

# Conditional DETR for Fast Training Convergence

Depu Meng<sup>1\*</sup> Xiaokang Chen<sup>2\*</sup> Zejia Fan<sup>2</sup> Gang Zeng<sup>2</sup>  
Houqiang Li<sup>1</sup> Yuhui Yuan<sup>3</sup> Lei Sun<sup>3</sup> Jingdong Wang<sup>3†</sup>

<sup>1</sup>University of Science and Technology of China

<sup>2</sup>Peking University

<sup>3</sup>Microsoft Research Asia

## Abstract

The recently-developed DETR approach applies the transformer encoder and decoder architecture to object detection and achieves promising performance. In this paper, we handle the critical issue, slow training convergence, and present a conditional cross-attention mechanism for fast DETR training. Our approach is motivated by that the cross-attention in DETR relies highly on the content embeddings for localizing the four extremities and predicting the box, which increases the need for high-quality content embeddings and thus the training difficulty.

Our approach, named conditional DETR, learns a conditional spatial query from the decoder embedding for decoder multi-head cross-attention. The benefit is that through the conditional spatial query, each cross-attention head is able to attend to a band containing a distinct region, e.g., one object extremity or a region inside the object box. This narrows down the spatial range for localizing the distinct regions for object classification and box regression, thus relaxing the dependence on the content embeddings and easing the training. Empirical results show that conditional DETR converges  $6.7\times$  faster for the backbones R50 and R101 and  $10\times$  faster for stronger backbones DC5-R50 and DC5-R101. Code is available at <https://github.com/Atten4Vis/ConditionalDETR>.

## 1. Introduction

The DEtection TRansformer (DETR) method [3] applies the transformer encoder and decoder architecture to object detection and achieves good performance. It effectively eliminates the need for many hand-crafted components, including non-maximum suppression and anchor generation.

The DETR approach suffers from slow convergence on training, and needs 500 training epochs to get good performance. The very recent work, deformable DETR [53], handles this issue by replacing the global dense attention

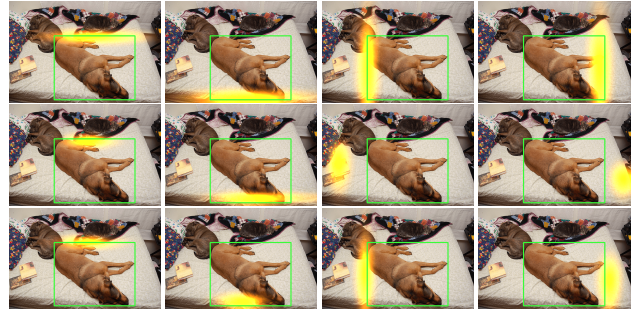


Figure 1. Comparison of spatial attention weight maps for our conditional DETR-R50 with 50 training epochs (the first row), the original DETR-R50 with 50 training epochs (the second row), and the original DETR-R50 with 500 training epochs (the third row). The maps for our conditional DETR and DETR trained with 500 epochs are able to highlight the four extremity regions satisfactorily. In contrast, the spatial attention weight maps responsible for the left and right edges (the third and fourth images in the second row) from DETR trained with 50 epochs cannot highlight the extremities satisfactorily. The green box is the ground-truth box.

(self-attention and cross-attention) with deformable attention that attends to a small set of key sampling points and using the high-resolution and multi-scale encoder. Instead, we still use the global dense attention and propose an improved decoder cross-attention mechanism for accelerating the training process.

Our approach is motivated by high dependence on content embeddings and minor contributions made by the spatial embeddings in cross-attention. The empirical results in DETR [3] show that if removing the positional embeddings in keys and the object queries from the second decoder layer and only using the content embeddings in keys and queries, the detection AP drops slightly<sup>1</sup>.

Figure 1 (the second row) shows that the spatial attention weight maps from the cross-attention in DETR trained with 50 epochs. One can see that two among the four maps do not correctly highlight the bands for the corresponding extremities, thus weak at shrinking the spatial range for the

<sup>\*</sup>Equal Contribution. This work was done when D. Meng, X. Chen, and Z. Fan were interns at Microsoft Research, Beijing, P.R. China

<sup>†</sup>Corresponding author.

<sup>1</sup>The minor AP drop 1.4 is reported on R50 with 300 epochs in Table 3 from [3]. We empirically got the consistent observation: the AP drops to 34.0 from 34.9 for 50 training epochs.

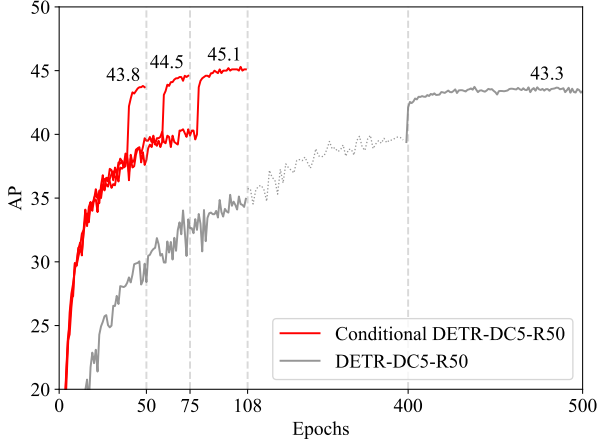


Figure 2. Convergence curves for conditional DETR-DC5-R50 and DETR-DC5-R50 on COCO 2017 val. The conditional DETR is trained for 50, 75, 108 epochs. Conditional DETR training is converged much faster than DETR.

content queries to precisely localize the extremities. The reasons are that (i) the spatial queries, i.e., object queries, only give the general attention weight map without exploiting the specific image information; and that (ii) due to short training the content queries are not strong enough to match the spatial keys well as they are also used to match the content keys. This increases the dependence on high-quality content embeddings, thus increasing the training difficulty.

