

# Evaluating the Adversarial Robustness of Detection Transformers

Amirhossein Nazeri\*, Chunheng Zhao\* and Pierluigi Pisu

**Abstract**—Robust object detection is critical for autonomous driving and mobile robotics, where accurate detection of vehicles, pedestrians, and obstacles is essential for ensuring safety. Despite the advancements in object detection transformers (DETRs), their robustness against adversarial attacks remains underexplored. This paper presents a comprehensive evaluation of DETR model and its variants under both white-box and black-box adversarial attacks, using the MS-COCO and KITTI datasets to cover general and autonomous driving scenarios. We extend prominent white-box attack methods (FGSM, PGD, and C&W) to assess DETR’s vulnerability, demonstrating that DETR models are significantly susceptible to adversarial attacks, similar to traditional CNN-based detectors. Our extensive transferability analysis reveals high intra-network transferability among DETR variants, but limited cross-network transferability to CNN-based models. Additionally, we propose a novel untargeted attack designed specifically for DETR, exploiting its intermediate loss functions to induce misclassification with minimal perturbations. Visualizations of self-attention feature maps provide insights into how adversarial attacks affect the internal representations of DETR models. These findings reveal critical vulnerabilities in detection transformers under standard adversarial attacks, emphasizing the need for future research to enhance the robustness of transformer-based object detectors in safety-critical applications.

## I. INTRODUCTION

Object detection, a fundamental task in computer vision, plays a critical role in the rapidly evolving fields of autonomous driving and mobile robotics [1]. This technology forms the backbone of perception systems that enable self-driving vehicles and robots to detect and locate other vehicles, pedestrians, traffic signs, and potential obstacles in real-time. The dual challenge of predicting both bounding boxes and corresponding classes for objects of interest within an image is particularly crucial in these domains, where the reliability and accuracy of such detection can have profound implications for safety and operational efficiency [2].

Inspired by accomplishments of transformers in Natural Language Processing (NLP) [3], the application of transformers in computer vision has gained substantial traction. Vision Transformer (ViT) [4] and its variants have shown remarkable results in image classification, outperforming traditional Convolutional Neural Networks (CNNs). Their self-attention mechanism effectively captures complex, long-range dependencies within an image. Extending this approach to object detection, Carion et al. [5] introduced Detection Transformers (DETR), which eliminate the need for

hand-crafted components such as anchor generation and non-maximum suppression. By leveraging a simple transformer encoder-decoder structure connected to a CNN backbone for feature extraction, DETR allows for end-to-end training of object detectors with reduced computational complexity and fewer trainable parameters [5].

Despite deep learning’s widespread success in computer vision, studies highlight the vulnerability of Deep Neural Networks (DNNs) to well-crafted adversarial inputs, which are subtly modified data that remain nearly imperceptible but cause misclassifications by DNNs. This vulnerability raises significant concerns about the reliability of DNNs in critical applications like autonomous vehicles, where misdetection of obstacles could potentially lead to accidents. Szegedy et al. [6] demonstrated this vulnerability by adding small crafted perturbations that deceive state-of-the-art (SOTA) DNNs into causing incorrect image classifications. Numerous adversarial attacks have since been developed, and several are reported highly effective in compromising object detection performance [7], [8], [9]. Recent studies [10], [11] have explored the vulnerability of ViT under adversarial attacks. However, the adversarial robustness of Detection Transformer (DETR) remains underexplored, with no comprehensive evaluation of their performance under standard white-box and black-box attacks, nor any comparison between results on general object detection datasets and application-specific datasets. This gap in the literature is particularly concerning given the increasing adoption of transformer-based object detectors in critical applications such as autonomous driving and robotics.

In this paper, we present the first comprehensive study on the adversarial robustness of DETR and its variants using three prominent adversarial attacks (FGSM, PGD, and C&W) over two SOTA object detection datasets: MS-COCO and KITTI. While PGD and C&W were originally designed for image classification, they have demonstrated SOTA performance on object detection models [12], with PGD even outperforming SOTA methods [13]. By focusing on DETR as the foundational transformer-based object detection model, we isolate and examine the vulnerabilities of its core encoder-decoder components. This approach allows us to avoid the complexities introduced by variants like Anchor-DETR [14] and Efficient-DETR [15], which primarily address training convergence and efficiency rather than adversarial robustness. In summary, we present our key findings as follows: (1) Our results uncover significant vulnerabilities in DETR even to basic attacks, aligning with observations from CNN-based object detection models and vision transformers. (2) Extensive transferability analysis shows good intra-network transferability but limited cross-

\*Equal Contribution

The authors are with the Department of Automotive Engineering, Clemson University, 4 Research Drive, Greenville, SC 29607, USA (e-mail: anazeri@clemson.edu; chunhez@clemson.edu; pisup@clemson.edu)

The manuscript is under review, all rights are preserved.

network transferability in black-box settings; (3) Adversarial attacks can significantly alter self-attention feature maps, indicating that the attention mechanism is not as robust to adversarial examples as we might expect. (4) Through simple yet non-trivial modifications to the loss function, we propose our novel attack designed for DETR, which achieves SOTA performance with less visible perturbations on the COCO dataset. The code will be made available upon publication.

## II. RELATED WORK

Since 2021, the robustness of vision transformers for image classification task has garnered significant attention, with numerous studies investigating their performance under various adversarial attacks [16], [17], and [18]. The first work in evaluation of transformer-based image classification robustness by Mahmood et al, [10] multiple classifiers including ViTs [4] and BiTs [19] were tested against a variety of white-box and black-box attacks such as FGSM [8], PGD [20], BPDA [21], and C&W [22]. They also explored the transferability of white-box adversarial examples between traditional CNN-based and transformer-based image classifiers. While the adversarial robustness of image classification transformers has been extensively studied in the literature, the robustness of object detection transformers remains less-explored due to their intricate architectures and task's complexity.

