

# MULTI-TASK CONSISTENCY-BASED DETECTION OF ADVERSARIAL ATTACKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep Neural Networks (DNNs) have found successful deployment in numerous vision perception systems. However, their susceptibility to adversarial attacks has prompted concerns regarding their practical applications, specifically in the context of autonomous driving. Existing research on defenses often suffers from cost inefficiency, rendering their deployment impractical for resource-constrained applications. In this work, we propose an efficient and effective adversarial attack detection scheme leveraging the multi-task perception within a complex vision system. Adversarial perturbations are detected by the inconsistencies between the inference outputs of multiple vision tasks, e.g., objection detection and instance segmentation. To this end, we developed a consistency score metric to measure the inconsistency between vision tasks. Next, we designed an approach to select the best model pairs for detecting this inconsistency effectively. Finally, we evaluated our defense by implementing PGD attacks across multiple vision models on the BDD100k validation dataset. The experimental results demonstrated that our defense achieved a ROC-AUC performance of 99.9% detection within the considered attacker model.

## 1 INTRODUCTION

The camera-based perception system is critical to enable automated driving (AD). Indeed, camera is the only sensor able to read traffic signs, identify lane markings or drivable areas, and see traffic light colors. To perform such perception tasks (e.g., object detection, classification, segmentation), a wide range of machine learning models were developed, each with its own objective and network architecture (Zou et al., 2023). For example, from an input image, 2D object detection models output bounding boxes, while semantic segmentation models output masks, or multi-object tracking models output track identifiers. The model outputs help to understand the scene and allow the automated vehicle to maneuver appropriately.

However, camera inputs can be maliciously manipulated to affect the performance of perception tasks, or even downstream tasks of automated vehicles (e.g., path planning, motion control). The idea of adversarial inputs (commonly called *adversarial examples*) is to add specially-crafted noise to images such that the underlying machine learning models do not perform as originally intended (Madry et al., 2018). Adversarial examples have been demonstrated in the form of full image perturbations or patches, realized digitally or physically, and with some high attack success rate and universality (Chow et al., 2020). Because of their low level of sophistication and effectiveness, it is key to deploy defenses to protect automated vehicles against such threats. Defenses range from preemptive techniques (e.g., adversarial training (Shafahi et al., 2019), certified robustness (Xiang et al., 2024)) to reactive techniques (e.g., real-time detection of perturbations (Xiang et al., 2022), image compression (Das et al., 2018)).

In this paper, we focus on reactive techniques, aiming at real-time detection of perturbations, because it does not require any adversarial data generation or additional training. Especially, we propose to leverage the output of multiple perception tasks to identify perturbations on every image prior to use by downstream tasks. Prior work showed the effectiveness of checking inconsistencies of edge extractions between outputs of semantic segmentation and depth estimation (Klingner et al., 2022), but with some limitations. Their inconsistency check only detects adversarial perturbations on the entire image and might show limited performance on local perturbations. Therefore, we propose

a consistency-based detection technique that is effective regardless of the perturbation’s location. As long as the perturbation causes inconsistent inference output across models, locally or globally, our defense can capture the inconsistency. Especially, we demonstrate the benefits of cross-model consistency by using 2D object detection and instance segmentation models. Indeed, 2D object detection models are commonly used in AD to detect obstacles or traffic signs, and then to convert 2D bounding boxes to 3D bounding boxes (Feng et al., 2020; Arnold et al., 2019). Instance segmentation is also used in AD to provide finer object boundaries (Zhou et al., 2020). Both model share the objective of detecting objects, and hence, can be used to identify inconsistencies.

Our **contributions** are as follows:

- We propose a lightweight consistency detector based on outputs from object detection and instance segmentation models.
- We develop a technique to select the optimal model pair, deriving requirements w.r.t model architecture.
- We define a metric to capture the consistency score between two models’ output.
- We generate and publish an adversarial BDD100k dataset to assess the effectiveness of our defense, and allow reproducibility and comparison of future defenses.

## 2 SYSTEM MODEL

### 2.1 VISION MULTI-TASK SYSTEM

Perception systems perform multiple vision tasks such as object detection, segmentation, and depth estimation. Because of its better generalization performance and efficiency (Guo et al., 2020b), one architecture considered for automated driving is Multi-Task Learning (MTL) (Miraliev et al., 2023). A common approach in MTL is to have a shared feature extractor and multiple task-specific heads (Caruana, 1997; Kokkinos, 2017; Lu et al., 2017). In this paper, because our detection method must work with MTL and non-MTL architecture, our architecture consists of one model per task. With this flexible approach, we can evaluate the performance of our detector when the tasks share (or not) the same backbone. Indeed, the attack success rate strongly correlates with the architecture similarity between tasks as highlighted by Xie et al. (2017). Interestingly, from a security perspective, it may be more robust to have an architecture with different backbone per task than a common backbone architecture for all tasks (like in the MTL architecture).

### 2.2 ATTACKER MODEL

We follow the same attacker model as defined by Xiang et al. (2022), where the attacker performs a white-box attack (i.e., has access to the model’s architecture and weights). We assume a model  $\mathbb{F}$  with an underlying data distribution  $\mathcal{D}$  over pairs consisting of image  $\mathbf{x}$  and its corresponding ground truth  $y$ .  $\mathcal{X}$  denotes the image space. The attacker adds the perturbation  $\delta$  to the genuine image  $\mathbf{x}$  to create an adversarial image ( $\mathbf{x}' = \mathbf{x} + \delta$ ) (with  $\|\delta\|_p \leq \epsilon$ , where  $\epsilon$  is the bound on the  $L_p$  norm perturbation) such as  $\mathbf{x}' \in \mathcal{A}(\mathbf{x}) \subset \mathcal{X}$ , where constraint  $\mathcal{A}$  defines the attacker’s capability. The goal of the attacker is to minimize the alteration of the genuine image  $\mathbf{x}$  while ensuring the attack succeed, and is formulated as:

$$\min \|\mathbf{x}' - \mathbf{x}\| \text{ s.t. } \mathbb{F}(\mathbf{x}') \neq \mathbb{F}(\mathbf{x}) \quad (1)$$

where,  $\mathbb{F}(\mathbf{x}') \neq \mathbb{F}(\mathbf{x})$  can be the removal or injection of bounding boxes/masks.

To achieve her goal, the attacker uses a projected gradient descent (PGD) attack (Madry et al., 2017).

$$\mathbf{x}'_{t+1} = \Pi_{\mathbf{x}+\mathcal{X}}(\mathbf{x}'_t + \alpha \text{sgn}(\nabla_{\mathbf{x}} L(\theta, \mathbf{x}', y))) \quad (2)$$

where,  $L(\theta, \mathbf{x}', y)$  is the global loss function defined as the sum of classification loss and localization loss ( $L = L_{cls} + L_{loc}$ ). Hence, the perturbation targets a misclassification or mislocalization.

## 3 MULTI-TASK CONSISTENCY

We first define the multi-task consistency score between model outputs across different vision tasks. In particular, we use object detection (OD) and instance segmentation (SEG) as example vision tasks in this paper. Then, we explain how to use the consistency score to detect adversarial perturbations.

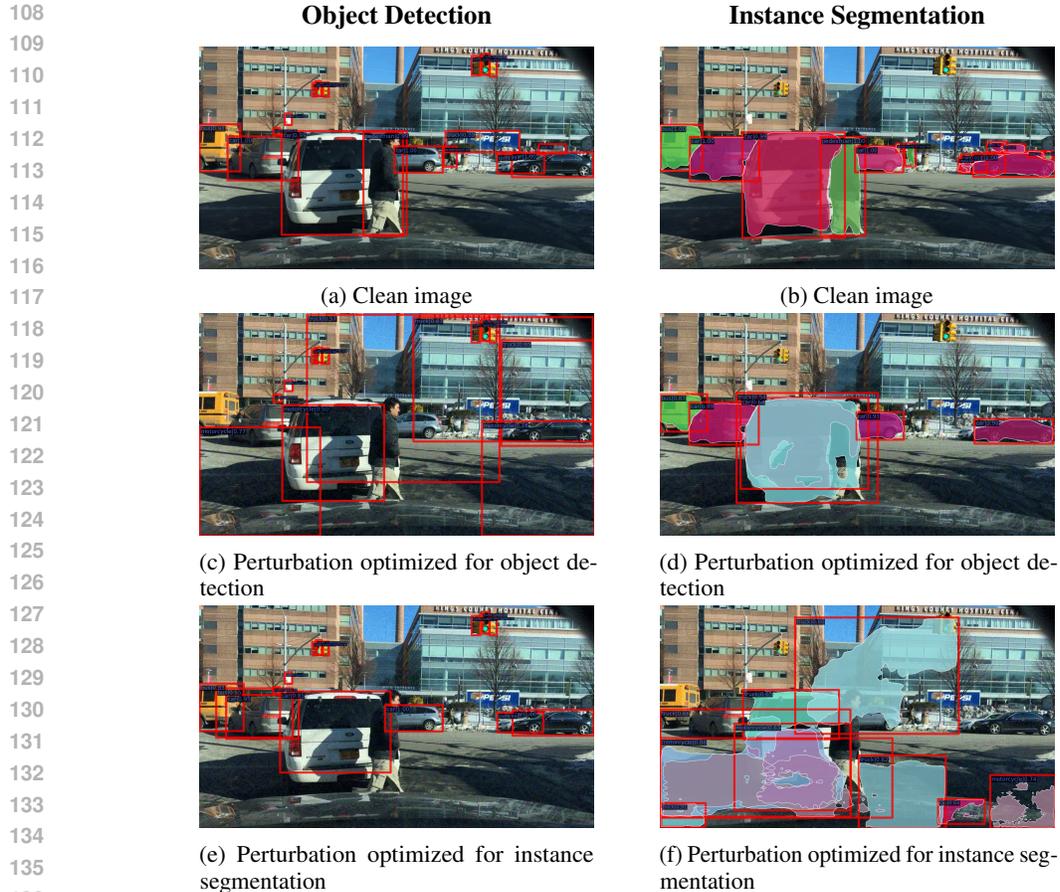


Figure 1: Impact of adversarial perturbation on the vision models

### 140 3.1 CONSISTENCY BETWEEN VISION TASKS

141 As shown in Figure 1, the inference outputs for object detection and instance segmentation on  
 142 clean images exhibit overall consistency. Indeed, the object bounding boxes match with the object  
 143 masks. However, on the perturbed images, discrepancies arise. For instance, in Figures 1c-1d,  
 144 the perturbation optimized for the object detection model successfully deceived the object detector,  
 145 leading to numerous false positive and false negative predictions. On the other hand, the same  
 146 perturbation did not fool the instance segmentation model, which accurately predicted the bounding  
 147 boxes and masks<sup>1</sup>. Similar impact is observed in Figure 1e-1f where the perturbation is optimized for  
 148 instance segmentation. In fact, we can identify two types of consistency between the model outputs:

- 149 • **Location Consistency:** refers to detecting an object at the same location within an input image  
 150 using both an object detection model and an instance segmentation model. It involves calculating  
 151 the Intersection over Union (IoU) between each detected object from both models. If the IoU  
 152 exceeds a predefined threshold (e.g., 50%), the object pair is considered *location consistent*.
- 153 • **Semantic Consistency:** goes beyond location and ensures that the labels of the object pair are  
 154 identical as well. In this paper, we consider a detection as *consistent* if both location and semantic  
 155 consistency are proven.

157 **Consistency Score.** In this work, we call *consistent detection (CD)* a matching pair of box and  
 158 mask (location and label wise). In order to measure the overall consistency of a single image,  
 159 Equation 3 defines the *Consistency Score*  $C_{\text{task}}$  as the ratio of total number of consistent detection  
 160 over the total number of detection from either model ( $N_{\text{task}}$ ).

161 <sup>1</sup>We note a slight impact on the foreground objects' masks.

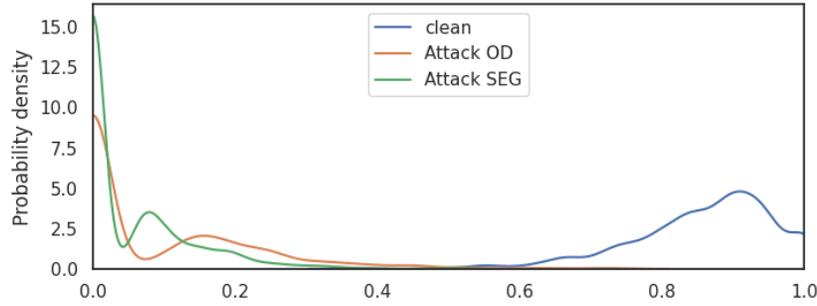


Figure 2: Empirical study of the consistency score distribution for (FRCNN R50, MRCNN R50) pair on BDD100k dataset. Blue line shows consistency scores for clean images. Orange line shows consistency scores for perturbed images (target OD). Green line shows consistency scores for perturbed images (target SEG). The clear divergence between distributions, confirms the ability of our detector to identify perturbations.

$$C_{\text{task}} = \frac{|CD|}{N_{\text{task}}} \quad \text{task} \in \{\text{det}, \text{seg}\} \quad (3)$$

Then, as in Equation 4, we define **consistency score**  $C$  as a harmonic mean of  $C_{\text{det}}$  and  $C_{\text{seg}}$  to measure the overall consistency of the inferences on input images by both models.

$$C = \frac{2 \cdot C_{\text{det}} \cdot C_{\text{seg}}}{C_{\text{det}} + C_{\text{seg}}} \quad (4)$$

**Empirical Study.** We performed an empirical study on the consistency score distribution for a clean image dataset and perturbed image datasets. The clean dataset consists of 1000 images from BDD100k validation dataset. The perturbed dataset is created by applying the PGD attack on each image in the clean dataset. Then, we calculated the consistency score of (FRCNN R50, MRCNN R50) model pair using Equation 4 on each image of the datasets and plot the distribution using kernel density estimation. From Figure 2, we observe that the model pair on clean images have higher consistency score, while perturbed images have much lower consistency score. This implies that we can distinguish between clean and perturbed images using the consistency score. We present other models distribution plots in Appendix A.4.

### 3.2 CONSISTENCY SCORE BASED ATTACK DETECTION

Inspired by the above observation, we propose a consistency score based adversarial attack detection scheme illustrated in Figure 3.

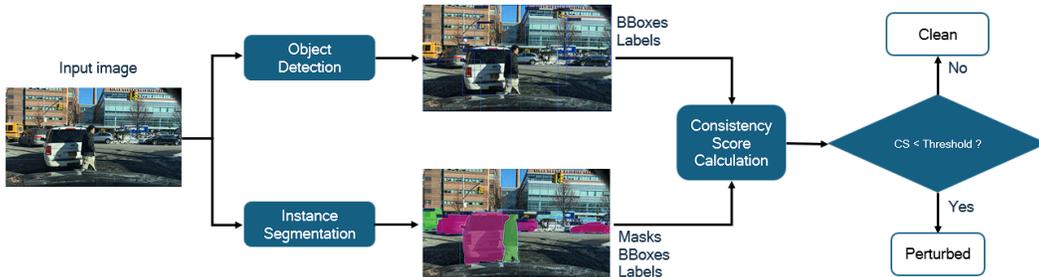


Figure 3: Pipeline of consistency score based adversarial perturbation detection

**Notations.** In order to formulate the problem, we denote the output of the object detection model as a set of annotations of detected objects  $S_{\text{det}} = \{(\text{BBox}_{\text{det},i}, \text{Label}_{\text{det},i}) | i = 1, \dots, K_{\text{det}}\}$  where  $\text{BBox}_{\text{det},i}$  is the bounding box coordinates for the  $i$ -th detection,  $\text{Label}_{\text{det},i}$  is its corresponding class label, and  $K_{\text{det}}$  is the total number of detection by the object detection model. Similarly, we denote the output of the instance segmentation model as  $S_{\text{seg}} = \{(\text{BBox}_{\text{seg},j}, \text{Label}_{\text{seg},j}) | j = 1, \dots, K_{\text{seg}}\}$ .