We present a conditional DETR approach, which learns a conditional spatial embedding for each query from the corresponding previous decoder output embedding, to form a so-called conditional spatial query for decoder multi-head cross-attention. The conditional spatial query is predicted by mapping the information for regressing the object box to the embedding space, the same to the space that the 2D coordinates of the keys are also mapped to.

We empirically observe that using the spatial queries and keys, each cross-attention head spatially attends to a band containing the object extremity or a region inside the object box (Figure 1, the first row). This shrinks the spatial range for the content queries to localize the effective regions for class and box prediction. As a result, the dependence on the content embeddings is relaxed and the training is easier. The experiments show that conditional DETR converges  $6.7\times$  faster for the backbones R50 and R101 and  $10\times$  faster for stronger backbones DC5-R50 and DC5-R101. Figure 2 gives the convergence curves for conditional DETR and the original DETR [3].

## 2. Related Work

**Anchor-based and anchor-free detection.** Most existing object detection approaches make predictions from initial guesses that are carefully designed. There are two main initial guesses: anchor boxes or object centers. The an-

chor box-based methods inherit the ideas from the proposal-based method, Fast R-CNN. Example methods include Faster R-CNN [9], SSD [26], YOLOv2 [31], YOLOv3 [32], YOLOv4 [1], RetinaNet [24], Cascade R-CNN [2], Libra R-CNN [29], TSD [35] and so on.

The anchor-free detectors predict the boxes at points near the object centers. Typical methods include YOLOv1 [30], CornerNet [21], ExtremeNet [50], CenterNet [49, 6], FCOS [39] and others [23, 28, 52, 19, 51, 22, 15, 46, 47].

**DETR and its variants.** DETR successfully applies transformers to object detection, effectively removing the need for many hand-designed components like non-maximum suppression or initial guess generation. The high computation complexity issue, caused by the global encoder self-attention, is handled in adaptive clustering transformer [48] and by sparse attentions in deformable DETR [53].

The other critical issue, slow training convergence, has been attracting a lot of recent research attention. The TSP (transformer-based set prediction) approach [37] eliminates the cross-attention modules and combines the FCOS and R-CNN-like detection heads. Deformable DETR [53] adopts deformable attention, which attends to sparse positions learned from the content embedding, to replace decoder cross-attention.

The spatially modulated co-attention (SMCA) approach [7], which is concurrent to our approach, is very close to our approach. It modulates the DETR multi-head global cross-attentions with Gaussian maps around a few (shifted) centers that are learned from the decoder embeddings, to focus more on a few regions inside the estimated box. In contrast, the proposed conditional DETR approach learns the conditional spatial queries from the decoder content embeddings, and predicts the spatial attention weight maps without human-crafting the attention attenuation, which highlight four extremities for box regression, and distinct regions inside the object for classification.

**Conditional and dynamic convolution.** The proposed conditional spatial query scheme is related to conditional convolutional kernel generation. Dynamic filter network [16] learns the convolutional kernels from the input, which is applied to instance segmentation in CondInst [38] and SOLOv2 [42] for learning instance-dependent convolutional kernels. CondConv [44] and dynamic convolution [4] mix convolutional kernels with the weights learned from the input. SENet [14], GEMNet [13] and Lite-HRNet [45] learn from the input the channel-wise weights.

These methods learn from the input the convolutional kernel weights and then apply the convolutions to the input. In contrast, the linear projection in our approach is learned from the decoder embeddings for representing the displacement and scaling information.

**Transformers.** The transformer [40] relies on the attention mechanism, self-attention and cross-attention, to

draw global dependencies between the input and the output. There are several works closely related to our approach. Gaussian transformer [11] and T-GSA (Transformer with Gaussian-weighted self-attention) [18], followed by SMCA [7], attenuate the attention weights according to the distance between target and context symbols with learned or human-crafted Gaussian variance. Similar to ours, TUPE [17] computes the attention weight also from the spatial attention weight and the content attention weight. Instead, our approach mainly focuses on the attention attenuation mechanism in a learnable form other than a Gaussian function, and potentially benefits speech enhancement [18] and natural language inference [11].

### 3. Conditional DETR

#### 3.1. Overview

**Pipeline.** The proposed approach follows detection transformer (DETR), an end-to-end object detector, and predicts all the objects at once without the need for NMS or anchor generation. The architecture consists of a CNN backbone, a transformer encoder, a transformer decoder, and object class and box position predictors. The transformer encoder aims to improve the content embeddings output from the CNN backbone. It is a stack of multiple encoder layers, where each layer mainly consists of a self-attention layer and a feed-forward layer.

The transformer decoder is a stack of decoder layers. Each decoder layer, illustrated in Figure 3, is composed of three main layers: (1) a self-attention layer for removing duplication prediction, which performs interactions between the embeddings, outputted from the previous decoder layer and used for class and box prediction, (2) a cross-attention layer, which aggregates the embeddings output from the encoder to refine the decoder embeddings for improving class and box prediction, and (3) a feed-forward layer.

**Box regression.** A candidate box is predicted from each decoder embedding as follows,

$$\mathbf{b} = \text{sigmoid}(\text{FFN}(\mathbf{f}) + [\mathbf{s}^\top \ 0 \ 0]^\top). \quad (1)$$

Here,  $\mathbf{f}$  is the decoder embedding.  $\mathbf{b}$  is a four-dimensional vector  $[b_{cx} \ b_{cy} \ b_w \ b_h]^\top$ , consisting of the box center, the box width and the box height.  $\text{sigmoid}()$  is used to normalize the prediction  $\mathbf{b}$  to the range  $[0, 1]$ .  $\text{FFN}()$  aims to predict the unnormalized 2D coordinate of the reference point, and is  $(0, 0)$  in the original DETR. In our approach, we consider two choices: learn the reference point  $\mathbf{s}$  as a parameter for each candidate box prediction, or generate it from the corresponding object query.

