

# Unsupervised Domain Adaptation of Object Detectors: A Survey

Poojan Oza, *Member, IEEE*, Vishwanath A. Sindagi, *Member, IEEE*, Vibashan VS, *Member, IEEE*, and Vishal M. Patel, *Senior Member, IEEE*

**Abstract**—Recent advances in deep learning have led to the development of accurate and efficient models for various computer vision applications such as classification, segmentation, and detection. However, learning highly accurate models relies on the availability of large-scale annotated datasets. Due to this, model performance drops drastically when evaluated on label-scarce datasets having visually distinct images, termed as domain adaptation problem. There are a plethora of works to adapt classification and segmentation models to label-scarce target dataset through unsupervised domain adaptation. Considering that detection is a fundamental task in computer vision, many recent works have focused on developing novel domain adaptive detection techniques. Here, we describe in detail the domain adaptation problem for detection and present an extensive survey of the various methods. Furthermore, we highlight strategies proposed and the associated shortcomings. Subsequently, we identify multiple aspects of the problem that are most promising for future research. We believe that this survey shall be valuable to the pattern recognition experts working in the fields of computer vision, biometrics, medical imaging, and autonomous navigation by introducing them to the problem, and familiarizing them with the current status of the progress while providing promising directions for future research.

**Index Terms**—Object detection, domain adaptation, unsupervised learning, transfer learning, deep learning.

## 1 INTRODUCTION

THE success of deep learning has been greatly beneficial for various fields such as natural language processing [1], [2], [3], robotics [4], [5], [6], computer vision [7], [8], [9], etc. This is especially evident in the case of computer vision, where majority of the progress can be largely attributed to the advancements in deep convolutional neural networks (DCNN) [7]. Owing to their learning capacity, DCNN models have achieved state-of-the-art performance in many vision tasks such as object classification ([9], [10], [11]), semantic segmentation ([12], [13], [14]), and object detection ([15], [16], [17]). This has led to DCNN's increased popularity in several real world applications as compared to the classical computer vision techniques. Specifically, deep learning based object detection has become an integral part of many real-world applications ranging from video security/surveillance, augmented reality, autonomous navigation, human computer interface, self-checkout convenience stores. Major advancements like Faster-RCNN [15], You Only Look Once (YOLO) [16] and Single Shot Multi-box Detector (SSD) [17] have resulted in significant improvements of detection performance and speed.

It is important to note that most DCNN models need to be trained in a supervised fashion, which has been made possible due to the availability of large datasets having thousands of images annotated with ground-truth labels [18], [19], [20]. However, one of the major drawbacks is the poor generalization capability of DCNN models to visually distinct images compared to the training images.

For instance, a detection model trained with a dataset collected in Rome may not necessarily perform well on images from Tokyo due to the changes in the appearance of scenes/objects and/or weather between them, as illustrated in Fig. 1(c). A similar example is shown for cases such as sunny to foggy weather (Fig. 1(a)), visible to thermal (Fig. 1(b)), and synthetic to real-world (Fig. 1(d)). Fig. 2 shows quantitatively the performance drop of different deep learning based object detectors that are trained on one particular dataset, when evaluated on different datasets. This problem where models, trained on one particular dataset (also known as source dataset), do not generalize well to a dataset that has a different distribution (also known as target dataset) is commonly referred to as *domain shift* or *distribution shift* in the literature [21], [22], [23].

A straightforward approach to solving this distributional shift problem is annotating the target dataset images with ground-truth detection labels. However, this might prove to be infeasible considering that the labor cost of the annotation process is prohibitively expensive for all visually distinct conditions. To circumvent this issue, many methods rely on the principles of unsupervised domain adaptation [21], [22] which involves training the DCNN model with both label-rich source dataset and label-scarce target dataset having visually distinct appearance. Techniques [23], [24], [25] for domain adaptive training with source and unlabeled target datasets have demonstrated improved generalization capabilities, resulting in improved performance on the visually distinct target domain. The unsupervised domain adaptation has been extensively studied for the task of classification [26], [27], [23], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], and semantic segmentation [24], [43], [44], [45], [46], [47], [48], [49], [50], [49], [51], [52], [53], [54], [46], [55], [56].

• Poojan Oza, Vishwanath Sindagi, Vibashan VS, and Vishal M. Patel are with the department of Electrical and Computer Engineering (ECE) at Johns Hopkins University, Baltimore, MD - 21218. E-mail: {poza2, vishwanathsindagi, vvishnu2, vpatel36}@jhu.edu.

Manuscript received Month XX, XXXX; revised Month XX, XXXX.

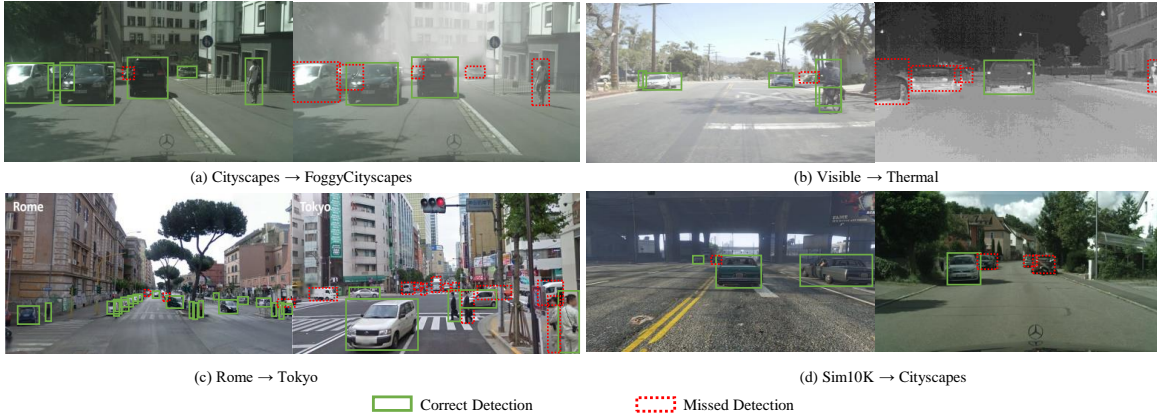


Fig. 1. Visualization of object detection results. Left: Source trained model on source domain, Right: Source trained model on the target domain. (a) The model, trained on Cityscapes dataset, performs well on the Cityscapes images. However, it fails to detect multiple cars when evaluated on FoggyCityscapes images, which has a domain shift due to fog. (b) Under visible to thermal domain shift, the model fails to detect person and cars in the thermal domain. (c) A detection model trained in Rome, when evaluated on another city such as Tokyo undergoes drastic performance reduction due to differences in scene appearances, weather, objects, etc. (d) In the case of Sim10K to Cityscapes, multiple cars are missed in the Cityscapes domain as the model was trained on images captured from the simulated virtual world. These examples show that the detection models generalize poorly under the domain shift/dataset distribution shift.

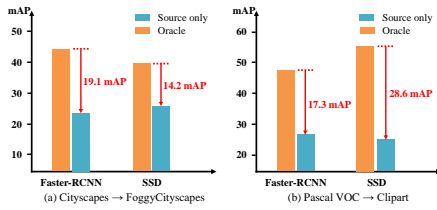


Fig. 2. Illustration of detector performance; In (a), the model is trained on Cityscapes and evaluated on FoggyCityscapes and in (b), the model is trained on Pascal VOC and evaluated on Clipart. We can observe a significant drop in the performance of the detector when there is a distribution shift in the training and test data.

