

Fast Convergence of DETR with Spatially Modulated Co-Attention

Peng Gao¹ Minghang Zheng³ Xiaogang Wang¹ Jifeng Dai² Hongsheng Li¹

¹Multimedia Laboratory, The Chinese University of Hong Kong

²SenseTime Research ³Peking University

1155102382@link.cuhk.edu.hk daijifeng@sensetime.com

{xgwang, hqli}@ee.cuhk.edu.hk

Abstract

The recently proposed Detection Transformer (DETR) model successfully applies Transformer to objects detection and achieves comparable performance with two-stage object detection frameworks, such as Faster-RCNN. However, DETR suffers from its slow convergence. Training DETR [4] from scratch needs 500 epochs to achieve a high accuracy. To accelerate its convergence, we propose a simple yet effective scheme for improving the DETR framework, namely Spatially Modulated Co-Attention (SMCA) mechanism. The core idea of SMCA is to conduct regression-aware co-attention in DETR by constraining co-attention responses to be high near initially estimated bounding box locations. Our proposed SMCA increases DETR’s convergence speed by replacing the original co-attention mechanism in the decoder while keeping other operations in DETR unchanged. Furthermore, by integrating multi-head and scale-selection attention designs into SMCA, our fully-fledged SMCA can achieve better performance compared to DETR with a dilated convolution-based backbone (45.6 mAP at 108 epochs vs. 43.3 mAP at 500 epochs). We perform extensive ablation studies on COCO dataset to validate the effectiveness of the proposed SMCA.

1. Introduction

The recently proposed DETR [4] has significantly simplified object detection pipeline by removing hand-crafted anchor [35] and non-maximum suppression (NMS) [2]. However, the convergence speed of DETR is slow compared with two-stage [11, 10, 35] or one-stage [27, 33, 25] detectors (500 vs. 40 epochs). Slow convergence of DETR increases the algorithm design cycle, makes it difficult for researchers to further extend this algorithm, and thus hinders its widespread usage.

In DETR, there are a series of object query vectors responsible for detecting objects at different spatial locations. Each object query interacts with the spatial visual features

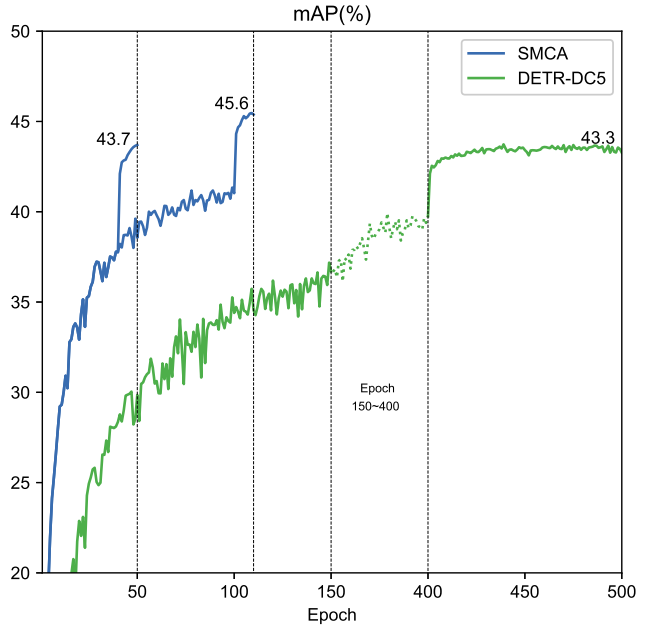


Figure 1. Comparison of convergence of DETR-DC5 trained for 500 epochs, and our proposed SMCA trained for 50 epochs and 108 epochs. The convergence speed of the proposed SMCA is much faster than the original DETR.

encoded by a Convolution Neural Network (CNN) [15] and adaptively collects information from spatial locations with a co-attention mechanism and then estimates the bounding box locations and object categories. However, in the decoder of DETR, the co-attended visual regions for each object query might be unrelated to the bounding box to be predicted by the query. Thus the decoder of DETR needs long training epochs to search for the properly co-attended visual regions to accurately identify the corresponding objects.

Motivated by this observation, we propose a novel module named Spatially Modulated Co-attention (SMCA), which is a plug-and-play module to replace the existing co-attention mechanism in DETR and achieves faster convergence and improved performance with very simple modi-

fications. The proposed SMCA dynamically predicts initial center and scale of the box corresponding to each object query to generate a 2D spatial Gaussian-like weight map. The weight map is element-wisely multiplied with the co-attention feature maps of object query and image features to more effectively aggregate query-related information from the visual feature map. In this way, the spatial weight map effectively modulates the search range of each object query’s co-attention to be properly around the initially estimated object center and scale. By leveraging the predicted Gaussian-distributed spatial prior, our SMCA can significantly speed up the training of DETR.

Although naively incorporating the spatially-modulated co-attention mechanism into DETR speeds up the convergence, the performance is worse compared with DETR (41.0 mAP at 50 epochs, 42.7 at 108 epochs vs. 43.3 mAP at 500 epochs). Motivated by the effectiveness of multi-head attention-based Transformer [40] and multi-scale feature [24] in previous research work, our SMCA is further augmented with the multi-scale visual feature encoding in the encoder and the multi-head attention in the decoder. For multi-scale visual feature encoding in the encoder, instead of naively rescaling and upsampling the multi-scale features from the CNN backbone to form a joint multi-scale feature map, Intra-scale and multi-scale self-attention mechanisms are introduced to directly and efficiently propagate information between the visual features of multiple scales. For the proposed multi-scale self-attention, visual features at all spatial locations of all scales interact with each other via self-attention. However, as the number of all spatial locations at all scales is quite large and leads to large computational cost, we introduce the intra-scale self-attention to alleviate the heavy computation. The properly combined intra-scale and multi-scale self-attention achieve efficient and discriminative multi-scale feature encoding. In the decoder, each object query can adaptively select features of proper scales via the proposed scale-selection attention. For the multiple co-attention heads in the decoder, all heads estimate head-specific object centers and scales to generate a series of different spatial weight maps for spatially modulating the co-attention features. Each of the multiple heads aggregates visual information from slightly different locations and thus improves the detection performance.