**Category prediction.** The classification score for each candidate box is also predicted from the decoder embedding through an FNN,  $\mathbf{e} = \text{FFN}(\mathbf{f})$ .

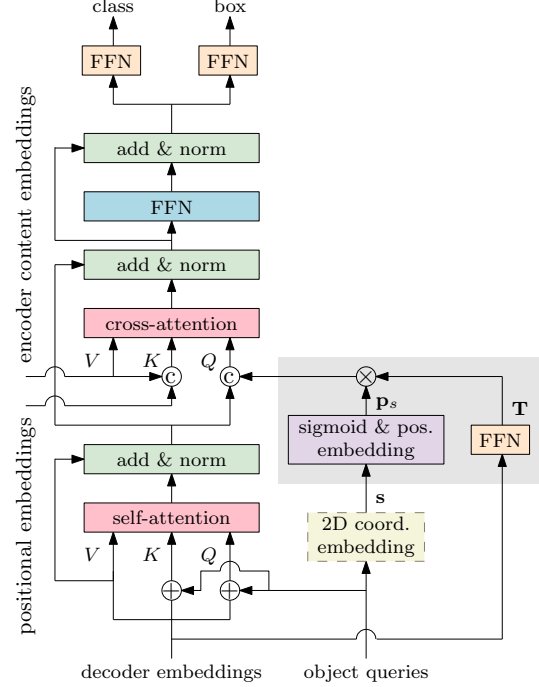


Figure 3. Illustrating one decoder layer in conditional DETR. The main difference from the original DETR [3] lies in the input queries and the input keys for cross-attention. The conditional spatial query is predicted from learnable 2D coordinates  $\mathbf{s}$  and the embeddings output from the previous decoder layer, through the operations depicted in the gray-shaded box. The 2D coordinate  $\mathbf{s}$  can be predicted from the object query (the dashed box), or simply learned as model parameters. The spatial query (key) and the content query (key) are concatenated as the query (key). The resulting cross-attention is called conditional cross-attention. Same as DETR [3], the decoder layer is repeated 6 times.

**Main work.** The cross-attention mechanism aims to *localize the distinct regions, four extremities for box detection and regions inside the box for object classification, and aggregates the corresponding embeddings*. We propose a conditional cross-attention mechanism with introducing conditional spatial queries for improving the localization capability and accelerating the training process.

#### 3.2. DETR Decoder Cross-Attention

The DETR decoder cross-attention mechanism takes three inputs: queries, keys and values. Each key is formed by adding a content key  $\mathbf{c}_k$  (the content embedding output from the encoder) and a spatial key  $\mathbf{p}_k$  (the positional embedding of the corresponding normalized 2D coordinate). The value is formed from the content embedding, same with the content key, output from the encoder.

In the original DETR approach, each query is formed by adding a content query  $\mathbf{c}_q$  (the embedding output from the decoder self-attention), and a spatial query  $\mathbf{p}_q$  (i.e., the object query  $\mathbf{o}_q$ ). In our implementation, there are  $N = 300$

object queries, and accordingly there are  $N$  queries<sup>2</sup>, each query outputting a candidate detect result in one decoder layer.

The attention weight is based on the dot-product between the query and the key, used for attention weight computation,

$$\begin{aligned} & (\mathbf{c}_q + \mathbf{p}_q)^\top (\mathbf{c}_k + \mathbf{p}_k) \\ &= \mathbf{c}_q^\top \mathbf{c}_k + \mathbf{c}_q^\top \mathbf{p}_k + \mathbf{p}_q^\top \mathbf{c}_k + \mathbf{p}_q^\top \mathbf{p}_k \\ &= \mathbf{c}_q^\top \mathbf{c}_k + \mathbf{c}_q^\top \mathbf{p}_k + \mathbf{o}_q^\top \mathbf{c}_k + \mathbf{o}_q^\top \mathbf{p}_k. \end{aligned} \quad (2)$$

### 3.3. Conditional Cross-Attention

The proposed conditional cross-attention mechanism forms the query by concatenating the content query  $\mathbf{c}_q$ , outputting from decoder self-attention, and the spatial query  $\mathbf{p}_q$ . Accordingly, the key is formed as the concatenation of the content key  $\mathbf{c}_k$  and the spatial key  $\mathbf{p}_k$ .

The cross-attention weights consist of two components, content attention weight and spatial attention weight. The two weights are from two dot-products, content and spatial dot-products,

$$\mathbf{c}_q^\top \mathbf{c}_k + \mathbf{p}_q^\top \mathbf{p}_k. \quad (3)$$

Different from the original DETR cross-attention, our mechanism separates the roles of content and spatial queries so that spatial and content queries focus on the spatial and content attention weights, respectively.

An additional important task is to compute the spatial query  $\mathbf{p}_q$  from the embedding  $\mathbf{f}$  of the previous decoder layer. We first identify that the spatial information of the distinct regions are determined by the two factors together, decoder embedding and reference point. We then show how to map them to the embedding space, forming the query  $\mathbf{p}_q$ , so that the spatial query lies in the same space the 2D coordinates of the keys are mapped to.

**The decoder embedding contains the displacements** of the distinct regions with respect to the reference point. The box prediction process in Equation 1 consists of two steps: (1) predicting the box with respect to the reference point in the unnormalized space, and (2) normalizing the predicted box to the range  $[0, 1]$ <sup>3</sup>.

Step (1) means that the decoder embedding  $\mathbf{f}$  contains the displacements of the four extremities (forming the box) with respect to the reference point  $\mathbf{s}$  in the unnormalized space. This implies that both the embedding  $\mathbf{f}$  and the reference point  $\mathbf{s}$  are necessary to determine the spatial information of the distinct regions, the four extremities as well as the region for predicting the classification score.

