ), for scaling both closer and farther, showing the limitations of each norm for the scaling task. While ARE is still relatively small for larger norms, standard deviation grows larger – meaning the perturbations can no longer scale the scene consistently with low error.

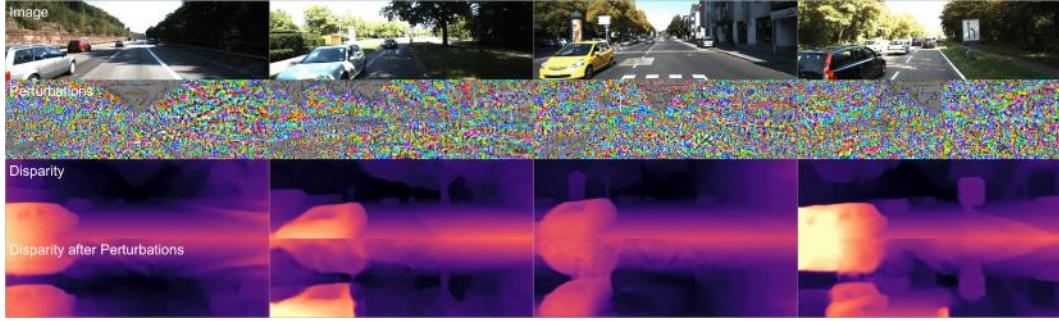


Figure 18: *Additional examples of failure cases for vertical flip.* While adversarial perturbations with  $\xi = 2 \times 10^{-2}$ , can fool Monodepth2 to predicted scenes that are scaled by large amounts, and horizontally flipped. They cannot cause the network to vertically flip the scene, leaving behind cars and roads as artifacts.

## G Limitations

Fig. 17 shows the absolute relative error (ARE) with respect to the target scaling factor for each upper norm. As we can see, for smaller norms of  $2 \times 10^{-3}$  and  $5 \times 10^{-3}$ , the perturbations are limited to scaling the scene by  $\approx 15\%$  and  $\approx 20\%$ , respectively. Scaling factors higher than such increases the ARE by  $\approx 4\%$  for every 5% increase in scaling factor, signaling the limit for these norms. For larger norms of  $1 \times 10^{-2}$  and  $2 \times 10^{-2}$ , the perturbations can afford to scale the scene by a much larger factor. For  $\xi = 1 \times 10^{-2}$ , perturbations can scale the scene by as much as 30% closer and farther less than 5% error. Whereas, for  $\xi = 1 \times 10^{-2}$ , perturbations can scale the scene up to  $\pm 45\%$  with less than 2% ARE. However, while large scaling still has a low ARE, the standard deviation for larger norms increases drastically showing that it can no longer consistently scale the scene.

While for smaller scales (e.g.  $\pm 5, \pm 10$ ) the ARE and amount of noise required is approximately the same (see Sec. 4.1, main text), suggesting similar difficulty levels. As we plot the errors for larger scales in Fig. 17, scaling the scene farther generally yields lower error than scaling the scene closer.

Fig. 18 shows additional examples of failure cases for vertical flip. While we have shown in the main paper as well as Sec. D and H that it is possible to manipulate the scene with small norm perturbations, we show here that perturbations cannot fool a network into vertically flipping the scene.

## H Additional Results on Outdoor Scenarios

In this section, we show (i) side by side visualization of the perturbations required to scale the scene, (ii) additional visualizations of perturbations to horizontally and vertically flip the scene, (iii) quantitative results on targeted attacks to semantic categories and (iv) qualitative results on targeted attacks to instances.

### H.1 Scaling the Scene



Figure 19: Visually imperceptible perturbations  $v(x)$ , with  $\xi = 2 \times 10^{-2}$ , can fool Monodepth2 to predicted scenes that are 5% or 10% closer and also 5% or 10% farther.

Here, we show qualitative results for the task of scaling the scene (Sec. 4.1, main text) by a factor of  $1 + \alpha$  where  $\alpha \in \{-0.10, -0.05, +0.05, +0.10\}$ . As seen in Fig. 19, the perturbations are successful in fooling state-of-the-art monocular depth prediction method, Monodepth2 [Godard et al., 2019], into predicting the scene 5% or 10% closer and also 5% or 10% farther. Additionally, the perturbations are concentrated in similar regions for scaling the scene 5% or 10% closer and for 5% or 10% farther as well. As noted in the main text, the amount of noise required for scaling the scene by  $\pm 5\%$  are approximately the same, as is the amount for scaling the scene by  $\pm 10\%$ . This is visible in Fig. 19.