**Step 1: Consistency Score Calculation.** Following Equation 3, the consistency score is calculated between the two tasks output. In Appendix A.1, we propose Algorithm 1 as an implementation of the Consistency Score Calculation module of Figure 3.

**Step 2: Adversarial determination.** With the consistency score generated for the input image, the next step is to decide if it is a clean or perturbed image. As shown in Figure 3, a threshold-based binary classification takes the consistency score as input. If the consistency score is lower than the predefined threshold, the input is labelled as “perturbed”. As shown in Figure 2, setting a high cut-off threshold (e.g., 0.75) would trigger false positives. Conversely, selecting a low threshold (e.g., 0.2) would trigger false negatives. Therefore, there is a trade-off between false positive rate and false negative rate. Implementers would have to pick the appropriate threshold using known techniques (Lan et al., 2020).

**Cross-task model selection.** When designing a multi-task consistency detector, it is important to select the appropriate models used for each task. Indeed, the two models<sup>2</sup> could share the same backbone and underlying structure, or only share the same backbone, or share similar backbone but with different layer depth. We aim at answering the question “What model architectures or parameters affect the ability to detect adversarial inputs via multi-task consistency?”. For example, should the feature extractors be different? if so, to what extent? Ghamizi et al. (2022) hinted that one should carefully select the auxiliary tasks added to reduce model vulnerability. Indeed, the addition of auxiliary tasks can have negative effects (e.g., larger model size, slower convergence of the common encoder layers, deterioration of clean performance). They raised the (still open) question of how to select the combination that yields the lowest vulnerability. One could think that picking the most adversarially robust backbone would be preferable. For example, when investigating ResNet50 and ResNet101 backbones, the only difference is that ResNet101 has 23 conv4\_x layers while ResNet50 has 6 (so a total of 51 additional convolution layers as the name indicates). This means that ResNet101 has larger receptive fields than ResNet50. As shown by Xiang et al. (2024), smaller receptive fields impose a bound on the number of features that can be corrupted, hence more adversarially robust. This could justify the use of ResNet50 backbone over ResNet101. However, in our context, we select models that, even if fooled by the attack, yield to inconsistent outputs. So, having two weak models could be acceptable as long as their outputs are inconsistent.

## 4 EXPERIMENTS

In this section, we outline the implementation details of the datasets, models, attack parameters, and evaluation metrics used to ensure reproducibility. We then present and discuss the experimental results of the multi-task consistency-based detector. Additionally, we offer recommendations for a cross-model strategy to select the best model pairs for the detector.

### 4.1 IMPLEMENTATION DETAILS

**Datasets.** Our evaluation relies on a set of genuine and adversarial datasets based on the BDD100k dataset. The description of the BDD100k dataset can be found in the Appendix( A.2).

**Models.** We use 11 existing models from BDD100k model zoo (Huang, 2021) and from mmdetection 2.0 framework (Chen et al., 2019): six models for object detection (OD) and five models for instance segmentation (SEG). All models are fine-tuned on the BDD100k dataset. Our selection of models aims to maximize the diversity of models for a given vision task to understand how it affects the performance of our defense. Indeed, an adversarial attack may transfer from one model to

<sup>2</sup>For the sake of conciseness, we restrict ourselves to a model pair. However, the system can be extended to larger tuple (see Section 5).

another if their architecture are similar. Therefore, we chose our models based on a set of criteria. The first one is the type of architecture (e.g., transformer or CNN). A second criteria is the depth of the backbone (ResNet50 versus ResNet101). The last criteria is to ensure a diversity of heads among the models (e.g., FRCNN versus RetinaNet).

**Attack.** We utilized 1,000 clean images from the BDD100k instance segmentation validation dataset for our attack. This dataset was selected due to its comprehensive annotations, which include both segmentation masks and bounding boxes, allowing us to fairly assess the impact on both object detection (OD) and segmentation (SEG) models. We then applied the PGD-40 attack (40 iterations with a perturbation strength  $\epsilon = 16/255$ ) to each of the eleven models. This resulted in 11 adversarial datasets: six from attacking the OD models and five from attacking the SEG models. We use the clean dataset alongside these 11 adversarial datasets to evaluate the performance of the models and our detection scheme.

**Evaluation Metrics.** To evaluate the prediction performance of the models on both the clean dataset and the eleven adversarial datasets, we utilize the mean Average Precision (mAP), a widely accepted metric for assessing computer vision models. For evaluating our detection scheme, we employ the receiver operating characteristic (ROC) curve, a popular metric that illustrates the performance of a classification model across all classification thresholds. The area under the curve (AUC) provides a measure of our adversarial attack detection performance.

## 4.2 EXPERIMENTAL EVALUATION

First, we study the transferability of the attack. Next, we assess the performance of our detector scheme in detecting perturbations. We then discuss the cross-model strategy. Finally, we compare our detector with one state-of-the-art defense.

### 4.2.1 PREDICTION PERFORMANCE UNDER ATTACK

Table 1 shows the mAP for each model across twelve test datasets. The table’s diagonal highlights that the attack is most effective on the target model for which the perturbation is optimized. For instance, the attack on the FRCNN R50 model decreases its mAP from 30.2 to 0.18. Comparable performance declines can also be noted for the other models under attack, which was expected given that our attack is a white-box PGD attack on the target models.

The perturbations demonstrate transferability across models with similar network architectures, regardless of the task. For instance, the adversarial dataset generated by attacking the OD model FRCNN R50 decreases the mAP of the SEG model MRCNN R50 from 19.8 to 1.5. Conversely, the attack on MRCNN R50 reduces the mAP of FRCNN R50 from 30.2 to 8.8. This indicates that the perturbation can transfer to different tasks or models with the same backbone architectures. Transferability is also observed in models that share the same backbone type but differ in depth. As shown in Table 1, attacks on models with an R50 backbone (see columns) can transfer to models with an R101 backbone (see rows), and vice versa. However, for models with the same baseline architecture but different backbones, such as FRCNN R50 and FRCNN SwinT, or RetinaNet R50

Table 1: Impact of the attack on the mAP of vision models

Task	Vision Models	Clean mAP ( $\uparrow$ )	Attack Object Detection ( $\downarrow$ )						Attack Segmentation ( $\downarrow$ )				
			F R50	F R101	F SwinT	R R50	R R101	R PVT	M R50	M R101	G R50	G R101	M2F SwinT
OD	F R50	30.2	<b>0.18</b>	5.7	18.4	0.34	3.6	11.8	8.8	9.2	8.8	10.0	24.7
	F R101	30.3	7.5	<b>0.17</b>	18.5	3.2	1.0	13.0	14.5	6.4	14.1	8.06	25.0
	F SwinT	31.8	17.0	16.5	<b>1.5</b>	11.5	12.0	14.8	20.7	18.6	20.7	19.0	22.4
	R R50	28.7	2.2	4.4	17.0	<b>0.01</b>	2.7	10.2	7.1	7.4	7.0	8.9	23.1
	R R101	29.2	7.7	2.2	17.9	2.63	<b>0.02</b>	11.9	14.0	5.6	13.5	7.1	24.3
	R PVT	29.8	12.8	13.0	18.4	7.5	8.9	<b>0.04</b>	18.2	15.0	17.8	15.0	24.8
SEG	M R50	19.8	1.5	2.6	10.1	0.6	2.6	7.1	<b>0.01</b>	1.6	0.5	2.2	13.4
	M R101	20.5	4.2	1.8	10.3	2.5	1.4	8.0	4.4	<b>0.01</b>	4.2	0.72	13.4
	G R50	20.1	1.7	3.1	10.7	0.62	2.9	6.9	0.73	1.8	<b>0.01</b>	2.1	13.3
	G R101	20.7	4.2	2.0	10.3	2.6	1.7	7.6	4.4	0.3	4.1	<b>0.01</b>	13.2
	M2F SwinT	21.0	9.4	9.1	7.1	7.3	7.8	9.7	9.7	7.9	10.1	9.1	<b>2.8</b>

Acronyms: Object Detection (OD), Instance Segmentation (SEG), FRCNN (F), RetinaNet (RN), MRCNN (M), GCNET (G), Mask2Former (M2F)

A bold value is the lowest mAP score among all targeted models for a given adversarial dataset.

An underlined value indicates the adversarial dataset successfully dropped the mAP score of the targeted model below 5 mAP.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

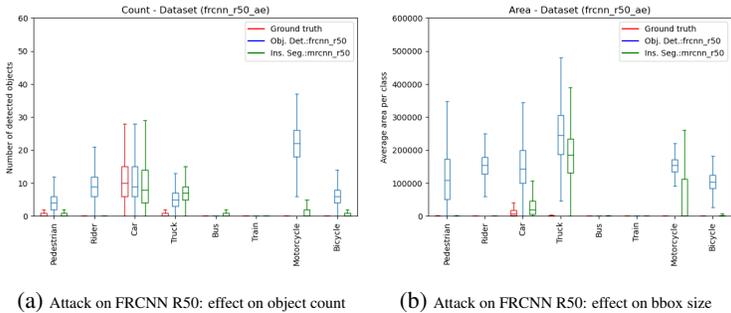


Figure 4: Perturbations optimized for FRCNN R50 transfer to MRCNN R50, impacting them differently in terms of number of detected objects and their sizes.

and RetinaNet PVT, the transferability is less evident. This indicates that the backbone plays a more crucial role in the transferability of the attack.

While perturbations can transfer between different models and tasks, their fine-grained impact varies significantly across models. This variation is evident in several aspects, such as the number of objects detected or the area of the bounding boxes. Figure 4 illustrates this using one model pairs: (FRCNN R50, MRCNN R50). Subfigure 4a shows the distribution of the number of objects for each category given the adversarial dataset optimized on FRCNN R50. The attack generates significantly more objects on FRCNN R50 than on MRCNN R50, especially for categories like rider and motorcycle. Subfigure 4b indicates that the attack also causes larger objects for FRCNN R50 in most categories, while for MRCNN R50, this effect is seen only in the truck and motorcycle categories. This demonstrates that even when perturbations transfer, they can lead to inconsistent impacts on different models<sup>3</sup>.

#### 4.2.2 PERTURBATION DETECTION PERFORMANCE

We evaluated the performance of our detector across 30 (6 OD × 5 SEG) model pairs. As previously noted in Figure 2, we aim for a model pair that exhibits a high consistency score (CS) on clean inputs while a lower score on adversarial inputs, facilitating the identification of perturbations. Figure 5 (top) illustrates the average CS of model pairs across the three datasets (clean, attack OD, attack SEG). The blue stars represent the average CS for clean inputs. Generally, the CS for clean inputs is high, especially for model pairs with similar baseline architectures (RCNN) and backbones (ResNet), which can extract consistent features from the clean inputs, resulting in consistent outputs. Model pairs with different architectures or backbones exhibit slightly lower CS due to their varying feature extraction capabilities, leading to inconsistent outputs. The red squares represent the CS for adversarial datasets optimized on OD models. As discussed in the previous section, similar architectures (RCNN and ResNet) result in high transferability but also high inconsistency, causing CS to drop as low as 0 for those model pairs. For model pairs with different architectures, the attack shows less transferability, and thus, higher consistency. Similar findings are observed when attacking SEG models (green circles). The full analysis can be found in Appendix A.3.

The AUC curves in Figure 5 (bottom) demonstrate that all model pairs achieve an AUC greater than 85%, with most exceeding 95%, when either model of the pairs is attacked. This highlights the exceptional performance of our consistency-based detector in identifying perturbations. Model pairs with similar backbone types (ResNet) and baseline architecture (RCNN) exhibit the highest performance, achieving an AUC of 99.9%. Again, it shows that, although the attack can easily transfer between these models, this transferability leads to distinct variations in the number, label, and size of the detected objects. These variations result in a higher level of inconsistency, which our detector can effectively identify. In contrast, model pairs with different backbones or baseline architectures exhibit low transferability and low inconsistency, resulting in a relatively lower AUC.

<sup>3</sup>See Appendix A.3 for full transferability analysis across all model pairs considered.

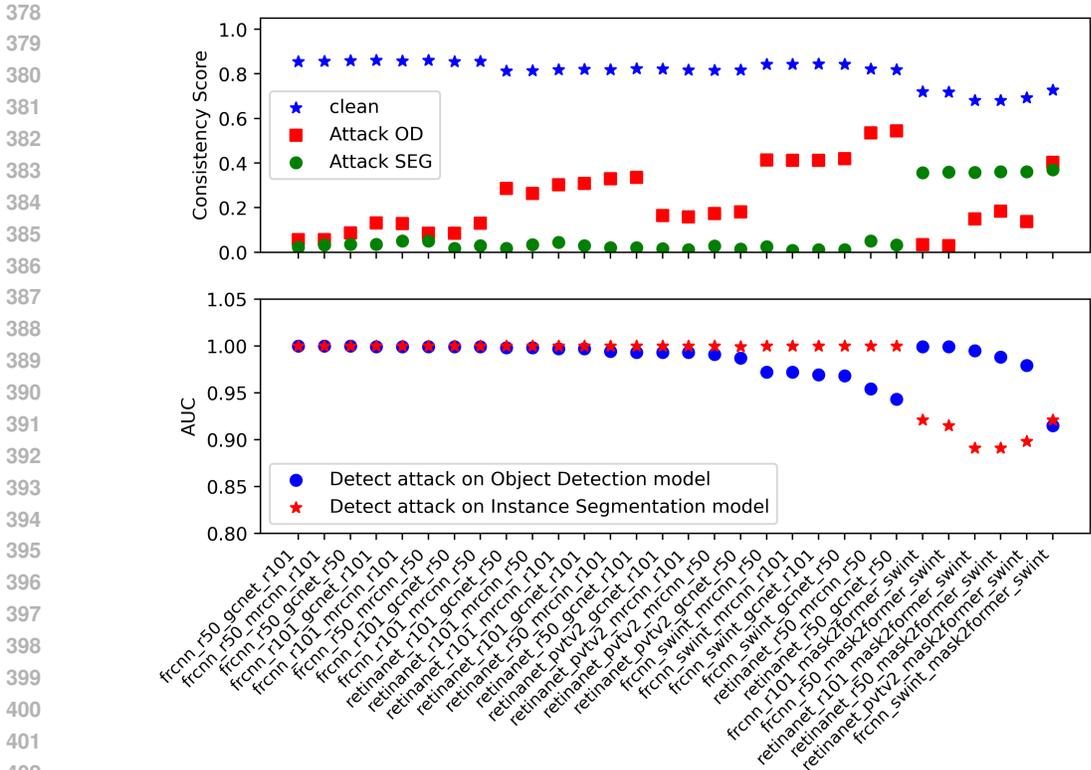


Figure 5: Top: Consistency Score for all model pairs. The lower the better the pair is for our detector. Bottom: The AUC for all model pairs. The higher the better the pair is our detector.

Table 2: Effect of perturbation strength  $\epsilon$  on mAP. Target model is FRCNN R50.

Model	Perturbation Strength					
	clean	1	2	4	8	16
<b>FRCNN R50</b>	30.2	28.4	25.4	18.7	4.2	0.18
<b>MRCNN R50</b>	19.8	18.9	17.0	13.4	6.2	1.5

Next, we are interested to learn how the perturbation strength of the attack can impact the prediction performance of the models and the detection performance of our detector. We evaluate the robustness of the models against attack size  $\epsilon \in \{1/255, 2/255, 4/255, 8/255, 16/255\}$ . Here, we use one of the best model pair (FRCNN R50, MRCNN R50) as an example. More results of other model pairs can be found in the Appendix A.4. As shown in Table 2, the mAP of both models decreases as the perturbation strength increases, which is expected. Conversely, as illustrated in Figure 6b, the detection performance in terms of AUC increases. This is the desired behavior because stronger perturbation leads to greater inconsistency (lower consistency score as in Figure 6a) between the outputs of model pairs, resulting in higher detection performance for our detector.

**Takeaway on multi-task architecture.** The empirical analysis indicates that our detector performs optimally when model pairs exhibit high inconsistency in their outputs. Under our attacker model, the most effective model pairs are those with similar architectures and backbones, as they demonstrate high adversarial transferability but also high inconsistency.