<sup>2</sup>For description simplicity and clearness, we drop the query, key, and value indices.

<sup>3</sup>The origin  $(0, 0)$  in the unnormalized space for the original DETR method is mapped to  $(0.5, 0.5)$  (the center in the image space) in the normalized space through the sigmoid function.

**Conditional spatial query prediction.** We predict the conditional spatial query from the embedding  $\mathbf{f}$  and the reference point  $\mathbf{s}$ ,

$$(\mathbf{s}, \mathbf{f}) \rightarrow \mathbf{p}_q, \quad (4)$$

so that it is aligned with the positional space which the normalized 2D coordinates of the keys are mapped to. The process is illustrated in the gray-shaded box area of Figure 3.

We normalize the reference point  $\mathbf{s}$  and then map it to a 256-dimensional sinusoidal positional embedding in the same way as the positional embedding for keys:

$$\mathbf{p}_s = \text{sinusoidal}(\text{sigmoid}(\mathbf{s})). \quad (5)$$

We then map the displacement information contained in the decoder embedding  $\mathbf{f}$  to a linear projection in the same space through an FFN consisting of learnable linear projection + ReLU + learnable linear projection:  $\mathbf{T} = \text{FFN}(\mathbf{f})$ .

The conditional spatial query is computed by transforming the reference point in the embedding space:  $\mathbf{p}_q = \mathbf{T}\mathbf{p}_s$ . We choose the simple and computationally-efficient projection matrix, a diagonal matrix. The 256 diagonal elements are denoted as a vector  $\lambda_q$ . The conditional spatial query is computed by the element-wise multiplication:

$$\mathbf{p}_q = \mathbf{T}\mathbf{p}_s = \lambda_q \odot \mathbf{p}_s. \quad (6)$$

**Multi-head cross-attention.** Following DETR [3], we adopt the standard multi-head cross-attention mechanism. Object detection usually needs to implicitly or explicitly localize the four object extremities for accurate box regression and localize the object region for accurate object classification. The multi-head mechanism is beneficial to disentangle the localization tasks.

We perform multi-head parallel attentions by projecting the queries, the keys, and the values  $M = 8$  times with learned linear projections to low dimensions. The spatial and content queries (keys) are separately projected to each head with different linear projections. The projections for values are the same as the original DETR and are only for the contents.

### 3.4. Visualization and Analysis

**Visualization.** Figure 4 visualizes the attention weight maps for each head: the spatial attention weight maps, the content attention weight maps, and the combined attention weight maps. The maps are soft-max normalized over the spatial dot-products  $\mathbf{p}_q^\top \mathbf{p}_k$ , the content dot-products  $\mathbf{c}_q^\top \mathbf{c}_k$ , and the combined dot-products  $\mathbf{c}_q^\top \mathbf{c}_k + \mathbf{p}_q^\top \mathbf{p}_k$ . We show 5 out of the 8 maps, and other three are the duplicates, corresponding to bottom and top extremities, and a small region inside the object box<sup>4</sup>.

<sup>4</sup>The duplicates might be different for models trained several times, but the detection performance is almost the same.



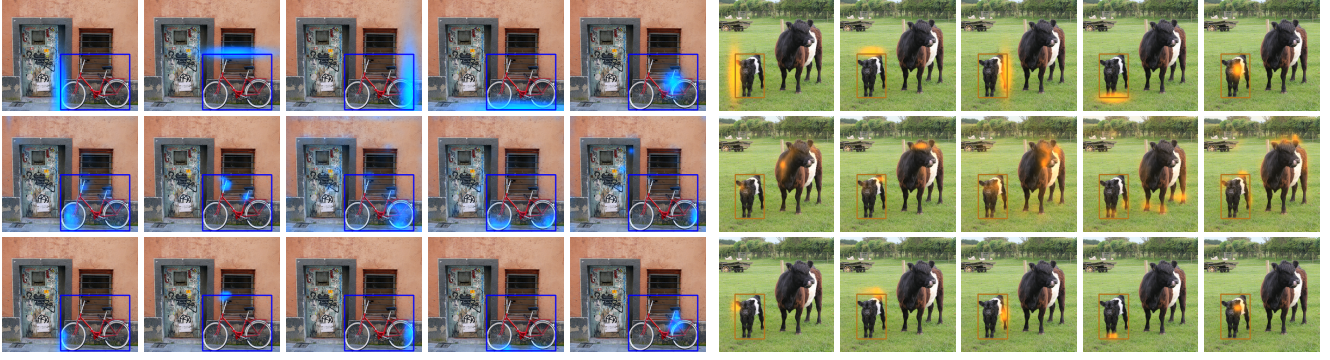


Figure 4. Illustrating the spatial attention weight maps (the first row), the content attention weight maps (the second row), and the combined attention weight maps (the third row) computed from our conditional DETR. The attention weight maps are from 5 heads out of the 8 heads and are responsible for the four extremities and a region inside the object box. The content attention weight maps for the four extremities highlight scattered regions inside the box (bicycle) or similar regions in two object instances (cow), and the corresponding combined attention weight maps highlight the extremity regions with the help of the spatial attention weight maps. The combined attention weight map for the region inside the object box mainly depends on the spatial attention weight map, which implies that the representation of a region inside the object might encode enough class information. The maps are from conditional DETR-R50 trained with 50 epochs.

We can see that the spatial attention weight map at each head is able to localize a distinct region, a region containing one extremity or a region inside the object box. It is interesting that each spatial attention weight map corresponding to an extremity highlights a spatial band that overlaps with the corresponding edge of the object box. The other spatial attention map for the region inside the object box merely highlights a small region whose representations might already encode enough information for object classification.

The content attention weight maps of the four heads corresponding to the four extremities highlight scattered regions in addition to the extremities. The combination of the spatial and content maps filters out other highlights and keeps extremity highlights for accurate box regression.