However, unlike classification (image-level prediction) and semantic segmentation (pixel-level prediction), the object detection task involves bounding-box localization and the bounding-box-level category prediction task. This poses unique challenges while addressing unsupervised domain adaptation for object detection models. This has sparked an interest in addressing unsupervised domain adaptation for object detection with many novel approaches proposed very recently [25], [57], [58], [59], [60], [61], [62], [63], [64], [65], [66], [67], [68], [69], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80]. A timeline of some of the key papers proposed in the recent past is shown in Fig. 3.

While there exist multiple survey papers that extensively review domain adaptation techniques both classification [22], [81], [82] and semantic segmentation [83], [84], to the best of our knowledge, there is no comprehensive survey of unsupervised domain adaptation of object detectors. Although Li *et al.* [85] attempt to review some of the domain adaptive object detection literature, their discussions are limited and lack a comprehensive comparison of existing methods. This motivates us to present a comprehensive literature analysis of all the domain adaptive object detection methods proposed in the past few years, along with detailed discussions and comprehensive comparison. The major contributions of this work are summarized as follows:

- 1) As opposed to the existing survey [85], we provide a more detailed and thorough discussion on all the ex-

isting works (to the best of our knowledge) on domain adaptive object detection. We also define a taxonomy of the various works in literature. Furthermore, we discuss the preliminaries of the relevant topics like object detection and domain adaptation in an attempt to make the paper self-sufficient and useful to the readers who are not familiar with these concepts.

- 2) We present a comprehensive comparison of the existing methods on all publicly available datasets used in the literature for unsupervised domain adaptive object detection with thorough discussions on the methods and detailed comparison of the performances along with their respective experimental settings.
- 3) Finally, we identify potential research directions that might prove beneficial for the researchers working in the area in order to further advance the state-of-the-art.

## 2 PRELIMINARIES

In this section, we provide an introduction to two of the important aspects related to domain adaptive object detection, i.e., object detection and unsupervised domain adaptation. We formally set up the problem and notations that are used throughout the paper.

### 2.1 Object detection

Over the years, deep convolutional neural network based object detectors have demonstrated exceptional improvements in performance on a variety of datasets and have become an integral part of various computer vision applications. There are a variety of surveys [86], [87], [88] on the topic covering wide range of techniques proposed over the past decade for object detection. The most popular frameworks for object detection are Faster-RCNN [15], You Only Look Once (YOLO) [16] and Single Shot Multi-box Detector (SSD) [17]. The majority of domain adaptive object detection works are based on the Faster-RCNN and a few others use SSD. Other recent frameworks include, Fully Convolutional One Stage (FCOS) Object Detection [89] and

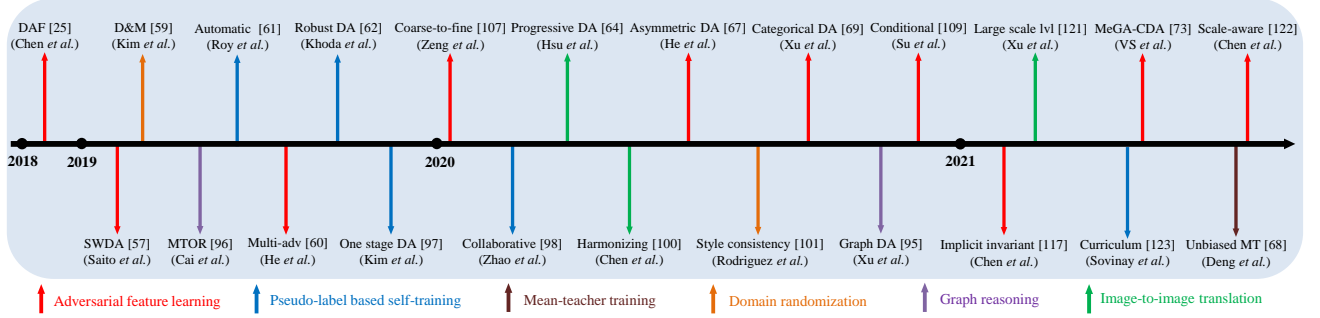


Fig. 3. A timeline of the key papers related to the domain adaptation of object detectors published over the recent years. These methods are broadly categorized into six classes: Adversarial feature learning, Pseudo-label based self-training, Graph-reasoning, Image-to-image translation, Domain randomization and Mean-teacher training.

DEtection TRansformer (DETR) [90]. However, these frameworks have been only scarcely used for the domain adaptive object detectors. In what follows, we briefly describe the widely used detection frameworks in the domain adaptive detection literature, i.e., Faster-RCNN and SSD.

### 2.1.1 Faster-RCNN

The Faster-RCNN framework, proposed by Ren *et al.* [15], follows a two-stage object detection approach and it consists of three major components: 1) a backbone CNN, 2) a Region Proposal Network (RPN), and 3) a Region-of-Interest (RoI) based classifier (RCN). Fig. 4(a) shows an overview of the Faster-RCNN architecture. Consider a dataset,  $\mathcal{D} = \{X^i, Y^i\}_{i=1}^N$ , having  $N$  images, with each image  $X^i$  with ground-truth annotation  $Y^i$ . Here, the ground-truth annotation  $Y^i$  denotes both bounding boxes and respective object categories in the corresponding image  $X^i$ . As shown in Fig. 4(a), an input image ( $X^i$ ) is forwarded through the backbone network resulting in a set of feature maps. These feature maps are then fed to the RPN network which generates a set of candidate object proposals. The RPN network relies on pre-defined anchor boxes of multiple sizes and aspect ratios in order to effectively learn to generate the candidate proposals. Subsequently, each proposal is then transformed into fixed-size features using RoI-pooling. Finally, the pooled features are then forwarded through the RCN, which predicts the category label for each candidate proposal in addition to refining its bounding box. For training the RPN candidate, a category-agnostic binary label (of being an object or not) is assigned to each anchor. The  $j^{th}$  anchor is assigned a label, denoted as  $y_b^j \in \{0, 1\}$ , as positive (or 1) if it has the highest Intersection over Union (IoU) overlap with one of the ground-truth boxes or if it has an IoU overlap higher than 0.7 with any of the ground-truth boxes in the corresponding image. Similarly, a negative label (or 0) is assigned to the anchor if IoU ratio is lower

than 0.3 for all ground-truth boxes. The RPN is then tasked to perform a binary classification to identify whether the candidate bounding box proposal corresponds to one of the objects in the image and learn the offset between the ground-truth bounding box, denoted as  $\mathbf{b}^j$ , and respective anchor box to get final bounding box prediction, denoted as  $\tilde{\mathbf{b}}^j$ . The offset learning is supervised with the help of a regression loss applied on the bounding box parameters. Both these losses are combined together to obtain the final loss for region proposal network as shown below:

$$\mathcal{L}_{rpn} = \frac{1}{N_b} \sum_j \mathcal{L}_{rpn}^{bce}(y_b^j, p_b^j) + \lambda_{rpn} \frac{1}{N_{bbox}} \sum_j p_b^j \mathcal{L}_{rpn}^{reg}(\mathbf{b}^j, \tilde{\mathbf{b}}^j), \quad (1)$$

where  $j$  is the index of an anchor box in the mini-batch and  $p_b^j$  is the probability assigned to the respective anchor box being an object. The loss,  $\mathcal{L}_{rpn}^{bce}$ , computes the smooth-L1 distance between the given ground truth bounding box and the predicted bounding box  $\tilde{\mathbf{b}}^j$ . Both  $\mathbf{b}^j$  and  $\tilde{\mathbf{b}}^j$  are vectors having four bounding box parameters, namely center x-coordinate, center y-coordinate, height and width to represent a bounding box. Also,  $bce$  denotes binary cross entropy loss and  $reg$  denotes regression loss, which is smooth L1 loss for Faster-RCNN [15]. Here,  $\mathcal{L}_{rpn}^{bce}$  is normalized with the size of mini-batch and  $\mathcal{L}_{rpn}^{reg}$  is normalized with number of bounding box locations  $N_{bbox}$ .