### H.2 Symmetrically Flipping the Scene

In Sec. 4.2 and Fig. 3 in the main text, we demonstrated that adversarial perturbation can cause a monocular depth prediction network to predict a horizontally or vertically flipped scene. Here, we show additional qualitative results on the horizontal and vertical flipping tasks in Fig. 20 and Fig. 21. We note that perturbations can cause the network to predict a horizontally flip scene, they have trouble fooling the network to predict a vertically flipped scene. This is unlike our findings in the indoor scenario (Sec. E) as seen in Fig. 14 and 13. Fooling the network to vertically flip the scene is in fact *easier* than fooling it to horizontally flip the scene. This confirms

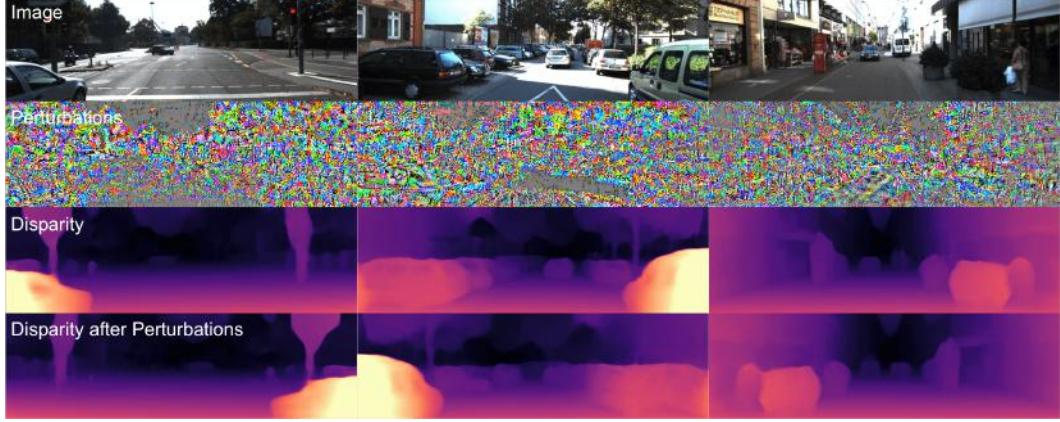


Figure 20: *Additional examples of horizontal flip.* Adversarial perturbations with  $\xi = 2 \times 10^{-2}$ , can fool Monodepth2 to predicted scenes that horizontally flipped.

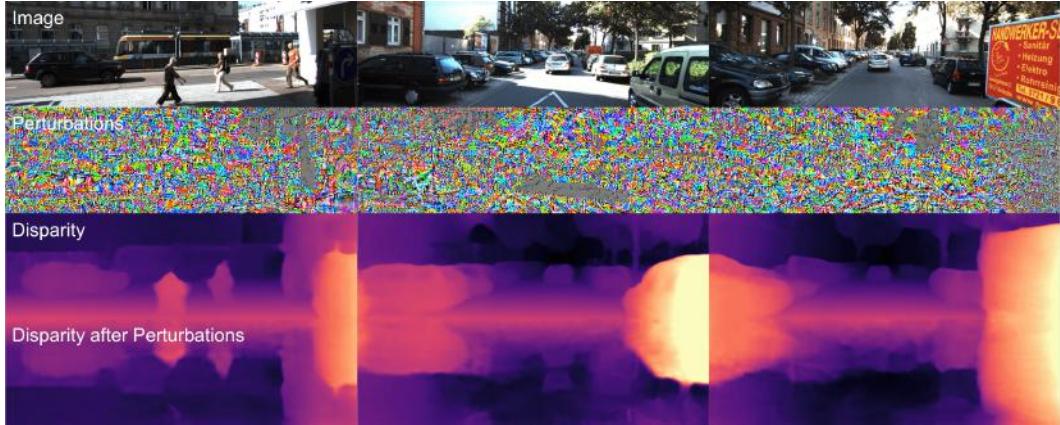


Figure 21: *Additional examples of vertical flip.* Adversarial perturbations with  $\xi = 2 \times 10^{-2}$ , can fool Monodepth2 to predicted scenes that vertically flipped. Even in these successful examples, there are still artifacts (ripples, wavy-ness) in the output.

the biases (roads on bottom, sky on top) that the network learned from the outdoor dataset that are not present in the indoor dataset.