Recently, Lovisotto et al, [23] demonstrated that global reasoning process in dot-product attention can be significantly vulnerable when subjected to well-crafted adversarial patch causing huge performance degradation of DETR object detector. However, a limited conclusion can be drawn on the adversarial robustness of DETR as the adversarial transferability and validation on DETR variants are missing. Moreover, digital patch attacks are of less interest due to their higher visibility, easier detectability, less transferability, and greater dependency on the patch location and input visibility that make them less practical for real-world scenarios. In 2023, Leng et al, [24] proposed an Object-Aware mechanism based on black-box FGSM attack on traditional object detectors and one detection transformer. However, the work did not provide a comprehensive study on the robustness of detection transformers, and lacked generalization on other adversarial attacks and comparing the results on a variety of datasets. Zhang et al. [25] evaluated robustness of various object detectors (e.g. including FasterRCNN [26], SSD [27], and Deformable DETR [28]) against a transferable physical attack crafted based on multi-scale attention maps. The paper's outcome regarding evaluation of robustness of DETR variant is limited to one dataset, one specific black-box attack being considered where the attack is transferred to Deformable-DETR and is not directly exploited the transformer model.

Although a limited number of studies strived to address the adversarial robustness of transformer-based object detection models, they primarily focus on black-box settings where attacks are not explicitly designed to compromise DETR models or use patch attacks, which differ significantly from

pixel-level attacks and belong to a separate category. To the best of our knowledge, our work presents the first comprehensive study that evaluates the adversarial robustness of DETR and its variants in both white-box and black-box settings. We conduct our analysis on a general object detection dataset (COCO [29]) and a scenario-specific dataset (KITTI [30]), exploring DETR's vulnerabilities across different applications. This comprehensive approach allows us to offer insights that are crucial for developing more robust transformer-based object detection models for real-world applications such as autonomous driving and mobile robotics.

## III. PRELIMINARIES

In this paper, we utilize DETR baseline variants with two CNN backbones, ResNet50 and ResNet101, and their dilation versions. The models we use are DETR-R50, DETR-R50-DC5, DETR-R101, and DETR-R101-DC5, as described in [5]. The DC5 variants include a dilated C5 stage in the last stage of the backbone and remove a stride from the first convolution of this stage. The final outputs of DETR transformer with respect to inputs  $x$  and network parameters  $\theta$  can be defined as:

$$O(\theta, x) = [O_p(\theta, x), O_b(\theta, x)] \\ = [\text{softmax}(P(\theta, x), \text{sigmoid}(B(\theta, x))] \quad (1)$$

where  $O(\theta, x)$  is the final outputs consisting of class probabilities  $O_p(\theta, x)$  and coordinates of bounding boxes  $O_b(\theta, x)$ .  $O_p(\theta, x)$  is computed by applying softmax functions to the logits  $P(\theta, x)$  while  $O_b(\theta, x)$  is computed by applying sigmoid functions to  $B(\theta, x)$ .  $P(\theta, x)$  and  $B(\theta, x)$  are outputs from Prediction feed-forward networks (FFNs), which are the linear modules on the top of transformer architecture and act as detection heads of DETR [5]. They are defined as:

$$\begin{cases} P(\theta, x) = \text{LL}(h_s) \\ B(\theta, x) = \text{MLP}(h_s) \end{cases} \quad (2)$$

where LL is a single linear layer network and MLP is a multi-layer perceptron network.  $h_s$  represents the hidden states from final decoder layer, while  $h_s^k$  represents the hidden states from the  $k^{th}$  decoder layer.

The labels are:

$$t = \{t_c, t_a\} \quad (3)$$

where  $t_c$  denotes the ground-truth classes and  $t_a$  denotes and ground-truth annotations.

## IV. ADVERSARIAL ATTACKS ON DETR

The primary goals of generating adversarial images are to (1) minimize the perturbation added to the original images and (2) induce misclassifications (i.e., altering object labels or detecting non-existent objects), subject to input data constraints. In this section, misclassifications are induced through untargeted attacks, where any incorrect class predictions are considered as successful attacks. We explore various methods for generating adversarial images, extend them to DETR, and evaluate their transferability across different models. Additionally, we introduce our proposed attack specific to DETR.

### A. White-box Attacks Extension to DETR

In this section, we select three prominent attacking methods originally developed for image classifiers and extend them to generate adversarial examples for DETR. These attacks assume that the adversary has white-box access to the deep learning model, meaning they have full knowledge of the model structure and weights, enabling computation of both outputs and gradients. FGSM [8] is selected for its simplicity and representativeness as a basic, one-step adversarial attack. Despite its simplicity, it effectively demonstrates the concept of adversarial attacks. For the stronger multi-step iterative attack, we select PGD [20]. C&W attack [22] is also considered as it performs extremely small perturbations while achieving high success rates. While newer attack methods exist [17], [7], we deliberately focus on these classic attacks to demonstrate that DETR models are susceptible even to basic and classic adversarial techniques, similar to CNN-based object detectors [13].

To extend these classification-based attacks to DETR, we focus on manipulating the model's class predictions. In this Sec. IV-A, we use only ground-truth classes  $t_c$  in our adaptations. While this approach doesn't explicitly account for bounding box predictions, it provides a foundation for understanding DETR's vulnerabilities to adversarial attacks.

The FGSM attack can be formulated as:

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x(-J(\theta, x, t_c))) \quad (4)$$

where  $x_{adv}$  is the generated adversarial image,  $\epsilon$  is the multiplier to control the perturbation size, and  $J(\theta, x, t_c)$  is a cross-entropy loss function dependent on the neural network parameters  $\theta$ , inputs  $x$  and targeted classes  $t_c$ . As an untargeted attack, we set  $J = -J$  to make the loss diverge.