Next, the RCN network is trained to perform classification of RoI-pooled features using cross entropy loss with  $K + 1$  class classification, denoted as  $\mathcal{L}_{rcn}^{ce}$ . Here,  $K$  denotes the number of categories in the dataset and an additional class to represent the background category. Additionally, the RCN is also tasked to predict the bounding box offset through regression loss similar to the RPN network

$$\mathcal{L}_{rcn} = \frac{1}{N_b} \sum_j \mathcal{L}_{rcn}^{ce}(y_c^j, \mathbf{p}^j) + \lambda_{rcn} \frac{1}{N_{bbox}} \sum_j p_b^j \mathcal{L}_{rcn}^{reg}(\mathbf{b}^j, \tilde{\mathbf{b}}^j), \quad (2)$$

where  $y_c^j \in \{1, \dots, K + 1\}$  denotes the ground-truth category label of  $j^{th}$  RoI-pooled feature and  $\mathbf{p}^j$  is the predicted probability vector denoting probabilities assigned for all  $K + 1$  categories. The loss  $\mathcal{L}_{rcn}^{ce}$  denotes cross entropy loss and  $\mathcal{L}_{rcn}^{reg}$  loss is same as the one used for the RPN network.

The overall loss function used to train the entire Faster-

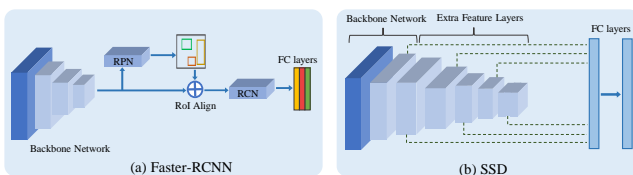


Fig. 4. Illustration of popular detection frameworks: (a) Faster-RCNN [15], (b) Single Shot Multi-Box Detector (SSD) [17].

RCNN network is trained is defined as:

$$\mathcal{L}_{det}^{rcnn} = \mathcal{L}_{rpn} + \mathcal{L}_{rcn}. \quad (3)$$

More details regarding the anchor boxes, bounding box regression losses, training procedure, and architecture can be found in the [15].

### 2.1.2 Single Shot Multi-Box Detector (SSD)

Liu *et al.* [17] proposed a single shot object detection framework which consists of forwarding the image through a single stage as opposed to two stages in the Faster-RCNN detector. Fig. 4(b) illustrates the SSD detection architecture. By following this approach, SSD eliminates the need for an object proposal stage, making it simpler and computationally efficient as compared to the Faster-RCNN approach. The SSD framework employs VGG16 as the backbone network which is used for extracting feature map of size  $H \times W$  from an input image  $X$ . For each feature map location, SSD discretizes the output space of the bounding boxes into a set of default bounding boxes. A convolutional layer is added that for each feature map location predicts a score for a category or offsets relative to the default box coordinates. The set of default boxes contain bounding boxes of multiple predefined aspect ratios and scales to match any object shape in the image better. Furthermore, SSD combines predictions from feature maps at multiple scales to better handle the object scales with respect to the image.

Once the model predictions are available, they are matched with the ground-truth box and category to perform an end-to-end training with regression and classification loss. The regression loss used in SSD is a smooth L1 loss, denoted here as  $L_1^s$ . For each location in the feature map, a default box is matched with a ground-truth bounding box. The final bounding box prediction is computed by adding the predicted offset to the default boxes and the regression loss is computed to correct the offsets based on the matched ground-truth bounding box. The matching strategy is to find default box which has best jaccard overlap with the ground-truth bounding box and then matching default boxes any ground-truth having jaccard overlap higher than 0.5. For a given  $i^{th}$  predicted bounding box  $\tilde{\mathbf{b}}^i$  and matched  $j^{th}$  ground-truth bounding box  $\mathbf{b}^j$ , the corresponding label is defined as  $y^{ij} \in \{0, 1\} = 1$ , the regression loss is given as:

$$\mathcal{L}_{reg} = \sum_{i=1}^{HW} \sum_{j=1}^{N_b} y^{ij} L_1^s(\tilde{\mathbf{b}}^i, \mathbf{b}^j), \quad (4)$$

where  $N_b$  denotes number of ground-truth bounding box per image. Both  $\tilde{\mathbf{b}}^i, \mathbf{b}^j$  are bounding box vector having center  $x-y$  location and height ( $h$ ) and width ( $w$ ). For each predicted bounding box, the classification loss is computed over  $K+1$  categories as shown below:

$$\mathcal{L}_{cls} = - \sum_i \sum_{c=1}^{K+1} \mathbf{y}_c^i \log(\mathbf{p}_c^i), \quad (5)$$

where  $\mathbf{y}_c^i \in \{1, \dots, K+1\}$  denotes one-hot vector indicating the category label respective predicted bounding boxes and  $\mathbf{p}_c^i$  is the corresponding prediction probability vector. Specifically,  $\mathbf{p}_c^i$  denotes the probability of  $i^{th}$  bounding box belonging to  $c^{th}$  category. As earlier, there are  $K$  categories in the dataset and  $K+1^{th}$  label denotes the background class. The final detection loss is a combination of both

regression and classification losses and is defined as follows:

$$\mathcal{L}_{det}^{ssd} = \mathcal{L}_{reg} + \mathcal{L}_{cls}. \quad (6)$$

In the case where there are no predicted bounding boxes that can be matched with one of the ground-truth bounding boxes, the regression loss is set to zero. More details regarding the default boxes, box matching algorithm bounding box regression losses, training procedure, and architecture details can be found in [17].

## 2.2 Domain adaptation

In the domain adaptation problem, we consider two domains, namely source and target, denoted as  $\mathcal{S}$  and  $\mathcal{T}$ , respectively. The source and target domains are assumed to have different data distributions, i.e.,  $P_{\mathcal{S}} \neq P_{\mathcal{T}}$ . Most domain adaptation formulations consider that the source dataset is label-rich, while the target dataset is label-scarce in nature [22]. Multiple variations of this formulation that are commonly studied in the literature include, *semi-supervised* [91], [92], [93], *weakly-supervised* [94], and *unsupervised* domain adaptation [21], [25], [57], [72], [73], [95], [96]. In the context of object detection, the *semi-supervised domain adaptation* formulation assumes that source domain is fully labeled with bounding box annotations and corresponding category labels and only a subset of the target domain samples are fully annotated with bounding box and respective category labels. *Weakly-supervised domain adaptation* formulation assumes that source domain is fully annotated and all target domain samples have binary annotations indicating the presence/absence of any category and no bounding box annotations. Lastly, the *unsupervised domain adaptation* formulation assumes that source domain is fully annotated while no annotations are available for target domain. Among these formulations, the unsupervised formulation is more practical and challenging. Further, the solutions obtained for this formulation can be easily adopted to address the semi-supervised and weakly-supervised domain adaptation tasks as well. For these reasons, we mainly focus on reviewing works that address unsupervised domain adaptation for object detection. In what follows, we formally define the unsupervised domain adaptation formulation and provide a brief overview of the same.