#### 4.2.3 COMPARISON TO OTHER DEFENSES

In this section, we compare our detector with the adversarial training method RobustDet, as proposed by Dong et al. (2022). For a fair comparison, we applied RobustDet to FRCNN R50 which resulted in a robust model named RobustFRCNN. More details about our implementation can be found in

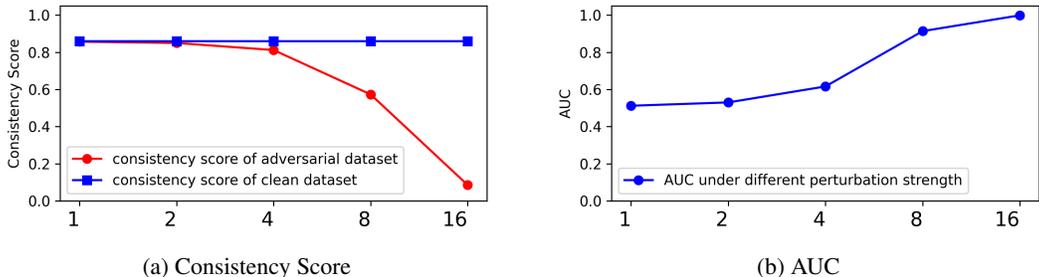
Figure 6: Impact of perturbation strength  $\epsilon$  on our detector (consistency score AUC)

Table 3: Impact of the attack on the mAP of regular and RobustFRCNN

Model	Clean				Attack			
	mAP	mAP <sub>small</sub>	mAP <sub>medium</sub>	mAP <sub>large</sub>	mAP	mAP <sub>small</sub>	mAP <sub>medium</sub>	mAP <sub>large</sub>
<b>FRCNN R50</b>	30.2	12.4	34.6	54.4	0.18	0.07	0.24	0.33
<b>RobustFRCNN</b>	19.8	8.1	22.4	36.7	6.2	2.4	7.3	11.9

Appendix B.2.1. First, we evaluated the prediction performance of two models based on FRCNN R50: the standard model and RobustFRCNN. Table 3 presents the performance of both models under clean and adversarial datasets. For the standard model, we used the same adversarial dataset as previously mentioned. For the robust model, we applied PGD attack using same attack parameters. Table 3 shows that the standard model experiences a significant performance drop due to the attack, compared to the robust model. Specifically, its mAP decreases from 30.2 to 0.18, while for the RobustFRCNN, it decreases from 19.8 to 6.2. Hence, RobustDet technique enhances the model’s adversarial robustness. It is worth noting that the clean mAP for the robust model is not high, indicating there is still room to adjust the training parameters to improve both its clean and adversarial performance.

Next, we compare our detector with RobustDet. Our detector functions as a binary classifier, determining whether an input is adversarial or not. In contrast, RobustDet, similar to adversarial training, enhances the model’s robustness. To ensure a fair comparison, we introduce the metric *Detection Rate*, which represents the true positive rate for a given adversarial dataset. Specifically, we utilize one of our best model pairs (FRCNN R50, MRCNN R50) for our detector and assess its detection rate on the adversarial dataset for FRCNN R50. For RobustFRCNN, we evaluate its performance by calculating the consistency score between its output and the ground truth annotations under adversarial conditions. A high consistency score indicates that RobustFRCNN successfully mitigates the perturbation, whereas a low score signifies failure. Therefore, the detection rate is the ratio of adversarial inputs with a consistency score above the consistency threshold.

As shown in Table 4, our consistency-based detector successfully identifies all adversarial inputs (99.9%) in the adversarial datasets, thanks to the high inconsistency between the outputs of the model pair. In contrast, RobustFRCNN performs poorly (mAP = 6.2), failing to ensure prediction outputs align with the ground truth, resulting in a very low detection rate (19%). On top of being less able to detect adversarial inputs, RobustFRCNN employs a dynamic convolution kernel that is four times the size of a regular convolution kernel, significantly increasing its model size. In comparison, our detector has a combined weight size of only 350MB for OD and SEG. Finally, our detector achieves faster inference speeds on the same hardware due to its less complex architecture. In summary, our detector demonstrates stronger performance compared to RobustDet.

Table 4: Performance of RobustFRCNN vs our detector

Defense	Detection Rate	Model Weight Size (MB)	Inference Speed (FPS)
<b>RobustFRCNN</b>	19	643	11
<b>Our detector</b>	<b>99.9</b>	<b>350</b>	<b>20</b>

## 5 OPEN CHALLENGES

**Stronger attacker model.** Our attacker model targets only one perception task. A stronger attacker could target both tasks. Prior work have demonstrated attacks fooling both semantic segmentation and object detection with some success (Xie et al., 2017). More generally, techniques were designed to improve the adversarial transferability cross-model or cross-task (Gu et al., 2023; Wei et al., 2024; Hu et al., 2024; Lu et al., 2020). We expect that an attacker that uses these techniques would be able to bypass our multi-task consistency detector. However, we noted that these techniques, despite being able to fool both tasks in silo, do not create cross-task consistent adversarial output (e.g., a fake bounding box in OD does not match the location or size of the fake instance segmentation mask).

**Generalization.** In this paper, we investigated object detection and instance segmentation models, finding the best model pairs to use in a multi-task consistency detector. We are interested in generalizing the approach to any combination of tasks. Especially, we would like to understand if the recommendations (about the model architectures) generalized across tasks.

**Tuple multi-task consistency.** We propose to extend the detector with more tasks and investigate how the detection rate correlates to the number of tasks. Though, Ghamizi et al. (2022) demonstrated that what matters the most is not the number of tasks or how they correlate, but how much the tasks individually impact the vulnerability of the model. Indeed, the more vulnerable the tasks in the model are, the less likely adding new tasks increases the robustness of the model; and adding a vulnerable task may actually decrease the robustness of the whole model. Thus, a comprehensive analysis is required to answer this challenge.

## 6 RELATED WORK

In recent years, many defenses were created to detect (Hendrycks & Gimpel, 2016; Liu et al., 2019; Tian et al., 2021; Sperl et al., 2020) or to improve the robustness of vision systems against adversarial perturbations (Hendrycks et al., 2019; Mađry et al., 2018). Especially, a strong emphasis has been put on the security of image classification task. Examples of defenses used in image classification include: use of additional detection networks (Liu et al., 2019), analysis of network output (Hendrycks & Gimpel, 2016; Tian et al., 2021), or use of certain activation patterns within the hidden layers (Sperl et al., 2020). These detection methods focus on the output structure or network topology of an image classifier and are thereby not transferable to more complex vision tasks.

As described earlier, multi-task learning (MTL) (Kendall et al., 2018) tackles a wide range of vision tasks in an efficient way. Mao et al. (2020) showed that MTL increases the adversarial robustness due to the increased difficulty of successfully attacking several tasks. Thus, subsequent work explored other task combinations (Xie et al., 2017; Klingner et al., 2020; Wang et al., 2020; Kumar et al., 2021), or compared the effectiveness of adding different auxiliary tasks (Ghamizi et al., 2022; Gurulingan et al., 2021; Haleta et al., 2021). While the positive effects of MTL on adversarial robustness are quite well-explored, we are the first to check the consistency between outputs from object detection and instance segmentation models, deriving recommendations to select best model pairs.

## 7 CONCLUSION

Vision models are paramount to many applications such as autonomous driving. Their robustness have been shown to be brittle under adversarial setting. From the observation that adversarial inputs yield different effects when fed to different models, we propose an adversarial perturbation detection method based on multi-task perception. We showed an example of our lightweight defense using instance segmentation and object detection tasks. We generated adversarial BDD100k datasets and demonstrated our consistency score can effectively detect perturbations. Then, we empirically identified the optimal model pairs, demonstrating that even if sharing the same backbone, the attack can be detected because of uncoordinated perturbations on both models. The optimal models pair had a 99.9% detection rate. Future work will focus on joint multi-task perturbations and assess the effectiveness of our defense against stronger attacker models.

## REFERENCES

- 540  
541  
542 Eduardo Arnold, Omar Y Al-Jarrah, Mehrdad Dianati, Saber Fallah, David Oxtoby, and Alex  
543 Mouzakitis. A survey on 3d object detection methods for autonomous driving applications. *IEEE*  
544 *Transactions on Intelligent Transportation Systems*, 20(10):3782–3795, 2019.
- 545 Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.
- 546  
547 Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen  
548 Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie  
549 Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang,  
550 Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark.  
551 *arXiv preprint arXiv:1906.07155*, 2019.
- 552 Pin-Chun Chen, Bo-Han Kung, and Jun-Cheng Chen. Class-aware robust adversarial training for  
553 object detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
554 10420–10429, 2021.
- 555 Ka-Ho Chow, Ling Liu, Margaret Loper, Juhyun Bae, Mehmet Emre Gursoy, Stacey Truex, Wenqi  
556 Wei, and Yanzhao Wu. Adversarial objectness gradient attacks in real-time object detection  
557 systems. In *IEEE International Conference on Trust, Privacy and Security in Intelligent Systems*  
558 *and Applications*, pp. 263–272. IEEE, 2020.
- 559 Nilaksh Das, Madhuri Shanbhogue, Shang-Tse Chen, Fred Hohman, Siwei Li, Li Chen, Michael E  
560 Kounavis, and Duen Horng Chau. Compression to the rescue: Defending from adversarial attacks  
561 across modalities. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2018.
- 562 Ziyi Dong, Pengxu Wei, and Liang Lin. Adversarially-aware robust object detector. In *European*  
563 *Conference on Computer Vision*, pp. 297–313. Springer, 2022.
- 564  
565 Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaeser, Fabian Timm,  
566 Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic seg-  
567 mentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on*  
568 *Intelligent Transportation Systems*, 22(3):1341–1360, 2020.
- 569 Salah Ghamizi, Maxime Cordy, Mike Papadakis, and Yves Le Traon. Adversarial robustness in  
570 multi-task learning: Promises and illusions. In *AAAI Conference on Artificial Intelligence*, pp.  
571 697–705, 2022.
- 572  
573 Jindong Gu, Xiaojun Jia, Pau de Jorge, Wenqain Yu, Xinwei Liu, Avery Ma, Yuan Xun, Anjun Hu,  
574 Ashkan Khakzar, Zhijiang Li, et al. A survey on transferability of adversarial examples across  
575 deep neural networks. *arXiv preprint arXiv:2310.17626*, 2023.
- 576  
577 Minghao Guo, Yuzhe Yang, Rui Xu, Ziwei Liu, and Dahua Lin. When nas meets robustness: In  
578 search of robust architectures against adversarial attacks. In *IEEE/CVF Conference on Computer*  
579 *Vision and Pattern Recognition*, pp. 631–640, 2020a.
- 580 Pengxin Guo, Yuancheng Xu, Baijiong Lin, and Yu Zhang. Multi-task adversarial attack. *arXiv*  
581 *preprint arXiv:2011.09824*, 2020b.
- 582  
583 Naresh Kumar Gurulingan, Elahe Arani, and Bahram Zonooz. Uninet: A unified scene understanding  
584 network and exploring multi-task relationships through the lens of adversarial attacks. In *IEEE/CVF*  
585 *International Conference on Computer Vision*, pp. 2239–2248, 2021.
- 586 Pavlo Haleta, Dmytro Likhomanov, and Oleksandra Sokol. Multitask adversarial attack with disper-  
587 sion amplification. *EURASIP Journal on Information Security*, 2021(1):10, 2021.
- 588  
589 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
590 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 591  
592 Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning  
593 can improve model robustness and uncertainty. *Advances in neural information processing systems*,  
32, 2019.

- 594 Anjun Hu, Jindong Gu, Francesco Pinto, Konstantinos Kamnitsas, and Philip Torr. As firm as their  
595 foundations: Can open-sourced foundation models be used to create adversarial examples for  
596 downstream tasks? *arXiv preprint arXiv:2403.12693*, 2024.
- 597 Thomas E. Huang. Bdd100k Model Zoo. <https://github.com/SysCV/bdd100k-models>,  
598 2021.
- 600 Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses  
601 for scene geometry and semantics. In *IEEE conference on computer vision and pattern recognition*,  
602 pp. 7482–7491, 2018.
- 603 Marvin Klingner, Andreas Bar, and Tim Fingscheidt. Improved noise and attack robustness for  
604 semantic segmentation by using multi-task training with self-supervised depth estimation. In  
605 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 320–321,  
606 2020.
- 607 Marvin Klingner, Varun Ravi Kumar, Senthil Yogamani, Andreas Bär, and Tim Fingscheidt. Detecting  
608 adversarial perturbations in multi-task perception. In *IEEE/RSJ International Conference on*  
609 *Intelligent Robots and Systems (IROS)*, pp. 13050–13057. IEEE, 2022.
- 611 Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and  
612 high-level vision using diverse datasets and limited memory. In *IEEE conference on computer*  
613 *vision and pattern recognition*, pp. 6129–6138, 2017.
- 614 Varun Ravi Kumar, Marvin Klingner, Senthil Yogamani, Stefan Milz, Tim Fingscheidt, and Patrick  
615 Mader. Syndistnet: Self-supervised monocular fisheye camera distance estimation synergized with  
616 semantic segmentation for autonomous driving. In *IEEE/CVF winter conference on applications*  
617 *of computer vision*, pp. 61–71, 2021.
- 618 Mindi Lan, Jun Luo, Senchun Chai, Ruiqi Chai, Chen Zhang, and Baihai Zhang. A novel industrial  
619 intrusion detection method based on threshold-optimized cnn-bilstm-attention using roc curve.  
620 In *Chinese Control Conference (CCC)*, pp. 7384–7389, 2020. doi: 10.23919/CCC50068.2020.  
621 9188872.
- 623 Jiayang Liu, Weiming Zhang, Yiwei Zhang, Dongdong Hou, Yujia Liu, Hongyue Zha, and Nenghai  
624 Yu. Detection based defense against adversarial examples from the steganalysis point of view. In  
625 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4825–4834, 2019.
- 626 Yantao Lu, Yunhan Jia, Jianyu Wang, Bai Li, Weiheng Chai, Lawrence Carin, and Senem Velipasalar.  
627 Enhancing cross-task black-box transferability of adversarial examples with dispersion reduction.  
628 In *IEEE/CVF conference on Computer Vision and Pattern Recognition*, pp. 940–949, 2020.
- 629 Yongxi Lu, Abhishek Kumar, Shuangfei Zhai, Yu Cheng, Tara Javidi, and Rogerio Feris. Fully-  
630 adaptive feature sharing in multi-task networks with applications in person attribute classification.  
631 In *IEEE conference on computer vision and pattern recognition*, pp. 5334–5343, 2017.
- 633 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
634 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,  
635 2017.
- 636 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
637 Towards deep learning models resistant to adversarial attacks. In *International Conference on*  
638 *Learning Representations*, 2018.
- 639 Aleksander Mađry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
640 Towards deep learning models resistant to adversarial attacks. *stat*, 1050:9, 2018.
- 642 Chengzhi Mao, Amogh Gupta, Vikram Nitin, Baishakhi Ray, Shuran Song, Junfeng Yang, and Carl  
643 Vondrick. Multitask learning strengthens adversarial robustness. In *European Conference on*  
644 *Computer Vision*, pp. 158–174. Springer, 2020.
- 645 Shokhrukh Miraliev, Shakhboz Abdigapporov, Vijay Kakani, and Hakil Kim. Real-time memory  
646 efficient multitask learning model for autonomous driving. *IEEE Transactions on Intelligent*  
647 *Vehicles*, 2023.