The PGD attack can be formulated as a multi-step FGSM:

$$x_{adv}^{i+1} = \Pi_{x+S}(x_{adv}^i + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, t_c))) \quad (5)$$

where  $\Pi_{x+S}$  projects perturbations into the set  $S$  to limit the perturbation size at each step  $t$  and  $L_\infty$  is used as the distance metric.

The C&W attack can be formulated as:

$$\begin{cases} \text{Minimize } \|x_{adv} - x\|_2^2 + c \cdot f(x_{adv}) \\ f = \text{Max}(\text{Max}\{P(\theta, x_{adv})_j : j \neq t\} - P(\theta, x_{adv})_t, -\kappa) \end{cases} \quad (6)$$

Here, a Change of Variables [22] method is introduced such that  $x_{adv}$  is not directly optimized; instead, a new variable  $w$  is optimized, formulated as  $x_{adv} = \frac{1}{2} \cdot (\tanh(w) + 1)$ . This transformation ensures that the pixel values of  $x_{adv}$  remain in a valid range while allowing unconstrained optimization of  $w$ . The parameter  $\kappa$  encourages the solver to find an adversarial instance  $x_{adv}$  that is classified as class  $t$  with high confidence. The constant  $c$  balances between the two loss terms.

### B. Adversarial Attacks Transferability (Black-box Attacks)

In this subsection, we explore the transferability of adversarial attacks on DETR, an important aspect that can reveal potential vulnerabilities in practical applications (e.g.,

autonomous driving). Adversarial examples generated on one network have shown to be effective on other networks, a phenomenon known as adversarial attack transferability. This raises significant concerns for real-world deployments, as it suggests that DNNs may be vulnerable to adversarial attacks even when there is no direct access to the target network (i.e., black-box attacks). In black-box scenarios, attackers may train their own substitute model, generate adversarial examples against this substitute, and then transfer these examples to a victim model. This approach can be effective even with limited information about the victim model.

In this paper, we investigate both intra-network and cross-network transferability to further explore the potential vulnerabilities of DETR. Intra-network transferability refers to the effectiveness of adversarial examples when generated and evaluated using DETR and its variants. This type of transferability helps us understand how robust different versions of DETR are to attacks generated on similar architectures. Additionally, we conduct a preliminary cross-network transferability analysis by evaluating how adversarial examples generated with DETR perform when applied to a CNN-based object detector, specifically Faster R-CNN [26]. To quantify the effectiveness of these transfers, we formally define untargeted transferability as follows:

$$TR_{m,n} = \frac{AP_{clean}^m - AP_{adv(n)}^m}{AP_{clean}^n - AP_{adv(n)}^n} \quad (7)$$

where  $TR_{m,n}$  represents the transfer rate of adversarial examples generated on model  $n$  and tested on model  $m$ .  $adv(n)$  denotes the adversarial examples generated using model  $n$ .  $AP_{adv(n)}^m$  denotes the average precision evaluated on model  $m$  using  $adv(n)$ , while  $AP_{adv(n)}^n$  denotes the average precision evaluated on model  $n$  using  $adv(n)$ .

### C. Our Attack on DETR

In this subsection, we introduce a novel untargeted attack specifically designed for DETR object detection models assuming the same white-box adversary as detailed in Sec. IV-A. Our approach improves upon existing methods by combining an initial one-step perturbation with a multi-step process based on modified C&W attacks, considering DETR's intermediate loss functions.

Our attack method, as outlined in Algorithm 1, consists of two main stages: (1) initialize the attack with a one-step slight perturbation; (2) apply a multi-step C&W process with modified loss functions to the slightly-perturbed images. By using slightly-perturbed images as inputs, our multi-step attack generates perturbations based on already perturbed images, tending to maintain similarity to these slightly-perturbed images rather than the pure clean images. This approach ensures that even if the iterative procedures fail to produce effective perturbations, the initial one-step attack can compensate and provide a baseline level of adversarial effect. The updated input images (i.e., slightly-perturbed images)

---

**Algorithm 1** Our Attack on DETR

---

**Input:** initial image  $x$ , ground-truth label  $t$ , steps  $m$ , perturbation constant  $c$ , initial perturbation constant  $\alpha$ .

**Output:** adversarial image  $x_{adv}$

```
1: Initialize:
2:  $x \leftarrow x + \alpha \cdot \nabla_x(-J_{cls}(\theta, x, t_c))$ 
3:  $w_0 = \text{zeros}(x)$ 
4: for  $i = 0$  to  $m - 1$  do
5:    $x_{adv} = \frac{1}{2} \cdot (\tanh(w_i) + 1)$ 
6:    $Loss_{dm} = L_2(x_{adv}, x)$ 
7:    $Loss_{cls} = c \cdot f(x_{adv})$ 
8:    $Loss_{bb}^{\{o,k\}} = -J_{bb}^{\{o,k\}}(\theta, x_{adv}, t_a)$ 
9:    $Loss_{iou}^{\{o,k\}} = -J_{iou}^{\{o,k\}}(\theta, x_{adv}, t_a)$ 
10:  Update  $w_i$  with gradient decent
     $w_i \leftarrow \nabla_{w_i}(Loss_{total} = Loss_{dm} + Loss_{cls} +$ 
     $Loss_{bb}^{\{o,k\}} + Loss_{iou}^{\{o,k\}})$ 
11:  if  $Loss_{total}$  does not converge then
12:    return:  $x_{adv} = \frac{1}{2} \cdot (\tanh(w_i) + 1)$ 
13:  end if
14: end for
15: return:  $x_{adv} = \frac{1}{2} \cdot (\tanh(w_i) + 1)$ 
```

---

are generated as follows:

$$\begin{cases} x = x + \alpha \cdot \nabla_x(-J_{cls}(\theta, x, t_c)) \\ J_{cls}(\theta, x, t_c) = J_{cls}^o(\theta, x, t_c) + \sum_{k=1}^N J_{cls}^k(\theta, x, t_c) \\ J_{cls}^{\{o,k\}}(\theta, x, t_c) = L_{CE}(\theta, x, t_c) \end{cases} \quad (8)$$

where  $J_{cls}^{\{o,k\}}(\theta, x, t)$  computes the cross-entropy loss between predicted and targeted classes;  $J_{cls}^o$  represents the classification loss from the final output layer and  $J_{cls}^k$  denotes the classification loss from the  $k^{th}$  decoder layer;  $N$  is the total number of decoder layers. Gradients  $\nabla_x$  are computed with respect to the negative loss  $-J_{cls}$ , aiming to minimize the confidence score of the targeted classes.

