

# Cross-domain Detection Transformer based on Spatial-aware and Semantic-aware Token Alignment

Jinhong Deng Xiaoyue Zhang Wen Li Lixin Duan  
University of Electronic Science and Technology of China

{jhdeng1997, liwenbnu, lxduan}@gmail.com, xzhangeo@connect.ust.hk

## Abstract

*Detection transformers like DETR [3] have recently shown promising performance on many object detection tasks, but the generalization ability of those methods is still quite challenging for cross-domain adaptation scenarios. To address the cross-domain issue, a straightforward way is to perform token alignment with adversarial training in transformers. However, its performance is often unsatisfactory as the tokens in detection transformers are quite diverse and represent different spatial and semantic information. In this paper, we propose a new method called Spatial-aware and Semantic-aware Token Alignment (SSTA) for cross-domain detection transformers. In particular, we take advantage of the characteristics of cross-attention as used in detection transformer and propose the spatial-aware token alignment (SpaTA) and the semantic-aware token alignment (SemTA) strategies to guide the token alignment across domains. For spatial-aware token alignment, we can extract the information from the cross-attention map (CAM) to align the distribution of tokens according to their attention to object queries. For semantic-aware token alignment, we inject the category information into the cross-attention map and construct domain embedding to guide the learning of a multi-class discriminator so as to model the category relationship and achieve category-level token alignment during the entire adaptation process. We conduct extensive experiments on several widely-used benchmarks, and the results clearly show the effectiveness of our proposed method over existing state-of-the-art baselines.*

## 1. Introduction

Object detection, as a fundamental task for visual understanding, has been one of the most attractive research problems in the computer vision community [2, 3, 10, 23, 26, 27, 35]. With the thriving of deep convolutional neural networks (CNN) [12, 19], many CNN-based object detection approaches (e.g., Faster RCNN [27] and FCOS [35])

have been proposed in the last decade. Recently, detection transformers (e.g., DETR [3]) have gained increasing attention from researchers. Based on the design of visual transformer, detection transformers remove the requirement of hand-designed components such as non-maximum suppression (NMS) and anchor generation in traditional CNN-based object detection methods, and at the same time, achieve new state-of-the-art performance in many object detection tasks [3, 24, 29, 39, 43, 47]. Despite the success of detection transformers, the cross-domain generalization ability remains a challenge when adapting a learned model to a novel domain (i.e., target domain). Usually, existing detection transformers often suffer from severe performance degradation due to domain discrepancy between the source and target domains [38].

However, addressing the domain shift issue for detection transformers is non-trivial. Researchers have proposed many ways to improve the cross-domain generalization ability for CNN-based object detectors. For example, a variety of studies for cross-domain object detection (CDOD) [5, 8, 25, 31, 46] are proposed to eliminate the domain discrepancy by aligning the feature distributions of the source and target via adversarial training. Similarly, for the cross-domain detection transformer, a potential and straightforward solution for the cross-domain detection transformer is to perform token alignment with adversarial training, since the visual features are often converted into tokens as the input to the transformer blocks. However, aligning the token distributions is difficult, especially when there exists a significant domain gap between domains.

Recent work [38] attempts to apply adversarial training strategies on tokens in transformers, but the improvements are still unsatisfactory. One of the major reasons is that tokens in detection transformers are quite diverse. In detection transformers (e.g., DETR), the tokens are passed through several multi-head self-attention layers to obtain new token embeddings for representing different spatial and semantic information. Then, object queries are introduced to probe useful tokens and leverage those tokens to predict the positions and categories of different objects. On the one

hand, since some tokens are more useful while less for others, it is desirable to take the importance of tokens into consideration in the cross-domain detection transformer. On the other hand, the semantic information embedded in tokens is also helpful for aligning the token distributions w.r.t. the corresponding category, which can ease the adversarial training process.

In this work, we propose a new cross-domain detection method named Spatial-aware and Semantic-aware Token Alignment (SSTA) under the transformer framework. In particular, we take advantage of the characteristics of cross-attention as used in the detection transformers and newly developed two strategies, *i.e.*, spatial-aware token alignment (SpaTA) and semantic-aware token alignment (SemTA) to guide the token alignment across domains. The cross-attention in the decoder of SSTA utilizes the object queries to aggregate information from encoder outputs (tokens). During this process, only a small part of them are attended to for detecting objects accurately. For spatial-aware token alignment, we can extract the information from the cross-attention map (CAM) to align the distribution of tokens according to their attention to object queries. For semantic-aware token alignment, we inject the category information into the cross-attention map and construct domain embeddings to guide the learning of a multi-class domain discriminator so as to model the category relationship and achieve category-level alignment during the entire adaptation process.

We have conducted extensive experiments on three domain adaptive benchmarks, including adverse weather, synthetic-to-real, and scene adaptation, where we achieve new state-of-the-art performance for cross-domain object detection. The experimental results show the effectiveness of our proposed method. We also show the usefulness of each component in our approach by conducting careful ablation studies. The contributions of our work are three-fold:

- We propose a novel approach named Spatial-aware and Semantic-aware Token Alignment (SSTA) for cross-domain object detection, under the transformer framework. To the best of our knowledge, we make the first attempt to explore the intrinsic cross-attention property for improving the cross-domain generalization ability of detection transformers.
- Two new modules, *i.e.*, token alignment (SpaTA) and semantic-aware token alignment (SemTA), are developed respectively to align the token distributions according to their attentions to object queries and to achieve the category-level alignment.
- We conduct extensive experiments on several widely-used benchmarks (*e.g.*, FoggyCityscapes, Sim10K and BDD100K), and promising results demonstrate the ef-

fectiveness of our proposed method over existing state-of-the-art baselines.