### H.3 Category Conditioned Scaling

In Sec. 5.1 and Fig. 5 of the main text, we showed category specific attacks to scale all objects belonging a given category to be a factor of  $1 + \alpha$  closer or farther where  $\alpha \in \{-0.10, -0.05, +0.05, +0.10\}$ . Here, we provide performance, measured in ARE, of adversarial perturbations crafted for each category. We use the same convention for grouping different classes into categories as the Cityscapes dataset [Cordts et al., 2016], with the exception of the ‘‘Human’’ category, which includes the bicycles that the bikers are riding.

Fig. 22 shows a comparative study between different categories. Not all categories are equally easy to be fooled by the perturbations, some are more robust to adversarial attacks than others. As seen in Fig. 22, each category exhibits a different level of robustness to adversarial noise – ‘‘Human’’ and ‘‘Traffic’’ categories are the hardest to fool, ‘‘Construction’’, ‘‘Vehicle’’ and ‘‘Flat’’ are more susceptible, and ‘‘Sky’’ and ‘‘Nature’’ are the easiest to attack. Plots are cropped at the maximum error across different categories to enable comparison of difficulty in fooling different categories. We note that attacking localized regions in the scene is considerably harder than attacking the entire scene. Fig. 12 shows that perturbations can attack the entire scene with small errors across various norms while Fig. 22 shows that, even with large norms, there are still errors ( $\approx 2\%$  to  $6\%$  ARE). We show visualizations for the ‘‘Construction’’, ‘‘Nature’’, and ‘‘Vehicle’’ categories in Fig. 23, 24, and 25 respectively.