After updating the input images, we leverage and customize C&W attacks with redesigned loss functions specifically targeting DETR architectures. The optimization process minimizes  $Loss_{total}$ , aiming to achieve two main objectives: minimizing perturbation and inducing misclassification. The total loss is composed of four components:

$$Loss_{total} = \omega_1 \cdot Loss_{dm} + \omega_2 \cdot Loss_{cls} + \omega_3 \cdot Loss_{bb}^{\{o,k\}} + \omega_4 \cdot Loss_{iou}^{\{o,k\}} \quad (9)$$

The optimal weight values  $\omega$  are determined using a grid search approach. Each component of the loss function is designed to achieve the following objectives:

**Reducing the difference between adversarial and input images:** The first objective aims to minimize perturbations by evaluating the  $L_2$  similarity between generated adversarial images  $x_{adv}$  and input images  $x$  (i.e., slightly-perturbed images).

$$Loss_{dm} = L_2(x_{adv}, x) = \|x_{adv} - x\|_2^2 \quad (10)$$

**Eliminating the detection of the targeted class.** The second objective focuses on inducing misclassifications. This objective is achieved through three sub-objectives:

- Maximizing the confidence score of untargeted classes.

$$Loss_{cls} = c \cdot f(x_{adv}) \quad (11)$$

where  $f(x_{adv})$  is defined as in Eq. (6), and  $\kappa = 0$  is set as in the original C&W attacks.

- Minimizing the regression score of bounding box coordinates of targeted classes:

$$Loss_{bb}^{\{o,k\}} = -(J_{bb}^o(\theta, x_{adv}, t_a) + \sum_{k=1}^N J_{bb}^k(\theta, x_{adv}, t_a)) \quad (12)$$

where  $J_{bb}^{\{o,k\}}(\theta, x_{adv}, t_a)$  computes the  $L_1$  regression loss from the output layer and  $k^{th}$  decoder layer, correspondingly.

- Minimizing the Intersection over Union (IoU) score of targeted classes:

$$Loss_{iou}^{\{o,k\}} = -(J_{iou}^o(\theta, x_{adv}, t_a) + \sum_{k=1}^N J_{iou}^k(\theta, x_{adv}, t_a)) \quad (13)$$

where  $J_{iou}^{\{o,k\}}(\theta, x_{adv}, t_a)$  computes the Generalized Intersection over Union (GIoU) loss from the output layer and  $k^{th}$  decoder layer, correspondingly.

Overall, the proposed attack considers not only the outputs from the last decoder layer  $h_s$  but also the outputs from the  $k^{th}$  intermediate decoder layers  $h_s^k$ , which is specifically designed for the transformer-based architecture of DETR. Unlike the attacks described in Sec. IV-A that only consider object classes, the last two loss components ( $loss_{bb}$  and  $loss_{iou}$ ) in this attack aim to cause the detector to incorrectly predict object locations, potentially resulting in misdetection. While our primary goal is to mislabel the targeted classes, incorrect bounding box placement can further degrade the detector's performance.

## V. EXPERIMENTS

### A. Experimental Setup

1) *White-box attacks setup:* **Datasets:** We employ two popular object detection benchmarks, the MS COCO 2017 [29] and KITTI Vision [30]. COCO is used as a general object detection benchmark with 80 categories, and our experiments are conducted on its validation set, which contains 5,000 images. KITTI is specifically used for autonomous driving or mobile robotics applications, containing 9 classes. For this experiment, we use a subset of it to build a validation set, which consists of 1871 images. These datasets are chosen to evaluate attacks on both general object detection tasks (COCO) and practical application scenarios (KITTI).

**Models:** We evaluate four DETR variants: DETR-R50, DETR-R50-DC5, DETR-R101, and DETR-R101-DC5. These models differ in their backbone architectures (ResNet-50 or ResNet-101) and the use of dilation technique (DC5). Further details can be found in [5].



**TABLE I.** White-box attacks against DETR models on COCO and KITTI. COCO evaluation metrics  $AP(IoU = 0.50 : 0.95)$  and  $AR(maxDets = 100)$  are reported.