## 2. Related Work

### 2.1. Object Detection

Object detection aims to recognize and localize one or multiple objects in a given image. Traditional object detection methods [2, 10, 23, 26, 27, 35] are based on convolutional neural networks (CNN) [12, 19, 33] and can be divided into two directions, one-stage, and two-stage methods. Two-stage methods [2, 10, 27] typically first generate some region proposals and then refine their classification and bounding boxes. In contrast to two-stage methods, one-stage methods [23, 26, 35] ignore the proposal generation stage and directly predict the category and coordinates of objects. Although these CNN-based detectors have achieved a remarkable breakthrough, they need many hand-designed components like removing duplicated detections by non-maximum suppression and anchor generation which explicitly encodes our prior knowledge about the task. Recently, Carion *et al.*, proposed DETR [3] that reaches an end-to-end object detection without anchor generation and any sophisticated post-process procedure. Many DETR-like models [24, 29, 39, 47] are proposed to further improve the performance of the DETR model in both convergence speed and accuracy. Among these works, one of the most representative works is Deformable DETR [47] which adopts deformable attention mechanism [6] into DETR and designs a multi-scale attention module so that it reduces the training time and improves detection performance significantly. Nevertheless, these methods suffer from severe performance degradation due to the domain discrepancy between the training and test domains. To address this problem, we present Spatial-aware and Semantic-aware Token Alignment (SSTA) to learn domain-invariant token representations. Following [38], we choose Deformable DETR [47] as the base detector for a fair comparison.

### 2.2. Cross-domain Object Detection

Cross-domain object detection (CDOD) aims to transfer the knowledge from the label-rich source domain to the label-scarce target domain by bridging the domain discrepancy between them. Previous works [1, 5, 8, 15, 17, 25, 31, 36, 46] can be roughly categorized into image translation, self-supervision, and adversarial training. Image translation methods [15, 17] adopt style transfer algorithms to enhance the image diversity so as to reduce the domain gap at the pixel level. Self-supervision approaches [1, 8, 25, 30] deploy the pseudo-labeling techniques to provide additional supervision signal for the target domain. Adversarial training methods [5, 31] align the feature distribution and eliminate the domain discrepancy to bridge the domain gap. Early

works align the features with diverse levels, *e.g.*, strong-weak alignment [31], global-instance level [5].

However, these methods are based on the Faster RCNN or FCOS, and the transferability of detection transformers remains a challenge. SFA [38] has developed a domain adaptive detection transformer to align domain query feature and token-wise feature and design an additional bipartite matching consistency loss to enhance the feature discriminability. Different from SFA [38], our SSTA takes advantage of the cross-attention map and leverages the spatial and semantic information to help the token distribution alignment. Our model follows the principle of giving minimal modification to the DETR model so that the inference has no extra overload. To the best of our knowledge, our method is the first domain adaptation work that takes advantage of the characteristics of cross-attention to improve the generalization ability of the DETR model.

### 3. Methodology

In the task of CDOD, we are given a source domain consisting of labeled images with object bounding boxes and their class labels and a target domain consisting of unlabeled images. Let us denote  $\mathcal{D}_s = \{(x_i^s, \mathbf{y}_i^s)\}_{i=1}^{N_s}$  drawn from distribution  $\mathcal{P}_s$  as the labeled source domain and  $\mathcal{D}_t = \{x_j^t\}_{j=1}^{N_t}$  drawn from distribution  $\mathcal{P}_t$  as the unlabeled target domain, where  $\mathcal{P}_s \neq \mathcal{P}_t$ . And  $\mathbf{y}_i^s = \{(\mathbf{b}_j^s, c_j^s)\}_{j=1}^m$ , where  $\mathbf{b}_j^s \in \mathbb{R}^4$  and  $c_j^s \in \{1, \dots, C\}$  are the bounding box and corresponding category for each object, and  $m$  is the total number of objects in an image  $x_i^s$ . Our goal is to learn an object detection model that performs well on the target domain.

In the following, we introduce the motivation of our proposed method in Sec. 3.1. And then, we first give the vanilla token alignment in Sec. 3.2 and describe the detailed design of spatial-aware token alignment (Sec. 3.3) and semantic-aware token alignment (Sec. 3.4). Lastly, we give the overall objective of the proposed method.

#### 3.1. Motivation

In this section, we give a brief preliminary to the DETR model. And then, we demonstrate the cross-domain challenges in DETR as well as our new solution.

**DEtection TRansformer (DETR):** DETR consists of CNN backbone, transformer encoder and transformer decoder. The image  $x \in \mathbb{R}^{3 \times H_0 \times W_0}$  are firstly fed into CNN backbone (*e.g.*, ResNet50 [12]) and to generate a lower-resolution feature map  $f \in \mathbb{R}^{C \times H \times W}$ , where  $C = 2048$ ,  $H = \frac{H_0}{32}$  and  $W = \frac{W_0}{32}$ . The encoder uses a  $1 \times 1$  convolution to reduce the channel  $C$  into a smaller dimension  $d$  and then collapse the spatial dimensions into one dimension, resulting token inputs  $z_c \in \mathbb{R}^{d \times N_k}$ , where  $N_k = WH$  is the length of sequence. The encoder layer adopts tokens  $z_c$  along with position embedding to make interac-

tion among tokens and outputs new tokens  $z_e \in \mathbb{R}^{d \times N_k}$  through standard architecture that consists of a multi-head self-attention and a feed forward network (FFN). The decoder comprises of multi-head self-attention and multi-head cross-attention mechanisms. Different with encoder, the decoder first deploys self-attention for  $N_q$  object queries and then uses cross-attention (*i.e.*, encoder-decoder attention) to aggregate features from the outputs of the encoder, resulting a sequence  $z_d \in \mathbb{R}^{d \times N_q}$ . Finally, the decoder will result  $N_q$  predictions. DETR utilizes Hungarian algorithm to find a bipartite matching between the sets of predictions and ground truth. The loss of DETR can be summarized as follows:

$$\mathcal{L}_{det} = \mathcal{L}_{cls} + \mathcal{L}_{reg}, \quad (1)$$

where the  $\mathcal{L}_{cls}$  is for classification and  $\mathcal{L}_{reg}$  is for bounding boxes regression.

DETR requires much longer training epochs (*i.e.*, 500) to converge than traditional detectors and has relatively low detection accuracy on small objects. Thus Deformable DETR [47] adopts efficient deformable attention module to replace the dense attention in DETR. The deformable attention mechanism can be naturally extended to aggregating multi-scale features, leading to fast convergence and high performance. Following [38], we choose Deformable DETR [47] as the base detector for a fair comparison. For more detail, please refer to [3, 47].