**Comparison to DETR.** Figure 1 shows the spatial attention weight maps of our conditional DETR (the first row) and the original DETR trained with 50 epochs (the second row). The maps of our approach are computed by soft-max normalizing the dot-products between spatial keys and queries,  $\mathbf{p}_q^\top \mathbf{p}_k$ . The maps for DETR are computed by soft-max normalizing the dot-products with the spatial keys,  $(\mathbf{o}_q + \mathbf{c}_q)^\top \mathbf{p}_k$ .

It can be seen that our spatial attention weight maps accurately localize the distinct regions, four extremities. In contrast, the maps from the original DETR with 50 epochs can not accurately localize two extremities, and 500 training epochs (the third row) make the content queries stronger, leading to accurate localization. This implies that it is really hard to learn the content query  $\mathbf{c}_q$  to serve as two roles<sup>5</sup>: match the content key and the spatial key simultaneously, and thus more training epochs are needed.

<sup>5</sup>Strictly speaking, the embedding output from decoder self-attention for more training epochs contains both spatial and content information. For discussion convenience, we still call it content query.

**Analysis.** The spatial attention weight maps shown in Figure 4 imply that the conditional spatial query, used to form the spatial query, have at least two effects. (i) Translate the highlight positions to the four extremities and the position inside the object box: interestingly the highlighted positions are spatially similarly distributed in the object box. (ii) Scale the spatial spread for the extremity highlights: large spread for large objects and small spread for small objects.

The two effects are realized in the spatial embedding space through applying the transformation  $\mathbf{T}$  over  $\mathbf{p}_s$  (further disentangled through image-independent linear projections contained in cross-attention and distributed to each head). This indicates that the transformation  $\mathbf{T}$  not only contains the displacements as discussed before, but also the object scale.

### 3.5. Implementation Details

**Architecture.** Our architecture is almost the same with the DETR architecture [3] and contains the CNN backbone, transformer encoder, transformer decoder, prediction feed-forward networks (FFNs) following each decoder layer (the last decoder layer and the 5 internal decoder layers) with parameters shared among the 6 prediction FFNs. The hyperparameters are the same as DETR.

The main architecture difference is that we introduce the conditional spatial embeddings as the spatial queries for conditional multi-head cross-attention and that the spatial query (key) and the content query (key) are combined through concatenation other than addition. In the first cross-attention layer there are no decoder content embeddings, we make simple changes based on the DETR implementation [3]: concatenate the positional embedding predicted from the object query (the positional embedding) into the original query (key).

Table 1. Comparison of conditional DETR with DETR on COCO 2017 val. Our conditional DETR approach for high-resolution backbones DC5-R50 and DC5-R101 is  $10\times$  faster than the original DETR, and for low-resolution backbones R50 and R101  $6.67\times$  faster. Conditional DETR is empirically superior to other two single-scale DETR variants. \*The results of deformable DETR are from the GitHub repository provided by the authors of deformable DETR [53].

Model	#epochs	GFLOPs	#params (M)	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
DETR-R50	500	86	41	42.0	62.4	44.2	20.5	45.8	61.1
DETR-R50	50	86	41	34.9	55.5	36.0	14.4	37.2	54.5
Conditional DETR-R50	50	90	44	40.9	61.8	43.3	20.8	44.6	59.2
Conditional DETR-R50	75	90	44	42.1	62.9	44.8	21.6	45.4	60.2
Conditional DETR-R50	108	90	44	43.0	64.0	45.7	22.7	46.7	61.5
DETR-DC5-R50	500	187	41	43.3	63.1	45.9	22.5	47.3	61.1
DETR-DC5-R50	50	187	41	36.7	57.6	38.2	15.4	39.8	56.3
Conditional DETR-DC5-R50	50	195	44	43.8	64.4	46.7	24.0	47.6	60.7
Conditional DETR-DC5-R50	75	195	44	44.5	65.2	47.3	24.4	48.1	62.1
Conditional DETR-DC5-R50	108	195	44	45.1	65.4	48.5	25.3	49.0	62.2
DETR-R101	500	152	60	43.5	63.8	46.4	21.9	48.0	61.8
DETR-R101	50	152	60	36.9	57.8	38.6	15.5	40.6	55.6
Conditional DETR-R101	50	156	63	42.8	63.7	46.0	21.7	46.6	60.9
Conditional DETR-R101	75	156	63	43.7	64.9	46.8	23.3	48.0	61.7
Conditional DETR-R101	108	156	63	44.5	65.6	47.5	23.6	48.4	63.6
DETR-DC5-R101	500	253	60	44.9	64.7	47.7	23.7	49.5	62.3
DETR-DC5-R101	50	253	60	38.6	59.7	40.7	17.2	42.2	57.4
Conditional DETR-DC5-R101	50	262	63	45.0	65.5	48.4	26.1	48.9	62.8
Conditional DETR-DC5-R101	75	262	63	45.6	66.5	48.8	25.5	49.7	63.3
Conditional DETR-DC5-R101	108	262	63	45.9	66.8	49.5	27.2	50.3	63.3
<i>Other single-scale DETR variants</i>									
Deformable DETR-R50-SS*	50	78	34	39.4	59.6	42.3	20.6	43.0	55.5
UP-DETR-R50 [5]	150	86	41	40.5	60.8	42.6	19.0	44.4	60.0
UP-DETR-R50 [5]	300	86	41	42.8	63.0	45.3	20.8	47.1	61.7
Deformable DETR-DC5-R50-SS*	50	128	34	41.5	61.8	44.9	24.1	45.3	56.0

**Reference points.** In the original DETR approach,  $\mathbf{s} = [0 \ 0]^\top$  is the same for all the decoder embeddings. We study two ways forming the reference points: regard the unnormalized 2D coordinates as learnable parameters, and the unnormalized 2D coordinate predicted from the object query  $\mathbf{o}_q$ . In the latter way that is similar to deformable DETR [53], the prediction unit is an FFN and consists of learnable linear projection + ReLU + learnable linear projection:  $\mathbf{s} = \text{FFN}(\mathbf{o}_q)$ . When used for forming the conditional spatial query, the 2D coordinates are normalized by the sigmoid function.