### 2.2.1 Unsupervised Domain adaptation

Let us denote the source dataset as,  $\mathcal{S} = \{X_s^i, Y_s^i\}_{i=1}^{N_s}$ , and it consists of  $N_s$  number of images. Here,  $X_s^i$  denotes  $i^{th}$  image and  $Y_s^i$  denotes the corresponding bounding box annotations with category label. Similarly, let us denote the target dataset as,  $\mathcal{T} = \{X_t^i\}_{i=1}^{N_t}$  having  $N_t$  number of target domain images with no ground-truth annotations. Ben *et al.* [21] proposed a framework to perform domain adaptation for the given setup, i.e., labeled source dataset and unlabeled target dataset, with theoretical upper bounds on the target performance. Ben *et al.* [21] designed a  $\mathcal{H}\Delta\mathcal{H}$ -distance to measure the divergence between two sets of samples that have different data distributions, as is the case for the domain adaptation problem. Let us consider an arbitrary source domain image  $X_s \in \mathcal{S}$  and an arbitrary target domain image  $X_t \in \mathcal{T}$ . Furthermore, let us consider a domain discriminator denoted as,  $D : X \rightarrow \{0, 1\}$ , that



takes in any image  $X \in \{\mathcal{S} \cup \mathcal{T}\}$  and predicts the domain of the input image. classifies source domain image  $X_s \in \mathcal{S}$  as label 0, and target domain image  $X_t \in \mathcal{T}$  as label 1. Considering  $\mathcal{H}$  to be a set of possible domain discriminators, the  $\mathcal{H}\Delta\mathcal{H}$ -distance can be defined as follows:

$$d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) = 2 \sup_{(D, D') \in \mathcal{H}^2} \left| \mathbf{E}_{X \sim \mathcal{S}} [D(X) \neq D'(X)] - \mathbf{E}_{X \sim \mathcal{T}} [D(X) \neq D'(X)] \right|, \quad (7)$$

where  $\mathbf{E}_{X \sim \mathcal{S}}$  and  $\mathbf{E}_{X \sim \mathcal{T}}$  denotes the expected domain classification errors over the source and target domain dataset, respectively. More precisely, the Eq. 7 measures the divergence by the disagreement of the hypothesis sampled from  $\mathcal{H}$ . The ideal joint hypothesis is defined as:

$$D^* = \underset{D \in \mathcal{H}}{\operatorname{argmin}} \operatorname{Err}_{\mathcal{S}}(D^*) + \operatorname{Err}_{\mathcal{T}}(D^*). \quad (8)$$

Here, the terms  $\operatorname{Err}_{\mathcal{S}}$  and  $\operatorname{Err}_{\mathcal{T}}$  denote the expected prediction errors on the source and target domain data samples, respectively. This distance is often used in the domain adaptation literature to measure the adaptability between any give source and target domain datasets. Ben *et al.* [21] present a theorem that further defines the upper bound on the target error as:

$$\forall D \in \mathcal{H}, \operatorname{Err}_{\mathcal{T}}(D) \leq \operatorname{Err}_{\mathcal{S}}(D) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \operatorname{Const}. \quad (9)$$

We can note from the Eq. 9, the target error is upper bounded by three terms, namely expected prediction error on the source domain, domain divergence denoted in Eq.7, and few constant terms. More details regarding both Eq. 7 and Eq. 9 are provided in [21]. A majority of the domain adaptation works in the literature rely on this formulation and focus on minimizing the upper bound on the target error by reducing the domain divergence between the source and target domain. In what follows, we discuss the different strategies used in the literature to address this specifically for the task of object detection.

### 3 METHODS

As discussed earlier, a variety of methods [25], [57], [58], [59], [60], [61], [62], [63], [64], [65], [66], [67], [68], [69], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80] have been proposed for the task of domain adaptive object detection. Based on a meticulous review of these approaches, we categorize them into the following classes.

1) **Adversarial feature learning:** This class of adaptation approach performs an adversarial training of object detector model with the help of a domain discriminator. The training follows a gradient reversal layer based feature learning proposed by Ganin *et al.* [23]. Specifically, the detector model is trained to produce features that fool the domain discriminator, while domain discriminator is tasked to correctly classify the features as source/target domain. This results in detector model producing domain invariant features which are useful to perform detection in target domain. Many methods in the literature utilize this strategy to adapt detectors to target domain [25], [57], [72], [73], [71].

2) **Pseudo-label based self-training:** Many works in the literature [94], [61], [62], [97], [98] utilize highly confident predictions by source-trained detector model to train the it on the target. Since confident predictions on target domain have higher chance of being correct, such training strategy progressively makes detector model better on the target.

3) **Image-to-image translation:** The image-to-image translation based strategy utilize an unpaired image-translation model to translate target image to a source-like image or vice versa. This reduces distribution shift in the visual domain and makes it easier for detector to perform well on the source-like target images. Many work in the literature utilize such approach for improving detector performance on the target domain [99], [64], [100], [101].

4) **Domain randomization:** Another interesting way to improve the detector performance on target domain is to devoid all source-style bias from the model. Domain randomization strategy creates multiple stylized version of source domain data to train the detector model such that the model is not biased towards any one style and generalizes better on the target domain. Some works in the literature such as [59], [101] follow this strategy.

5) **Mean-teacher training:** Mean-teacher is an effective way to utilize unlabeled data to improve model generalization by progressively training a detector model in a student-teacher framework. This motivated a few works in the literature [68], [96] to explore mean-teacher training to adapt detector model by utilizing unlabeled target domain.

6) **Graph reasoning:** Some works in the literature [96], [95], [120] exploit the inter-object and intra-object relationships that exist in detection dataset. These object relations are modeled through graphs which help the detector on target domain by training to enforce the same object relations.

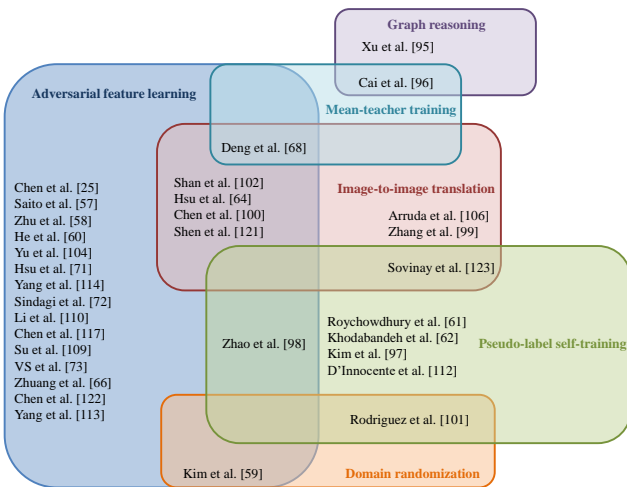


Fig. 5. The figure provides a Venn diagram of domain adaptive object detection methods. Each method falls into a class of adaptation approach (discussed in Sec. 3), listed in the set representing respective techniques. Some of the work fall into more than one class of approaches, listed in the overlap region of set of respective class of adaptation techniques.

Fig. 5 illustrates the various categories of approaches. A comprehensive list of the existing approaches is presented in

TABLE 1  
List of different domain adaptive object detection approaches.