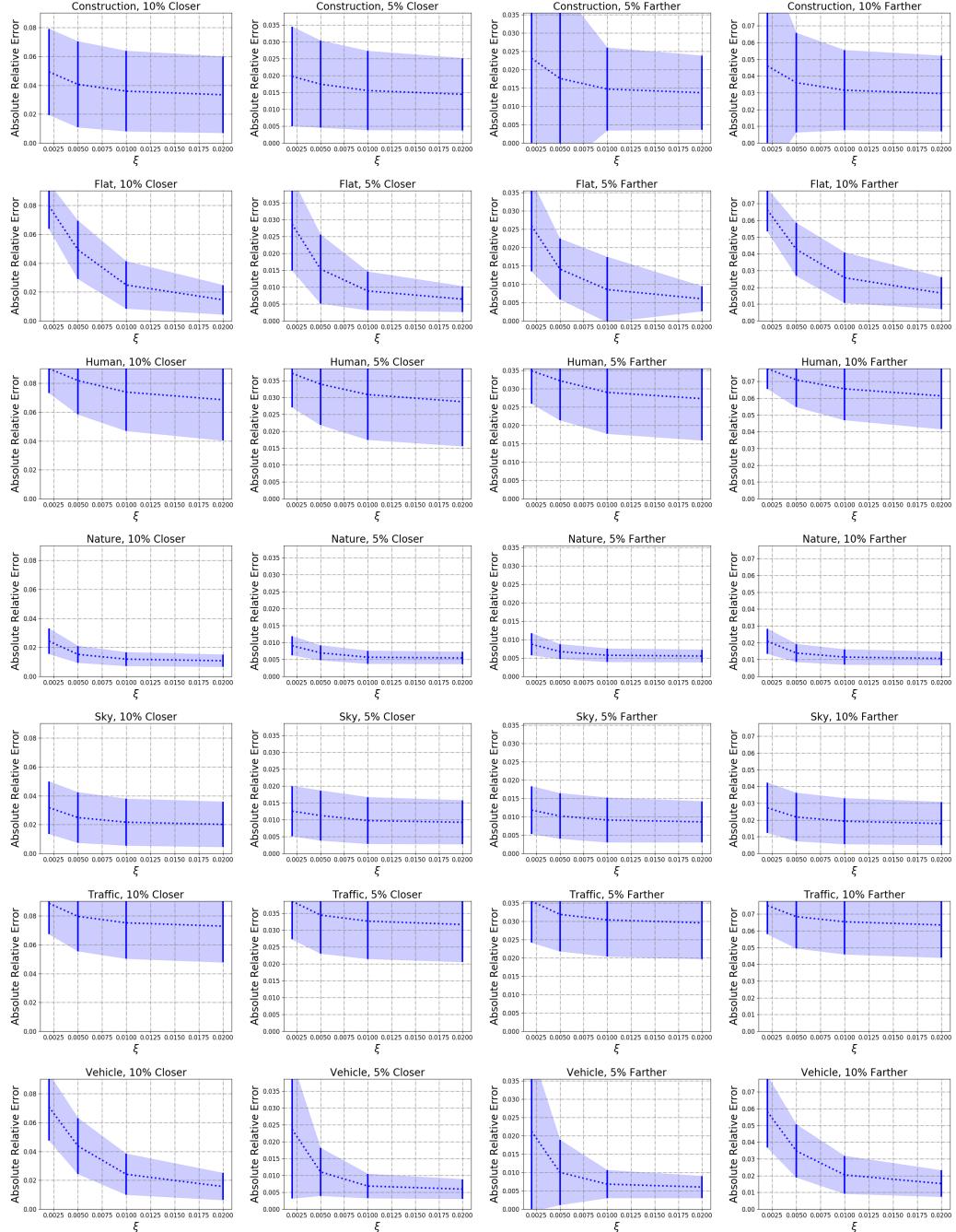


Figure 22: ARE for scaling different categories closer and farther. From top to bottom: “Construction”, “Flat”, “Human”, “Nature”, “Sky”, “Traffic”, “Vehicle”. From left to right: 10% closer, 5% closer, 5% farther, 10% farther. Y-axis is kept the same for the same scale, for making the comparison across categories possible. It is easier to fool the network to predict vehicle and nature categories closer and farther than is to fool human and traffic categories.

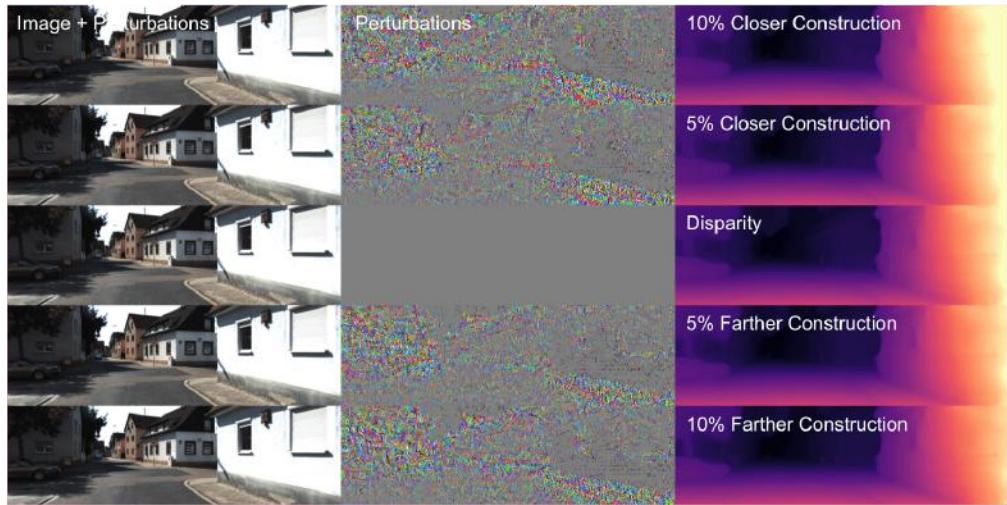


Figure 23: Examples of targeted attacks on regions belonging to “Construction” category.