Our SMCA is motivated by the following research directions. DRAW [12] proposed a differential read-and-write operator with dynamically predicted Gaussian sampling points for image generation. Gaussian Transformer [13] has been proposed for accelerating natural language inference with Gaussian prior. Different from Gaussian Transformer, SMCA predicts a dynamically spatial weight map to tackle the dynamic search range of the objects. Deformable DETR [46] achieved fast convergence of DETR with learnable sparse sampling. Compared with Deformable DETR,

our proposed SMCA explores another direction for fast convergence of DETR by exploring dynamic Gaussian-like spatial prior. Besides, SMCA can accelerate the training of DETR by only replacing co-attention in the decoder. Deformable DETR replaces the Transformer with deformable attention for both the encoder and decoder, which explores local information rather than global information. SMCA demonstrates that exploring global information can also result in the fast convergence of DETR. Besides the above-mentioned methods, SMCA is also motivated by feature pyramids and dynamic modulation, which will be introduced in related work.

We summarize our contributions below:

- We propose a novel Spatial Modulated Co-Attention (SMCA), which can accelerate the convergence of DETR by conducting location-constrained object regression. SMCA is a plug-and-play module in the original DETR. The basic version of SMCA without multi-scale features and multi-head attention can already achieve 41.0 mAP at 50 epochs and 42.7 mAP at 108 epochs. It takes 265 V100 GPU hours to train the basic version of SMCA for 50 epochs.
- Our full SMCA further integrates multi-scale features and multi-head spatial modulation, which can further significantly improve and surpass DETR with much fewer training iterations. SMCA can achieve 43.7 mAP at 50 epochs and 45.6 mAP at 108 epochs, while DETR-DC5 achieves 43.3 mAP at 500 epochs. It takes 600 V100 GPU hours to train the full SMCA for 50 epochs.
- We perform extensive ablation studies on COCO 2017 dataset to validate the proposed SMCA module and the network design.

2. Related Work

2.1. Object Detection

Motivated by the success of deep learning on image classification [22, 15], deep learning has been successfully applied to object detection [11]. Deep learning-based object detection frameworks can be categorized into two-stage, one-stage, and end-to-end ones.

For two-stage object detectors including RCNN [11], Fast RCNN [10] and Faster RCNN [35], the region proposal layer generates a few regions from dense sliding windows first, and the ROI align [14] layer then extracts fine-grained features and perform classification over the pooled features. For one-stage detectors such as YOLO [33] and SSD [27], they conduct object classification and location estimation directly over dense sliding windows. Both two-stage and one-stage methods need complicated post-processing to generate the final bounding box predictions.

Recently, another branch of object detection methods [37, 36, 34, 4] beyond one-stage and two-stage ones has gained popularity. They directly supervise bounding box predictions end-to-end with Hungarian bipartite matching. However, DETR [4] suffered from slow convergence compared with two-stage and one-stage object detectors. Deformable DETR [46] accelerates the convergence speed of DETR via learnable sparse sampling coupled with multi-scale deformable encoder. TSP [38] analyzed the possible causes of slow convergence in DETR and identify co-attention and bipartite matching are two main causes. It then combined RCNN- or FCOS-based methods with DETR. TSP-RCNN and TSP-FCOS achieve fast convergence with better performance. Deformable DETR, TSP-RCNN and TSP-FCOS only explored local information while our SMCA explores global information with a self-attention and co-attention mechanism. Adaptive Clustering Transformer (ACT) [45] proposed a run-time pruning of attention on DETR’s encoder by LSH approximate clustering. Different from ACT, we accelerate the converging speed while ACT targets at acceleration of inference without re-training. UP-DETR [5] propose a novel self-supervised loss to enhance the convergence speed and performance of DETR.

Loss balancing and multi-scale information has been actively studied in object detection. There usually exist imbalance between positive and negative samples. Thus the gradient of negative samples would dominate the training process. Focal loss [25] proposed an improved version of cross entropy loss to attenuate the gradients generated by negative samples in object detection. Feature Pyramid Network (FPN) [24] and its variants [20] proposed a bottom-up and top-down way to generate multi-scale features for achieving better performance for object detection. Different from the multi-scale features generated from FPN, SMCA adopts a simple cascade of intra-scale and multi-scale self-attention modules to conduct information exchange between features at different positions and scales.

2.2. Transformer

CNN [23] and LSTM [16] can be used for modeling sequential data. CNN processes input sequences with a weight-shared sliding window manner. LSTM processes inputs with a recurrence mechanism controlled by several dynamically predicted gating functions. Transformer [40] introduces a new architecture beyond CNN and LSTM by performing information exchange between all pairs of input using key-query value attention. Transformer has achieved success on machine translation, after which Transformer has been adopted in different fields, including model pre-training [6, 31, 32, 3], visual recognition [30, 7], and multi-modality fusion [44, 8, 29]. Transformer has quadratic complexity for information exchange between all pairs of

inputs, which is difficult to scale up for longer input sequences. Many methods have been proposed to tackle this problem. Reformer [21] proposed a reversible FFN and clustering self-attention. Linformer [41] and FastTransformer [19] proposed to remove the softmax in the transformer and perform matrix multiplication between query and value first to obtain a linear-complexity transformer. LongFormer [1] perform self-attention within a local window instead of the whole input sequence. Transformer has been utilized in DETR to enhance the features by performing feature exchange between different positions and object query. In SMCA, intra-scale and multi-scale self-attention has been utilized for information exchange inside and outside each scale. In this paper, our SMCA is based on the original Transformer. We will explore memory-efficient transformers in SMCA in future work.

2.3. Dynamic Modulation

Dynamic modulation has been actively studied in different research fields of deep learning. In LSTM, a dynamic gate would be predicted to control the temporal information flow. Recent attention mechanism can be seen as a variant of dynamic modulation. Look-Attend-Tell [43] applied dynamic modulation in image captioning using attention. At each time step, an extra attention map is predicted and a weighted summation over the residual features and predict the word for the current step. The attention patterns in [43] can be interpreted, where the model is looking at. Dynamic filter [18] generates a dynamic convolution kernel from a prediction network and apply the predicted convolution over features in a sliding window fashion. Motivated by the dynamic filter, QGHC [9] adopted a dynamic group-wise filter to guide the information aggregation in the visual branch using language guided convolution. Lightweight convolution [42] used dynamic predicted depth-wise filters in machine translation and surpass the performance of Transformer. SE-Net [17] successfully applies channel-wise attention to modulate deep features for image recognition. Motivated by the dynamic modulation mechanism in previous research, we design a simple scale-selection attention to dynamically select the corresponding scale for each object query.

3. Spatially Modulated Co-Attention

3.1. Overview

In this section, we will first revisit the basic design of DETR [4] and then introduce the basic version of SMCA. After introducing SMCA, we will introduce how to integrate multi-head and scale-selection attention mechanisms into SMCA. The overall pipeline of SMCA is illustrated in Figure 2.