**Cross-domain Challenges in DETR:** To improve the generalization ability of detection transformer, a potential solution is to perform token alignment with adversarial learning. Recent work [38] also attempts to apply adversarial training strategies on tokens in transformers, but the improvements are still unsatisfactory. One of the main reasons is that the tokens in detection transformer are quite diverse. In detection transformers (*e.g.*, DETR), the tokens are passed through several multi-head self-attention layers to obtain new token embeddings for representing different spatial and semantic information. Then, object queries are introduced to probe useful tokens and leverage those tokens to predict the positions and categories of different objects. On the one hand, since some tokens are more useful while less for others, it is desirable to take the importance of tokens into consideration in the cross-domain detection transformer. On the other hand, the semantic information embedded in tokens is also helpful for aligning the token distributions of the corresponding category. This would ease the adversarial training when aligning the token distributions between domains.

To this end, we propose the spatial-aware token alignment (SpaTA) and the semantic-aware token alignment (SemTA) strategies to guide the token alignment across domains by leveraging the characteristics of cross-attention in detection transformer. As shown in Fig. 1, the proposed spatial-aware and the semantic-aware to-

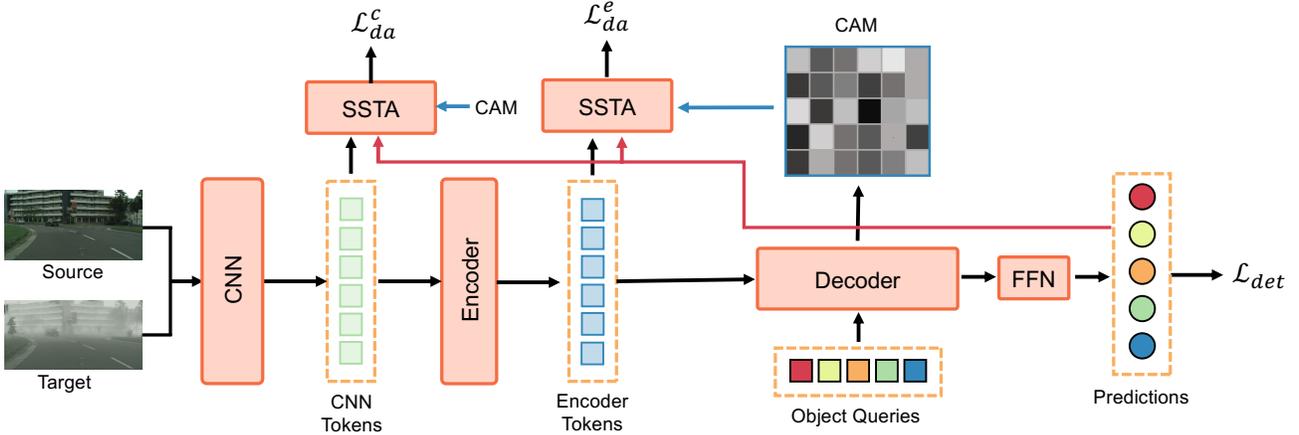


Figure 1. The overview of our method. We design a new Spatial-aware and Semantic-aware Token Alignment (SSTA) module to align CNN token and encoder token distribution across two domains. We take advantage of the characteristics of the cross-attention in the decoder and feed the cross-attention map (CAM) and the predictions of the detection head (FFN) to improve the token alignment. The details of the SSTA module are shown in Fig. 2.

ken alignment (SSTA) module adopts the cross-attention map (CAM) and predictions of the decoder to align the distributions of tokens from the CNN and encoder. The detail will be presented below.

### 3.2. Vanilla Token Alignment

Before we dive into the design of our SSTA module, we first introduce the vanilla token alignment. The existing adversarial methods [5, 31, 44] usually take a discriminator to reduce domain discrepancy via aligning feature distribution between domains. The discriminator tries to distinguish which domain the features come from, while the feature extractor aims to confuse features and deceive the discriminator in a minimax manner. It can be placed at a certain layer or multiple layers of feature extractor. In practice, a gradient reverse layer (GRL) [9] is used to connect the discriminator and feature extractor and flips the gradients when it flows through the feature extractor, leading to an end-to-end learning instead of sophisticated multi-stage iterative optimization like [11]. To bridge the domain gap, a naive solution is to simply align the distribution of tokens where the domain discriminator tries to recognize each token. Formally, the adversarial objective of vanilla token alignment can be defined as follows:

$$\mathcal{L}_{ta} = - \sum_{i=1}^{N_q} \{d \log(D(z_i)) + (1-d) \log(1-D(z_i))\}, \quad (2)$$

where  $N_q$  is the length of sequence,  $z_i$  is the  $i$ -th token representation and can be from CNN backbone or transformer encoder, and  $d$  is the domain label with  $d = 1$  for the source and  $d = 0$  for the target. When the above adversarial learning loss being optimized, the sign of gradi-

ent back-propagated from discriminator to feature extractor will be inverted by GRL, thus making the feature extractor learn domain-invariant representations.

The overall objective of vanilla token alignment can be formulated as:

$$\mathcal{L} = \mathcal{L}_{det} + \lambda \cdot (\mathcal{L}_{ta}^c + \mathcal{L}_{ta}^e), \quad (3)$$

where  $\lambda$  is the trade-off parameter, and  $\mathcal{L}_{ta}^c$  and  $\mathcal{L}_{ta}^e$  are the vanilla token alignment loss for the CNN and encoder tokens.

### 3.3. Spatial-aware Token Alignment

As the analysis in Sec. 3.1, object queries are introduced to probe useful tokens and leverage those tokens to predict the positions and categories of different objects. In other words, tokens contribute differently to the detection results. Simply aligning the token distribution between domains has unsatisfactory improvements, as tokens in detection transformer have different importances to object detection task. If we consider the tokens equally contributing to the adversarial training, we will overlook matching the distribution of critical tokens that may contain essential instances and global context for accurately predicting the positions and categories of different objects. Consequently, the efforts to reduce the domain gap will eventually meet difficulties, making the alignment less effective.