**Loss function.** We follow DETR [3] to find an optimal bipartite matching [20] between the predicted and ground-truth objects using the Hungarian algorithm, and then form the loss function for computing and back-propagate the gradients. We use the same way with deformable DETR [53] to formulate the loss: the same matching cost function, the same loss function with 300 object queries, and the same trade-off parameters; The classification loss function is focal loss [24], and the box regression loss (including L1 and GIoU [34] loss) is the same as DETR [3].

## 4. Experiments

### 4.1. Setting

**Dataset.** We perform the experiments on the COCO 2017 [25] detection dataset. The dataset contains about 118K training images and 5K validation (val) images.

**Training.** We follow the DETR training protocol [3]. The backbone is the ImageNet-pretrained model from TORCHVISION with batchnorm layers fixed, and the transformer parameters are initialized using the Xavier initialization scheme [10]. The weight decay is set to be  $10^{-4}$ . The AdamW [27] optimizer is used. The learning rates for the backbone and the transformer are initially set to be  $10^{-5}$  and  $10^{-4}$ , respectively. The dropout rate in transformer is 0.1. The learning rate is dropped by a factor of 10 after 40 epochs for 50 training epochs, after 60 epochs for 75 training epochs, and after 80 epochs for 108 training epochs.

We use the augmentation scheme same as DETR [3]: resize the input image such that the short side is at least 480 and at most 800 pixels and the long size is at most 1333 pixels; randomly crop the image such that a training image is cropped with probability 0.5 to a random rectangular patch.

Table 2. Results for multi-scale and higher-resolution DETR variants. We do not expect that our approach performs on par as our approach (single-scale,  $16\times$  resolution) does not use a strong multi-scale or  $8\times$  resolution encoder. Surprisingly, the AP scores of our approach with DC5-R50 and DC5-R101 are close to the two multi-scale and higher-resolution DETR variants.

Model	#epochs	GFLOPs	#params (M)	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Faster RCNN-FPN-R50 [33]	36	180	42	40.2	61.0	43.8	24.2	43.5	52.0
Faster RCNN-FPN-R50 [33]	108	180	42	42.0	62.1	45.5	26.6	45.5	53.4
Deformable DETR-R50 [53]	50	173	40	43.8	62.6	47.7	26.4	47.1	58.0
TSP-FCOS-R50 [37]	36	189	—	43.1	62.3	47.0	26.6	46.8	55.9
TSP-RCNN-R50 [37]	36	188	—	43.8	63.3	48.3	28.6	46.9	55.7
TSP-RCNN-R50 [37]	96	188	—	45.0	64.5	49.6	29.7	47.7	58.0
Conditional DETR-DC5-R50	50	195	44	43.8	64.4	46.7	24.0	47.6	60.7
Conditional DETR-DC5-R50	108	195	44	45.1	65.4	48.5	25.3	49.0	62.2
Faster RCNN-FPN-R101 [33]	36	246	60	42.0	62.5	45.9	25.2	45.6	54.6
Faster RCNN-FPN-R101 [33]	108	246	60	44.0	63.9	47.8	27.2	48.1	56.0
TSP-FCOS-R101 [37]	36	255	—	44.4	63.8	48.2	27.7	48.6	57.3
TSP-RCNN-R101 [37]	36	254	—	44.8	63.8	49.2	29.0	47.9	57.1
TSP-RCNN-R101 [37]	96	254	—	46.5	66.0	51.2	29.9	49.7	59.2
Conditional DETR-DC5-R101	50	262	63	45.0	65.5	48.4	26.1	48.9	62.8
Conditional DETR-DC5-R101	108	262	63	45.9	66.8	49.5	27.2	50.3	63.3

**Evaluation.** We use the standard COCO evaluation. We report the average precision (AP), and the AP scores at 0.50, 0.75 and for the small, medium, and large objects.

## 4.2. Results

**Comparison to DETR.** We compare the proposed conditional DETR to the original DETR [3]. We follow [3] and report the results over four backbones: ResNet-50 [12], ResNet-101, and their  $16\times$ -resolution extensions DC5-ResNet-50 and DC5-ResNet-101.

The corresponding DETR models are named as DETR-R50, DETR-R101, DETR-DC5-R50, and DETR-DC5-R101, respectively. Our models are named as conditional DETR-R50, conditional DETR-R101, conditional DETR-DC5-R50, and conditional DETR-DC5-R101, respectively.

Table 1 presents the results from DETR and conditional DETR. DETR with 50 training epochs performs much worse than 500 training epochs. Conditional DETR with 50 training epochs for R50 and R101 as the backbones performs slightly worse than DETR with 500 training epochs. Conditional DETR with 50 training epochs for DC5-R50 and DC5-R101 performs similarly as DETR with 500 training epochs. Conditional DETR for the four backbones with 75/108 training epochs performs better than DETR with 500 training epochs. In summary, conditional DETR for high-resolution backbones DC5-R50 and DC5-R101 is  $10\times$  faster than the original DETR, and for low-resolution backbones R50 and R101  $6.67\times$  faster. In other words, conditional DETR performs better for stronger backbones with better performance.

In addition, we report the results of single-scale DETR extensions: deformable DETR-SS [53] and UP-DETR [5] in Table 1. Our results over R50 and DC5-R50 are bet-

ter than deformable DETR-SS: 40.9 vs. 39.4 and 43.8 vs. 41.5. The comparison might not be fully fair as for example parameter and computation complexities are different, but it implies that the conditional cross-attention mechanism is beneficial. Compared to UP-DETR-R50, our results with fewer training epochs are obviously better.