3.2. A Revisit of DETR

End-to-end object DETection with TRansformer (DETR) [4] formulates object detection as a set prediction problem. A Convolution Neural Network (CNN) [15] extracts visual feature maps $f \in \mathbb{R}^{C \times H \times W}$ from an image $I \in \mathbb{R}^{3 \times H_0 \times W_0}$, where H, W and H_0, W_0 are the height/width of the input image and the visual feature map, respectively.

The visual features augmented with position embedding f_{pe} would be fed into the encoder of the Transformer. Self-attention would be applied to f_{pe} to generate the key, query, and value features K, Q, V to exchange information between features at all spatial positions. To increase the feature diversity, such features would be split into multiple groups along the channel dimension for the multi-head self-attention. The multi-head normalized dot-product attention is conducted as

$$\begin{aligned} E_i &= \text{Softmax}(K_i^T Q_i / \sqrt{d}) V_i, \\ E &= \text{Concat}(E_1, \dots, E_H), \end{aligned} \quad (1)$$

where K_i, Q_i, V_i denote the i th feature group of the key, query, and value features. There are H groups for each type of features, and the output encoder features E is then further transformed and input into the decoder of the Transformer.

Given the visual feature E encoded from the encoder, DETR performs co-attention between object queries $O_q \in \mathbb{R}^{N \times C}$ and the visual features $E \in \mathbb{R}^{L \times C}$, where N denotes the number of pre-specified object queries and L is the number of the spatial visual features.

$$\begin{aligned} Q &= \text{FC}(O_q), \quad K, V = \text{FC}(E) \\ C_i &= \text{Softmax}(K_i^T Q_i / \sqrt{d}) V_i, \\ C &= \text{Concat}(C_1, \dots, C_H), \end{aligned} \quad (2)$$

where FC denotes a single-layer linear transformation, and C_i denotes the co-attended feature for the object query O_q from the i th co-attention head. The decoder's output features of each object query is then further transformed by a Multi-Layer Perceptron (MLP) to output class score and box location for each object.

Given box and class prediction, the Hungarian algorithm is applied between predictions and ground-truth box annotations to identify the learning targets of each object query.

3.3. Spatially Modulated Co-Attention

The original co-attention in DETR is unaware of the predicted bounding boxes and thus requires many iterations to generate the proper attention map for each object query. The core idea of our SMCA is to combine the learnable co-attention maps with handcrafted query spatial priors, which constrain the attended features to be around the object queries' initial estimations and thus to be more related

to the final object predictions. SMCA module is illustrated in the Figure 2 in orange.

Dynamic spatial weight maps. Each object query first dynamically predicts the center and scale of its responsible object, which are then used to generate a 2D Gaussian-like spatial weight map. The center of the Gaussian-like distribution are parameterized in the normalized coordinates of $[0, 1] \times [0, 1]$. The initial prediction of the normalized center $c_h^{\text{norm}}, c_w^{\text{norm}}$ and scale s_h, s_w of the Gaussian-like distribution for object query O_q is formulated as

$$\begin{aligned} c_h^{\text{norm}}, c_w^{\text{norm}} &= \text{sigmoid}(\text{MLP}(O_q)), \\ s_h, s_w &= \text{FC}(O_q), \end{aligned} \quad (3)$$

where the object query O_q is projected to obtain normalized prediction center in the two dimensions $c_h^{\text{norm}}, c_w^{\text{norm}}$ with a 2-layer MLP followed by a sigmoid activation function. The predicted center is then unnormalized to obtain the center coordinates c_h, c_w in the original image. O_q would also dynamically estimate the object scales s_h, s_w along the two dimensions to create the 2D Gaussian-like weight map, which is then used to re-weight the co-attention map to emphasize features around the predicted object location.

Objects in natural images show diverse scales and height/width ratios. The design of predicting width- and height-independent s_h, s_w can better tackle the complex object aspect ratios in real-world scenarios. For large or small objects, SMCA dynamically generates s_h, s_w of different values, so that the modulated co-attention map by the spatial weight map G can aggregate sufficient information from all parts of large objects or suppress background clutters for the small objects. After predicting the object center c_w, c_h and scale s_w, s_h , SMCA generates the Gaussian-like weight map as

$$G(i, j) = \exp\left(-\frac{(i - c_w)^2}{\beta s_w^2} - \frac{(j - c_h)^2}{\beta s_h^2}\right), \quad (4)$$

where $(i, j) \in [0, W] \times [0, H]$ is the spatial indices of the weight map G , and β is a hyper-parameter to modulate the bandwidth of the Gaussian-like distribution. In general, the weight map G assigns high importance to spatial locations near the center and low importance to positions far from the center. β can be manually tuned with a handcrafted scheme to ensure G covering a large spatial range at the beginning of training so that the network can receive more informative gradients.

Spatially-modulated co-attention. Given the dynamically generated spatial prior G , we modulate the co-attention maps C_i between object query O_Q and self-attention encoded feature E with the spatial prior G . For each co-attention map C_i generated with the dot-product attention (Eq. (2)), we modulate the co-attention maps C_i with the