Motivated by this, we propose a spatial-aware token alignment (SpaTA) module to discover instance-related tokens and emphasize their alignment by assigning higher weights to these tokens for adversarial training according to their attention to the object queries. Formally, we can

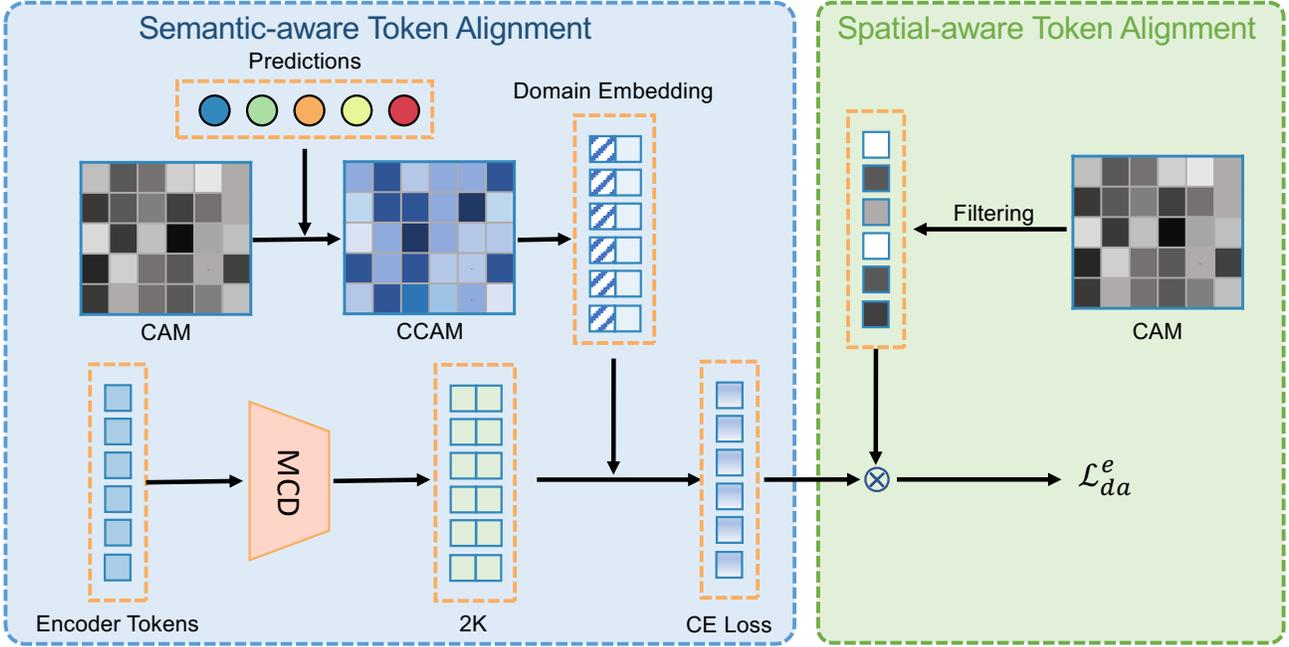


Figure 2. The overview of our Semantic-aware and Spatial-aware Token Alignment (SSTA) module. Take the SSTA module for encoder tokens as an example. The proposed SSTA module takes the tokens as the input and jointly utilizes Semantic-aware Token Alignment (SemTA) and Spatial-aware Token Alignment (SpaTA) to respectively align token distributions. SemTA affiliates the predictions of the detection head into the cross-attention map (CAM) and obtains a category cross-attention map (CCAM), which can be used to construct domain embedding to guide the learning of a multi-class discriminator (MCD) to achieve category-level token alignment. The SpaTA utilizes the CAM to give different weights to the adversarial learning of tokens according to their attention to object queries.

obtain the objective as follows:

$$\mathcal{L}_{spa} = \sum_{i=1}^{N_k} (1 + \mathcal{W}^i) \cdot \mathcal{L}_{ta}^i, \quad (4)$$

where  $\mathcal{W}^i$  is the weight for  $i$ -th token, intuitively, the more important the token should be assigned higher weights. As shown in the right part of Fig. 2, we utilize cross-attention map (CAM) as an alternative to providing the weights, as object queries probe features by giving different weights to tokens via the cross-attention mechanism.

However, the CAM cannot be directly obtained in deformable attention because of its special design. To this end, the key factor is determining how to obtain the CAM. We scatter and accumulate the cross-attention in the decoder from each object query to discrete token positions in the sequence. The deformable attention applies bilinear interpolation to obtain values from the surrounding position, as attention offset in deformable attention is fractional. Therefore, we also apply bilinear interpolation to obtain CAM. Specifically, let  $r$ ,  $\Delta r$ ,  $A$ , and  $v$  be one of the reference points of the decoder, corresponding offsets, attention weights, and values, respectively.

For the attention to each token, we can obtain CAM of

$i$ -th query as follows:

$$\mathcal{M}_i = \frac{1}{N_d} \sum_{l=1}^{N_d} \sum_{(A_l, r, \Delta r)} A_l \cdot \mathcal{B}(t, r + \Delta r), \quad (5)$$

where  $N_d$  is the number of decoder layer,  $\mathcal{B}(\cdot, \cdot)$  is the bilinear interpolation operation, and  $t$  enumerates all integral spatial locations of tokens. We provide more details in our Supplementary materials. After obtaining the CAM, we filter out some attentions that are less than a given threshold.

In summary, the important weight for tokens can be obtained via:

$$\mathcal{W} = \mathcal{M} \odot \mathbb{1}(\mathcal{M} \geq \tau(\mathcal{M})), \quad (6)$$

where  $\mathcal{M}$  is the average of CAM for all the queries and  $\tau(\mathcal{M}) = \text{mean}(\mathcal{M})$  is an adaptive threshold for each sample  $x$ .