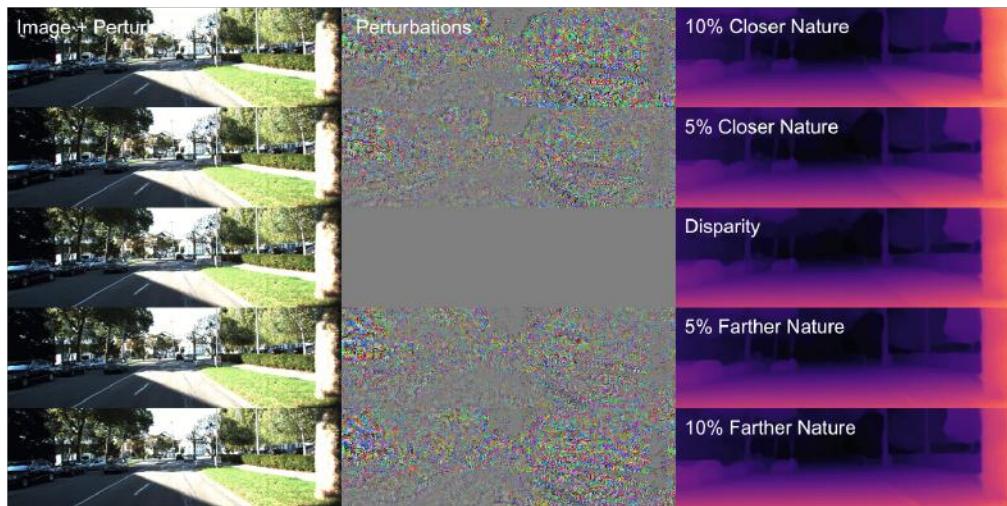


Figure 24: Examples of targeted attacks on regions belonging to “Nature” category.

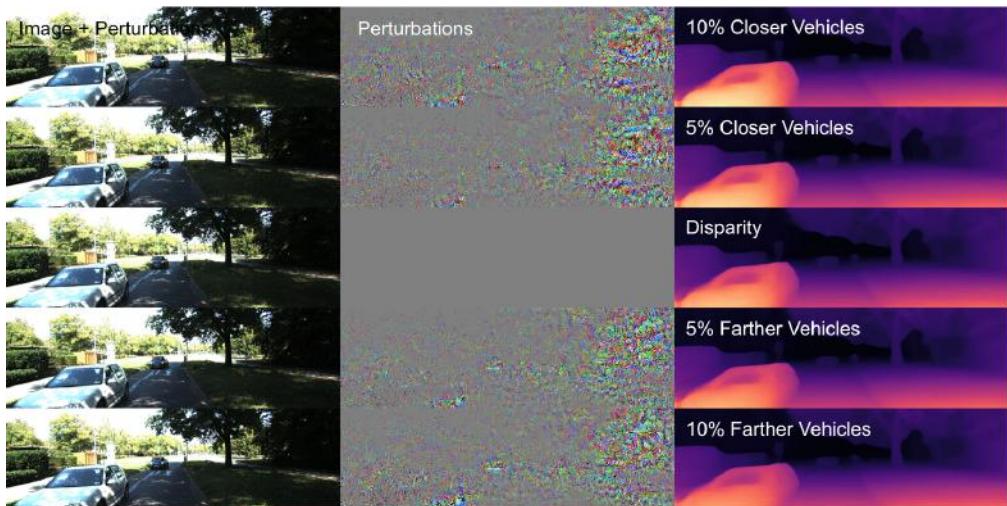


Figure 25: Examples of targeted attacks on regions belonging to “Vehicle” category.