- 648 Zhuang Qian, Kaizhu Huang, Qiu-Feng Wang, and Xu-Yao Zhang. A survey of robust adversarial  
649 training in pattern recognition: Fundamental, theory, and methodologies. *Pattern Recognition*, 131:  
650 108889, 2022.
- 651 Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer,  
652 Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in neural  
653 information processing systems*, 32, 2019.
- 654 Philip Sperl, Ching-Yu Kao, Peng Chen, Xiao Lei, and Konstantin Böttinger. Dla: dense-layer-  
655 analysis for adversarial example detection. In *IEEE European Symposium on Security and Privacy*,  
656 pp. 198–215. IEEE, 2020.
- 657 Jinyu Tian, Jiantao Zhou, Yuanman Li, and Jia Duan. Detecting adversarial examples from sensitivity  
658 inconsistency of spatial-transform domain. In *AAAI Conference on Artificial Intelligence*, pp.  
659 9877–9885, 2021.
- 660 Derui Wang, Chaoran Li, Sheng Wen, Surya Nepal, and Yang Xiang. Defending against adversarial  
661 attack towards deep neural networks via collaborative multi-task training. *IEEE Transactions on  
662 Dependable and Secure Computing*, 19(2):953–965, 2020.
- 663 Zhipeng Wei, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Adaptive cross-modal transferable  
664 adversarial attacks from images to videos. *IEEE Transactions on Pattern Analysis and Machine  
665 Intelligence*, 46(5):3772–3783, 2024. doi: 10.1109/TPAMI.2023.3347835.
- 666 Chong Xiang, Saeed Mahloujifar, and Prateek Mittal. PatchCleanser: Certifiably robust defense  
667 against adversarial patches for any image classifier. In *USENIX Security Symposium*, pp. 2065–  
668 2082, 2022.
- 669 Chong Xiang, Tong Wu, Sihui Dai, Jonathan Petit, Suman Jana, and Prateek Mittal. {PatchCURE}:  
670 Improving certifiable robustness, model utility, and computation efficiency of adversarial patch  
671 defenses. In *USENIX Security Symposium*, pp. 3675–3692, 2024.
- 672 Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial  
673 examples for semantic segmentation and object detection. In *IEEE international conference on  
674 computer vision*, pp. 1369–1378, 2017.
- 675 Shaokai Ye, Kaidi Xu, Sijia Liu, Hao Cheng, Jan-Henrik Lambrechts, Huan Zhang, Aojun Zhou,  
676 Kaisheng Ma, Yanzhi Wang, and Xue Lin. Adversarial robustness vs. model compression, or both?  
677 In *IEEE/CVF International Conference on Computer Vision*, pp. 111–120, 2019.
- 678 Haichao Zhang and Jianyu Wang. Towards adversarially robust object detection. In *IEEE/CVF  
679 International Conference on Computer Vision*, pp. 421–430, 2019.
- 680 Dingfu Zhou, Jin Fang, Xibin Song, Liu Liu, Junbo Yin, Yuchao Dai, Hongdong Li, and Ruigang  
681 Yang. Joint 3d instance segmentation and object detection for autonomous driving. In *IEEE/CVF  
682 Conference on Computer Vision and Pattern Recognition*, pp. 1839–1849, 2020.
- 683 Zhengxia Zou, Keyan Chen, Zhenwei Shi, Yuhong Guo, and Jieping Ye. Object detection in 20 years:  
684 A survey. *Proceedings of the IEEE*, 111(3):257–276, 2023.
- 685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A APPENDIX

### A.1 ALGORITHM

Algorithm 1 describes the different steps involved in the computation of the consistency score. The function *IoU* function is the function *box\_iou* defined in the *Torchvision* library.

---

#### Algorithm 1 Consistency Score Calculation

---

**Input:**  $\mathcal{S}_{det}$  // A set of pairs of bounding boxes and labels from object detection  
**Input:**  $\mathcal{S}_{seg}$  // A set of pairs of bounding boxes and labels from instance segmentation  
**Output:** Consistency Score ( $C$ )  
 $|CD| = 0$  // Number of pairs (box and mask)  
 $IoU = calc\_iou(\mathcal{S}_{det}, \mathcal{S}_{seg})$  // IoU score and label similarity for between pairs of  $\mathcal{S}_{det}$  and  $\mathcal{S}_{seg}$   
 $IoU = prune(IoU, threshold)$  // Prune each box with all IoU scores below threshold  
 $n\_box_{seg}, n\_box_{det} = get\_number(IoU)$  // Get remaining number of boxes for each task  
 $|CD| = compute\_n\_pairs(n\_box_{seg}, n\_box_{det})$  // Get total number of pairs  
 $C = compute\_c(|CD|, len(\mathcal{S}_{det}), len(\mathcal{S}_{seg}))$

---

### A.2 DATASET: BDD100K

The BDD100K dataset is a public dataset of driving scenes, which contains 100k frames and annotations for 10 vision tasks. Compared with other driving datasets, the BDD100k dataset has a diversity of geography, environment, and weather. Therefore, we use the BDD100k as the benchmark dataset to train the models and evaluate our detection. In particular, we use the 100k subfolder for object detection task, which is split to 70k training, 10k validation and 20k testing images. We also use the 10k subfolder for instance segmentation task, which is split to 7k training, 1k validation and 2k testing images.

### A.3 ADVERSARIAL TRANSFERABILITY

This section presents the distinct impact of the attack on OD and SEG models across all model pairs, focusing on the number and size of the detected objects. As previously mentioned, the attack exhibits high transferability between models with similar architectures and backbones, but it also leads to significant inconsistencies in the model outputs. For instance, Figure 8 shows the attack transfer between FRCNN R50 and MRCNN R50, but the number and size of hallucinated objects across the categories vary. For model pairs with different baseline architectures or backbones, such as FRCNN R50 and MASK2FORMER SwinT in Figure 10, the adversarial dataset optimized on MASK2FORMER does not fool FRCNN R50, whose outputs remain close to the ground truth in terms of number and size. Similarly, the adversarial dataset optimized on the FRCNN model does not fool MASK2FORMER whose object areas are close to ground truth. Although the number of objects is very large, this is due to the poor performance of MASK2FORMER, which predicts a large number of objects even on clean inputs, as shown in Figure 7.

This observation also supports the conclusion in Section 4.2.2. The consistency scores for model pairs with similar architectures are low because the attack transfers between them, resulting in distinct impacts, and thus, high inconsistency. For model pairs with different architectures, the attack does not transfer well, leading to low inconsistency. However, when one model in these pairs performs very poorly even on a clean dataset, it will output many hallucinated objects despite the attack not transferring to it, still resulting in high inconsistency, as seen with FRCNN R50 and MASK2FORMER SwinT.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