Method	Detection framework	Type	Publication	Year
Faster-RCNN in the wild [25]	Faster-RCNN	Adversarial feature learning	Chen <i>et al.</i> , CVPR	2018
Cross-domain weakly supervised adaptation [94]	SSD	Pseudo-label based self-training	Inoue <i>et al.</i> , CVPR	2018
Strong weak distribution alignment [57]	Faster-RCNN	Adversarial feature learning	Saito <i>et al.</i> , CVPR	2019
Selective cross-domain alignment [58]	Faster-RCNN	Adversarial feature learning	Zhu <i>et al.</i> , CVPR	2019
Diversify and match [59]	Faster-RCNN	Domain randomization, Adversarial feature learning	Kim <i>et al.</i> , CVPR	2019
Automatic adaptation from unlabeled videos [61]	Faster-RCNN	Pseudo-label based self-training	Roychowdhury <i>et al.</i> , CVPR	2019
Mean teacher with object relations [96]	Faster-RCNN	Graph-reasoning, Mean-teacher training	Cai <i>et al.</i> , CVPR	2019
Multi-adversarial adaptation [60]	Faster-RCNN	Adversarial feature learning	He <i>et al.</i> , ICCV	2019
Robust learning from noisy labels [62]	Faster-RCNN	Pseudo-label based self-training	Khodabandeh <i>et al.</i> , ICCV	2019
Self-training for one-stage detector [97]	SSD	Pseudo-label based self-training	Kim <i>et al.</i> , ICCV	2019
Multi-level adaptation [63]	Faster-RCNN	Adversarial feature learning	Xie <i>et al.</i> , ICCV Workshop	2019
Pixel and feature adaptation [102]	Faster-RCNN	Image-to-image translation, Adversarial feature learning	Shan <i>et al.</i> , Neurocomputing	2019
Adapting from synthesis to reality [103]	SSD	Adversarial feature learning	Xu <i>et al.</i> , IEEE Access	2019
Improving localization [104]	Faster-RCNN	Adversarial feature learning	Yu <i>et al.</i> , IEEE Access	2019
Cycle-consistent adaptation [99]	Faster-RCNN	Image-to-image translation	Zhang <i>et al.</i> , IEEE Access	2019
Cross-domain scene text [105]	Faster-RCNN	Adversarial feature learning	Chen <i>et al.</i> , ICNIP	2019
Cross domain detection image translation [106]	Faster-RCNN	Image-to-image translation	Arruda <i>et al.</i> , IJCNN	2019
Graph-induced prototype alignment [95]	Faster-RCNN	Graph-reasoning	Xu <i>et al.</i> , CVPR	2020
Coarse-to-fine adaptation [107]	Faster-RCNN	Adversarial feature learning	Zheng <i>et al.</i> , CVPR	2020
Harmonizing transferability and discriminability [100]	Faster-RCNN	Image-to-image translation, Adversarial feature learning	Chen <i>et al.</i> , CVPR	2020
Cross-domain document object detection [108]	Faster-RCNN	Adversarial feature learning	Li <i>et al.</i> , CVPR	2020
Categorical regularization [69]	Faster-RCNN	Adversarial feature learning	Xu <i>et al.</i> , CVPR	2020
Prior-based detector adaptation [72]	Faster-RCNN	Adversarial feature learning	Sindagi <i>et al.</i> , ECCV	2020
Every pixel matters [71]	FCOS [89]	Adversarial feature learning	Hsu <i>et al.</i> , ECCV	2020
Collaborative training [98]	Faster-RCNN	Pseudo-label based self-training, Adversarial feature learning	Zhao <i>et al.</i> , ECCV	2020
Conditional normalization network [109]	Faster-RCNN	Adversarial feature learning	Su <i>et al.</i> , ECCV	2020
Spatial attention pyramid network [110]	Faster-RCNN	Adversarial feature learning	Li <i>et al.</i> , ECCV	2020
Asymmetric tri-way training [67]	Faster-RCNN	Adversarial feature learning	He <i>et al.</i> , ECCV	2020
Dual multi-label prediction [111]	Faster-RCNN	Adversarial feature learning	Zhao <i>et al.</i> , ECCV	2020
One-shot cross-domain adaptation [112]	Faster-RCNN	Pseudo-label based self-training	D'Innocente <i>et al.</i> , ECCV	2020
Progressive adaptation [64]	Faster-RCNN	Image-to-image translation, Adversarial feature learning	Hsu <i>et al.</i> , WACV	2020
Multi-scale robust discrimination [70]	Faster-RCNN	Adversarial feature learning	Pan <i>et al.</i> , WACV	2020
Object detection via style consistency [101]	SSD	Domain randomization, Pseudo-label based self-training	Rodriguez <i>et al.</i> , BMVC	2020
Image-instance full alignment network [66]	Faster-RCNN	Adversarial feature learning	Zhuang <i>et al.</i> , AAAI	2020
Free lunch for source-free adaptation [74]	Faster-RCNN	Pseudo-label based self-training	Li <i>et al.</i> , AAAI	2020
Forward-backward cyclic adaptation [113]	Faster-RCNN	Adversarial feature learning	Yang <i>et al.</i> , ACCV	2020
Domain invariant region proposal [114]	Faster-RCNN	Adversarial feature learning	Yang <i>et al.</i> , ICME	2020
Uncertainty-aware distributional alignment [115]	Faster-RCNN	Adversarial feature learning	Nguyen <i>et al.</i> , ICM	2020
Region proposal oriented adaptation [116]	Faster-RCNN	Adversarial feature learning	Alqasir <i>et al.</i> , ACIVS	2020
Cross-device OCT lesion detection [113]	Faster-RCNN	Adversarial feature learning	Yang <i>et al.</i> , ISBI	2020
Memory-guided category-wise adaptation [73]	Faster-RCNN	Adversarial feature learning	VS <i>et al.</i> , CVPR	2021
Implicit invariant one-stage network [117]	SSD	Adversarial feature learning	Chen <i>et al.</i> , CVPR	2021
Unbiased mean-teacher [68]	Faster-RCNN	Mean-teacher training, Image-to-image translation, Adversarial feature learning	Deng <i>et al.</i> , CVPR	2021
Domain-specific suppression [118]	Faster-RCNN	Adversarial feature learning	Wang <i>et al.</i> , CVPR	2021
Augmented feature alignment [119]	Faster-RCNN	Image-to-image translation, Adversarial feature learning	Wang <i>et al.</i> , TIP	2021
Instance-invariant progressive disentanglement [120]	Faster-RCNN	Adversarial feature learning, Graph-reasoning	Wu <i>et al.</i> , TPAMI	2021
Large-scale instance-level image-to-image translation [121]	Faster-RCNN	Image-to-image translation, Adversarial feature learning	Shen <i>et al.</i> , IJCV	2021
Scale-Aware Domain Adaptive Faster RCNN [122]	Faster-RCNN	Adversarial feature learning	Chen <i>et al.</i> , IJCV	2021
Curriculum self-paced learning [123]	Faster-RCNN	Pseudo-label based self-training, Image-to-image translation	Sovinyay <i>et al.</i> , CVIU	2021
Adaptive transformer-based detector [78]	DETR [78]	Adversarial feature learning	Zhang <i>et al.</i> , archived	—

Table 1. In what follows, we discuss the key papers related to the respective categories in the aforementioned list.