### 3.4. Semantic-aware Token Alignment

Although we have discovered the critical tokens to emphasize their alignment and avoid the influence of noise tokens, the model still has the risk of misalignment during the adaptation process [36, 37]. The semantic information of tokens is helpful for aligning the token distributions of the corresponding category, so that the model can avoid the

Table 1. Average precisions (%) of different methods on Cityscapes→FoggyCityscapes.

| Method                        | Detector        | person      | rider       | car         | truck       | bus         | train       | mcycle      | bicycle     | mAP         |
|-------------------------------|-----------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Faster RCNN [27] (Source)     | Faster RCNN     | 26.9        | 38.2        | 35.6        | 18.3        | 32.4        | 9.6         | 25.8        | 28.6        | 26.9        |
| DA-Faster [5]                 |                 | 25.0        | 31.0        | 40.5        | 22.1        | 35.3        | 20.2        | 20.0        | 27.1        | 27.6        |
| SWDA [31]                     |                 | 29.9        | 42.3        | 43.5        | 24.5        | 36.2        | 32.6        | 30.0        | 35.3        | 34.3        |
| CFDA [45]                     |                 | 43.2        | 37.4        | 52.1        | 34.7        | 34.0        | 46.9        | 29.9        | 30.8        | 38.6        |
| UMT [8]                       |                 | 33.0        | 46.7        | 48.6        | 34.1        | 56.5        | 46.8        | 30.4        | 37.4        | 41.7        |
| MeGA [36]                     |                 | 37.7        | 49.0        | 52.4        | 25.4        | 49.2        | 46.9        | 34.5        | 39.0        | 41.8        |
| ICCR-VDD [40]                 |                 | 33.4        | 44.0        | 51.7        | <b>33.9</b> | <b>52.0</b> | 34.7        | 34.2        | 36.8        | 40.0        |
| ViSGA [28]                    |                 | 38.8        | 45.9        | 57.2        | 29.9        | 50.2        | <b>51.9</b> | 31.9        | 40.9        | 43.3        |
| DIDN [22]                     |                 | 38.3        | 44.4        | 51.8        | 28.7        | 53.3        | 34.7        | 32.4        | 40.4        | 40.5        |
| FCOS [35] (Source)            | FCOS            | 36.9        | 36.3        | 44.1        | 18.6        | 29.3        | 8.4         | 20.3        | 31.9        | 28.2        |
| EPM [14]                      |                 | 41.9        | 38.7        | 56.7        | 22.6        | 41.5        | 26.8        | 24.6        | 35.5        | 36.0        |
| SCAN [21]                     |                 | 41.7        | 43.9        | 57.3        | 28.7        | 48.6        | 48.7        | 31.0        | 37.3        | 42.1        |
| KTNet [34]                    |                 | 46.4        | 43.2        | 60.6        | 25.8        | 41.2        | 40.4        | 30.7        | 38.8        | 40.9        |
| SSAL [25]                     |                 | 45.1        | 47.4        | 59.4        | 24.5        | 50.0        | 25.7        | 26.0        | 38.7        | 39.6        |
| Deformable DETR [47] (Source) | Deformable DETR | 38.6        | 40.6        | 45.8        | 11.6        | 28.9        | 1.7         | 18.9        | 39.1        | 28.1        |
| SFA [38]                      |                 | 46.5        | 48.6        | 62.6        | 25.1        | 46.2        | 29.4        | 28.3        | 44.0        | 41.3        |
| SSTA (Ours)                   |                 | <b>50.5</b> | <b>53.0</b> | <b>67.2</b> | 24.7        | 47.7        | 33.0        | <b>36.7</b> | <b>46.6</b> | <b>44.9</b> |

class misalignment. For example, the “car” and the “truck” instances are forced to be very close in the feature space, deteriorating the model discriminant ability. Therefore, we propose to utilize a multi-class discriminator [37] (MCD) to capture the category information during adversarial training so that it realizes category-level token alignment. The multi-class discriminator contains not only domain information but also category relationship. Concretely, we remold the single-class discriminator to a multi-classes discriminator that outputs  $2K$  logits, where  $K = C + 1$ ,  $K$  for the source domain, and others for the target domain. The domain embedding  $\mathbf{d} \in \mathbb{R}^{2K \times 1}$  of the source and target are  $[0; \mathbf{s}]$  and  $[\mathbf{s}; 0]$ , respectively, where  $\mathbf{s} \in \mathbb{R}^{K \times 1}$  is the domain knowledge and  $\mathbf{0} \in \mathbb{R}^{K \times 1}$  is all-zero vector. The objective of semantic-aware token alignment can be written as follows:

$$\mathcal{L}_{sem}^i = - \sum_{k=1}^{2K} \mathbf{d}_k \cdot \log(\hat{D}(z_i)_k), \quad (7)$$

where  $\hat{D}$  is the multi-class domain discriminator. The key factor is determining how to obtain the domain knowledge  $\mathbf{s}$  to build domain embedding for these tokens. As illustrated in the left part of Fig. 2, we also utilize CAM to extract domain knowledge by injecting the category information into it. In particular, we affiliate the predictions of the detection head into the CAM and obtain a category cross-attention map (CCAM) which can be formally defined as follows:

$$\tilde{\mathcal{M}}_k = \frac{1}{N_q^k} \sum_i^{N_q} \mathbb{1}(\hat{y}_i = k) \cdot \mathcal{M}_i, \quad (8)$$

where  $\tilde{\mathcal{M}}^k \in \mathbb{R}^{N_k}$  refers to CCAM  $\tilde{\mathcal{M}} \in \mathbb{R}^{N_k \times K}$  for category  $k$ ,  $N_q^k$  is the number of queries that belong to category  $k$ . The  $\hat{y}_i$  is the category prediction from detection head

for  $i$ -th query and  $\mathbb{1}(\cdot)$  is the indicator function where if  $\cdot$  is true then equals 1, otherwise 0. The  $\mathbf{s}$  can be obtained after apply softmax function to the CCAM  $\tilde{\mathcal{M}}$ . Finally, we can obtain our domain adaptation loss by replacing the  $\mathcal{L}_{ta}^i$  by the semantic-aware token alignment in Eq. (4):

$$\mathcal{L}_{da} = - \sum_{i=1}^{N_k} (1 + \mathcal{W}^i) \cdot \mathcal{L}_{sem}^i, \quad (9)$$

### 3.5. Overall Objective

In summary, the overall objective includes the detection loss of Deformable DETR [47] on the source domain and domain adaptation loss for the CNN and encoder tokens. In summary, the overall objective can be defined as:

$$\mathcal{L} = \mathcal{L}_{det} + \lambda \cdot (\mathcal{L}_{da}^c + \mathcal{L}_{da}^e), \quad (10)$$

where  $\lambda$  is the trade-off parameter,  $\mathcal{L}_{da}^c$  and  $\mathcal{L}_{da}^e$  are the domain adaptation loss for the CNN and encoder tokens, respectively.