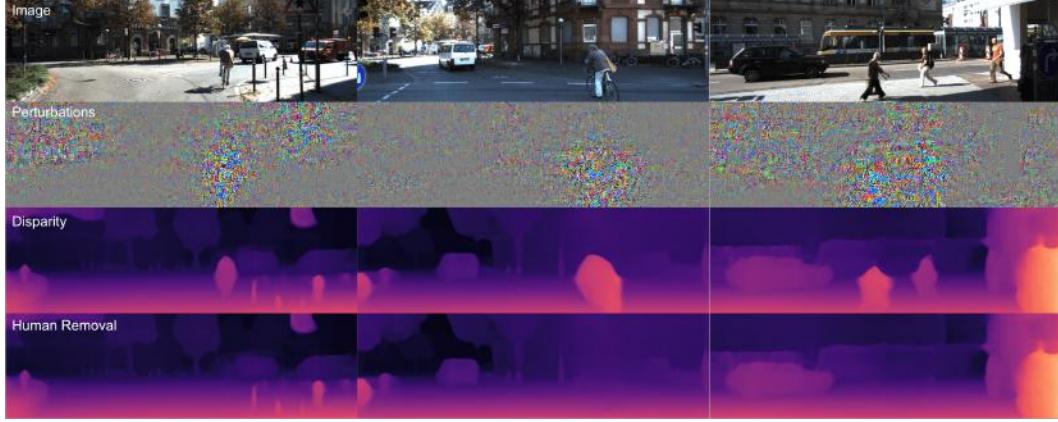


Figure 26: **Additional examples of human removal.** Targeted adversarial perturbations can remove humans from the predicted scene. Rightmost panel shows that we can target multiple humans and remove them from the scene without affecting the remaining pedestrian on the right.

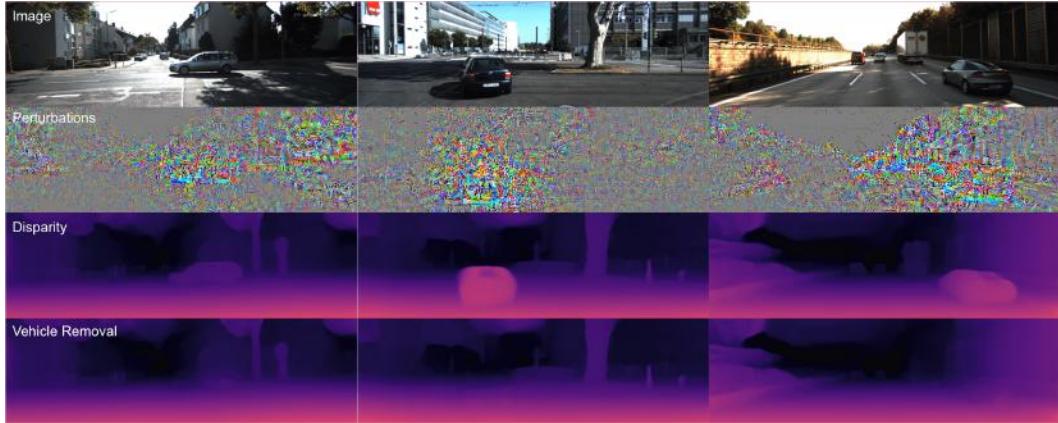


Figure 27: **Examples of vehicle removal.** Targeted adversarial perturbations can remove vehicles from the predicted scene. Rightmost panel shows that we can target a truck and a car on the right side and remove them. Note that the cars in the center still remain.

#### H.4 Instance Conditioned Targeted Attacks

In Sec. 5.2 and Fig. 6 in the main text, we show that, when given instance segmentation, adversarial perturbations can target specific instances and remove them from the scene and thus causing unforeseen consequences. Fig. 26 shows additional examples of removing humans from the scene and Fig. 27 demonstrates that it is possible to remove vehicles from the scene as well. In the rightmost panel of Fig. 26, we show that it is possible to remove *some* pedestrians from the scene without affecting others. Similarly, in the rightmost panel of Fig. 27, we removed a truck and a car on the right side and left the cars in the center untouched – leaving this as *still* a plausible highway driving scenario.

In Sec. 5.4 and Fig. 7 in the main text, we show that perturbations can move an instance to another location in the scene (requires removing the instance from its original location and creating it in the new location). In this section, we give more visuals for the perturbations used for moving an instance (e.g. vehicle, pedestrian) horizontally or vertically in the image space while keeping the rest of the scene unchanged.

Fig. 28 shows that perturbations can fool a network to move the target instance by  $\approx 8\%$  across the image in the left and right directions. Furthermore, Fig. 29 shows that perturbations can move select instances by  $\approx 42\%$  in the upward direction, creating the illusion that there are “flying vehicles” in the scene. We note that in both cases, the perturbations are concentrated on the instance and the region to which the instance is moved. For example, when moving a vehicle right or left, the corresponding perturbations also move right or left.