### 3.1 Adversarial feature learning

#### 3.1.1 Adversarial training through gradient reversal

The adversarial feature learning is built on the theory proposed by Ben *et al.* [21] (see Sec. 2.2.1 for details). Specifically, the overall strategy involves minimizing the upper bound given in Eq. 9 by directly minimizing the  $\mathcal{H}\Delta\mathcal{H}$ -distance. As we can notice from  $\mathcal{H}\Delta\mathcal{H}$ -distance given in Eq. 7, this distance is inversely proportional to the error rate of the domain classifier  $D$ . The goal in a domain adaptation scenario is to reduce this distance, i.e., increase the domain classifier

error. Ganin *et al.* [23] exploited this and proposed a novel gradient reversal approach to train any neural network model for domain adaptation. The overall goal is to achieve a domain invariant feature space of a backbone neural network with the help of a neural network-based domain classifier. Suppose we denote a domain classifier network as  $D$  and the backbone feature extractor network as  $F$ . In that case, the feature extractor network also tries to increase the domain classifier loss. The network  $F$  tries to minimize the task-specific loss (classification/segmentation/detection loss) and maximize the domain classification loss in the overall training pipeline. The network  $D$  is trained to minimize domain classification loss. In addition to the task-specific loss, an additional loss involving domain classifi-

cation is added. This loss is termed as adversarial loss [23] and it can be written as:

$$\max_F \min_{D \in \mathcal{H}} \{E_S(D) + E_T(D)\}, \quad (10)$$

where  $\mathcal{H}$  denotes the hypothesis space for the domain classifier and  $F$  is the feature extractor network.  $E_S(D)$  and  $E_T(D)$  denote the expected domain classification error over source and target domain, respectively. Eq. 10 is implemented with the help of a gradient reversal layer which is applied before the input to the domain classifier as shown in Fig. 6. The gradient reversal layer during feed-forward acts as an identity function and the gradients are multiplied with  $-1$  during backpropagation. In effect, this forces feature extractor  $F$  to maximize the domain classification loss while minimizing the task-specific loss resulting in the domain invariant feature space as proven by Ben *et al.* [21]. All the methods are utilizing this strategy to adapt a detector model using labeled source and unlabeled target domain fall under the *adversarial feature learning* category.

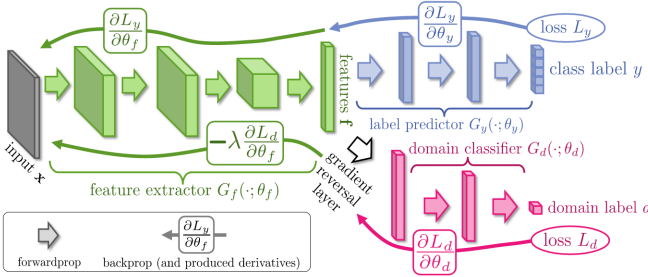


Fig. 6. Domain adaptation by backpropagation as proposed in [23] with the example of classification task.

### 3.1.2 Domain adaptive detection via adversarial training

Chen *et al.* [25] were among the first to formulate and address the problem of domain adaptive object detection. Fig. 7 illustrates an overview of their approach. Given the problem setup discussed in Sec. 2.2.1 with a source  $\mathcal{S}$  and target  $\mathcal{T}$  domain, the method utilizes Faster-RCNN detection framework and proposes to utilize gradient reversal training at multiple stages of the detection framework. Specifically, given a backbone network  $F_b$ , RCN network  $F_{rcnn}$ , and RPN network  $F_{rpn}$  in the detection model, they apply adversarial training at both the *image-level* features which are extracted from backbone  $F_b$  and the *instance-level* features that are extracted from the RCN network,  $F_{rcnn}$ . Let the features extracted from the backbone be denoted as  $F_b(X) \in \mathbb{R}^{C \times H \times W}$ , where  $C$  denotes the number of channels,  $H$  and  $W$  denote the height and width of the feature map, respectively. Furthermore, the discriminator networks for adversarial training at image-level and instance-level are denoted as  $D_{img}$  and  $D_{inst}$ . The image-level domain classification loss used to perform adversarial training at the image-level is then defined as:

$$\mathcal{L}_{img} = - \sum_{i,h,w} [y_d^i \log D_{img}(F_b(X^i)^{(h,w)}) + (1 - y_d^i) \log(1 - D_{img}(F_b(X^i)^{(h,w)}))], \quad (11)$$

where  $i$  indicates  $i^{th}$  image in the batch,  $h \in \{1, \dots, H\}$ ,  $w \in \{1, \dots, W\}$ , and  $F_b(X^i)^{(h,w)}$  denotes the feature of size  $1 \times C$  at location  $(h, w)$  in the feature map. The output of the discriminator network  $D_{img}$  is a probability score

indicating whether the given image is from source domain or target domain. Here,  $y_d^i \in \{0, 1\}$  denotes the domain label which is 0, when  $X^i \in \mathcal{S}$  and 1, when  $X^i \in \mathcal{T}$ . Similarly, the instance-level domain classification loss used to perform adversarial training at the instance level is defined as:

$$\mathcal{L}_{inst} = - \sum_{i,j} [y_d^i \log D_{inst}(f_{ij}^{pooled}) + (1 - y_d^i) \log(1 - D_{inst}(f_{ij}^{pooled}))], \quad (12)$$

where  $f_{ij}^{pooled}$  indicates the RoI-pooled feature of size  $1 \times d$  from the  $j^{th}$  proposal of the image  $X^i$ . Since, both image and instance domain discriminators are trained independently and for any image  $X^i$  should result in same domain prediction, Chen *et al.* [25] introduce a regularization that enforces consistency across predictions of the domain discriminators. This consistency regularization is defined as:

$$\mathcal{L}_{consistency} = \sum_{i,j} \left\| \frac{1}{N_{act}^i} \sum_{h,w} D_{img}(F_b(X^i)^{(h,w)}) - D_{inst}(f_{ij}^{pooled}) \right\|_2, \quad (13)$$

where  $N_{act}^i$  indicates number of activations in the feature map  $F_b(X^i)$  of a given image  $X^i$  and  $\|\cdot\|_2$  indicates  $L_2$ -norm. Let us denote the overall detection model as  $F = \{F_b, F_{rpn}, F_{rcnn}\}$ . Combining all these loss functions, the final training objective is given as:

$$\max_{D_{inst}, D_{img}} \min_F \mathcal{L}_{det}^{rcnn} - \lambda (\mathcal{L}_{img} + \mathcal{L}_{inst}), \quad (14)$$

$$\min_{F, D_{inst}, D_{img}} \lambda \mathcal{L}_{consistency}, \quad (15)$$

where  $\lambda$  is a trade-off parameter used to balance adversarial and consistency loss. The detection and the discriminator networks are trained with final objective, Eq. 14 and Eq. 15. To summarize, the detection network  $F$  aims to minimize the detection loss and maximize the image-level and instance-level domain classification loss. The discriminator networks aim to minimize the domain classification loss. Note that both the detector and the discriminators aim to minimize the consistency loss. The detection loss is applied only to the source data since target data does not have any label annotations. The domain classification loss and consistency regularization are applied on both labeled source and unlabeled target data.