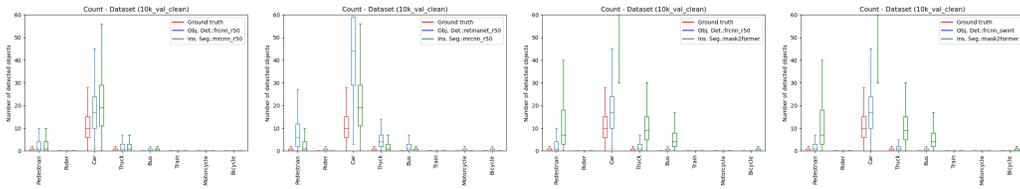


Figure 7: Clean images on model pairs

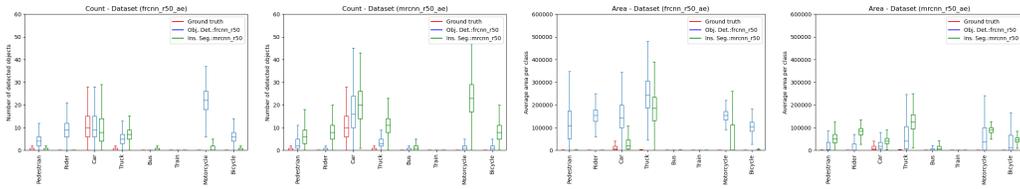


Figure 8: OD\_frncn\_r50\_SEG\_mrcnn\_r50

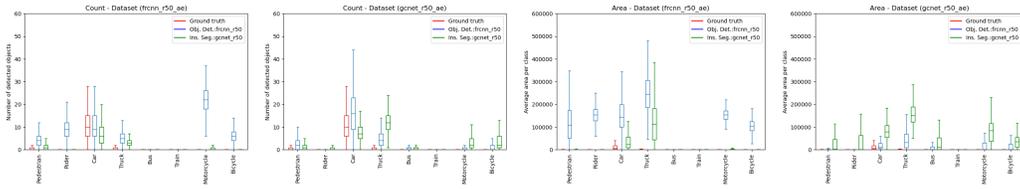


Figure 9: OD\_frncn\_r50\_SEG\_gcnet\_r50

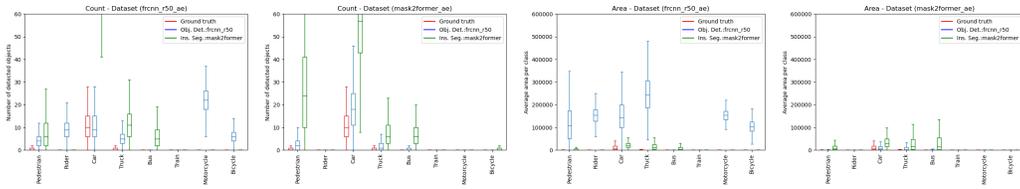


Figure 10: OD\_frncn\_r50\_SEG\_mask2former

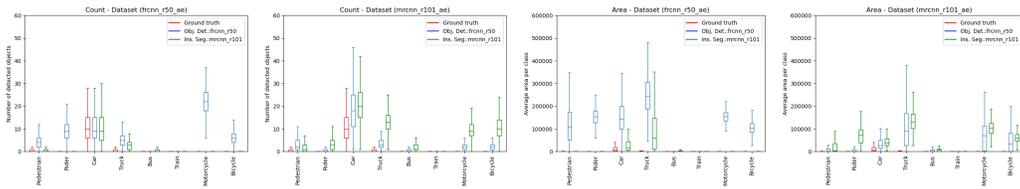


Figure 11: OD\_frncn\_r50\_SEG\_mrcnn\_r101

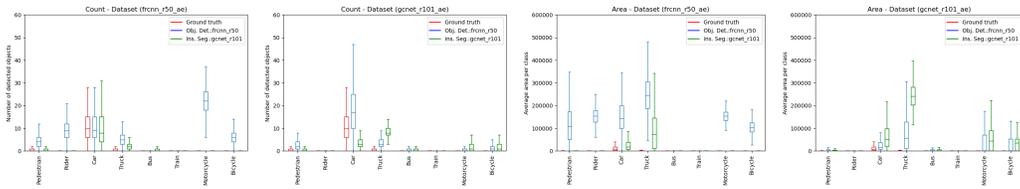


Figure 12: OD\_frncn\_r50\_SEG\_gcnet\_r101

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

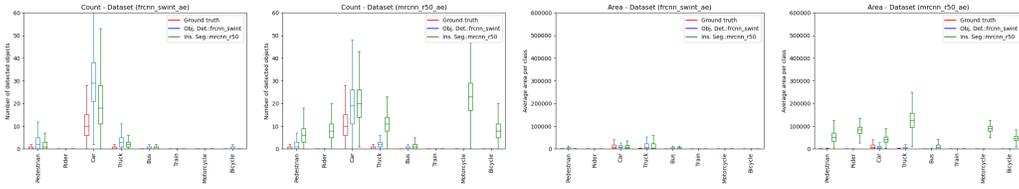


Figure 13: OD\_frncn\_swint\_SEG\_mrcnn\_r50

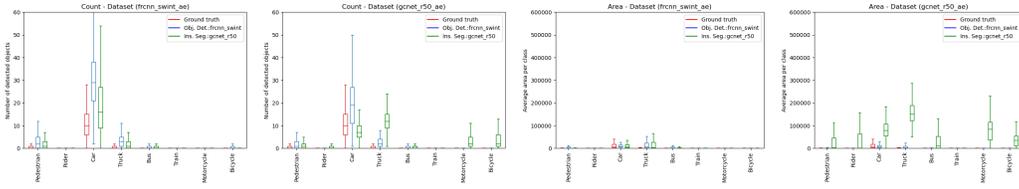


Figure 14: OD\_frncn\_swint\_SEG\_gcnet\_r50

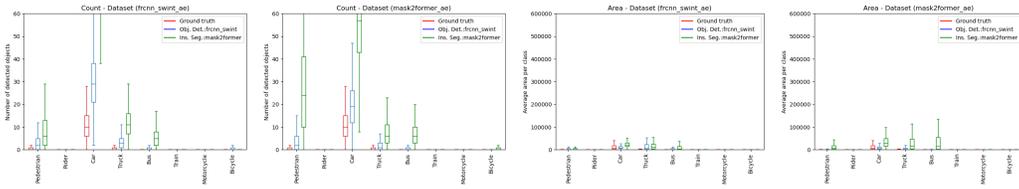


Figure 15: OD\_frncn\_swint\_SEG\_mask2former

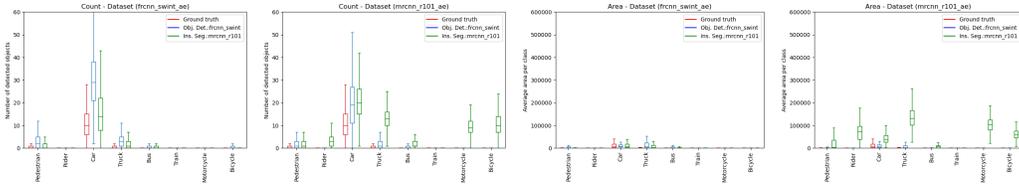


Figure 16: OD\_frncn\_swint\_SEG\_mrcnn\_r101

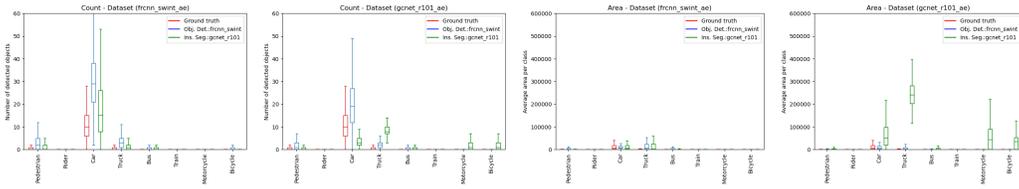


Figure 17: OD\_frncn\_swint\_SEG\_gcnet\_r101

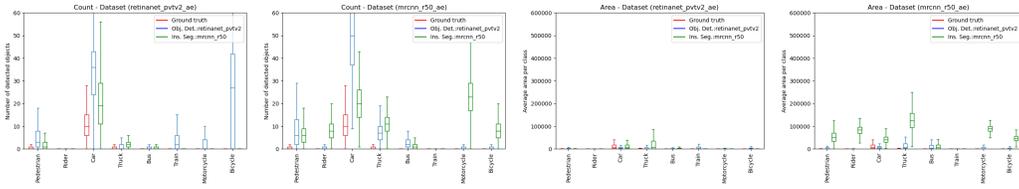


Figure 18: OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r50

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

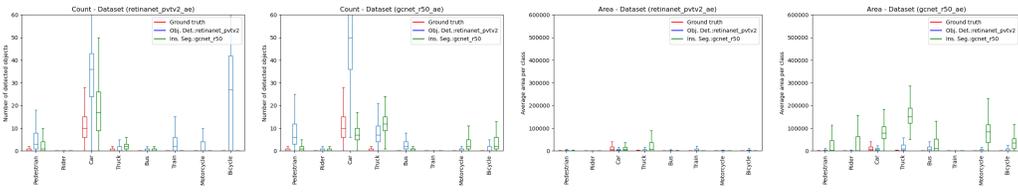


Figure 19: OD\_retinanet\_pvtv2\_SEG\_gcnet\_r50

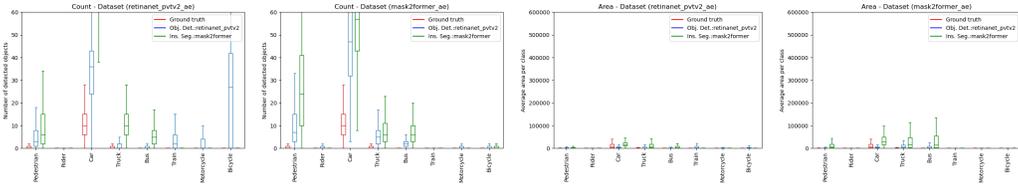


Figure 20: OD\_retinanet\_pvtv2\_SEG\_mask2former

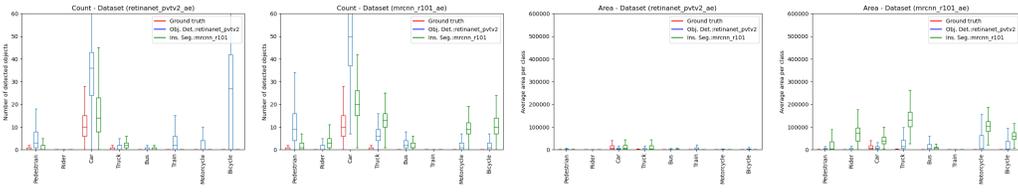


Figure 21: OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r101

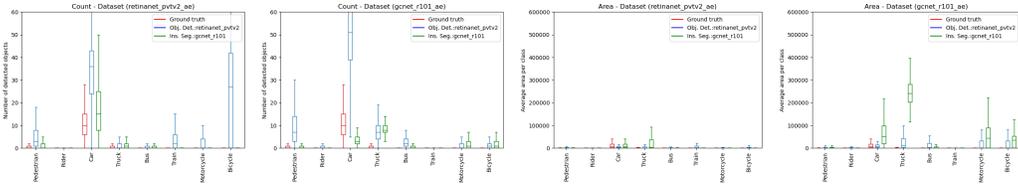


Figure 22: OD\_retinanet\_pvtv2\_SEG\_gcnet\_r101

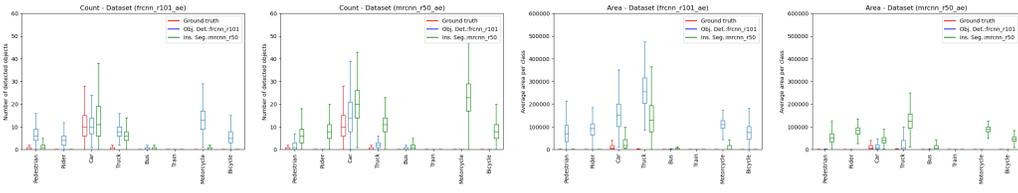


Figure 23: OD\_frcnn\_r101\_SEG\_mrcnn\_r50

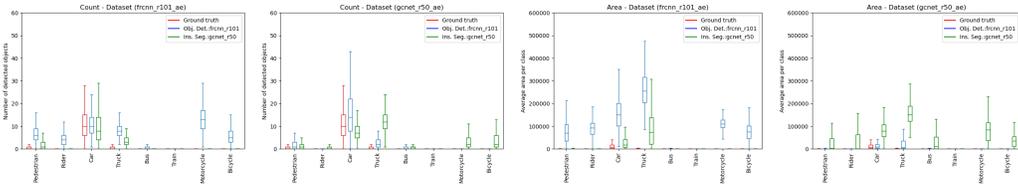


Figure 24: OD\_frcnn\_r101\_SEG\_gcnet\_r50

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

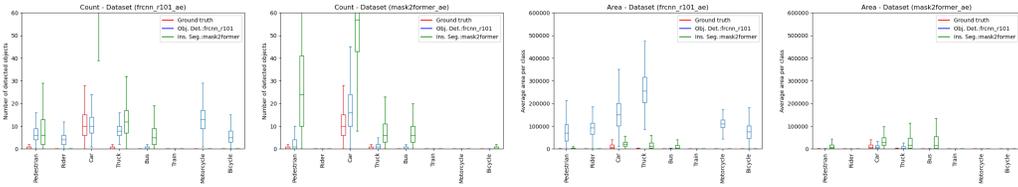


Figure 25: OD\_frncn\_r101\_SEG\_mask2former

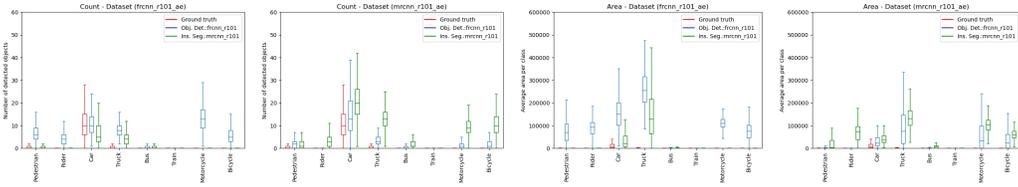


Figure 26: OD\_frncn\_r101\_SEG\_mrcnn\_r101

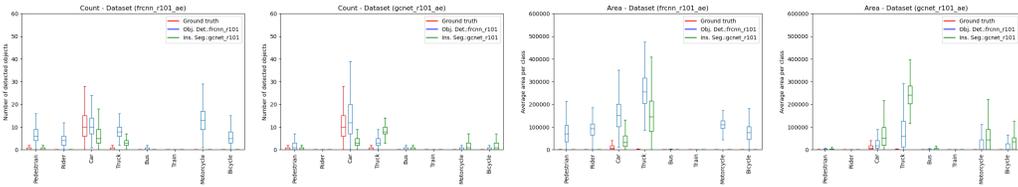


Figure 27: OD\_frncn\_r101\_SEG\_gcnet\_r101

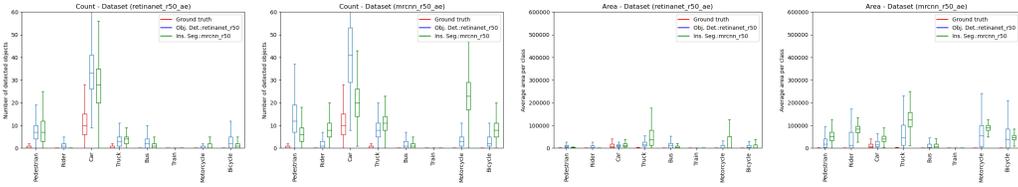


Figure 28: OD\_retinanet\_r50\_SEG\_mrcnn\_r50

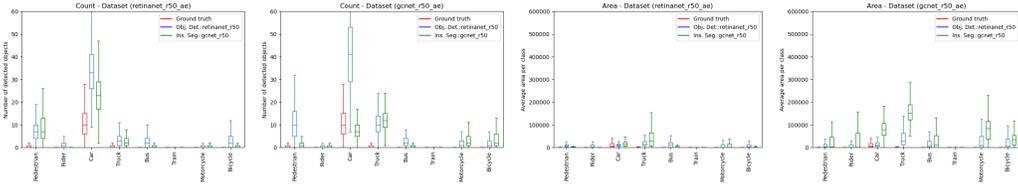


Figure 29: OD\_retinanet\_r50\_SEG\_gcnet\_r50

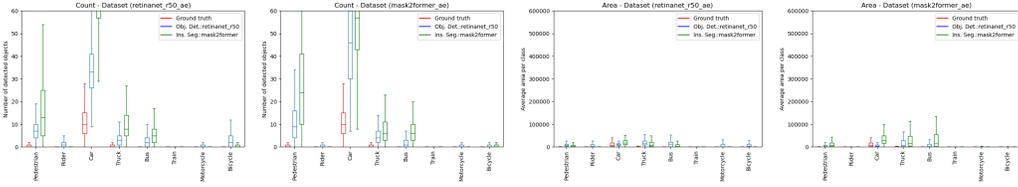


Figure 30: OD\_retinanet\_r50\_SEG\_mask2former

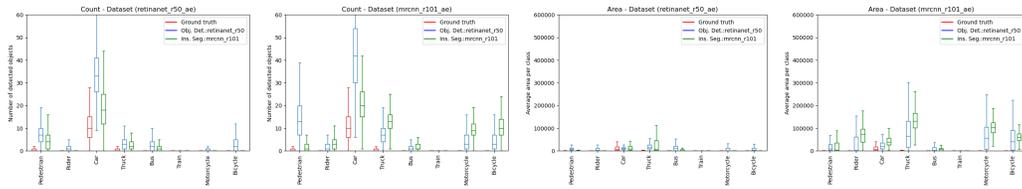


Figure 31: OD\_retinanet\_r50\_SEG\_mrcnn\_r101

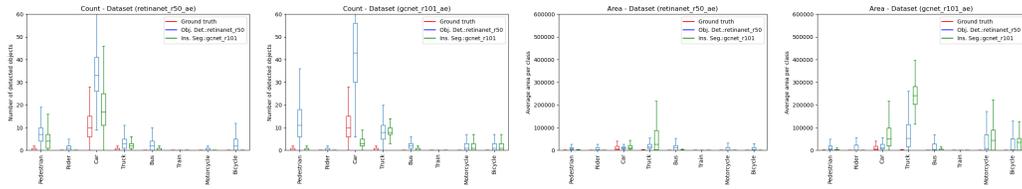


Figure 32: OD\_retinanet\_r50\_SEG\_gcnet\_r101

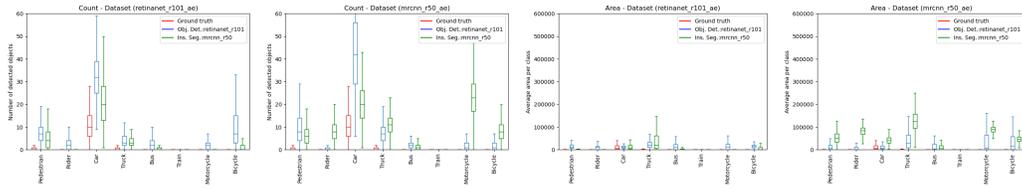


Figure 33: OD\_retinanet\_r101\_SEG\_mrcnn\_r50

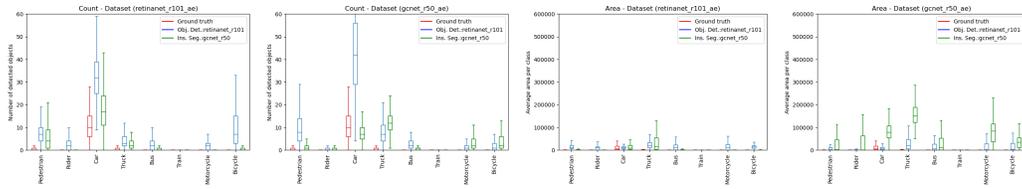


Figure 34: OD\_retinanet\_r101\_SEG\_gcnet\_r50

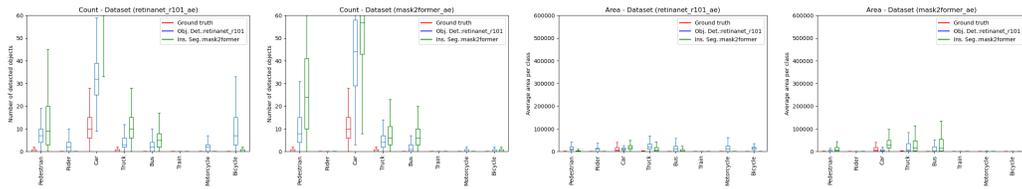


Figure 35: OD\_retinanet\_r101\_SEG\_mask2former

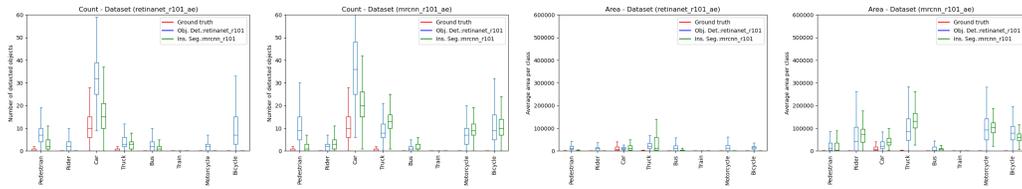


Figure 36: OD\_retinanet\_r101\_SEG\_mrcnn\_r101

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

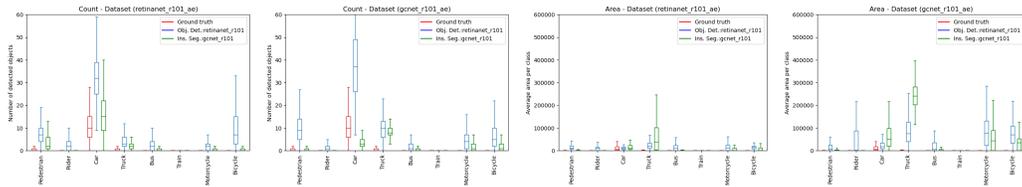


Figure 37: OD\_retinanet\_r101\_SEG\_gcnet\_r101

## A.4 DETECTOR PERFORMANCE

**Distribution of the consistency score.** This part contains additional results for CS distribution for all model pairs. As previous results showed, the model pairs with similar architecture results in distinct CS distributions between clean inputs and adversarial inputs, e.g., in Figure 38. This is desired for our detector to identify the perturbation. In contrast, the CS distributions for FRCNN R50 and MASK2FORMER SwinT is more difficult to distinguish, particularly when attacking MASK2FORMER SwinT model. This results in a relatively low AUC for this model pair as seen in Figure 5.