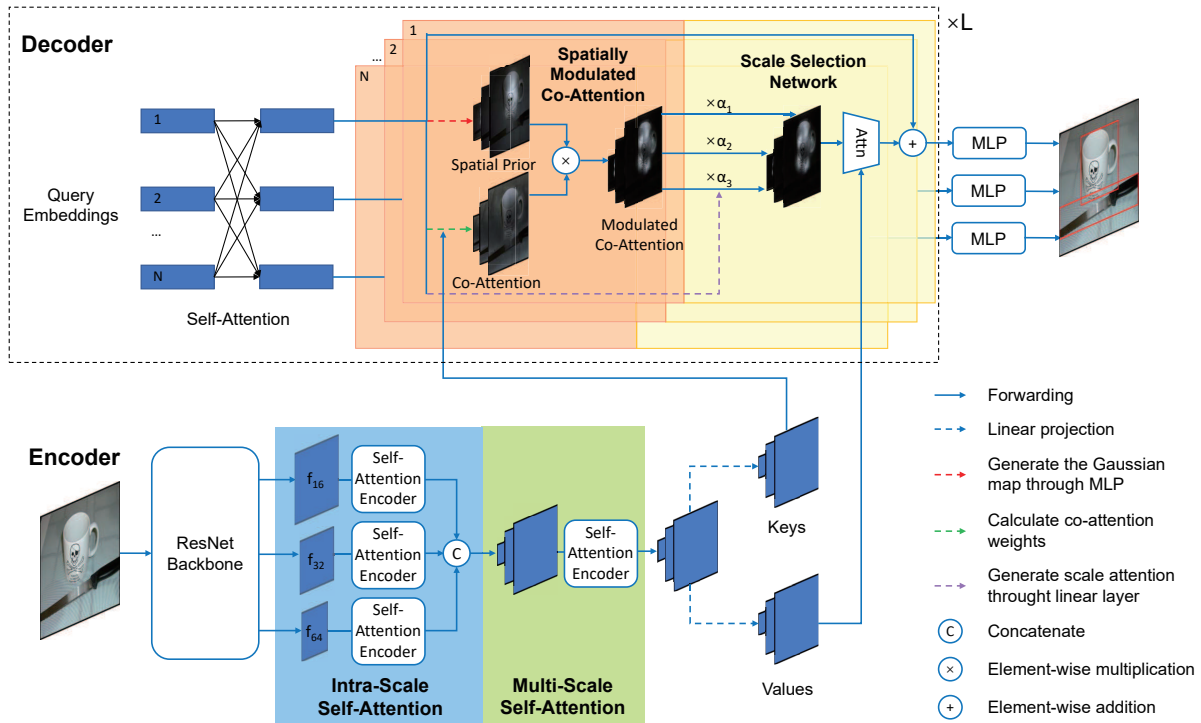


Figure 2. The overall pipeline of Spatially Modulated Co-Attention (SMCA) with intra-scale self-attention, multi-scale self-attention, spatial modulation, and scale-selection attention modules. Each object query performs spatially modulated co-attention and then predicts the target bounding boxes and their object categories. N stands for the number of object queries. L stands for the layers of decoder.

spatial weight map G , where G is shared for all co-attention heads in the basic version of our SMCA,

$$C_i = \text{softmax}(K_i^T Q_i / \sqrt{d} + \log G) V_i. \quad (5)$$

Our SMCA performs element-wise addition between the logarithm of the spatial map G and the dot-product co-attention $K_h^T Q_h / \sqrt{d}$ followed by softmax normalization over all spatial locations. By doing so, the decoder co-attention would weight more around the predicted bounding box locations, which can limit the search space of the spatial patterns of the co-attention and thus increases the convergence speed. The Gaussian-like weight map is illustrated in Figure 2, which constrains the co-attention to focus more on regions near the predicted bounding box location and thus significantly increases the convergence speed of DETR. In the basic version of SMCA, co-attention maps C_i of the multiple attention heads share the same Gaussian-like weight map G .

SMCA with multi-head modulation. We also investigate to modulate co-attention features differently for different co-attention heads. Each head starts from a head-shared center $[c_w, c_h]$, similar to that of the basic version of SMCA, and then predicts a head-specific center off-

set $[\Delta c_{w,i}, \Delta c_{h,i}]$ and head-specific scales $s_{w,i}, s_{h,i}$. The Gaussian-like spatial weight map G_i can thus be generated based on the head-specific center $[c_w + \Delta c_{w,i}, c_h + \Delta c_{h,i}]$ and scales $s_{w,i}, s_{h,i}$. The co-attention feature maps C_1, \dots, C_H can be obtained as

$$C_i = \text{softmax}(K_i^T Q_i / \sqrt{d} + \log G_i) V_i \quad \text{for } i = 1, \dots, H. \quad (6)$$

Different from Eq. (5) that shares $\log G$ for all attention heads, the above Eq. (6) modulates co-attention maps by head-specific spatial weight maps $\log G_i$. The multiple spatial weight maps can emphasize diverse context and improve the detection accuracy.

SMCA with multi-scale visual features. Feature pyramid is popular in object detection frameworks and generally leads to significant improvements over single-scale feature encoding. Motivated by the feature pyramid network [24] in previous works, we also integrate multi-scale features into SMCA. The basic version of SMCA conducts co-attention between object queries and single-scale feature maps. As objects naturally have different scales, we can further improve the framework by replacing single-scale feature encoding with multi-scale feature encoding in the encoder of

the Transformer.

Given an image, the CNN extracts the multi-scale visual features with downsampling rates 16, 32, 64 to obtain multi-scale features f_{16}, f_{32}, f_{64} , respectively. The multi-scale features are directly obtained from the CNN backbone and Feature Pyramid Network is not used to save the computational cost. For multi-scale self-attention encoding in the encoder, features at all locations of different scales are treated equally. The self-attention mechanism propagates and aggregates information between all feature pixels of different scales. However, the number of feature pixels of all scales is quite large and the multi-scale self-attention operation is therefore computationally costly. To tackle the issue, we introduce the intra-scale self-attention encoding as an auxiliary operator to assist the multi-scale self-attention encoding. Specifically, dot-product attention is used to propagate and aggregate features only between feature pixels within each scale. The weights of the Transformer block (with self-attention and feedforward sub-networks) are shared across different scales. Our empirical study shows that parameter sharing across scales enhances the generalization capability of intra-scale self-attention encoding. For the final design of the encoder in SMCA, it adopts 2 blocks of intra-scale self-attention encoding, followed by 1 block of multi-scale self-attention, and another 2 blocks of intra-scale self-attention. The design has a very similar detection performance to that of 5 blocks of multi-scale self-attention encoding but has a much smaller computational footprint.

Given the encoded multi-scale features E_{16}, E_{32}, E_{64} with downsampling rates of 16, 32, 64, a naive solution for the decoder to perform co-attention would be first rescaling and concatenating the multi-scale features to form a single-scale feature map, and then conducting co-attention between object query and the resulting feature map. However, we notice that some queries might only require information from a specific scale but not always from all the scales. For example, the information for small objects is missing in low-resolution feature map E_{64} . Thus the object queries responsible for small objects should more effectively acquire information only from high-resolution feature maps. On the other hand, traditional methods, such as FPN, assigns each bounding box explicitly to the feature map of a specific scale. Different from FPN [24], we propose to automatically select scales for each box using learnable scale-attention attention. Each object query generates scale-selection attention weights as

$$\alpha_{16}, \alpha_{32}, \alpha_{64} = \text{Softmax}(\text{FC}(O_q)), \quad (7)$$

where $\alpha_{16}, \alpha_{32}, \alpha_{64}$ stand for the importance of selecting f_{16}, f_{32}, f_{64} . To conduct co-attention between the object query O_q and the multi-scale visual features E_{16}, E_{32}, E_{64} , we first obtain the multi-scale key and value

features $K_{i,16}, K_{i,32}, K_{i,64}$ and $V_{i,16}, V_{i,32}, V_{i,64}$ for attention head i , respectively, from E_{16}, E_{32}, E_{64} with separate linear projections. To conduct co-attention for each head i between O_q and key/value features of each scale $j \in \{16, 32, 64\}$, the spatially-modulated co-attention in Eq. (5) is adaptively weighted and aggregated by the scale-selection weights $\alpha_{16}, \alpha_{32}, \alpha_{64}$ as

$$C_{i,j} = \text{Softmax}(K_{i,j}^T Q_i / \sqrt{d} + \log G_i) V_{i,j} \odot \alpha_j, \quad (8)$$

$$C_i = \sum_{\text{all } j} C_{i,j}, \quad \text{for } j \in \{16, 32, 64\}, \quad (9)$$

where $C_{i,j}$ stands for the co-attention features between the i th co-attention head between query and visual features of scale j . $C_{i,j}$'s are weightedly aggregated according to the scaled attention weights α_j obtained in Eq. (7). With such a scale-selection attention mechanism, the scale most related to each object query is softly selected while the visual features from other scales are suppressed.

