

# Investigating the Robustness and Properties of Detection Transformers (DETR) Toward Difficult Images

Zhao Ning Zou

The Australian National University  
Canberra, ACT

ZhaoNing.Zou@anu.edu.au

Yuhang Zhang

The Australian National University  
Canberra, ACT

Yuhang.Zhang@anu.edu.au

Robert Wijaya

The Australian National University  
Canberra, ACT

Robert.Wijaya@anu.edu.au

## Abstract

*Transformer-based object detectors (DETR) have shown significant performance across machine vision tasks, ultimately in object detection. This detector is based on a self-attention mechanism along with the transformer encoder-decoder architecture to capture the global context in the image. The critical issue to be addressed is how this model architecture can handle different image nuisances, such as occlusion and adversarial perturbations. We studied this issue by measuring the performance of DETR with different experiments and benchmarking the network with convolutional neural network (CNN) based detectors like YOLO and Faster-RCNN. We found that DETR performs well when it comes to resistance to interference from information loss in occlusion images. Despite that, we found that the adversarial stickers put on the image require the network to produce a new unnecessary set of keys, queries, and values, which in most cases, results in a misdirection of the network. DETR also performed poorer than YOLOv5 in the image corruption benchmark. Furthermore, we found that DETR depends heavily on the main query when making a prediction, which leads to imbalanced contributions between queries since the main query receives most of the gradient flow.*

## 1. Introduction

With the continuous development of deep learning, computer vision has reached a new stage, and target detection, as one of the very vital core directions, has also received attention and many applications based on object detection algorithms.

Before the concept of deep learning was introduced, object detection was mostly based on manual feature extraction. However, as manual feature extraction methods often failed to meet various features in the targets, traditional target detection algorithms could not meet people's needs. After the rise of deep learning, neural networks can automatically learn powerful feature extraction and fitting capabilities from large amounts of data. Thus many DL-based object detectors with excellent performance have emerged. These detectors can be broadly classified into three categories: two-stage object detection, one-stage object detection, and transformer-based object detection.

Faster R-CNN [1] is the most popular two-stage detector nowadays. It first generates a proposal for the object bounding box in the image through a network. It then extracts features from each candidate box, and uses them for object classification and bounding box regression tasks to obtain the final bounding box. On the other hand, the YOLO [2] series model is well-known as a one-stage detector, which discards the anchor frame setting in two-stage and extracts the prediction frames directly from the image. Furthermore, with the increasing popularity of transformers applied in computer vision tasks, new transformer-based object detectors, such as DETR [9], have also emerged. Instead of using anchor frames and NMS, DETR uses an encoder-decoder structure to classify each object in the image.

DETR is an end-to-end target detection network proposed by Facebook in 2020. Compared to traditional RNNs, DETR uses multiple self-attentive structures, and the parallel computing used therein allows DETR to extract relevance efficiently in context. It is also one of the best-performing target detection methods available.

This paper investigated the properties of DETR and com-









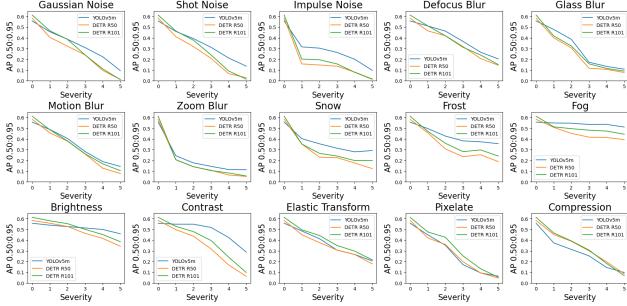


Figure 8. The benchmark results for three object detectors.

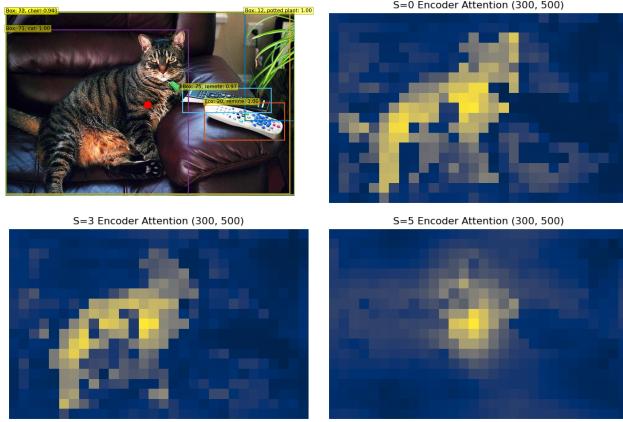


Figure 9. The detection result from DETR on original image and encoder attention at (300, 450) as the impulse noise became more severe.

#### 4.4. Query Properties

First, we evaluated the contribution of each query in object detection. Both DETR R50 and DETR R101 models were tested on the MSCOCO validation dataset. We collected all predictions with confidence larger than 0.8 and the corresponding query ID to compute the query frequency, as shown in the middle plots in Figure 10. We found an interesting phenomenon that both the DETR R50 and DETR R101 model have one main query (71 in DETR R50 and 68 in DETR R101) that detect lots of objects. Each main query can account for 7.5% of the total predictions. Intuitive thinking was that the main queries are responsible for detecting "person", which is the most common object in the dataset. Therefore, we investigated the relations between main queries and classes. We divided the main queries' frequency by the total frequency in each class to visualize their contribution, as shown in the right plots. Both main queries have shown very similar distribution over object categories. None of them has a high contribution ( $< 40\%$ ) in the "person" category, and both of them has a high contribution ( $< 40\%$ ) in the "airplane", "train", "cat" and "bear" category. This indicated that the main query has particular

preferences on object classes, but they are not detecting the most common object, "person".