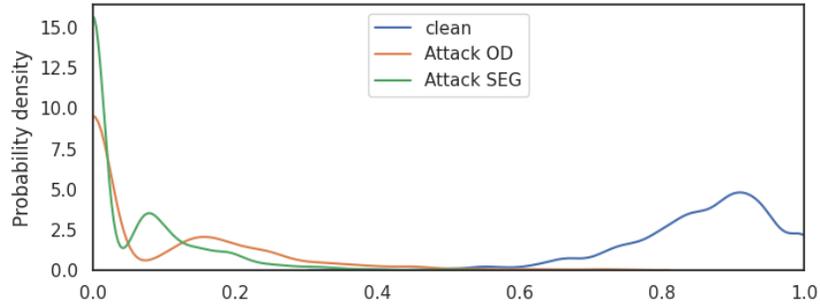


Figure 38: OD\_frcnn\_r50\_SEG\_mrcnn\_r50

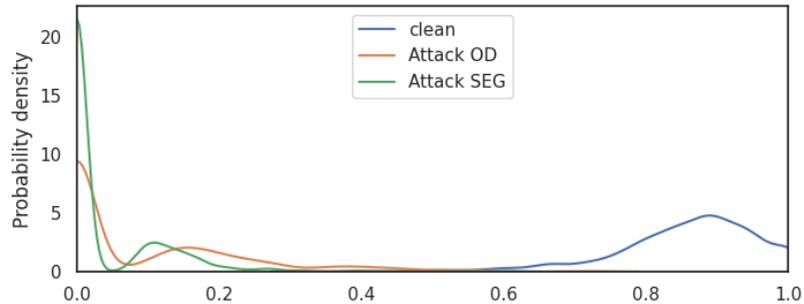


Figure 39: OD\_frcnn\_r50\_SEG\_gcnet\_r50

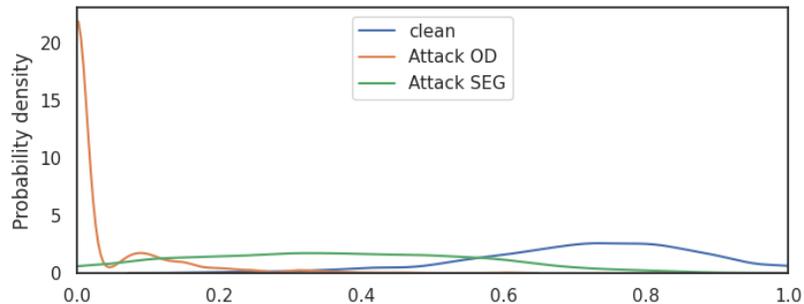


Figure 40: OD\_frcnn\_r50\_SEG\_mask2former

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

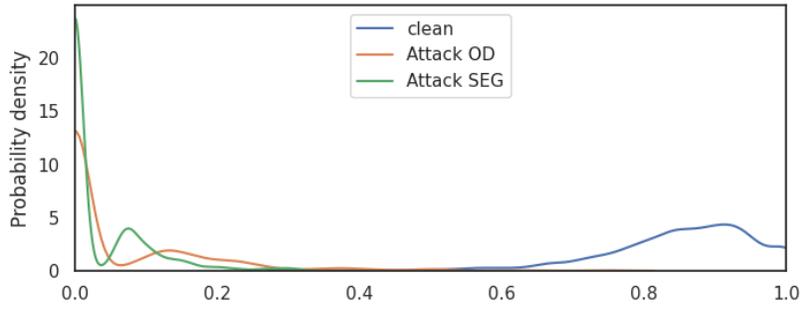


Figure 41: OD\_frcnn\_r50\_SEG\_mrcnn\_r101

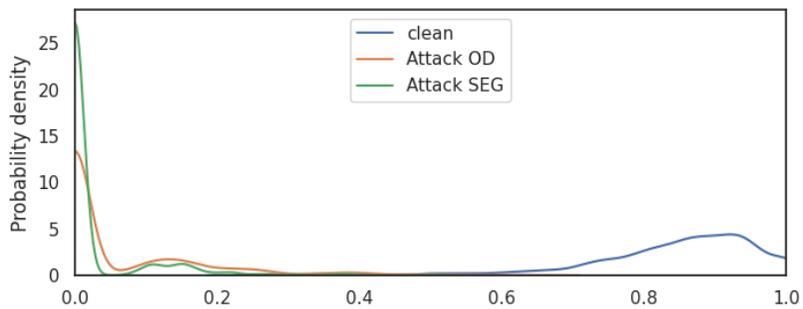


Figure 42: OD\_frcnn\_r50\_SEG\_gcnet\_r101

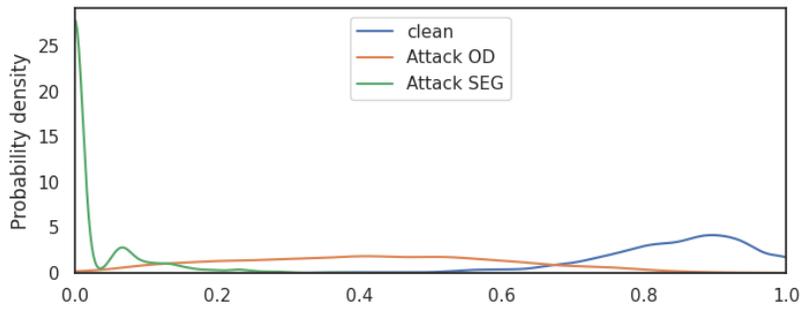


Figure 43: OD\_frcnn\_swint\_SEG\_mrcnn\_r50

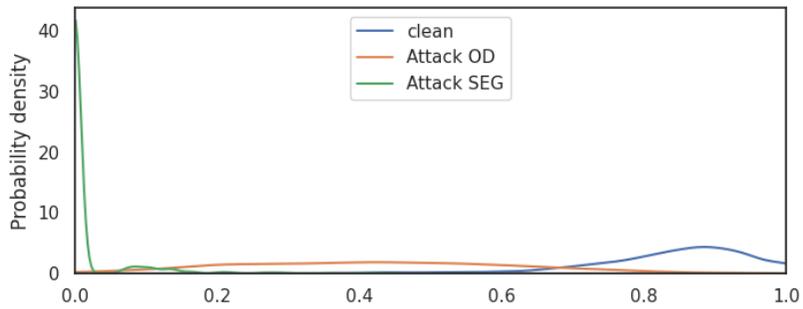


Figure 44: OD\_frcnn\_swint\_SEG\_gcnet\_r50

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

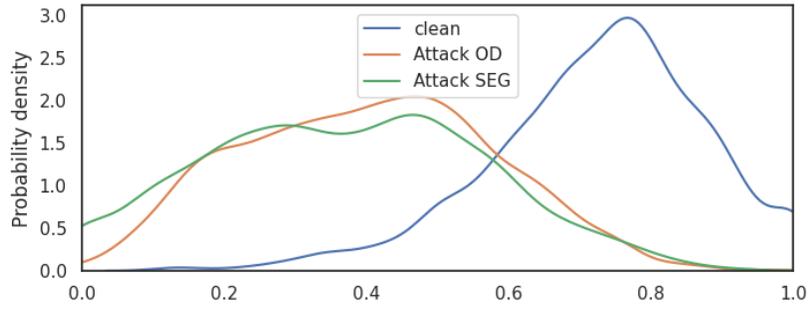


Figure 45: OD\_frncn\_swint\_SEG\_mask2former

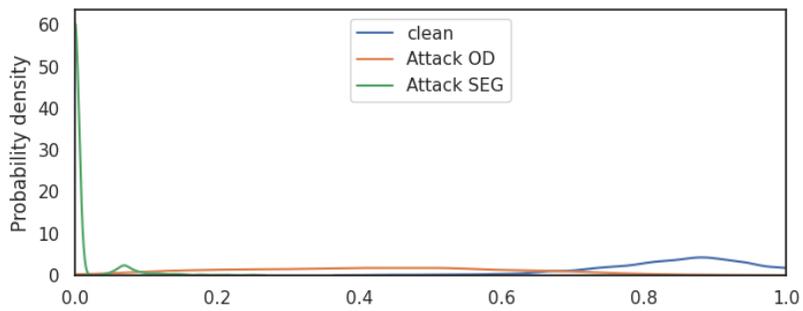


Figure 46: OD\_frncn\_swint\_SEG\_mrcnn\_r101

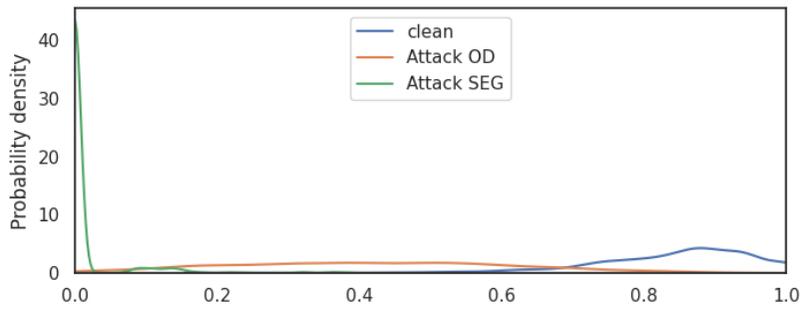


Figure 47: OD\_frncn\_swint\_SEG\_gcnet\_r101

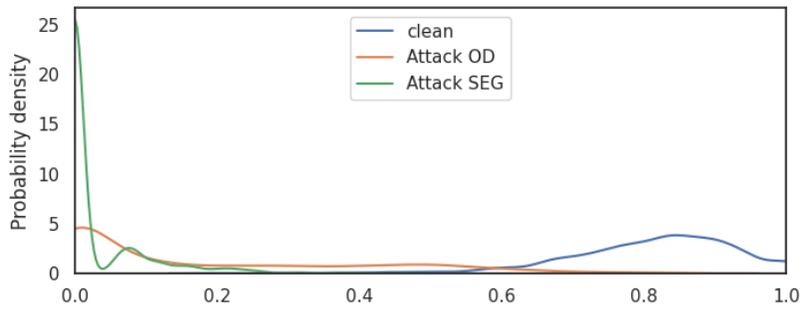


Figure 48: OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r50

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

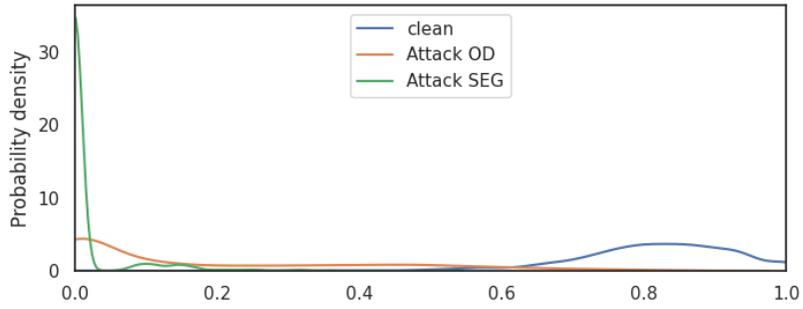


Figure 49: OD\_retinanet\_pvtv2\_SEG\_gcnet\_r50

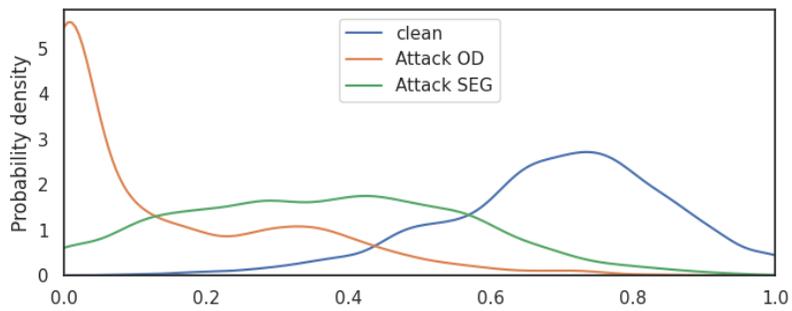


Figure 50: OD\_retinanet\_pvtv2\_SEG\_mask2former

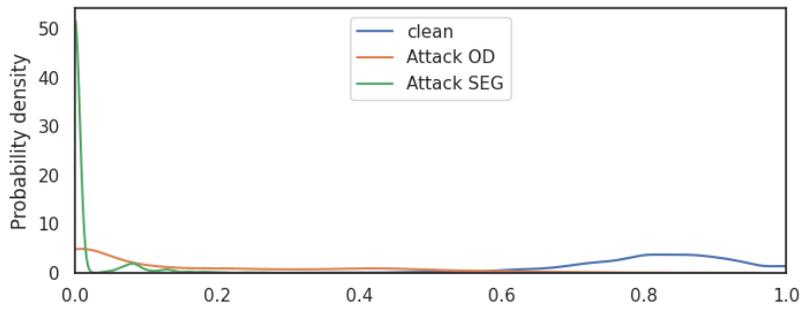


Figure 51: OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r101

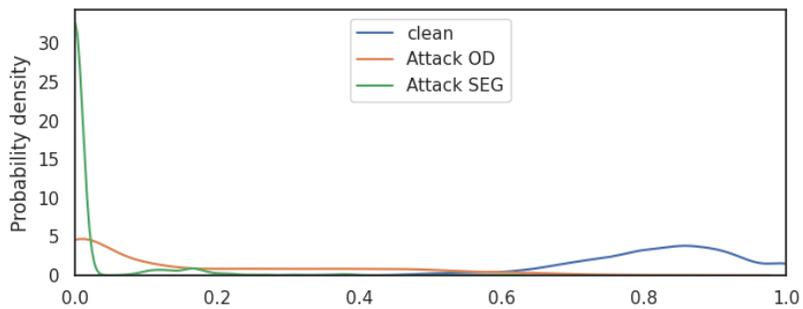


Figure 52: OD\_retinanet\_pvtv2\_SEG\_gcnet\_r101

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

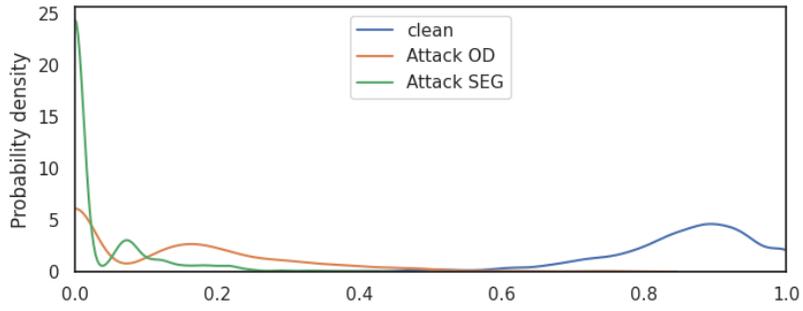


Figure 53: OD\_frcnn\_r101\_SEG\_mrcnn\_r50

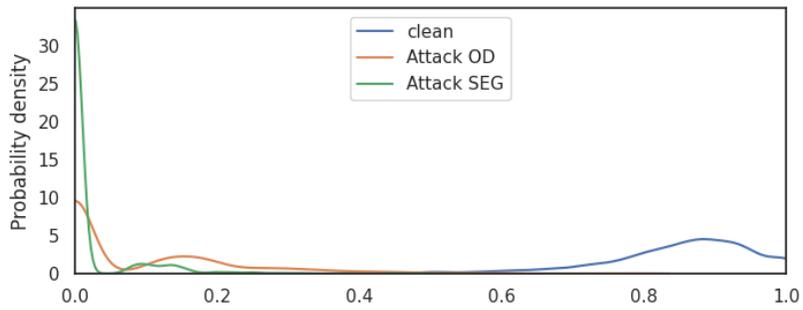


Figure 54: OD\_frcnn\_r101\_SEG\_gcnet\_r50

**Perturbation strength.** Perturbation strength affects the performance of detector using any model pair. As the perturbation strength increases, it results in stronger impact on the target model and cause higher inconsistency between model pairs.