Equipped with intra-inter-scale attention and scale selection attention mechanisms, our full SMCA can better tackle object detection than the basic version.

SMCA box prediction. After conducting co-attention between the object query O_q and the encoded image features, we can obtain the updated features $D \in \mathbb{R}^{N \times C}$ for object query O_q . In the original DETR, a 3-layer MLP and a linear layer are used to predict the bounding box and classification confidence. We denote the prediction as

$$\text{Box} = \text{Sigmoid}(\text{MLP}(D)), \quad (10)$$

$$\text{Score} = \text{FC}(D), \quad (11)$$

where ‘‘Box’’ stands for the center, height, width of the predicted box in the normalized coordinate system, and ‘‘Score’’ stands for the classification prediction. In SMCA, co-attention is constrained to be around the initially predicted object center $[c_h^{\text{norm}}, c_w^{\text{norm}}]$. We then use the initial center as a prior for constraining bounding box prediction, which is denoted as

$$\begin{aligned} \widehat{\text{Box}} &= \text{MLP}(D), \\ \widehat{\text{Box}}[:2] &= \widehat{\text{Box}}[:2] + [c_h^{\text{norm}}, c_w^{\text{norm}}], \\ \text{Box} &= \text{Sigmoid}(\widehat{\text{Box}}), \end{aligned} \quad (12)$$

where $\widehat{\text{Box}}$ stand for the box prediction, and $[c_h^{\text{norm}}, c_w^{\text{norm}}]$ represents the center of initial object prediction before the sigmoid function. In Eq. (12), we add the center of predicted box with the center of initial spatial prior $[c_h^{\text{norm}}, c_w^{\text{norm}}]$ before the sigmoid function. This procedure ensures that the bounding box prediction is highly related to the highlighted co-attention regions in SMCA.

Method	Epochs	time(s)	GFLOPs	mAP	AP _S	AP _M	AP _L
DETR	500	0.038	86	42.0	20.5	45.8	61.1
DETR-DC5	500	0.079	187	43.3	22.5	47.3	61.1
SMCA w/o multi-scale	50	0.043	86	41.0	21.9	44.3	59.1
SMCA w/o multi-scale	108	0.043	86	42.7	22.8	46.1	60.0
SMCA	50	0.100	152	43.7	24.2	47.0	60.4
SMCA	108	0.100	152	45.6	25.9	49.3	62.6

Table 1. Comparison with DETR model over training epochs, mAP, inference time and GFLOPs.

4. Experiments

4.1. Experiment setup

Dataset. We validate our proposed SMCA over COCO 2017 [26] dataset. Specifically, we train on COCO 2017 training dataset and validate on the validation dataset, which contains 118k and 5k images, respectively. We report mAP for performance evaluation following previous research [4].

Implementation details. We follow the experiment setup in the original DETR [4]. We denote the features extracted by ResNet-50 [15] as SMCA-R50. Different from DETR, we use 300 object queries instead of 100 and replace the original cross-entropy classification loss with focal loss [25]. To better tackle the positive-negative imbalance problem in foreground/background classification. The initial probability of focal loss is set as 0.01 to stabilize the training process.

We report the performance trained for 50 epochs and the learning rate decreases to 1/10 of its original value at the 40th epoch. The learning rate is set as 10^{-4} for the Transformer encoder-encoder and 10^{-5} for the pre-trained ResNet backbone and optimized by AdamW optimizer [28].

For multi-scale feature encoding, we use downsampling ratios of 16, 32, 64 by default. For bipartite matching [37, 4], the coefficients of classification loss, L1 distance loss, GIoU loss is set as 2, 5, 2, respectively. After bounding box assignment via bipartite matching, SMCA is trained by minimizing the classification loss, bounding box L1 loss, and GIoU loss with coefficients 2, 5, 2, respectively. For Transformer layers [40], we use post-norm similar to those in previous approaches [4]. We use random crop for data augmentation with the largest width or height set as 1333 for all experiments following [4]. All models are trained on 8 V100 GPUs with 1 image per GPU.

4.2. Comparison with DETR

SMCA shares the same architecture with DETR except for the proposed new co-attention modulation in the decoder and an extra linear network for generating the spatial modulation prior. The increase of computational cost of SMCA

and training time of each epoch are marginal. For SMCA with single-scale features (denoted as “SMCA w/o multi-scale”), we keep the dimension of self-attention to be 256 and the intermediate dimension of FFN to be 2048. For SMCA with multi-scale features, we set the intermediate dimension of FFN to be 1024 and use 5 layers of intra-scale and multi-scale self-attention in the encoder to have similar amount of parameters and fair comparison with DETR. As shown in Table 1, the performance of “SMCA w/o multi-scale” reaches 41.0 mAP with single-scale features and 43.7 mAP with multi-scale features at 50 epochs. Given longer training procedure, mAP of SMCA increases from 41.0 to 42.7 with single-scale features and from 43.7 to 45.6 with multi-scale features. “SMCA w/o multi-scale” can achieve better AP_s and AP_M compared with DETR. SMCA can achieve better overall performance on objects of all scales by adopting multi-scale information and the proposed spatial modulation. The convergence speed of SMCA is 10 times faster than DETR-based methods.

Given the significant increase of convergence speed and performance, the FLOPs and the increase of inference time of SMCA are marginal. With single-scale features, the inference time increases from $0.038s \rightarrow 0.041s$ and FLOPs increase by 0.06G. With multi-scale features, the inference speed increase from $0.079s \rightarrow 0.100s$, while the GFLOPs actually decrease because our multi-scale SMCA only uses 5 layers of self-attention layers for the encoder. Thin layers in the Transformer and convolution without dilation in the last stage of ResNet backbone achieve similar efficiency as the original dilated DETR model.

4.3. Ablation Study

To validate different components of our proposed SMCA, we perform ablation studies on the importance of the proposed spatial modulation, multi-head vs. head-shared modulation, and multi-scale encoding and scale-selection attention in comparison with the baseline DETR.