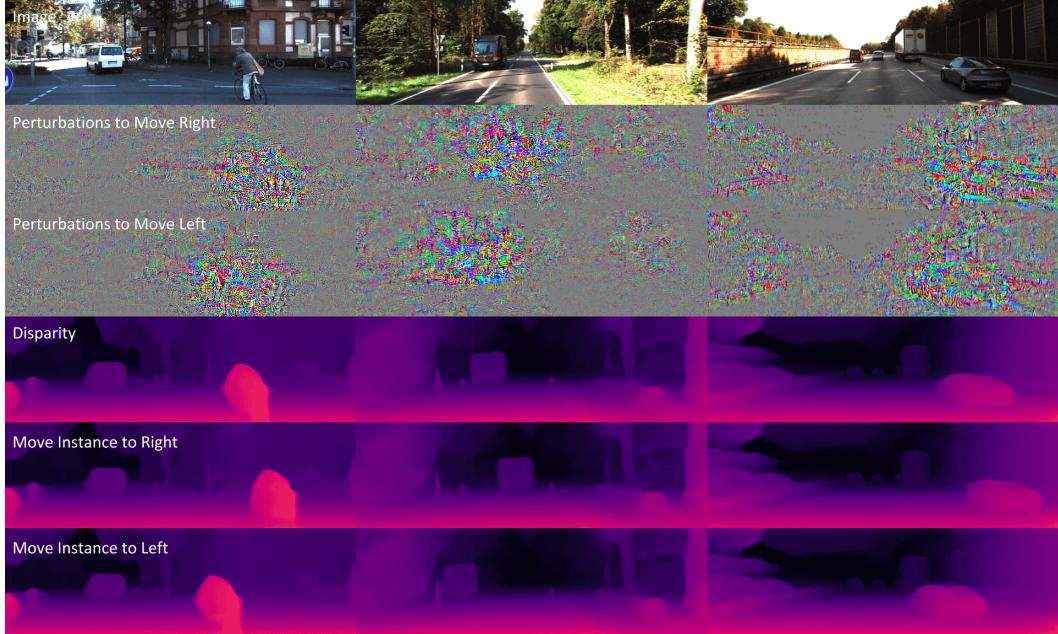


Figure 28: **Move instance horizontally.** Selected instance is moved by  $\approx 8\%$  in left and right directions while rest of the scene is preserved. Noise and disparity for both directions are given. We note that the noise is concentrated around the targeted instance and the region to which the instance is moved.

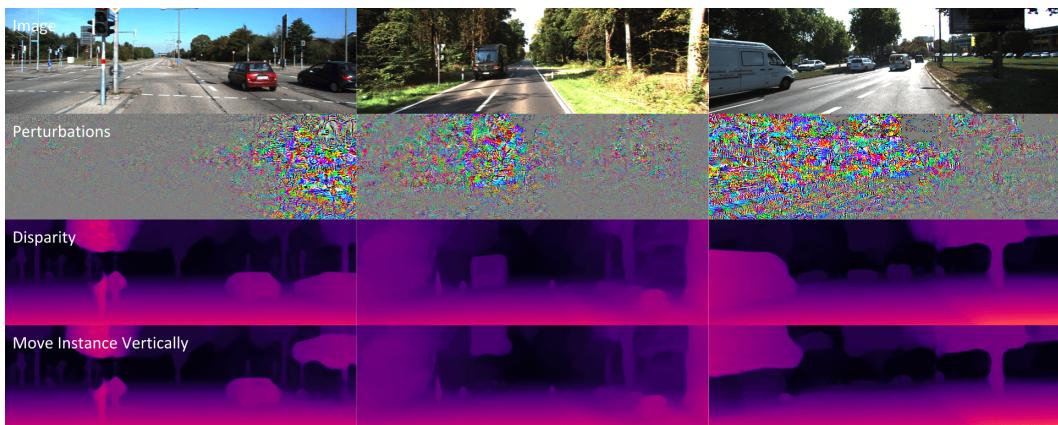


Figure 29: **Flying vehicles.** Selected vehicle is moved by  $\approx 42\%$  in the vertical direction while rest of the scene is preserved. We note that the noise is generally concentrated around the targeted instance and the region to which the instance is moved.