**Comparison to multi-scale and higher-resolution DETR variants.** We focus on accelerating the DETR training, without addressing the issue of high computational complexity in the encoder. We do not expect that our approach achieves on par with DETR variants w/ multi-scale attention and  $8\times$ -resolution encoders, e.g., TSP-FCOS and TSP-RCNN [37] and deformable DETR [53], which are able to reduce the encoder computational complexity and improve the performance due to multi-scale and higher-resolution.

The comparisons in Table 2 surprisingly show that our approach on DC5-R50 ( $16\times$ ) performs same as deformable DETR-R50 (multi-scale,  $8\times$ ). Considering that the AP of the single-scale deformable DETR-DC5-R50-SS is 41.5 (lower than ours 43.8) (Table 1), one can see that deformable DETR benefits a lot from the multi-scale and higher-resolution encoder that potentially benefit our approach, which is currently not our focus and left as our future work.

The performance of our approach is also on par with TSP-FCOS and TSP-RCNN. The two methods contain a transformer encoder over a small number of selected positions/regions (feature of interest in TSP-FCOS and region proposals in TSP-RCNN) without using the transformer decoder, are extensions of FCOS [39] and Faster RCNN [33]. It should be noted that position/region selection removes unnecessary computation in self-attention and reduces computation complexity dramatically.

Table 3. Ablation study for the ways forming the conditional spatial query. CSQ = our proposed conditional spatial query scheme. Please see the first two paragraphs in Section 5.3 for the meanings of CSQ variants. Our proposed CSQ manner performs better. The backbone ResNet-50 is adopted.

Exp.	CSQ-C	CSQ-T	CSQ-P	CSQ-I	CSQ
GFLOPs	89.3	89.5	89.3	89.5	89.5
AP	37.1	37.6	37.8	40.2	40.9

### 4.3. Ablations

**Reference points.** We compare three ways of forming reference points  $\mathbf{s}$ : (i)  $\mathbf{s} = (0, 0)$ , same to the original DETR, (ii) learn  $\mathbf{s}$  as model parameters and each prediction is associated with different reference points, and (iii) predict each reference point  $\mathbf{s}$  from the corresponding object query. We conducted the experiments with ResNet-50 as the backbone. The AP scores are 36.8, 40.7, and 40.9, suggesting that (ii) and (iii) perform on par and better than (i).

**The effect of the way forming the conditional spatial query.** We empirically study how the transformation  $\lambda_q$  and the positional embedding  $\mathbf{p}_s$  of the reference point, used to form the conditional spatial query  $\mathbf{p}_q = \lambda_q \odot \mathbf{p}_s$ , make contributions to the detection performance.

We report the results of our conditional DETR, and the other ways forming the spatial query with: (i) CSQ-P - only the positional embedding  $\mathbf{p}_s$ , (ii) CSQ-T - only the transformation  $\lambda_q$ , (iii) CSQ-C - the decoder content embedding  $\mathbf{f}$ , and (iv) CSQ-I - the element-wise product of the transformation predicted from the decoder self-attention output  $\mathbf{c}_q$  and the positional embedding  $\mathbf{p}_s$ . The studies in Table 3 imply that our proposed way (CSQ) performs overall the best, validating our analysis about the transformation predicted from the decoder embedding and the positional embedding of the reference point in Section 3.3.

**Focal loss and offset regression with respect to learned reference point.** Our approach follows deformable DETR [53]: use the focal loss with 300 object queries to form the classification loss and predict the box center by regressing the offset with respect to the reference point. We report how the two schemes affect the DETR performance in Table 4. One can see that separately using the focal loss or center offset regression without learning reference points leads to a slight AP gain and combining them together leads to a larger AP gain. Conditional cross-attention in our approach built on the basis of focal loss and offset regression brings a major gain 4.0.

**The effect of linear projections  $\mathbf{T}$  forming the transformation.** Predicting the conditional spatial query needs to learn the linear projection  $\mathbf{T}$  from the decoder embedding (see Equation 6). We empirically study how the linear projection forms affect the performance. The linear projection forms include: an *identity* matrix that means not to learn the linear projection, a *single* scalar, a *block* diagonal ma-

Table 4. The empirical results about the focal loss (FL), offset regression (OR) for box center prediction, and our conditional spatial query (CSQ). The backbone ResNet-50 is adopted.

OR	FL	CSQ	GFLOPs	AP
			85.5	34.9
✓			85.5	35.0
	✓		88.1	35.3
✓	✓		88.1	36.9
✓	✓	✓	89.5	40.9

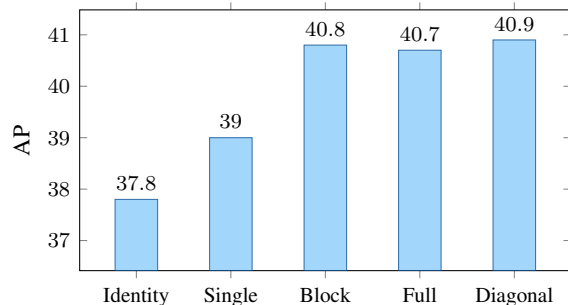


Figure 5. The empirical results for different forms of linear projections that are used to compute the spatial queries for conditional multi-head cross-attention. Diagonal (ours), Full, and Block perform on par. The backbone ResNet-50 is adopted.

trix meaning that each head has a learned  $32 \times 32$  linear projection matrix, a *full* matrix without constraints, and a *diagonal* matrix. Figure 5 presents the results. It is interesting that a single-scalar helps improve the performance, maybe due to narrowing down the spatial range to the object area. Other three forms, *block* diagonal, *full*, and *diagonal* (ours), perform on par.