Saito *et al.* [57] argue that directly applying the gradient reversal at multiple levels in the backbone network is not necessarily optimal. Since shallower convolutional layers of the backbone capture local information, directly applying the adversarial loss would be beneficial in learning domain

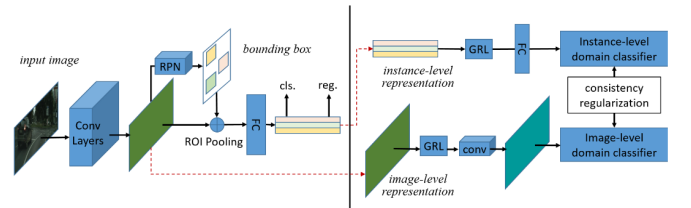


Fig. 7. Domain adaptive object detection in the wild [25] applies adversarial feature learning at image-level and instance-level with gradient reversal applied to detection backbone feature map and pooled features of the Faster-RCNN model. To regularize the adaptation, a consistency regularization is applied to make both image-level and instance-level domain classifier in sync with each other.

invariant local features between source and target domain data. This type of alignment is to as *strong alignment* in their work. Further, they reason that since the deeper layers in the backbone capture global information, performing a similar strong alignment would not be optimal. For example, when source and target domain data are sampled from different countries/cities, the number of objects in a scene, object co-occurrence, scene layout, etc can be very different. Consider the case when the source domain image contains only one object, whereas the target image contains multiple objects, performing strong alignment would likely increase the risk of misalignment. To tackle this issue, [57] modified the adversarial loss at the global-level (later convolutional layers of backbone network) by replacing the traditionally used binary cross-entropy loss with focal loss [124]. This strategy in [57] is referred to as *weak-alignment*. The overall approach then utilizes both the local *strong alignment* and global *weak-alignment* to reduce the domain gap between source and target domain data, resulting in increased detection performance for target images. Let the detection backbone network  $F_b$  be divided into two sub-networks that are cascaded together, namely global feature extractor  $F_g$  and local feature extractor  $F_l$  such that,  $F_b(X) = F_g(F_l(X))$ . The adversarial loss for strong alignment is defined as:

$$\begin{aligned}\mathcal{L}_{local_s} &= \frac{1}{N_s H W} \sum_{i=1}^{N_s} \sum_{w=1}^W \sum_{h=1}^H \|D_l(F_l(X_s^i))^{(h,w)}\|^2, \\ \mathcal{L}_{local_t} &= \frac{1}{N_t H W} \sum_{i=1}^{N_t} \sum_{w=1}^W \sum_{h=1}^H \|1 - D_l(F_l(X_t^i))^{(h,w)}\|^2, \\ \mathcal{L}_{local} &= \frac{1}{2} (\mathcal{L}_{local_s} + \mathcal{L}_{local_t}),\end{aligned}\quad (16)$$

where  $D_l$  denotes the local domain discriminator,  $N_s$  and  $N_t$  are the number of source and target examples respectively,  $\mathcal{L}_{local_s}$  and  $\mathcal{L}_{local_t}$  denote the source and target domain classification loss respectively, and  $F_l(X_s^i), F_l(X_t^i) \in \mathbb{R}^{C \times H \times W}$  are source and target local feature maps respectively. Denoting the global discriminator network as  $D_g$ , the weak alignment loss is defined as:

$$\begin{aligned}\mathcal{L}_{global_s} &= -\frac{1}{N_s} \sum_{i=1}^{N_s} \text{FL}(D_g(F_b(X_s^i))), \\ \mathcal{L}_{global_t} &= -\frac{1}{N_t} \sum_{i=1}^{N_t} \text{FL}(1 - D_g(F_b(X_t^i))), \\ \mathcal{L}_{global} &= \frac{1}{2} (\mathcal{L}_{global_s} + \mathcal{L}_{global_t}),\end{aligned}\quad (17)$$

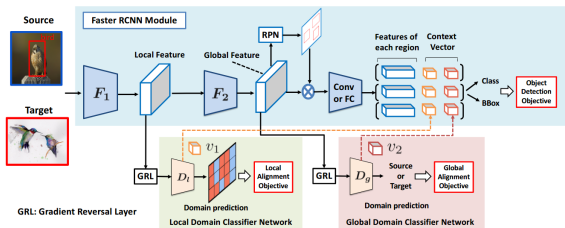


Fig. 8. Strong weak alignment for adaptive object detection [57]. One discriminator is applied at local-level (early layers of backbone) to perform a strong feature alignment via gradient reversal layer and another discriminator performs a weak alignment on the global-level (final layers of backbone). Furthermore, the features learned by the discriminator are fed to the classification network to provide domain context to improve the performance further.

where  $\mathcal{L}_{global_s}$  and  $\mathcal{L}_{global_t}$  are weak global alignment loss applied on source and target domain images, respectively.  $F_b(\cdot)$  extracts the global feature, and  $D_g(\cdot)$  denotes the probability of the global feature being from the source domain or target domain. FL denotes the focal loss defined as:

$$\text{FL}(p) = -(1 - p)^\gamma \log(p), \quad (18)$$

where  $\gamma$  is a focal loss parameter that controls the weight on hard-to-classify examples [124]. More specifically, the value of  $\gamma > 1$  will assign low loss values for the easy samples and high loss for hard samples, thereby focusing on the hard samples while training. As a result, while performing adversarial training with gradient reversal layer at the global level, hard target samples will be given more focus and easy to classify target examples will not be forced to align with the source domain. The paper demonstrated adaptation in the Faster-RCNN framework with this strong and weak alignment strategy and specifically showed the significance of weak alignment at the global level. Additionally, to further improve the performance of local and global domain discriminators,  $D_l$  and  $D_g$  are concatenated with RCN features to add domain-specific context information for object classification, as shown in Fig. 8. The final objective of the network is defined as:

$$\max_{D_g, D_l} \min_F \mathcal{L}_{det} - \lambda (\mathcal{L}_{global} + \mathcal{L}_{local}), \quad (19)$$

where  $\lambda$  is a trade-off hyper-parameter and  $F$  denotes the entire detection network.

A major drawback of the methods discussed so far is that they try to utilize the entire feature map to perform the alignment. However, a more optimal approach would be to perform the feature alignment on regions corresponding to objects in detection. Zhu *et al.* [58] base their method on this observation and selectively align the features of source and target domain data by mining the regions that are discriminative. For this, the authors exploited the region proposal network of the Faster-RCNN detection framework. Their method is divided into two parts, namely “where to look” and “how to align”, as illustrated in Fig. 9. In the “where to look” stage, region proposals generated by the RPN network are mined to find groups within the feature maps. To overcome the noisy proposals of the RPN network, *K-means* clustering is performed using the center coordinates of the proposals. The cluster centers obtained through the *K-means* are used as grouped regions. Based on these mined groups, a fixed number of features are reassigned to each group (i.e. cluster). Instead of performing alignment with gradient reversal layer on these grouped features, a generative adversarial network [125] based strategy is used to perform indirect feature alignment through generation. A similar strategy has been shown to work well in the case of classification [41], [126]. Specifically, they use the features in a group to reconstruct the corresponding patch of the original image by performing within-domain (i.e.  $s \rightarrow s, t \rightarrow t$ ) and cross-domain (i.e.  $t \rightarrow s, s \rightarrow t$ ) patch reconstructions. Denoting domain-specific generators as  $G_s$  and  $G_t$  and domains specific discriminators as  $D_s$  and  $D_t$  corresponding to the source and target domain, respectively, the patch reconstruction adversarial joint loss is defined as:

$$\mathcal{L}_{adv}^{joint} = \mathcal{L}_{D_s} + \mathcal{L}_{D_t} + \mathcal{L}_{G_s} + \mathcal{L}_{G_t} + \mathcal{L}_F. \quad (20)$$