The baseline DETR model. We choose DETR with ResNet-50 backbone as our baseline model. It is trained for 50 epochs with the learning rate dropping to 1/10 of the

Method		mAP	AP50	AP75
Baseline	DETR-R50	34.8	56.2	36.9
Head-shared Spatial Modulation	+Indep. (bs8)	40.2	61.4	42.7
	+Indep. (bs16)	40.2	61.3	42.9
	+Indep. (bs32)	39.9	61.0	42.4
Multi-head Spatial Modulation	+Fixed	38.5	60.7	40.2
	+Single	40.4	61.8	43.3
	+Indep.	41.0	62.2	43.6

Table 2. Ablation study on the importance of spatial modulation, multi-head mechanism. mAP, AP50, and AP75 are reported on COCO 2017 validation set.

Method		mAP	Params (M)
SMCA		41.0	41.0
SMCA (2Intra-Multi-2Intra)		43.7	39.5
SMCA w/o SSA (2Intra-Multi-2Intra)		42.6	39.5
3Intra		42.9	37.9
3Multi		43.3	37.9
5Intra		43.3	39.5
Weight Share	Shared FFN	43.0	42.2
	Shared SA	42.8	44.7
	No Share	42.3	47.3

Table 3. Ablation study on the importance of combining intra-scale and multi-scale propagation, and the weight sharing for intra-scale self-attention. ‘‘Shared FFN’’ stands for only sharing weights of the feedforward network of intra-scale self-attention. ‘‘Shared SA’’ stands for sharing the weights of the self-attention network. ‘‘No share’’ stands for no weight sharing in intra-scale self attention.

original value at the 40th epoch. Different from the original DETR, we increase the object query from 100 to 300 and replace the original cross entropy loss with focal loss. As shown in Table 2, the baseline DETR model can achieve an mAP of 34.8 at 50 epochs.

Head-shared spatially modulated co-attention. Based on the baseline DETR, we first test adding a head-shared spatial modulation as specified in Eq. (5) by keeping factors including the learning rate, training schedule, self-attention parameters, and coefficients of the loss to be the same as the baseline. The spatial weight map is generated based on the predicted height and width shared for all heads contain height- and width-independent scale prediction to better tackle the scale variance problem. We denote the method as ‘‘Head-shared Spatial Modulation + Indep.’’ in Table 3. The performance increase from 34.8 to 40.2 compared with baseline DETR. The large performance gain (+5.4) validates the effectiveness of SMCA, which not only acceler-

ates the convergence speed of DETR but also improve its performance by a large margin. We further test the performance of head-shared spatial modulation with different batch sizes of 8, 16, and 32 as shown in Table 3. The results show that our SMCA is insensitive to different batch sizes.

Multi-head vs. head-shared spatially modulated co-attention. For spatial modulation with multiple heads of separate predictable scales, all heads in Transformer are modulated by different spatial weight maps G_i following Eq. (6). All heads start from the same object center and predict offsets w.r.t. the common center and head-specific scales. The design of multi-head spatial modulation for co-attention enables the model to learn diverse attention patterns simultaneously. After switching from head-shared spatial modulation to multi-head spatial modulation (denoted as ‘‘Multi-head Spatial Modulation + Indep.’’ in Table 2), the performance increases from 40.2 to 41.0 compared with the head-shared modulated co-attention in SMCA. The importance of multi-head mechanism has also been discussed in Transformer [40]. From visualization in Figure 3, we observe that the multi-head modulation naturally focuses on different parts of the objects to be predicted by the object queries.

Design of multi-head spatial modulation for co-attention.

We test whether the width and height scales of the spatial weight maps should be manually set, shared, or independently predicted. As shown in Table 2, we test fixed-scale Gaussian-like spatial map (only predicting the center and fixing the scale of the Gaussian-like distribution to be the constant 1). The fixed-scale spatial modulation results in a 38.5 mAP (denoted as ‘‘+Fixed’’), which has +3.7 gain over the baseline DETR-R50 and validates the effectiveness of predicting centers for spatial modulation to constrain the co-attention. As objects in natural images have varying sizes, scales can be predicted to adapt to objects of different size. Thus we allow the scale to be a single predictable variable as in Eq. (3). If such a single predictable scale for spatial modulation (denoted as ‘‘+Single’’), SMCA can achieve 40.4 mAP and is +1.9 compared with the above fixed-scale modulation. By further predicting independent scales for height and width, our SMCA can achieve 41.0 mAP (denoted as ‘‘+Indep.’’), which is +0.6 higher compared with the SMCA with a single predictable scale. The results demonstrate the importance of predicting height and width scales for the proposed spatial modulation. As visualized by the co-attention patterns in Figure 3, we observe that independent spatial modulation can generate more accurate and compact co-attention patterns compared with fixed-scale and shared-scale spatial modulation.

Multi-scale feature encoding and scale-selection attention. The above SMCA only conducts co-attention between single-scale feature maps and the object query. As

objects in natural images exist in different scales, we conduct multi-scale feature encoding in the encoder via adopting 2 layers of intra-scale self-attention, followed by 1 layer of multi-scale self-attention, and then another 2 layers of intra-scale self-attention. We denote the above design as “SMCA (2Intra-Multi-2Intra)”. As shown in Table 3, we start from SMCA with a single-scale visual feature map, which achieves 41.0 mAP. After integrating multi-scale features with the 2intra-multi-2intra self-attention design, the performance can be enhanced from 41.0 to 43.7. As we introduce 3 convolutions to project features output from ResNet-50 to 256 dimensions, we make the hidden dimension of FFN decrease from 2048 to 1024 and the number of encoder layer decrease from 6 to 5 to make the parameter comparable to other models. To validate the effectiveness of scale-selection attention (SSA), we perform ablation studies on SMCA without integrating SSA (denoted as “SMCA w/o SSA”). As shown in Table 3, SMCA w/o SSA decreases the performance from 43.7 to 42.6.

After validating the effectiveness of the proposed multi-scale feature encoding and scale-selection attention module, we further validate the effectiveness of the design of 2intra-multi-2intra-scale self-attention. By switching the 2intra-multi-2intra design to simply stacking 5 intra-scale self-attention layers, the performance drops from 43.7 to 43.3, due to the lack of cross-scale information exchange. 5 layers of intra-scale self-attention (denoted as “5Intra”) encoder achieves better performance than 3Intra self-attention, which validates the effectiveness of a deeper intra-scale self-attention encoder. A 3-layer multi-scale (denoted as “3Multi”) self-attention encoder achieves better performance than a 3-layer intra-scale (3Intra) self-attention encoder. It demonstrates that enabling multi-scale information exchange leads to better performance than only conducting intra-scale information exchange alone. However, the large increase of FLOPs by replacing intra-scale with multi-scale self-attention encoder makes us choose a combination of intra-scale and multi-scale self-attention encoders, namely, the design of 2intra-inter-2intra. In the previously mentioned multi-scale encoder, we share both Transformer and FFN weights for features from intra-scale self-attention layers, which reduces the number of parameters and learns common patterns of multi-scale features. It increases the generalization of the proposed SMCA and achieves a better performance of 43.7 with fewer parameters.