COCO											
Models	Metric	Clean	FGSM			PGD		C&W			Ours
			$\epsilon = 0.03$	$\epsilon = 0.05$	$\epsilon = 0.1$	$\epsilon = 0.03$	$\epsilon = 0.1$	$c = 1$	$c = 3$	$c = 5$	
Detr-R50	$AP$	0.420	0.327	0.293	0.188	0.095	0.070	0.135	0.105	0.087	0.084
	$AR$	0.574	0.492	0.459	0.343	0.209	0.165	0.275	0.238	0.209	0.203
Detr-R50-DC5	$AP$	0.433	0.341	0.309	0.203	0.094	0.073	0.111	0.099	0.081	0.047
	$AR$	0.594	0.512	0.484	0.360	0.209	0.170	0.243	0.231	0.199	0.139
Detr-R101	$AP$	0.435	0.336	0.300	0.195	0.086	0.060	0.137	0.112	0.088	0.063
	$AR$	0.590	0.506	0.470	0.349	0.207	0.168	0.291	0.251	0.218	0.166
Detr-R101-DC5	$AP$	0.449	0.354	0.322	0.205	0.090	0.063	0.130	0.107	0.087	0.034
	$AR$	0.604	0.525	0.492	0.363	0.210	0.164	0.172	0.246	0.210	0.112
KITTI											
Models	Metric	Clean	FGSM			PGD		C&W			Ours
			$\epsilon = 0.03$	$\epsilon = 0.05$	$\epsilon = 0.1$	$\epsilon = 0.03$	$\epsilon = 0.1$	$c = 1$	$c = 3$	$c = 5$	
Detr-R50	$AP$	0.378	0.324	0.162	0.086	0.044	0.034	0.056	0.052	0.049	0.075
	$AR$	0.561	0.528	0.319	0.192	0.129	0.112	0.146	0.137	0.129	0.194
Detr-R50-DC5	$AP$	0.352	0.209	0.176	0.097	0.040	0.031	0.064	0.050	0.034	0.064
	$AR$	0.544	0.398	0.351	0.226	0.122	0.092	0.167	0.130	0.110	0.179
Detr-R101	$AP$	0.367	0.219	0.167	0.079	0.036	0.023	0.071	0.048	0.038	0.061
	$AR$	0.581	0.400	0.326	0.155	0.099	0.072	0.160	0.130	0.109	0.172
Detr-R101-DC5	$AP$	0.383	0.233	0.177	0.089	0.048	0.032	0.067	0.049	0.040	0.065
	$AR$	0.616	0.407	0.324	0.172	0.116	0.092	0.152	0.128	0.114	0.174

**FGSM:**  $\epsilon$  is set to 0.03, 0.05 and 0.1.

**PGD:**  $\epsilon$  is set to 0.03, and 0.1, with a total of 10 iterations. The  $L_\infty$  bounds are set to  $\pm 10/255$ .

**C&W:**  $c$  is set to 1, 3 and 5, with a total of 200 iterations.

**Ours:**  $\alpha$  is set to 0.3, and  $c$  is set to 0.8, with 200 iterations.

The perturbation size is limited based on perceptible levels to ensure the adversarial examples remain visually similar to the original images. In addition to standard object detection metrics like Average Precision ( $AP$ ) and Average Recall ( $AR$ ), we use another metric for assessing adversarial robustness, the Robustness Score ( $RS$ ). This metric evaluates a model's robustness against specific adversarial attacks, with higher  $RS$  values indicating stronger resistance to the adversarial attack.  $RS$  is defined as:  $RS = AP_{adv}/AP_{clean}$ , where  $AP_{adv}$  and  $AP_{clean}$  are average precision score with adversarial images and clean images, respectively.

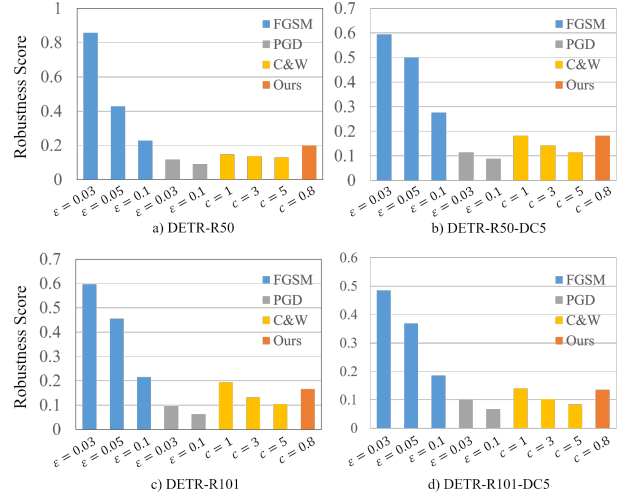


Fig. 2: Robustness Score of DETR variants on KITTI.

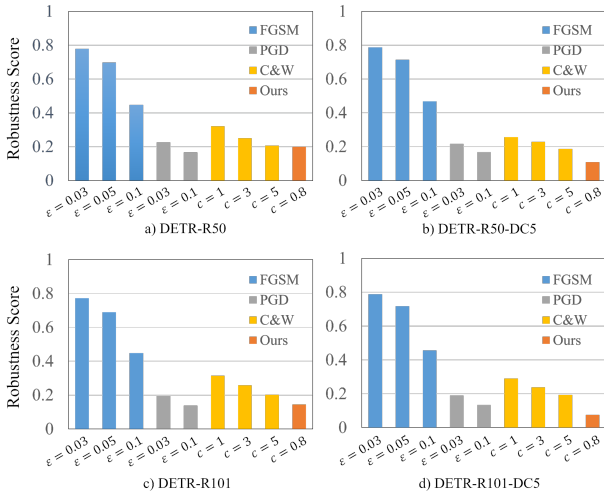


Fig. 1: Robustness Score of DETR variants on COCO.

2) *Transferability analysis setup:* For the transferability analysis, we use the same datasets, models, and attack parameters as described in Sec. V-A.1. For intra-network transferability, the same four DETR variants are used. As cross-network transferability is not the main focus in this paper, only one Faster R-CNN model (Faster R-CNN-R50-FPN) [26] is used as Faster R-CNN is one of the most representative CNN-based object detectors. This allows us to evaluate transferability across distinct object detection models due to the architectural differences between transformers and traditional CNNs.

To evaluate intra-network and cross-network transferability on the KITTI dataset, we retrain the four DETR variants along with the Faster R-CNN model. Since the KITTI object detection test set is unannotated, we split the original training set into new train and validation sets using a 3:1

**TABLE II.** Intra-network and cross-network transferability results on COCO and KITTI dataset. The models in different rows represent the model used to generate the adversarial examples. The models in different columns represent the model used to evaluate the adversarial examples. *TR* is reported as the primary metric.