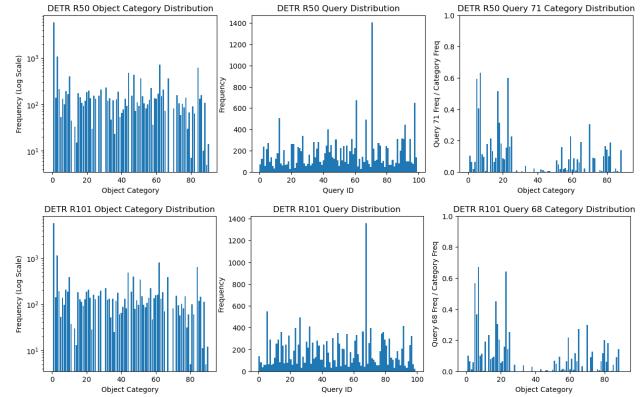


Figure 10. Left columns: Class distribution from model's prediction. Middle columns: Frequency of each query in predictions. Right columns: Main queries' contribution in detecting each object category.

Apart from class relations, we also analyzed bounding boxes' location and size predictions from main queries, shown in Figure 11. Both queries prefer to detect medium to large objects in the center region. However, there are also a large amount of small to medium-sized boxes circulating the center.

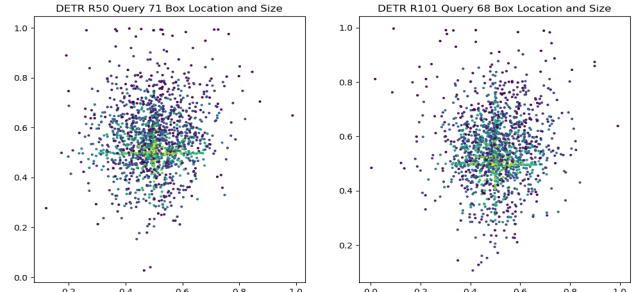


Figure 11. The predicted bounding boxes center location and size of main queries. Yellow colored points refers to boxes with a large size.

To analyze the importance of the main query, we masked its outputs and evaluated the precision change, as shown in Table 3. The main query indeed has a massive influence on the model's performance. Without predictions from the main query, the average precision dropped by around 7 points. This shows that DETR is heavy depending on the predictions from the main query. It will make the model more vulnerable to attacks.



it comes to resistance to interference from information loss in occlusion images. Nevertheless, in the case of the sticker patch, it manages to produce a new set of keys, queries, and values in the network which in most cases results in the misdirection of the network. This result in increasing the attention weight to the corresponding adversarial patch key token, making the network misguided the attention to the adversarial stickers. For the benchmark, the experiment indicates that DETR precision on corrupted images is generally lower than YOLOv5. We also found that DETR depends heavily on the predictions from the main query. This impact the main query to receive most of the gradient flow, which leads to imbalanced contributions among all queries.

## References

- [1] Shaoqing, R. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks arXiv:1506.01497 1
- [2] Alexey, B. YOLOv4: Optimal Speed and Accuracy of Object Detection arXiv:2004.10934 1
- [3] Sharif, M. (2016) Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (ACM)*: 1528–1540. 2
- [4] Evtimov, I. (2017) Robust physical-world attacks on deep learning models. *arXiv preprint arXiv:1707.08945* 3
- [5] Muzammal, N. (2021) Intriguing Properties of Vision Transformers arXiv:2105.10497 2
- [6] kaziwa, S. (2021) Occlusion Handling in Generic Object Detection arXiv:2101.08845 2
- [7] Dan, H. (2019) Benchmarking Neural Network Robustness to Common Corruptions and Surface Variantions arXiv:1807.01697 2
- [8] Claudio M. (2020) Benchmarking Robustness in Object Detection: Autonomous Driving when Winter is Coming arXiv:1907.0748 2
- [9] Nicolas C. (2020) End-to-End Object Detection with Transformers *arXiv preprint arXiv:2005.12872* 1, 3
- [10] Lin, T. (2014) Microsoft coco: Common objects in context. In *European conference on computer vision (ECCV)*: 740–755. 4
- [11] Aldahdooh, A. (2021) Reveal of vision transformers robustness against adversarial attacks. *arXiv preprint arXiv:2106.03734* 2
- [12] Benz, P. (2021) Adversarial robustness comparison of vision transformer and mlp-mixer to cnns. *arXiv preprint arXiv:2110.02797* 2
- [13] Bhojanapalli, S. (2021) Understanding robustness of transformers for image classification. *arXiv preprint arXiv:2103.14586* 2
- [14] Fu, Y. (2022) Patch-fool: Are vision transformers always robust against adversarial perturbations? In *International Conference on Learning Representations (ICLR)* 2
- [15] Gu, J. (2022) Are vision transformers robust to patch-wise perturbations?, In *European conference on computer vision (ECCV)* 2
- [16] Mahmood, K. (2021) On the robustness of vision transformers to adversarial examples. *arXiv preprint arXiv:2104.02610* 2
- [17] Naseer, M. (2021) On improving adversarial transferability of vision transformers. *arXiv preprint arXiv:2106.04169* 2
- [18] Shao, R. (2021) On the adversarial robustness of visual transformers. *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2
- [19] Wei, Z. (2021) Towards transferable adversarial attacks on vision transformers. *arXiv preprint arXiv:2109.04176* 2
- [20] Zem L. (2021) Swin Transformer: Hierarchical Vision Transformer using Shifted Window arXiv:2103.14030 7