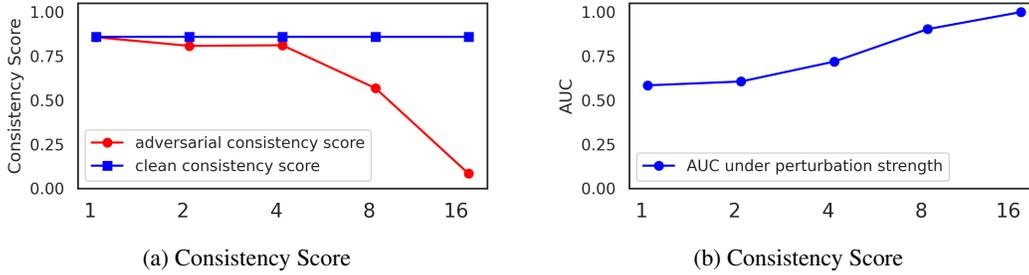


Figure 55: Impact of perturbation strength for OD\_frcnn\_r50\_SEG\_mrcnn\_r50

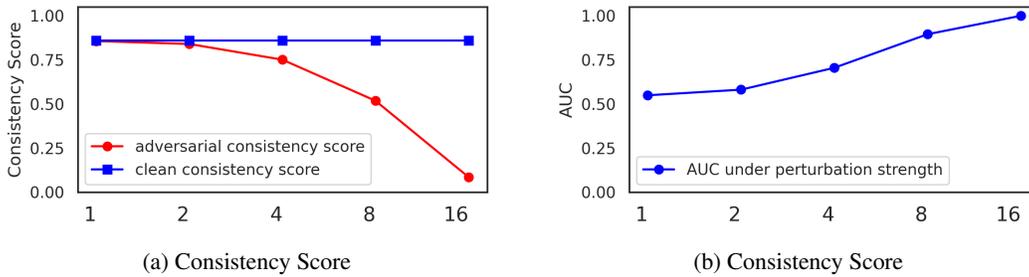


Figure 56: Impact of perturbation strength for OD\_frcnn\_r50\_SEG\_gcnet\_r50

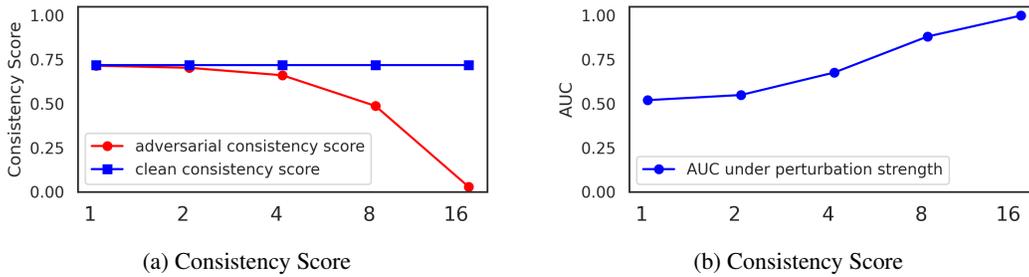


Figure 57: Impact of perturbation strength for OD\_frcnn\_r50\_SEG\_mask2former

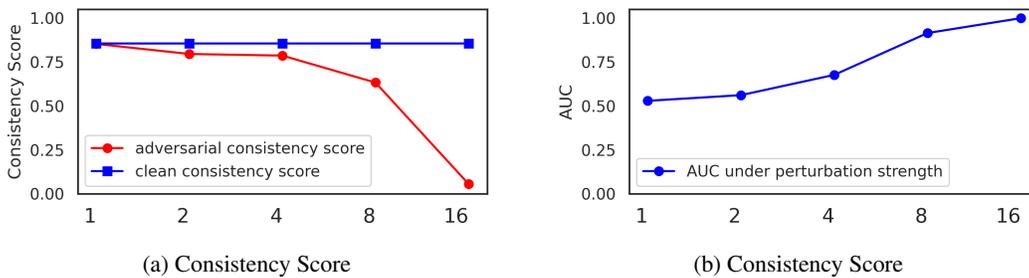


Figure 58: Impact of perturbation strength for OD\_frcnn\_r50\_SEG\_mrcnn\_r101

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

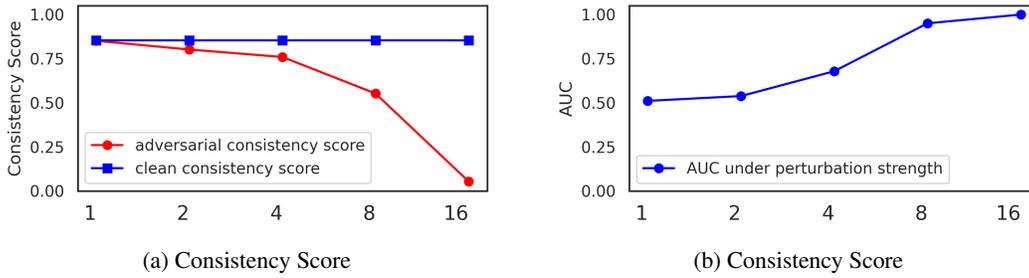


Figure 59: Impact of perturbation strength for OD\_frcnn\_r50\_SEG\_gcnet\_r101

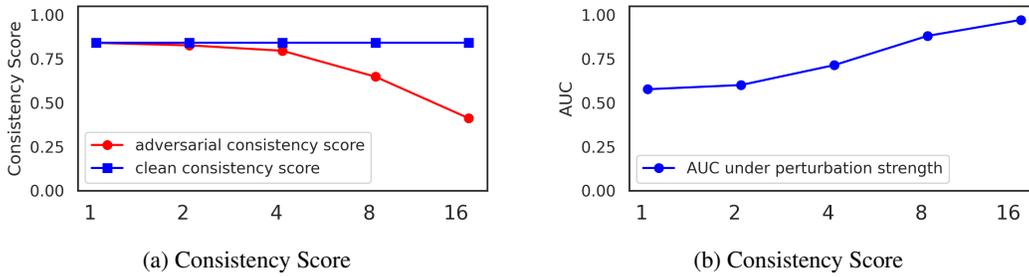


Figure 60: Impact of perturbation strength for OD\_frcnn\_swint\_SEG\_mrcnn\_r50

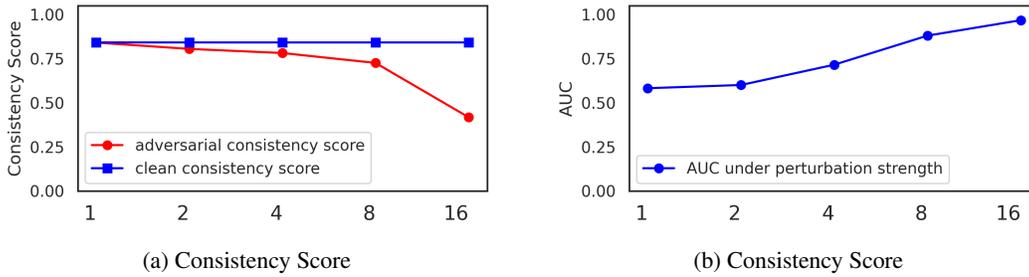


Figure 61: Impact of perturbation strength for OD\_frcnn\_swint\_SEG\_gcnet\_r50

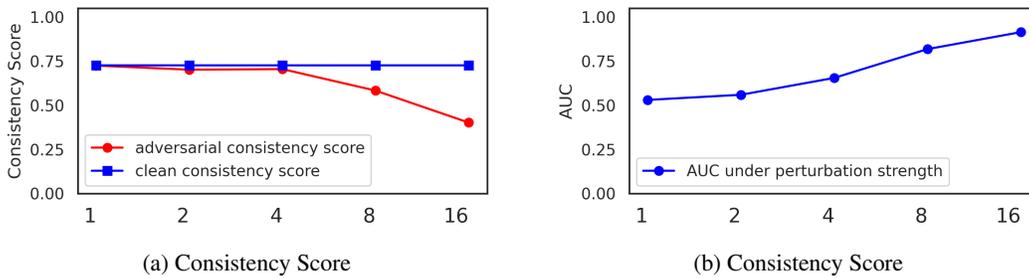


Figure 62: Impact of perturbation strength for OD\_frcnn\_swint\_SEG\_mask2former

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

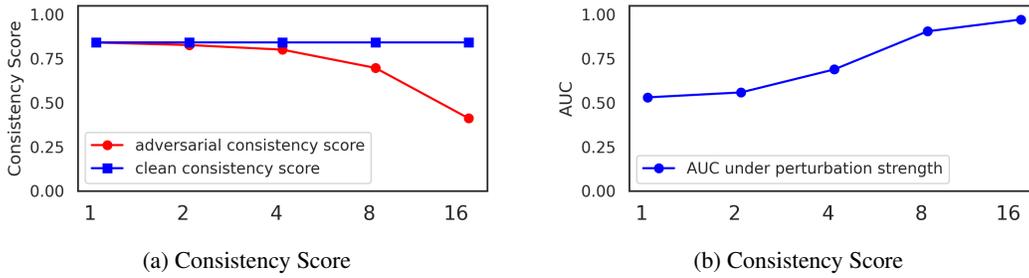


Figure 63: Impact of perturbation strength for OD\_frncn\_swint\_SEG\_mrcnn\_r101

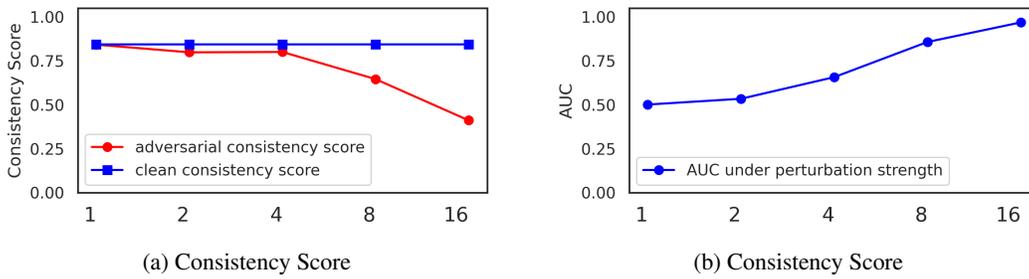


Figure 64: Impact of perturbation strength for OD\_frncn\_swint\_SEG\_gcnet\_r101

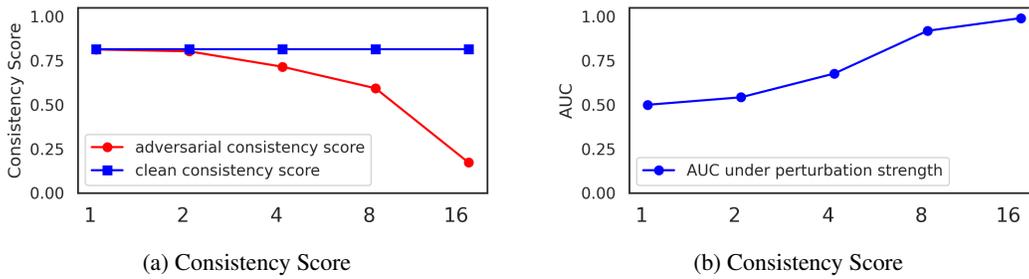


Figure 65: Impact of perturbation strength for OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r50

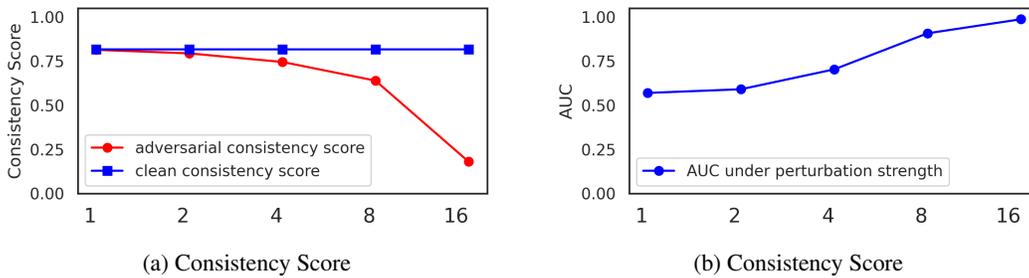


Figure 66: Impact of perturbation strength for OD\_retinanet\_pvtv2\_SEG\_gcnet\_r50

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

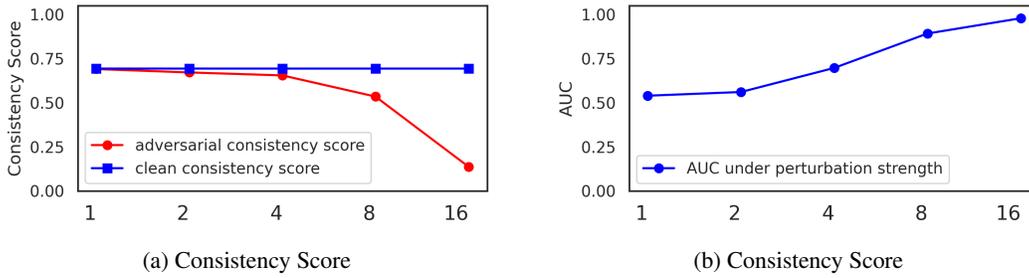


Figure 67: Impact of perturbation strength for OD\_retinanet\_pvtv2\_SEG\_mask2former

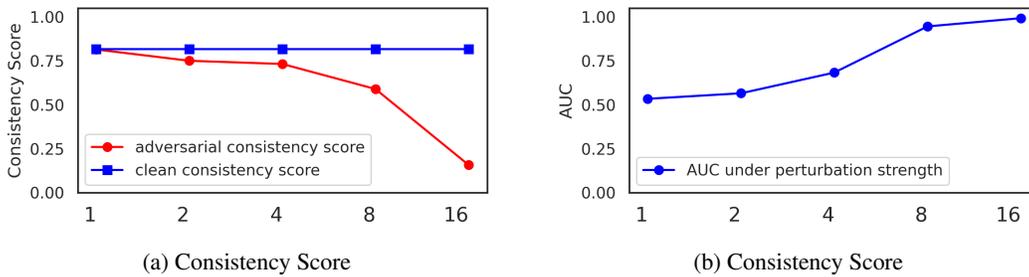


Figure 68: Impact of perturbation strength for OD\_retinanet\_pvtv2\_SEG\_mrcnn\_r101

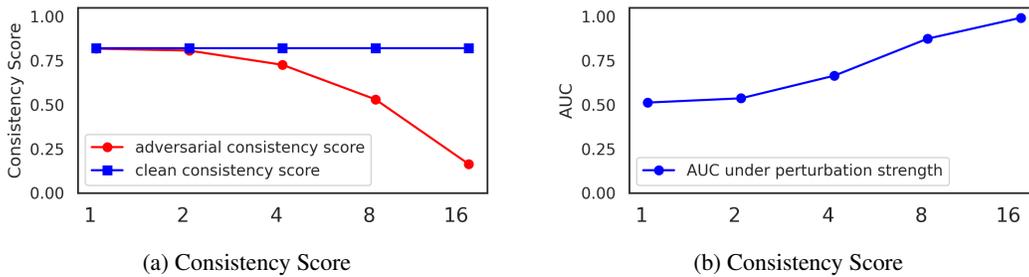


Figure 69: Impact of perturbation strength for OD\_retinanet\_pvtv2\_SEG\_gcnet\_r101

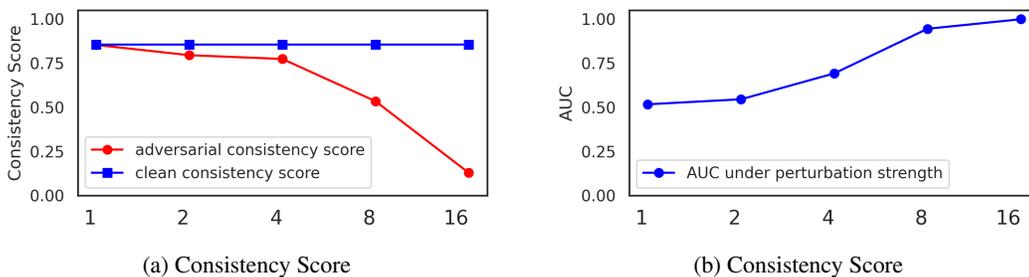


Figure 70: Impact of perturbation strength for OD\_frncnn\_r101\_SEG\_mrcnn\_r50

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

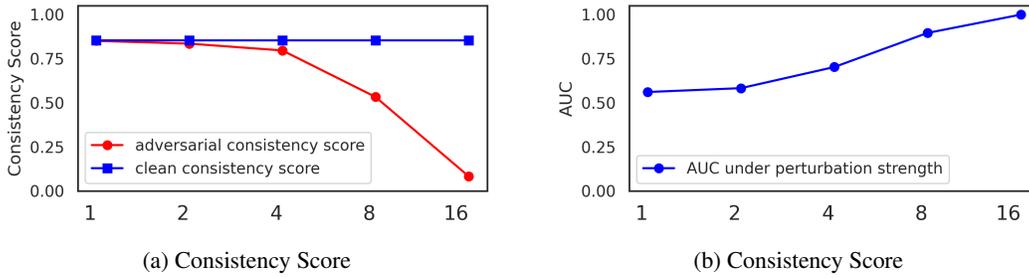


Figure 71: Impact of perturbation strength for OD\_frcnn\_r101\_SEG\_gcnet\_r50

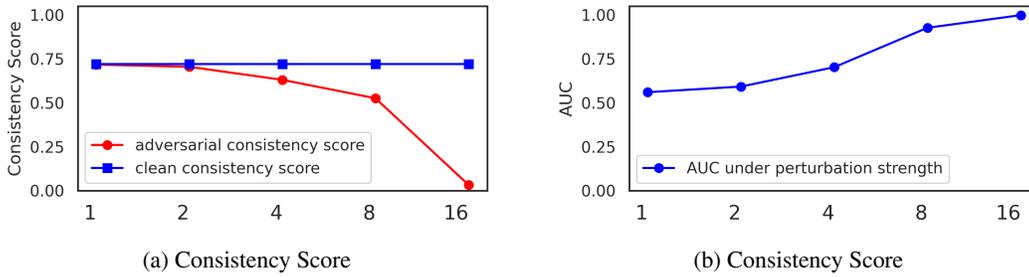


Figure 72: Impact of perturbation strength for OD\_frcnn\_r101\_SEG\_mask2former

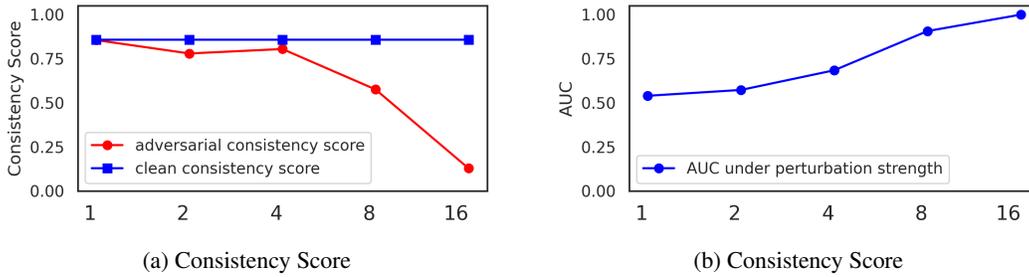


Figure 73: Impact of perturbation strength for OD\_frcnn\_r101\_SEG\_mrcnn\_r101

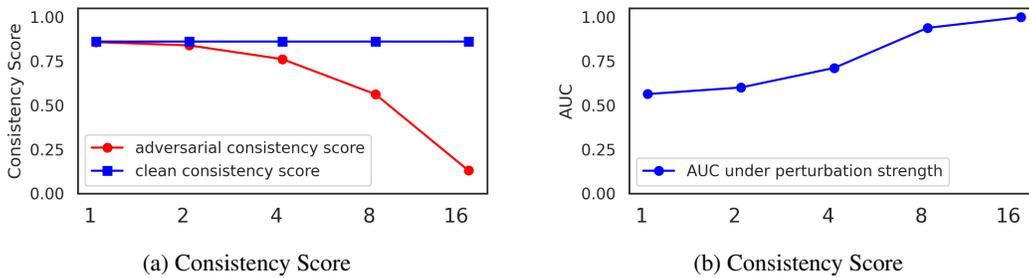


Figure 74: Impact of perturbation strength for OD\_frcnn\_r101\_SEG\_gcnet\_r101

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

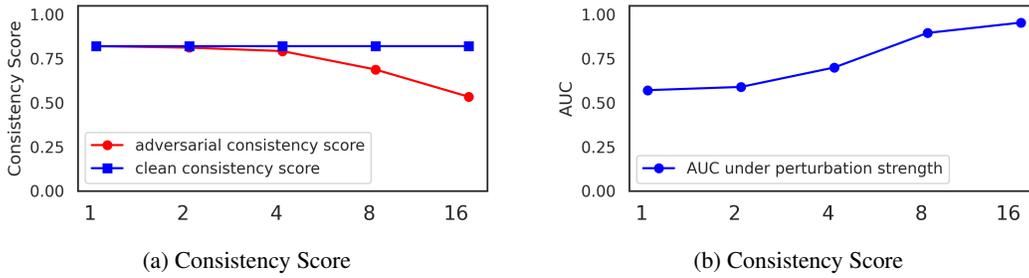


Figure 75: Impact of perturbation strength for OD\_retinanet\_r50\_SEG\_mrcnn\_r50