COCO						
Models	Attacks	Detr-R50	Detr-R50-DC5	Detr-R101	Detr-R101-DC5	Faster R-CNN
Detr-R50	FGSM( $\epsilon = 0.05$ )	100%	96.8%	71.8%	65.3%	66.9%
	PGD( $\epsilon = 0.03$ )	100%	109.0%	106.2%	109.3%	54.8%
	C&W( $c = 3$ )	100%	99.1%	89.5%	80%	66.7%
Detr-R50-DC5	FGSM( $\epsilon = 0.05$ )	107.5%	100%	78.3%	70.8%	70.0%
	PGD( $\epsilon = 0.03$ )	101.5%	100%	101.8%	104.7%	51.6%
	C&W( $c = 3$ )	94.2%	100%	85.5%	87.3%	69.1%
Detr-R101	FGSM( $\epsilon = 0.05$ )	86.4%	78.8%	100%	90.9%	70.5%
	PGD( $\epsilon = 0.03$ )	96.8%	100.1%	100%	107.4%	53.3%
	C&W( $c = 3$ )	82.5%	56.1%	100%	86.9%	59.0%
Detr-R101-DC5	FGSM( $\epsilon = 0.05$ )	89.7%	82.5%	104.0%	100%	71.4%
	PGD( $\epsilon = 0.03$ )	93.6%	96.1%	100.3%	100%	51.0%
	C&W( $c = 3$ )	79.8%	81.0%	91.2%	100%	67.8%
KITTI						
Models	Attacks	Detr-R50	Detr-R50-DC5	Detr-R101	Detr-R101-DC5	Faster R-CNN
Detr-R50	FGSM( $\epsilon = 0.05$ )	100%	77.9%	73.6%	59.6%	84.1%
	PGD( $\epsilon = 0.03$ )	100%	92.8%	95.8%	98.8%	77.6%
	C&W( $c = 3$ )	100%	88.7%	88.1%	87.5%	72.2%
Detr-R50-DC5	FGSM( $\epsilon = 0.05$ )	116.6%	100%	89.9%	76.3%	105.9%
	PGD( $\epsilon = 0.03$ )	110.2%	100%	104.2%	108.6%	85.9%
	C&W( $c = 3$ )	107.3%	100%	99.0%	96.0%	78.8%
Detr-R101	FGSM( $\epsilon = 0.05$ )	81.8%	62.6%	100%	93.4%	83.8%
	PGD( $\epsilon = 0.03$ )	101.5%	92.4%	100%	104.8%	83.7%
	C&W( $c = 3$ )	94.7%	86.3%	100%	100%	76.0%
Detr-R101-DC5	FGSM( $\epsilon = 0.05$ )	84.5%	69.5%	106%	100%	79.5%
	PGD( $\epsilon = 0.03$ )	99.7%	89.8%	98.5%	100%	82.6%
	C&W( $c = 3$ )	93.1%	85.3%	97.3%	100%	72.5%

ratio. The four DETR variants are retrained for 25 epochs, with a transformer learning rate of  $1e-5$ , following the recommended recipe by [5]. Faster R-CNN is retrained for 25 epochs, with a learning rate of  $1e-3$ .

### B. Quantitative Evaluation

1) *White-box attacks analysis:* We evaluate the robustness of four DETR models against four adversarial attacks on both COCO and KITTI datasets. The results are summarized in Tab. I, Fig. 1, and Fig. 2, showing Average Precision (*AP*), Average Recall (*AR*) and Robustness Score (*RS*) for each model-attack combination. In summary, our results indicate that DETR models, similar to vision transformers and CNN-based object detectors, remain vulnerable to adversarial attacks. FGSM, as the most basic attacking algorithm, shows limited success rate against all four DETR models. PGD, as an iterative attack, proves to be the most effective, significantly decreasing *AP* and *RS*. For instance, PGD attack on DETR-R50 achieves the lowest *AP* across different attacks ( $AP = 0.070$  on COCO,  $AP = 0.034$  on KITTI). However, PGD generates relatively large perturbation sizes, as visualized in Fig. 3 and Fig. 4. C&W attack achieves a balance between performance degradation and perturbation size. These results demonstrate that DETR models have significant vulnerabilities even against classic attacks.

In terms of our attack, it achieves SOTA performance comparable to PGD attacks with smaller perturbations on the COCO dataset. For example, our attack has  $AP = 0.047$  and  $AP = 0.034$  on DETR-R50-DC5 and DETR-R101-DC5 while PGD has  $AP = 0.073$  and  $AP = 0.060$

**TABLE III.** Transferability results of our attack on COCO dataset. *TR* is reported as the primary metric.

Models	Detr				Faster R-CNN
	R50	R50-DC5	R101	R101-DC5	
Detr-R50	100%	99.5%	95.8%	87.1%	43.5%

correspondingly. On the KITTI dataset, our attack achieves results comparable to C&W attacks. While not as effective as PGD attacks, our method produces much smaller perturbations, as shown in Fig. 3 and Fig. 4. With simple yet non-trivial modifications, our attack further explores DETR’s adversarial vulnerabilities, indicating the need for developing more robust DETR models. Notably, the effects of dilation and differences between backbones are not obvious on both datasets.