## 4. Experiments

Following [38], we train the model with labeled source data and unlabeled target data and test on the target data. We conduct extensive experiments on three CDOD scenarios. The detection results are evaluated with mean Average Precision (mAP) under the threshold of 0.5.

### 4.1. Experimental Setup

**Datasets:** Cityscapes dataset was collected for the scenes understanding of road and street. It comprises 2,975 and 500 images for training and validation, respectively. It contains 8 categories: *person*, *rider*, *car*, *truck*, *bus*, *train*, *motorbike*, and *bicycle*. FoggyCityscapes [32] dataset is the

foggy version of Cityscapes and generated using the depth information provided by Cityscapes. Thus, it shares the common annotations with Cityscapes. It contains three levels for foggy weather, including 0.01, 0.15, and 0.02. In experiments, we choose the worst foggy weather (*i.e.*, 0.02). Sim10K [16] dataset is a synthetic dataset rendered by the gaming engine Grand Theft Auto V (GTAV). This dataset contains 10,000 images with 58,701 bounding boxes with the category of “car”. BDD100K [42] dataset is a large-scale autonomous driving and contains 100k images with six types of weather, six different scenes, and three categories for the time of day. We extract the subset of daytime, resulting in 36,728 training and 5,258 validation images.

Following existing works [5,44], we evaluate our method on three benchmark settings:

- **Weather Adaptation:** We take Cityscapes as the source domain and FoggyCityscape as the target domain, and the model is trained on the train set of Cityscapes and FoggyCityscape and evaluated on the validation split of FoggyCityscapes.
- **Syn2Real:** We explore the adaptation of Sim10K to Cityscapes, we train the model using all the images of Sim10K and the train split of Cityscapes, and report mAP on the validation split of Cityscapes with “car” category.
- **Scene Adaptation:** We use Cityscapes as the source domain dataset and BDD100K containing distinct scenes as a large unlabeled target domain dataset. We evaluate the model on the validation set of BDD100K.

**Implementation Details:** Following the default setting in SFA [38], we adopt Deformable DETR [47] as base detector, which contains ResNet-50 [12] backbone pre-trained on ImageNet [7], six transformer encoders, six transformer decoders and multiple prediction heads. We adopt Adam [18] optimizer to update parameters. For Cityscapes to FoggyCityscapes, we first train the model with a learning rate  $2 \times 10^{-4}$  for 40 epochs, then decay the learning rate to  $2 \times 10^{-5}$  for 10 more epochs. And the trade-off parameter  $\lambda$  is set to 1.0. For Sim10K to Cityscapes and Cityscapes to BDD100K, we set the initial learning rate and the trade-off parameter  $\lambda$  to  $5 \times 10^{-5}$  and 0.01 respectively. We pre-train models on source data to obtain reliable CAM. All the experiments are conducted using four V100 GPUs with batch size of 16, *i.e.*, each GPU contains 2 source images and 2 target images. We implement our method with the PyTorch deep learning framework. The source code of our method will be released soon.

## 4.2. Results

We conduct extensive experiments and validate the effectiveness of our method by comparing various state-of-the-

Table 2. Average precisions (%) of different methods on SIM10K→Cityscapes.

| Method                       | Detector        | AP on Car   |
|------------------------------|-----------------|-------------|
| DA-Faster [5]                | Faster RCNN     | 39.0        |
| SCDA [46]                    |                 | 43.0        |
| SWDA [31]                    |                 | 40.1        |
| MAF [13]                     |                 | 41.1        |
| HTCN [4]                     |                 | 42.5        |
| SAP [20]                     |                 | 44.9        |
| UMT [8]                      |                 | 43.1        |
| ViSGA [28]                   |                 | 49.3        |
| EPM [14]                     | FCOS            | 49.0        |
| KTNet [34]                   |                 | 50.7        |
| SCAN [21]                    |                 | 52.6        |
| SSAL [25]                    |                 | 51.8        |
| Deformable DETR [47](Source) | Deformable DETR | 47.4        |
| SFA [38]                     |                 | 52.6        |
| SSTA (Ours)                  |                 | <b>57.7</b> |

art CDOD methods, mainly including three kinds of methods: 1) two-stage detector Faster RCNN 2) one-stage detector FCOS, 3) Deformable DETR. For all the methods, we report the results from the original papers. To validate the effectiveness of our proposed method, we also report the results of the Source model where the model is only trained on the source domain and directly evaluated on the target domain.

### Weather Adaptation (Cityscapes → FoggyCityscapes):

We show the adaptation results in Table 1. We can observe that our proposed method outperforms the previous state-of-the-art approaches by a large margin, reaching 44.9% in terms of mAP. Specifically, Deformable DETR (Source) achieves 28.1% in terms of mAP, which shows that Deformable DETR has a decent generalization but still suffers from the distribution discrepancy across domains. Both SFA [38] and our SSTA improve the Source baseline. However, our SSTA improved by 3.6% in terms of mAP compared with the counterpart SFA [38]. This demonstrates that our method by leveraging intrinsic cross-attention to conduct spatial-aware and semantic-aware token alignment can effectively improve the generalization ability of detection transformer on the target domain.

### Syn2Real (Sim10K → Cityscapes):

The results of synthetic-to-real adaptation are presented in Table 2. Our proposed method SSTA reaches the highest mAP (57.7%) that exceeds all compared state-of-the-art methods, including the two-stage, one-stage, and DETR works, by a large margin, that is 5.1% in terms of mAP over best-performing one-stage detector SCAN [21] and DETR counterpart SFA [38]. These results verify the effectiveness of our SSTA.