Visualization of SMCA. We provide visualization of co-attention weight maps by SMCA. As shown in Figure 3, we compare the detection result of fixed-scale SCMA, single-scale SMCA, and independent-scale SMCA (default SMCA). From the visualization, we can see independent-scale SMCA can better tackle objects of large aspect ratios. Different spatial modulation heads focus on different parts of the object to aggregate diverse information for final ob-

ject recognition. Finally, we show the co-attention map of the original DETR co-attention. Our SMCA can better focus on features around the object of interest, for which the query needs to estimate, while DETR’s co-attention maps show sparse patterns and are unrelated to the object it aims to predict.

4.4. Overall Performance Comparison

In Table 4, we compare our proposed SMCA with other object detection frameworks on COCO 2017 validation set. DETR [4] uses an end-to-end Transformer for object detection. DETR-R50 and DETR-DC5-R50 stand for DETR with ResNet-50 and DETR with dilated ResNet-50 backbone. Compared with DETR, our SMCA can achieve fast convergence and better performance in terms of detection of the small, medium, and large objects. Faster RCNN [35] with FPN [24] is a two-stage approach for object detection. Our method can achieve better mAP than Faster RCNN-FPN-R50 at 109 epochs (45.6 vs 42.0 AP). As Faster RCNN uses ROI-Align and feature pyramid with downsampled {8, 16, 32, 64} features, Faster RCNN is superior at detecting small objects (26.6 vs 25.9 mAP). Thanks to the multi-scale self-attention mechanism that can propagate information between features at all scales and positions, our SMCA is better for localizing large objects (62.6 vs 53.4 AP).

Deformable DETR [46] replaces the original self-attention of DETR with local deformable attention for both the encoder and the decoder. It achieves faster convergence compared with the original DETR. Exploring local information in Deformable DETR results in fast convergence at the cost of degraded performance for large objects. Compared with DETR, the AP_L of Deformable DETR drops from 61.1 to 58.0. Our SMCA explores a new approach for fast convergence of the DETR by performing spatially modulated co-attention. As SMCA constrains co-attention near dynamically estimated object locations, SMCA achieves faster convergence by reducing the search space in co-attention. As SMCA uses global self-attention for information exchange between all scales and positions, our SMCA can achieve better performance for large objects compared with Deformable DETR. Deformable DETR uses downsampled 8, 16, 32, 64 multi-scale features and 8 sampling points for deformable attention. Our SMCA only uses downsampled 16, 32, 64 features and 1 center point for dynamic Gaussian-like spatial prior. SCMA achieves comparable mAP with Deformable DETR at 50 epochs (43.7 vs. 43.8 AP). As SMCA focuses more on global information and deformable DETR focuses more on local features, SMCA is better at detecting large objects (60.4 vs 59.0 AP) while inferior at detecting small objects (24.2 vs 26.4 AP).

UP-DETR [5] explores unsupervised learning for DETR. UP-DETR can achieve fast convergence and better performance compared with the original DETR due to the ex-

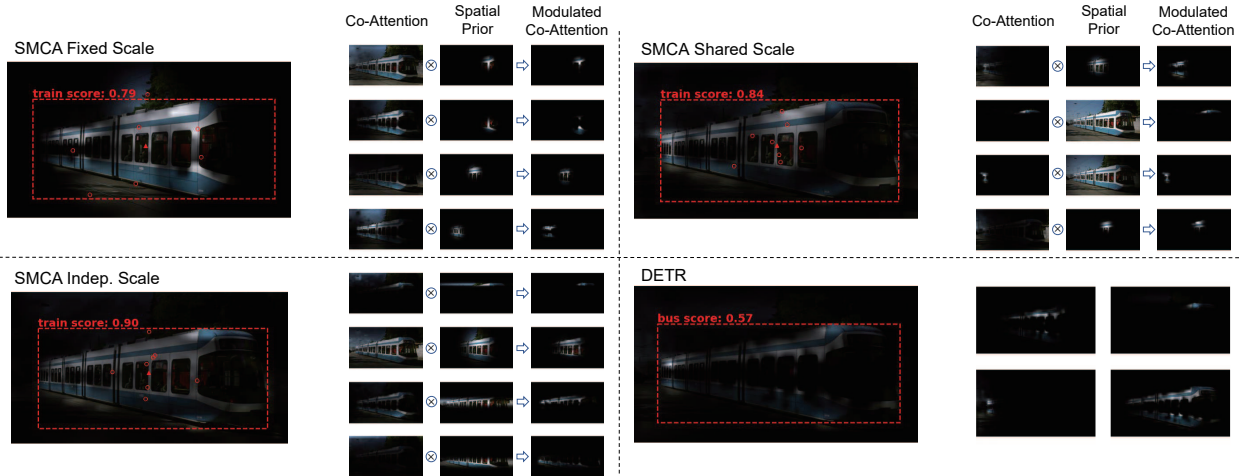


Figure 3. Visualization of co-attention of SMCA with fixed-scale, single-scale, independent-scale spatial modulation, and co-attention of DETR. The larger images show the average co-attention of 8 heads. Small images show the attention pattern of each head. In the head-specific modulation of co-attention of SMCA, we visualize the process of spatial modulation. Red circles in SMCA variants stand for the head-specific offset starting from the same red rectangular center.