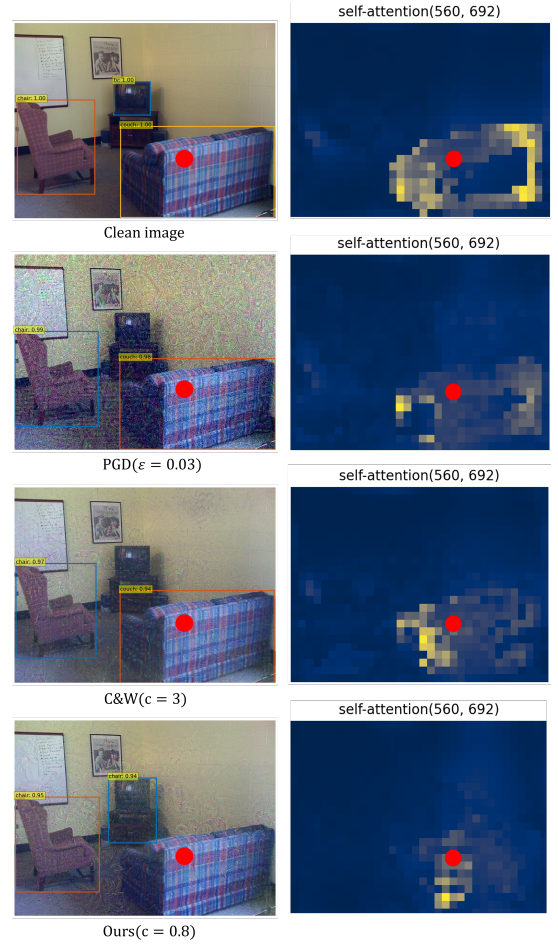


Fig. 3: Object detection on a sample image from COCO dataset with self-attention feature maps from the last encoder.

2) *Transferability analysis:* Transferability results for FGSM, PGD and C&W attacks are presented in Tab. II, while transferability results for our attack are shown in Tab. III. In summary, attacks generated using DETR models can transfer across DETR variants but show limited transferability to Faster R-CNN models. For intra-network transferability, results demonstrate that generated adversarial examples can transfer within the DETR family, although effectiveness

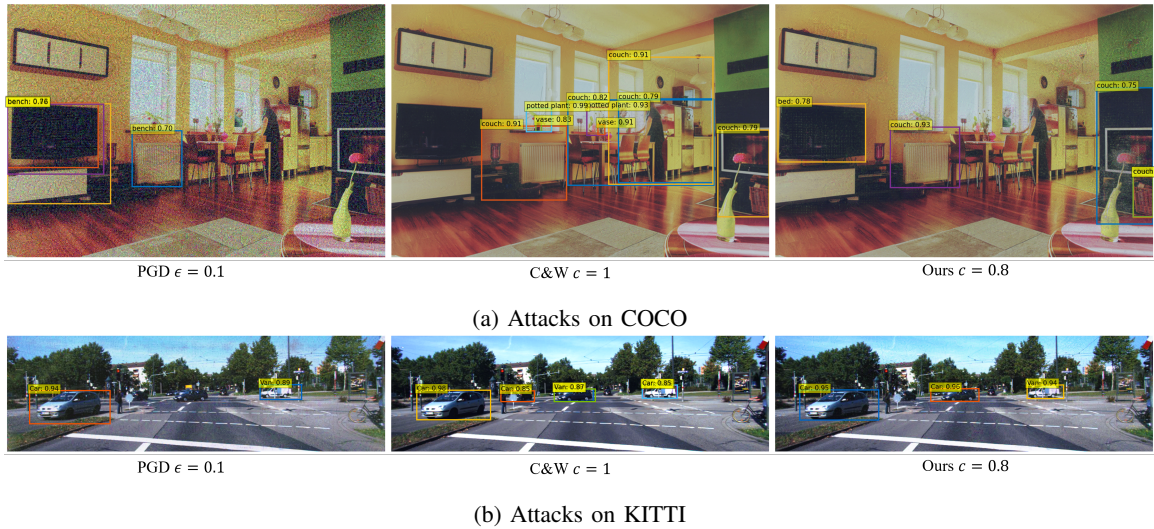


Fig. 4: Object detection results on sample images under various adversarial attacks.

decreases with increasing model complexity. For example, attacks generated on DETR-R50 show decreasing performance when transferred to DETR-R50-DC5, DETR-R101, and DETR-R101-DC5. Interestingly, attacks generated on models with dilation transfer well to their base models. For instance, FGSM attacks from DETR-R50-DC5 and DETR-R101-DC5 achieve better results on DETR-R50 (107.5% transfer rate) and DETR-R101 (104.0% transfer rate), respectively. For cross-network transferability, all the attacks generated on DETR variants show limited transferability to Faster R-CNN. This suggests a potential ensemble defense technique against adversarial attacks on DETR models. The KITTI dataset shows slightly better transferability to Faster R-CNN compared to COCO, possibly due to Faster R-CNN’s high clean  $AP$  (0.553) on KITTI after retraining, suggesting a potential trade-off between accuracy and robustness.

Among all attacks, PGD demonstrates the best intra-network transferability, while FGSM shows the best cross-network transferability. This might be due to FGSM’s one-step characteristic, which may rely less on DETR-specific gradients, allowing better transfer to non-DETR models. We argue that although FGSM is regarded as the simplest attack in the literature, it has potential benefits while conducting black-box attacks. For brevity, we conduct transferability analysis of our attack only on the COCO dataset with DETR-R50, as shown in Tab. III. Our attack shows good transferability across DETR variants but limited cross-network transferability. This is expected as our attack is specifically designed for DETR, leveraging its intermediate loss functions, which differ significantly from those in Faster R-CNN.

### C. Qualitative Evaluation

We provide a visual comparison of different attacks on sample images from the COCO and KITTI dataset as shown in Fig. 3 and Fig. 4. Our attack achieves a small perturbation level comparable to C&W attacks, despite employing a two-stage perturbation generation procedure. Notably, when

compared to PGD attacks, our method produces significantly less visible perturbations while maintaining similar performance on the COCO dataset. In contrast to C&W attacks, our approach provides superior performance with a similar perturbation level. Overall, our attack offers a balance between perturbation size and attack success rate.