### Scene Adaptation (Cityscapes → BDD100K):

The quantitative results are shown in Table 3. According to Table 3, our method SSTA achieves the new state-of-the-art results of 29.5% in terms of mAP, which surpasses the previous

Table 3. Average precisions (%) of different methods on Cityscapes  $\rightarrow$  BDD100K.

| Methods                       | Detector        | person      | rider       | car         | truck       | bus         | mcycle      | bicycle     | mAP         |
|-------------------------------|-----------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Faster R-CNN (Source)         | Faster RCNN     | 28.8        | 25.4        | 44.1        | 17.9        | 16.1        | 13.9        | 22.4        | 24.1        |
| DA-Faster [5]                 |                 | 28.9        | 27.4        | 44.2        | 19.1        | 18.0        | 14.2        | 22.4        | 24.9        |
| SWDA [31]                     |                 | 29.5        | 29.9        | 44.8        | 20.2        | 20.7        | 15.2        | 23.1        | 26.2        |
| SCDA [46]                     |                 | 29.3        | 29.2        | 44.4        | 20.3        | 19.6        | 14.8        | 23.2        | 25.8        |
| ECR [41]                      |                 | 32.8        | 29.3        | 45.8        | <b>22.7</b> | 20.6        | 14.9        | 25.5        | 27.4        |
| FCOS [35] (Source)            | FCOS            | 38.6        | 24.8        | 54.5        | 17.2        | 16.3        | 15.0        | 18.3        | 26.4        |
| EPM [14]                      |                 | 39.6        | 26.8        | 55.8        | 18.8        | 19.1        | 14.5        | 20.1        | 27.8        |
| Deformable DETR [47] (Source) | Deformable DETR | 38.4        | 27.1        | 56.1        | 14.6        | 12.3        | <b>16.3</b> | 20.7        | 26.5        |
| SFA [38]                      |                 | <b>40.2</b> | 27.6        | 57.5        | 19.1        | <b>23.4</b> | 15.4        | 19.2        | 28.9        |
| SSTA (Ours)                   |                 | 39.4        | <b>31.9</b> | <b>59.4</b> | 16.3        | 17.7        | 15.3        | <b>26.2</b> | <b>29.5</b> |

Table 4. Ablation studies of SSTA on Cityscapes  $\rightarrow$  FoggyCityscapes. TA indicates token alignment.

| Method                        | TA | SpaTA | SemTA | mAP (%) | $\Delta$ |
|-------------------------------|----|-------|-------|---------|----------|
| Deformable DETR [47] (Source) | -  | -     | -     | 28.1    | -        |
| Proposed                      | ✓  |       |       | 41.3    | 13.2↑    |
|                               |    | ✓     |       | 42.5    | 14.4↑    |
|                               |    |       | ✓     | 43.9    | 15.8↑    |
| SSTA                          |    | ✓     | ✓     | 44.9    | 16.8↑    |

Table 5. Average precisions (%) w.r.t. different values of  $\lambda$  on Cityscapes  $\rightarrow$  FoggyCityscapes.

| $\lambda$ | 0.0  | 0.1  | 0.5  | 1.0  | 1.5  | 2.0  |
|-----------|------|------|------|------|------|------|
| SSTA      | 28.1 | 42.1 | 44.3 | 44.9 | 44.9 | 44.6 |

works. This again demonstrates the generalization of our method.

**Ablation Studies:** To further verify the effectiveness of our proposed method, we have conducted detailed ablation studies by isolating each component of our SSTA. The experimental results are shown in Table 4. In particular, our SpaTA significantly boosts the baseline, leading to 14.4% mAP improvements compared with Source model (28.1%). This implies that the CAM can provide sufficient information to discover critical tokens, and emphasizing their contributions to distribution alignment will significantly improve the generalization ability of Deformable DETR. Moreover, SemTA also improves the accuracy of Deformable DETR, achieving 43.9% in terms of mAP. These improvements mainly come from our SemTA considering category information during token alignment and thus avoiding class misalignment. By synergizing SpaTA and SemTA together, we obtain 44.9% in terms of mAP, which shows their complementary to each other.

**Parameter Analysis:** We also investigate the influence of the trade-off parameter  $\lambda$  which is used to balance the weight between the source detection loss  $\mathcal{L}_{det}$  and the domain adaptation loss. Table 5 summarizes the experimental results on Cityscapes  $\rightarrow$  FoggyCityscapes. Note that when  $\lambda = 0$ , the method degenerates to the Source model. According to Table 5, we can conclude that our proposed

SSTA consistently improve the generalization ability of Deformable DETR in a wide range of  $\lambda$ , and  $\lambda = 1.0$  and  $\lambda = 1.5$  are the bests among them.

## 5. Conclusion

Detection transformers (*e.g.*, DETR) have shown promising results for object detection, when training and test images come from the same domain. However, they usually do not work well for cross-domain problems. In this work, we tackle cross-domain object detection by proposing a novel approach named Semantic-aware and Spatial-aware Token Alignment (SSTA) under the transformer framework. In SSTA, two new modules *i.e.*, spatial-aware token alignment (SpaTA) and semantic-aware token alignment (SemTA), are developed to guide the token alignment across domains. Promising results on benchmark datasets demonstrate the effectiveness of our method.

**Limitation:** Although our method outperforms existing cross-domain object detection works, it still faces challenges in detecting objects of rare classes. For example, the “truck” and “train” classes in Table 1 have relatively low AP compared with other classes (*e.g.*, “car”). We conjecture that this is caused by the label shift between the source and target domains. In the future, we will study how to improve the detection performance of our SSTA for these classes.