## 5. Conclusion

We present a simple conditional cross-attention mechanism. The key is to learn a spatial query from the corresponding reference point and decoder embedding. The spatial query contains the spatial information mined for the class and box prediction in the previous decoder layer, and leads to spatial attention weight maps highlighting the bands containing extremities and small regions inside the object box. This shrinks the spatial range for the content query to localize the distinct regions, thus relaxing the dependence on the content query and reducing the training difficulty. In the future, we will study the proposed conditional cross-attention mechanism for human pose estimation [8, 41, 36] and line segment detection [43].

**Acknowledgments.** We thank the anonymous reviewers for their insightful comments and suggestions on our manuscript. This work is supported by the National Key Research and Development Program of China (2017YFB1002601), National Natural Science Foundation of China (61375022, 61403005, 61632003, 61836011, 62021001), Beijing Advanced Innovation Center for Intelligent Robots and Systems (2018IRS11), and PEK-SenseTime Joint Laboratory of Machine Vision.



## References

- [1] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *CoRR*, abs/2004.10934, 2020. 2
- [2] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: delving into high quality object detection. In *CVPR*, 2018. 2
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 2, 3, 4, 5, 6, 7
- [4] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *CVPR*, 2020. 2
- [5] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. UP-DETR: unsupervised pre-training for object detection with transformers. *CoRR*, abs/2011.09094, 2020. 6, 7
- [6] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *ICCV*, 2019. 2
- [7] Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence of DETR with spatially modulated co-attention. *CoRR*, abs/2101.07448, 2021. 2, 3
- [8] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose estimation via disentangled keypoint regression. In *CVPR*, pages 14676–14686, June 2021. 8
- [9] Ross B. Girshick. Fast R-CNN. In *ICCV*, 2015. 2
- [10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 6
- [11] Maosheng Guo, Yu Zhang, and Ting Liu. Gaussian transformer: A lightweight approach for natural language inference. In *AAAI*, 2019. 3
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 7
- [13] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Andrea Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. In *NeurIPS*, 2018. 2
- [14] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 2
- [15] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. Densebox: Unifying landmark localization with end to end object detection. *CoRR*, abs/1509.04874, 2015. 2
- [16] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc Van Gool. Dynamic filter networks. In *NeurIPS*, 2016. 2
- [17] Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. *CoRR*, abs/2006.15595, 2020. 3
- [18] Jaeyoung Kim, Mostafa El-Khamy, and Jungwon Lee. T-GSA: transformer with gaussian-weighted self-attention for speech enhancement. In *ICASSP*, 2020. 3
- [19] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, and Jianbo Shi. Foveabox: Beyond anchor-based object detector. *CoRR*, abs/1904.03797, 2019. 2
- [20] Harold W. Kuhn. The hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 1995. 6
- [21] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, 2018. 2
- [22] Hei Law, Yun Teng, Olga Russakovsky, and Jia Deng. Cornernet-lite: Efficient keypoint based object detection. In *BMVC*. BMVA Press, 2020. 2
- [23] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Scale-aware trident networks for object detection. In *ICCV*, pages 6054–6063, 2019. 2
- [24] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *TPAMI*, 2020. 2, 6
- [25] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014. 6
- [26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. In *ECCV*, 2016. 2
- [27] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. In *ICLR*, 2017. 6
- [28] Xin Lu, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. Grid R-CNN. In *CVPR*, 2019. 2
- [29] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra R-CNN: towards balanced learning for object detection. In *CVPR*, 2019. 2
- [30] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 2
- [31] Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. In *CVPR*, 2017. 2
- [32] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018. 2
- [33] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *TPAMI*, 2017. 7
- [34] Hamid Rezaatoughi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian D. Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, 2019. 6
- [35] Guanglu Song, Yu Liu, and Xiaogang Wang. Revisiting the sibling head in object detector. In *CVPR*, 2020. 2
- [36] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, pages 5693–5703, 2019. 8
- [37] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris Kitani. Rethinking transformer-based set prediction for object detection. *CoRR*, abs/2011.10881, 2020. 2, 7
- [38] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *ECCV*, 2020. 2
- [39] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: fully convolutional one-stage object detection. In *ICCV*, 2019. 2, 7
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2
- [41] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui

- Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *TPAMI*, 2019. 8
- [42] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. In *NeurIPS*, 2020. 2
- [43] Yifan Xu, Weijian Xu, David Cheung, and Zhuowen Tu. Line segment detection using transformers without edges. In *CVPR*, pages 4257–4266, June 2021. 8
- [44] Brandon Yang, Gabriel Bender, Quoc V. Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. In *NeurIPS*, 2019. 2
- [45] Changqian Yu, Bin Xiao, Changxin Gao, Lu Yuan, Lei Zhang, Nong Sang, and Jingdong Wang. Lite-hrnet: A lightweight high-resolution network. In *CVPR*, pages 10440–10450, June 2021. 2
- [46] Jiahui Yu, Yuning Jiang, Zhangyang Wang, Zhimin Cao, and Thomas S. Huang. Unitbox: An advanced object detection network. In *MM*, 2016. 2
- [47] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z. Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *CVPR*, 2020. 2
- [48] Minghang Zheng, Peng Gao, Xiaogang Wang, Hongsheng Li, and Hao Dong. End-to-end object detection with adaptive clustering transformer. *CoRR*, abs/2011.09315, 2020. 2
- [49] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *CoRR*, abs/1904.07850, 2019. 2
- [50] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krähenbühl. Bottom-up object detection by grouping extreme and center points. In *CVPR*, 2019. 2
- [51] Chenchen Zhu, Fangyi Chen, Zhiqiang Shen, and Marios Savvides. Soft anchor-point object detection. In *ECCV*, 2020. 2
- [52] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *CVPR*, 2019. 2
- [53] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: deformable transformers for end-to-end object detection. *CoRR*, abs/2010.04159, 2020. 1, 2, 6, 7, 8