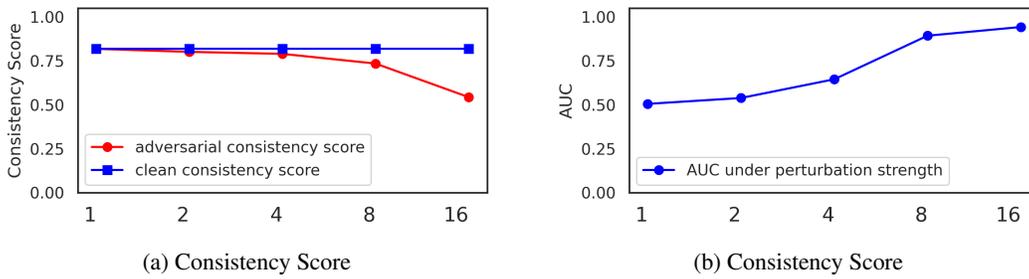


Figure 76: Impact of perturbation strength for OD\_retinanet\_r50\_SEG\_gcnet\_r50

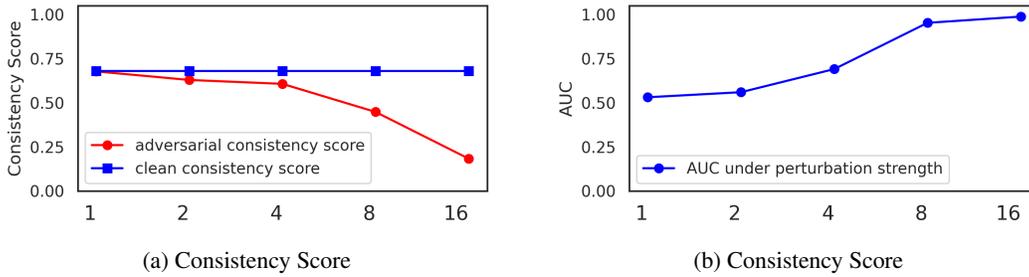


Figure 77: Impact of perturbation strength for OD\_retinanet\_r50\_SEG\_mask2former

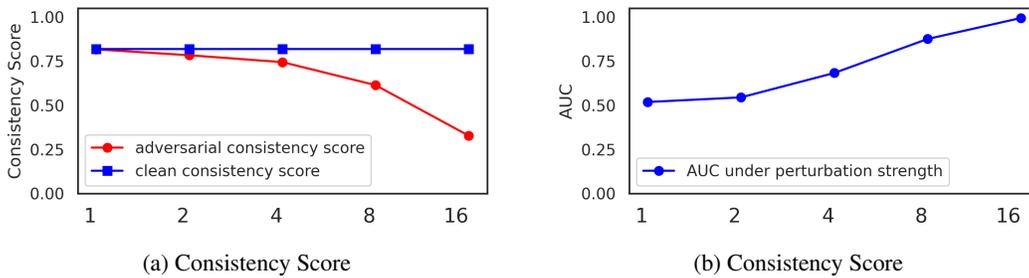


Figure 78: Impact of perturbation strength for OD\_retinanet\_r50\_SEG\_mrcnn\_r101

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

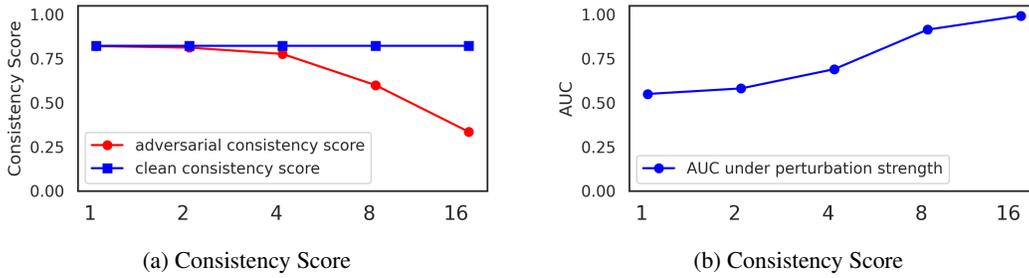


Figure 79: Impact of perturbation strength for OD\_retinanet\_r50\_SEG\_gcnet\_r101

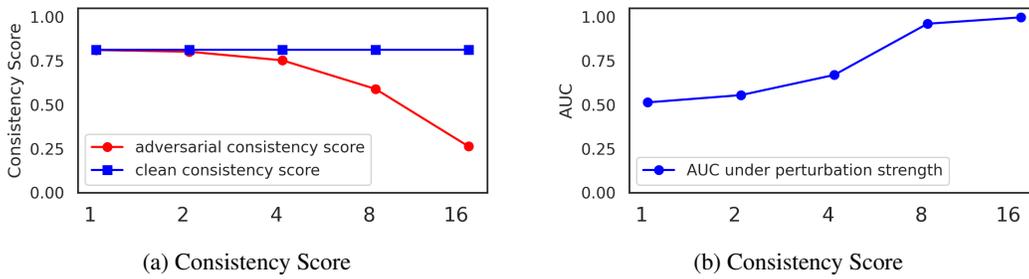


Figure 80: Impact of perturbation strength for OD\_retinanet\_r101\_SEG\_mrcnn\_r50

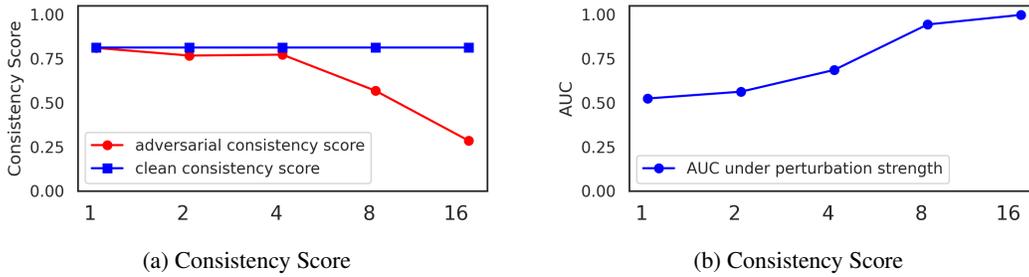


Figure 81: Impact of perturbation strength for OD\_retinanet\_r101\_SEG\_gcnet\_r50

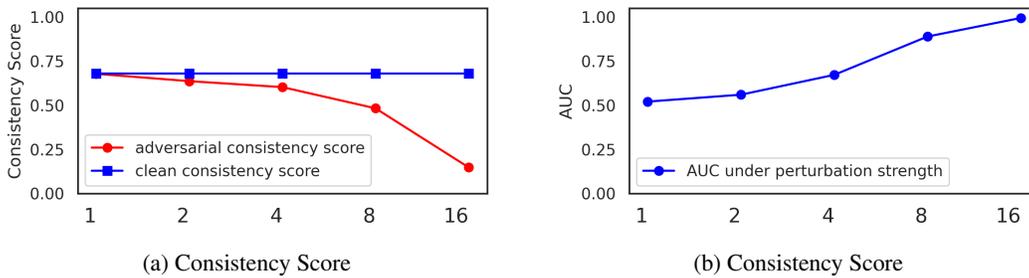


Figure 82: Impact of perturbation strength for OD\_retinanet\_r101\_SEG\_mask2former

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

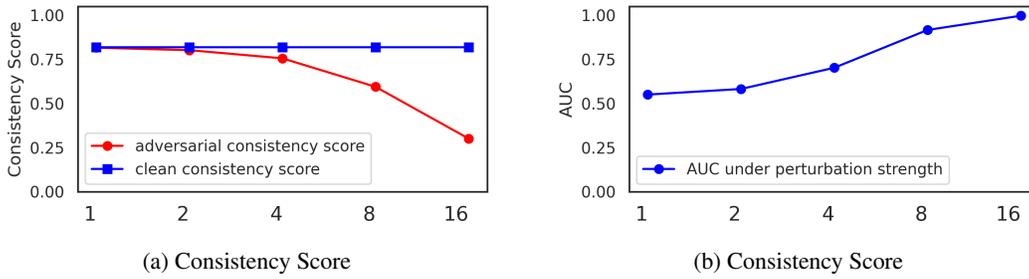


Figure 83: Impact of perturbation strength for OD\_retinanet\_r101\_SEG\_mrcnn\_r101

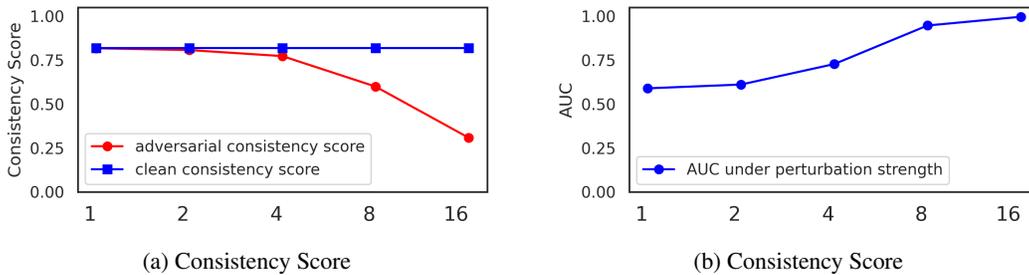


Figure 84: Impact of perturbation strength for OD\_retinanet\_r101\_SEG\_gcnet\_r101

## 1782 B OTHER DEFENSES

1783  
1784 To assess the benefits of our multi-task consistency detector, it is important to compare it against other  
1785 known defenses. Defenses usually fall into three categories: data-based, model-based, detection-  
1786 based. Data-based defenses use data augmentation at training to improve adversarial robustness.  
1787 Adversarial training is a common data-based defense (Qian et al., 2022). Model-based focuses on  
1788 selecting a specific network architecture that provides intrinsic adversarial robustness (Ye et al.,  
1789 2019; Guo et al., 2020a; Dong et al., 2022). Finally, detection-based defenses do not require data  
1790 augmentation or model changes, but add a processing (on the input or the output of the model) in  
1791 order to detect adversarial inputs. PatchCleanser (Xiang et al., 2022) is one example of a double  
1792 masking technique used to detect presence of adversarial patches in images.

### 1793 1794 B.1 ADVERSARIAL TRAINING

1795  
1796 The concept of adversarial training (AT) is to train a model on a dataset containing both genuine  
1797 and adversarial examples in order to build resilience against perturbations. Previous work such as  
1798 MTD (Zhang & Wang, 2019) and Class-Wise Adversarial Training (CWAT) (Chen et al., 2021)  
1799 defined loss functions to train the model to accurately localize and classify objects in an image despite  
1800 the presence of adversarial noise. Unfortunately, all AT schemes demonstrated a drop in model  
1801 accuracy, which is not desirable.

### 1802 1803 B.2 MODEL-BASED DEFENSE: ROBUST NETWORK

#### 1804 1805 B.2.1 ADVERSARIALLY-AWARE ROBUST OBJECT DETECTOR (ROBUSTDET)

1806  
1807 Dong et al. (2022) proposed a counter-proposal to adversarial training by modifying the model  
1808 architecture. The proposal, named RobustDet, aims to modify an existing backbone (e.g., SSD)  
1809 by adding three security components: an adversarial image discriminator (AID), an "adversarially-  
1810 aware convolution" (AAconv), and a consistent features with reconstruction (CFR). The AID is a  
1811 discriminator that outputs a probability vector based on the category of the image. For instance, if the  
1812 AID discriminates the image as genuine, then the AID will output the probability vector for a genuine  
1813 image. Otherwise, if the image is adversarial, then the AID will output the probability vector for an  
1814 adversarial image. For the training phase, the author formulated a dedicated loss function for the AID  
1815 to generate a probability vector specific to the category of the image (genuine or adversarial). This  
1816 probability vector will serve as an input for the next module: *AAconv*. Unlike in adversarial training,  
1817 *AAconv* aims to use specific weights for the model based on the category of the image. To achieve  
1818 this goal, *AAconv* uses the concept of dynamic convolution to generate different convolution kernels  
1819 based on the category of the image. The generation of those convolution kernels is possible thanks to  
1820 the (genuine or adversarial) probability vector provided by the *AID*. The probability vector serves  
1821 as the weights to generate convolution kernels. This approach allows to have dedicated weights for  
1822 genuine images and adversarial images instead of having a single set of weights for both categories  
of images (like in AT). Lastly, the CFR reconstructs the adversarial image into a clean image.

1823 Looking at their mAP evaluation, RobustDet has higher mAP scores than adversarial training methods  
1824 such as MTD and CWAT on both genuine and adversarial datasets. However, RobustDet still has at  
1825 best a 20 mAP score difference between the genuine dataset and the adversarial dataset. This issue  
1826 means RobustDet do not completely mitigate the mAP loss caused by adversarial examples.

#### 1827 1828 B.2.2 APPLICATION OF ROBUSTDET ON FASTER RCNN

1829  
1830 Dong et al. (2022) evaluated RobustDet on the object detection model SSD with a VGG16 backbone,  
1831 which was trained on the COCO dataset. However, in our paper, the models evaluated are trained on  
1832 BDD100k dataset and we do not use SSD. To fairly compare the performance of our consistency-  
1833 based detector with RobustDet, we applied the *AAconv* technique to one of the models, namely Faster  
1834 RCNN with ResNet50 backbone (FRCNN R50 in short).

1835 As previously explained, the core idea of *AAconv* is to replace any regular convolution kernel in a  
model network with a weighted sum of a set of dynamic convolution kernels, expressed as:

$$\dot{\theta}^{AA_{conv}} = \sum_{i=1}^M \theta_i^{AA_{conv}} \cdot \pi_i$$

where  $\theta_i^{AA_{conv}}$  is the  $i$ -th kernel in the set of  $M$  dynamic kernels, and  $\pi_i$  is its corresponding weight generated by *AID*. To clarify, each convolution kernel in the original network will be replaced by a unique set of dynamic kernels whose parameters are determined during the training phase. Thus, for each convolution layer in original Faster RCNN, we replace it with a dynamic convolution layer as defined by Dong et al. (2022). Regarding *AID*, we use the same network architecture Resnet18 as in the original paper. The performance of our Robust FRCNN is presented in Section 4.2.3.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889