To gain deeper insights into the impact of these attacks on DETR’s internal representations, we examine the self-attention feature maps from the last encoder layer for each attack, as illustrated in Figure 3. The clean image’s feature map clearly highlights the entire object shape, indicating that DETR correctly focuses on relevant image regions for object detection. However, both PGD and C&W attacks cause a noticeable shrinkage of the highlighted regions in the feature maps, suggesting successful disruption of DETR’s attention mechanism. Remarkably, our proposed attack results in a feature map with the least salient features, indicating that it most significantly impacts the model’s predictions. The ability of our attack to significantly alter the self-attention feature map while maintaining small perturbations underscores its effectiveness in exploiting DETR’s architectural vulnerabilities. These results lead us to conclude that the attention mechanism in DETR, contrary to expectations, does not provide sufficient protection against adversarial attacks.

## VI. CONCLUSION

The expanding use of detection transformers in critical applications, such as autonomous driving and robotics, raises significant concerns about their security and reliability. In this paper, we conducted the first comprehensive study on the adversarial vulnerability of DETR and its variants. We extended one basic adversarial attack (FGSM) and two classic yet strong attacks (PGD and C&W) to DETR models. Our results revealed substantial vulnerabilities in DETR models even against basic and classic attacks, consistent with the vulnerabilities observed in CNN-based object detection models and transformer-based image classification models.

In addition, we found strong transferability of adversarial examples generated by DETR models across different DETR variants, indicating shared vulnerabilities within similar architectures. However, the transferability to CNN-based object detectors like Faster R-CNN is limited. This highlights that employing an ensemble of diverse models may help mitigate the impact of adversarial attacks on DETR models.

We also introduced a novel adversarial attack leveraging the intermediate loss functions of DETR. Our attack is simple yet effective, causing significant performance degradation with less visible perturbations. Our overall results, including insights from self-attention map features, exhibited that the attention mechanism cannot provide additional protection against adversarial attacks in object detection tasks, contrary to expectations. These findings have important implications for the deployment of DETR models in safety-critical applications such as autonomous driving and robotics, where robustness to adversarial attacks is crucial. Future work could focus on exploring attacks transferability to single-stage object detector (e.g., YOLO) and other DETR variants. In addition, analyzing the impact of attacks on autonomous driving systems with DETR-based perception modules would offer valuable insights for real-world applications.

#### ACKNOWLEDGEMENT

This work is funded by the National Science Foundation CNS No. 2200457. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

#### REFERENCES

- [1] A. Balasubramanian and S. Pasricha, "Object detection in autonomous vehicles: Status and open challenges," *arXiv preprint arXiv:2201.07706*, 2022.
- [2] C. Michaelis, B. Mitzkus, R. Geirhos, E. Rusak, O. Bringmann, A. S. Ecker, M. Bethge, and W. Brendel, "Benchmarking robustness in object detection: Autonomous driving when winter is coming," *arXiv preprint arXiv:1907.07484*, 2019.
- [3] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [4] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [5] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.
- [6] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.
- [7] Z. Zhang and T. Wu, "Learning ordered top-k adversarial attacks via adversarial distillation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 776–777.
- [8] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [9] N. Papernot, P. McDaniel, A. Sinha, and M. Wellman, "Towards the science of security and privacy in machine learning," *arXiv preprint arXiv:1611.03814*, 2016.
- [10] K. Mahmood, R. Mahmood, and M. Van Dijk, "On the robustness of vision transformers to adversarial examples," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 7838–7847.
- [11] R. Shao, Z. Shi, J. Yi, P.-Y. Chen, and C.-J. Hsieh, "On the adversarial robustness of vision transformers," *arXiv preprint arXiv:2103.15670*, 2021.
- [12] T. Du, S. Ji, B. Wang, S. He, J. Li, B. Li, T. Wei, Y. Jia, R. Beyah, and T. Wang, "Detects ec: Evaluating the robustness of object detection models to adversarial attacks," *International Journal of Intelligent Systems*, vol. 37, no. 9, pp. 6463–6492, 2022.
- [13] Y. Wang, K. Wang, Z. Zhu, and F.-Y. Wang, "Adversarial attacks on faster r-cnn object detector," *Neurocomputing*, vol. 382, pp. 87–95, 2020.
- [14] Y. Wang, X. Zhang, T. Yang, and J. Sun, "Anchor detr: Query design for transformer-based detector," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 36, no. 3, 2022, pp. 2567–2575.
- [15] Z. Yao, J. Ai, B. Li, and C. Zhang, "Efficient detr: improving end-to-end object detector with dense prior," *arXiv preprint arXiv:2104.01318*, 2021.
- [16] J. Zhang, Y. Huang, W. Wu, and M. R. Lyu, "Transferable adversarial attacks on vision transformers with token gradient regularization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16 415–16 424.
- [17] Z. Wei, J. Chen, M. Goldblum, Z. Wu, T. Goldstein, and Y.-G. Jiang, "Towards transferable adversarial attacks on vision transformers," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 2668–2676.
- [18] Y. Wang, J. Wang, Z. Yin, R. Gong, J. Wang, A. Liu, and X. Liu, "Generating transferable adversarial examples against vision transformers," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 5181–5190.
- [19] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby, "Big transfer (bit): General visual representation learning," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*. Springer, 2020, pp. 491–507.
- [20] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
- [21] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," in *International conference on machine learning*. PMLR, 2018, pp. 274–283.
- [22] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.
- [23] G. Lovisotto, N. Finnie, M. Munoz, C. K. Mummadi, and J. H. Metzen, "Give me your attention: Dot-product attention considered harmful for adversarial patch robustness," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 234–15 243.
- [24] Z. Leng, Z. Cheng, P. Wei, and J. Chen, "Object-aware transfer-based black-box adversarial attack on object detector," in *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*. Springer, 2023, pp. 278–289.
- [25] Y. Zhang, Z. Gong, Y. Zhang, Y. Li, K. Bin, J. Qi, W. Xue, and P. Zhong, "Transferable physical attack against object detection with separable attention," *arXiv preprint arXiv:2205.09592*, 2022.
- [26] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [27] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 21–37.
- [28] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.
- [29] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [30] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.