

# End-to-End Object Detection with Transformers

Nicolas Carion<sup>1,2</sup><sup>[0000-0002-2308-9680]</sup>, Francisco Massa<sup>2</sup><sup>[000-0003-0697-6664]</sup>,  
Gabriel Synnaeve<sup>2</sup><sup>[0000-0003-1715-3356]</sup>, Nicolas Usunier<sup>2</sup><sup>[0000-0002-9324-1457]</sup>,  
Alexander Kirillov<sup>2</sup><sup>[0000-0003-3169-3199]</sup>, and Sergey  
Zagoruyko<sup>2</sup><sup>[0000-0001-9684-5240]</sup>

<sup>1</sup> Paris Dauphine University

<sup>2</sup> Facebook AI

{alcinos, fmassa, gab, usunier, akirillov, szagoruyko}@fb.com

**Abstract.** We present a new method that views object detection as a direct set prediction problem. Our approach streamlines the detection pipeline, effectively removing the need for many hand-designed components like a non-maximum suppression procedure or anchor generation that explicitly encode our prior knowledge about the task. The main ingredients of the new framework, called DETection TRansformer or DETR, are a set-based global loss that forces unique predictions via bipartite matching, and a transformer encoder-decoder architecture. Given a fixed small set of learned object queries, DETR reasons about the relations of the objects and the global image context to directly output the final set of predictions in parallel. The new model is conceptually simple and does not require a specialized library, unlike many other modern detectors. DETR demonstrates accuracy and run-time performance on par with the well-established and highly-optimized Faster R-CNN baseline on the challenging COCO object detection dataset. Moreover, DETR can be easily generalized to produce panoptic segmentation in a unified manner. We show that it significantly outperforms competitive baselines. Training code and pretrained models are available at <https://github.com/facebookresearch/detr>.

## 1 Introduction

The goal of object detection is to predict a set of bounding boxes and category labels for each object of interest. Modern detectors address this set prediction task in an indirect way, by defining surrogate regression and classification problems on a large set of proposals [36,5], anchors [22], or window centers [52,45]. Their performances are significantly influenced by postprocessing steps to collapse near-duplicate predictions, by the design of the anchor sets and by the heuristics that assign target boxes to anchors [51]. To simplify these pipelines, we propose a direct set prediction approach to bypass the surrogate tasks. This end-to-end philosophy has led to significant advances in complex structured prediction tasks such as machine translation or speech recognition, but not yet in object detection: previous attempts [42,15,4,38] either add other forms of prior 解决的问题