## References

- [1] Qi Cai, Yingwei Pan, Chong-Wah Ngo, Xinmei Tian, Lingyu Duan, and Ting Yao. Exploring object relation in mean teacher for cross-domain detection. In *CVPR*, pages 11457–11466, 2019. [2](#)
- [2] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. [1](#), [2](#)
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229. Springer, 2020. [1](#), [2](#), [3](#)
- [4] Chaoqi Chen, Zebiao Zheng, Xinghao Ding, Yue Huang, and Qi Dou. Harmonizing transferability and discriminability for adapting object detectors. In *CVPR*, pages 8869–8878, 2020. [7](#)
- [5] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *CVPR*, pages 3339–3348, 2018. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#)
- [6] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017. [2](#)
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. [7](#)
- [8] Jinhong Deng, Wen Li, Yuhua Chen, and Lixin Duan. Unbiased mean teacher for cross-domain object detection. In *CVPR*, pages 4091–4101, 2021. [1](#), [2](#), [6](#), [7](#)
- [9] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016. [4](#)
- [10] Ross Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015. [1](#), [2](#)
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. [4](#)
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. [1](#), [2](#), [3](#), [7](#)
- [13] Zhenwei He and Lei Zhang. Multi-adversarial faster-rcnn for unrestricted object detection. In *ICCV*, pages 6668–6677, 2019. [7](#)
- [14] Cheng-Chun Hsu, Yi-Hsuan Tsai, Yen-Yu Lin, and Ming-Hsuan Yang. Every pixel matters: Center-aware feature alignment for domain adaptive object detector. In *ECCV*, pages 733–748. Springer, 2020. [6](#), [7](#), [8](#)
- [15] Han-Kai Hsu, Chun-Han Yao, Yi-Hsuan Tsai, Wei-Chih Hung, Hung-Yu Tseng, Maneesh Singh, and Ming-Hsuan Yang. Progressive domain adaptation for object detection. In *WACV*, pages 749–757, 2020. [2](#)
- [16] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? *arXiv preprint arXiv:1610.01983*, 2016. [7](#)
- [17] Taekyung Kim, Minki Jeong, Seunghyeon Kim, Seokeon Choi, and Changick Kim. Diversify and match: A domain adaptive representation learning paradigm for object detection. In *CVPR*, pages 12456–12465, 2019. [2](#)
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. [7](#)
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 25, 2012. [1](#), [2](#)
- [20] Congcong Li, Dawei Du, Libo Zhang, Longyin Wen, Tiejian Luo, Yanjun Wu, and Pengfei Zhu. Spatial attention pyramid network for unsupervised domain adaptation. In *ECCV*, pages 481–497. Springer, 2020. [7](#)
- [21] Wuyang Li, Xinyu Liu, Xiwen Yao, and Yixuan Yuan. Scan: Cross domain object detection with semantic conditioned adaptation. In *AAAI*, 2022. [6](#), [7](#)
- [22] Chuang Lin, Zehuan Yuan, Sicheng Zhao, Peize Sun, Changhu Wang, and Jianfei Cai. Domain-invariant disentangled network for generalizable object detection. In *ICCV*, pages 8771–8780, 2021. [6](#)
- [23] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. [1](#), [2](#)
- [24] Depu Meng, Xiaokang Chen, Zejjia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *ICCV*, pages 3651–3660, 2021. [1](#), [2](#)
- [25] Muhammad Akhtar Munir, Muhammad Haris Khan, M Sarfraz, and Mohsen Ali. Ssal: Synergizing between self-training and adversarial learning for domain adaptive object detection. *NeurIPS*, 34, 2021. [1](#), [2](#), [6](#), [7](#)
- [26] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, pages 779–788, 2016. [1](#), [2](#)
- [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS*, 28, 2015. [1](#), [2](#), [6](#)
- [28] Farzaneh Rezaeianaran, Rakshith Shetty, Rahaf Aljundi, Daniel Olmeda Reino, Shanshan Zhang, and Bernt Schiele. Seeking similarities over differences: Similarity-based domain alignment for adaptive object detection. In *ICCV*, pages 9204–9213, 2021. [6](#), [7](#)
- [29] Byungseok Roh, JaeWoong Shin, Wuhyun Shin, and Sae-hoon Kim. Sparse detr: Efficient end-to-end object detection with learnable sparsity. *ICLR*, 2022. [1](#), [2](#)
- [30] Aruni RoyChowdhury, Prithvijit Chakrabarty, Ashish Singh, SouYoung Jin, Huaizu Jiang, Liangliang Cao, and Erik Learned-Miller. Automatic adaptation of object detectors to new domains using self-training. In *CVPR*, pages 780–790, 2019. [2](#)
- [31] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Strong-weak distribution alignment for adaptive object detection. In *CVPR*, pages 6956–6965, 2019. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#)

- [32] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 126(9):973–992, 2018. 6
- [33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [34] Kun Tian, Chenghao Zhang, Ying Wang, Shiming Xiang, and Chunhong Pan. Knowledge mining and transferring for domain adaptive object detection. In *ICCV*, pages 9133–9142, 2021. 6, 7
- [35] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *ICCV*, pages 9627–9636, 2019. 1, 2, 6, 8
- [36] Vibashan VS, Vikram Gupta, Poojan Oza, Vishwanath A Sindagi, and Vishal M Patel. Mega-cda: Memory guided attention for category-aware unsupervised domain adaptive object detection. In *CVPR*, pages 4516–4526, 2021. 2, 5, 6
- [37] Haoran Wang, Tong Shen, Wei Zhang, Ling-Yu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. In *ECCV*, pages 642–659. Springer, 2020. 5, 6
- [38] Wen Wang, Yang Cao, Jing Zhang, Fengxiang He, Zheng-Jun Zha, Yonggang Wen, and Dacheng Tao. Exploring sequence feature alignment for domain adaptive detection transformers. In *ACM MM*, pages 1730–1738, 2021. 1, 2, 3, 6, 7, 8
- [39] Yingming Wang, Xiangyu Zhang, Tong Yang, and Jian Sun. Anchor detr: Query design for transformer-based detector. *AAAI*, 2022. 1, 2
- [40] Aming Wu, Rui Liu, Yahong Han, Linchao Zhu, and Yi Yang. Vector-decomposed disentanglement for domain-invariant object detection. In *ICCV*, pages 9342–9351, 2021. 6
- [41] Chang-Dong Xu, Xing-Ran Zhao, Xin Jin, and Xiu-Shen Wei. Exploring categorical regularization for domain adaptive object detection. In *CVPR*, pages 11724–11733, 2020. 8
- [42] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *CVPR*, June 2020. 7
- [43] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 1
- [44] Jingyi Zhang, Jiaying Huang, Zhipeng Luo, Gongjie Zhang, and Shijian Lu. Da-detr: Domain adaptive detection transformer by hybrid attention. *arXiv preprint arXiv:2103.17084*, 2021. 4, 7
- [45] Yangtao Zheng, Di Huang, Songtao Liu, and Yunhong Wang. Cross-domain object detection through coarse-to-fine feature adaptation. In *CVPR*, pages 13766–13775, 2020. 6
- [46] Xinge Zhu, Jiangmiao Pang, Ceyuan Yang, Jianping Shi, and Dahua Lin. Adapting object detectors via selective cross-domain alignment. In *CVPR*, pages 687–696, 2019. 1, 2, 7, 8
- [47] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *ICLR*, 2020. 1, 2, 3, 6, 7, 8