Model	Epochs	GFLOPs	Params (M)	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
DETR-R50 [4]	500	86	41	42.0	62.4	44.2	20.5	45.8	61.1
DETR-DC5-R50 [4]	500	187	41	43.3	63.1	45.9	22.5	47.3	61.1
Faster RCNN-FPN-R50 [4]	36	180	42	40.2	61.0	43.8	24.2	43.5	52.0
Faster RCNN-FPN-R50++ [4]	108	180	42	42.0	62.1	45.5	26.6	45.4	53.4
Deformable DETR-R50 (Single-scale) [46]	50	78	34	39.7	60.1	42.4	21.2	44.3	56.0
Deformable DETR-R50 (50 epochs) [46]	50	173	40	43.8	62.6	47.7	26.4	47.1	58.0
Deformable DETR-R50 (150 epochs) [46]	150	173	40	45.3	*	*	*	*	*
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
TSP-FCOS-R50 [38]	36	189	*	43.1	62.3	47.0	26.6	46.8	55.9
TSP-RCNN-R50 [38]	36	188	*	43.8	63.3	48.3	28.6	46.9	55.7
TSP-RCNN+-R50 [38]	96	188	*	45.0	64.5	49.6	29.7	47.7	58.0
SMCA-R50	50	152	40	43.7	63.6	47.2	24.2	47.0	60.4
SMCA-R50	108	152	40	45.6	65.5	49.1	25.9	49.3	62.6
DETR-R101 [4]	500	152	60	43.5	63.8	46.4	21.9	48.0	61.8
DETR-DC5-R101 [4]	500	253	60	44.9	64.7	47.7	23.7	49.5	62.3
Faster RCNN-FPN-R101 [4]	36	256	60	42.0	62.1	45.5	26.6	45.4	53.4
Faster RCNN-FPN-R101+ [4]	108	246	60	44.0	63.9	47.8	27.2	48.1	56.0
TSP-FCOS-R101 [38]	36	255	*	44.4	63.8	48.2	27.7	48.6	57.3
TSP-RCNN-R101 [38]	36	254	*	44.8	63.8	49.2	29.0	47.9	57.1
TSP-RCNN+-R101 [38]	96	254	*	46.5	66.0	51.2	29.9	49.7	59.2
SMCA-R101	50	218	58	44.4	65.2	48.0	24.3	48.5	61.0

Table 4. Comparison with DETR-like object detectors on COCO 2017 validation set.

exploitation of unsupervised auxiliary tasks. The convergence speed and performance of SMCA is better than UP-DETR (45.6 at 108 epochs vs. 42.8 at 300 epochs). TSP-FCOS and TSP-RCNN [38] combines DETR’s Hungarian matching with FCOS [39] and RCNN [35] detectors, which results in

faster convergence and better performance than DETR. As TSP-FCOS and TSP-RCNN inherit the structure of FCOS and RCNN that uses local-region features for bounding box detection, they are strong at small objects but weak at large ones, similar to above mentioned deformable DETR

and Faster RCNN-FPN. For short training schedules, TSP-RCNN and GMCA-R50 achieve comparable mAP (43.8 at 38 epochs vs 43.7 at 50 epochs), which are better than 43.1 at 38 epochs by TSP-FCOS. For long training schedules, SMCA can achieve better performance than TSP-RCNN (45.6 at 108 epochs vs 45.0 at 96 epochs). We observe similar trends by replacing ResNet-50 backbone with ResNet-101 backbone as shown in the lower half part of Table 4.

5. Conclusion

DETR [4] proposed an end-to-end solution for object detection beyond previous two-stage [35] and one-stage approaches [33]. By integrating the Spatially Modulated Co-attention (SMCA) into DETR, the original 500 epochs training schedule can be reduced to 108 epochs and mAP increases from 43.4 to 45.6 under comparable inference cost. SMCA demonstrates the potential power of exploring global information for achieving high-quality object detection. In the future, we will explore the application of SMCA in more scenarios beyond object detection, such as general visual representation learning. We will also explore flexible fusions of local and global features for faster and more robust object detection.

References

- [1] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020. 3
- [2] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. Soft-nms—improving object detection with one line of code. In *Proceedings of the IEEE international conference on computer vision*, pages 5561–5569, 2017. 1
- [3] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 3
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. 1, 3, 4, 7, 9, 10, 11
- [5] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. *arXiv preprint arXiv:2011.09094*, 2020. 3, 9, 10
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3
- [8] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6639–6648, 2019. 3
- [9] Peng Gao, Hongsheng Li, Shuang Li, Pan Lu, Yikang Li, Steven CH Hoi, and Xiaogang Wang. Question-guided hybrid convolution for visual question answering. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 469–485, 2018. 3
- [10] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 1, 2
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Region-based convolutional networks for accurate object detection and segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 38(1):142–158, 2015. 1, 2
- [12] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623*, 2015. 2
- [13] Maosheng Guo, Yu Zhang, and Ting Liu. Gaussian transformer: a lightweight approach for natural language inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6489–6496, 2019. 2
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2, 4, 7
- [16] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 3
- [17] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 3
- [18] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *Advances in neural information processing systems*, pages 667–675, 2016. 3
- [19] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020. 3
- [20] Seung-Wook Kim, Hyong-Keun Kook, Jee-Young Sun, Mun-Cheon Kang, and Sung-Jea Ko. Parallel feature pyramid network for object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 234–250, 2018. 3
- [21] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020. 3
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 2

- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [3](#)
- [24] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. [2](#), [3](#), [5](#), [6](#), [9](#)
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. [1](#), [3](#), [7](#)
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [7](#)
- [27] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. [1](#), [2](#)
- [28] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018. [7](#)
- [29] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vlb- bert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23, 2019. [3](#)
- [30] Niki Parmar, Prajit Ramachandran, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. 2019. [3](#)
- [31] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018. [3](#)
- [32] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. [3](#)
- [33] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. [1](#), [2](#), [11](#)
- [34] Mengye Ren and Richard S Zemel. End-to-end instance segmentation with recurrent attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6656–6664, 2017. [3](#)
- [35] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016. [1](#), [2](#), [9](#), [10](#), [11](#)
- [36] Amaia Salvador, Miriam Bellver, Victor Campos, Manel Baradad, Ferran Marques, Jordi Torres, and Xavier Giro-i Nieto. Recurrent neural networks for semantic instance segmentation. *arXiv preprint arXiv:1712.00617*, 2017. [3](#)
- [37] Russell Stewart, Mykhaylo Andriluka, and Andrew Y Ng. End-to-end people detection in crowded scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2325–2333, 2016. [3](#), [7](#)
- [38] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris Kitani. Rethinking transformer-based set prediction for object detection. *arXiv preprint arXiv:2011.10881*, 2020. [3](#), [10](#)
- [39] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 9627–9636, 2019. [10](#)
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30:5998–6008, 2017. [2](#), [3](#), [7](#), [8](#)
- [41] Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. [3](#)
- [42] Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. *arXiv preprint arXiv:1901.10430*, 2019. [3](#)
- [43] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015. [3](#)
- [44] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6281–6290, 2019. [3](#)
- [45] Minghang Zheng, Peng Gao, Xiaogang Wang, Hongsheng Li, and Hao Dong. End-to-end object detection with adaptive clustering transformer. *arXiv preprint arXiv:2011.09315*, 2020. [3](#)
- [46] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [2](#), [3](#), [9](#), [10](#)