Here, the loss functions  $\mathcal{L}_{D_s}$  and  $\mathcal{L}_{D_t}$  are within domain



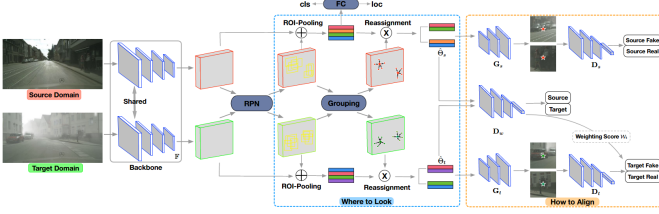


Fig. 9. Selective cross-domain alignment method [58]. First, the region of interest-based grouping strategy is used to identify the discriminative regions. Next, the discriminative regions are used to perform a weighted alignment using generative adversarial training.

losses that try to minimize the patch real/fake classification loss by predicting real input as real and fake input as fake. In contrast, for the generator network losses  $\mathcal{L}_{G_s}$  and  $\mathcal{L}_{G_t}$  then try to fool the discriminator networks by forcing them to identify fake as real and real as a fake within the same domain. The detection model  $F$  aims to minimize the cross-domain patch reconstruction loss, i.e., classifying a fake source as a real target and a fake target as a real source. Similar strategies are used in [41], [126] for classification task. By minimizing the cross-domain patch reconstructions, the network  $F$  will learn a domain invariant feature space resulting in a reduced gap between source and target domain. This strategy closely follows the gradient reversal training explained earlier. However, instead of directly applying adversarial loss on the feature space, the loss is applied to the reconstructed patch images. Additionally, a weight estimation network  $D_w$  is also added to balance source and target domain adversarial losses. The weight estimation network is trained with a binary cross-entropy loss. Denoting the source and target domain RoI-pooled feature groups as  $\tilde{\mathbf{f}}_s$  and  $\tilde{\mathbf{f}}_t$  respectively, the weight estimation loss is defined as:

$$\mathcal{L}_w = \log(D_w(\tilde{\mathbf{f}}_s)) + \log(1 - D_w(\tilde{\mathbf{f}}_t)). \quad (21)$$

The primary benefit of using weight estimation network is that it weighs the alignment loss based on how similar the target patches look to the source domain patches. If we were to denote the output of the weight estimation network for the target domain pooled features as  $\lambda_t$ , the final objective for the method can be written as:

$$\min_{F, G_s, G_t, D_w} \max_{D_t, D_s} \mathcal{L}_{det}^{frcnn} + \lambda_t(\mathcal{L}_{D_t} + \mathcal{L}_{G_t} + \mathcal{L}_F) + (\mathcal{L}_{D_s} + \mathcal{L}_{G_s}) + \mathcal{L}_w. \quad (22)$$

The network training is performed in a stage-wise manner by updating supervised detection loss, discriminator loss, weight loss, cross domain and within domain generation loss separately.

### 3.1.3 Weighted adversarial feature alignment

Most of the work in the domain adaptive detection literature that are based on adversarial feature learning focus on improving the gradient reversal training in order to obtain better feature alignment. Various methods such as [107], [69], [72], [71], [111], [73] build on the approaches discussed earlier and broadly follow the strategy of introducing a module that can control the gradient reversal information flow through loss weighting in addition to using a regularization technique to complement the proposed weighting module. For example, Zheng *et al.* [107] applies domain classifier at multiple levels of a detection

backbone network to perform adversarial feature learning with gradient reversal layer, as shown in Fig. 10. These adversarial losses are then multiplied by weights extracted from the backbone. In order to obtain the weights, the final feature map of the backbone is averaged over the channel dimension to obtain an attention map highlighting regions that potentially have an object. These weights are then used to modulate the adversarial loss. This strategy is referred to as Attention-based Region Transfer (ART). Furthermore, with the help of source ground-truth bounding boxes and predicted proposals for the target domain, class-specific prototypes are learned through feature averaging. At each step, both source and target prototypes of each class are aligned through prototype similarity loss. This strategy is referred to as Prototype Similarity Alignment (PSA). The ART loss helps gradient reversal training remove any noisy information coming from non-object regions resulting in better alignment, whereas prototype alignment maintains the semantic consistency while adapting to the target domain.

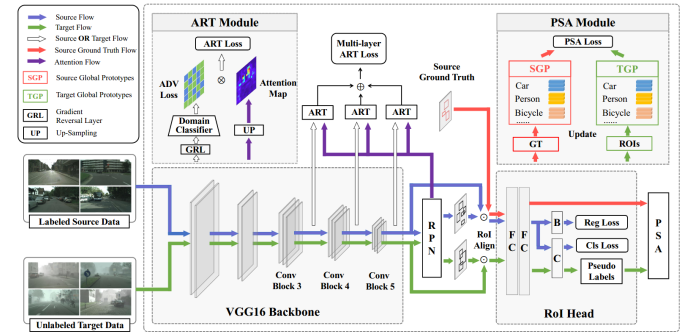


Fig. 10. The coarse-to-fine adaptation [107] proposes a multi-level alignment by applying gradient reversal-based adversarial loss at multiple layers in the detection backbone network. An attention mechanism is introduced that weights the adversarial loss based on the activation strength at any particular location in the respective feature map. This loss weighting mechanism is called as Attention-based Regional Transfer (ART) in the figure. Furthermore, a Prototype Similarity Alignment (PSA) loss is introduced that aligns the object prototype features of both source and target domain to regularize the alignment.

Similarly, He *et al.* [60] proposed a Multi-Adversarial Faster-RCNN (MAF) approach that extends the work by [25]. The first change in MAF is to apply image-level alignment at multiple layers of the backbone network. Additionally, the feature maps are reduced in the channel dimension to align the aggregated information rather than individual components. MAF also includes instance-level alignment but unlike [25], it weights the instance adversarial loss with the prediction probabilities produced by the RCN network. Fig. 11 illustrates the overall approach of MAF.

Xu *et al.* [69] a proposed categorical regularization strategy that can be combined with adversarial feature learning based approaches to further enhance the feature alignment between source and target domain. The categorical regularization strategy utilizes instance-level annotations from the labeled source domain data and adds a multi-label classification loss in addition to the detection losses as illustrated in Fig. 12. Such a strategy helps extract weak localization of objects in the feature maps through a multi-label classifier. The weak localization map is then used to gate image-level domain classification loss to block unnecessary infor-

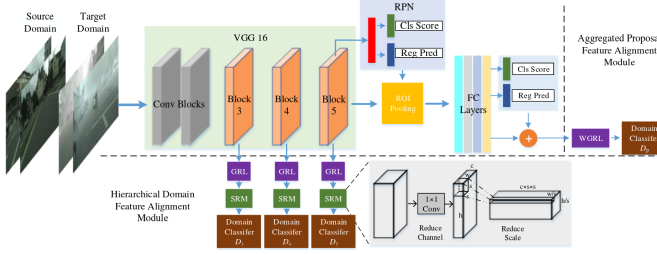


Fig. 11. Multi-adversarial Faster-RCNN [60] extends the adversarial training strategy of [25] by performing image-level gradient reversal training at multiple layers of the detector backbone network. Additionally, it performs hierarchical alignment where the feature maps across the channel dimension are reduced before being forwarded through the domain classifier. Furthermore, it also performs instance-level gradient reversal training and the adversarial loss is weighted with the prediction probabilities obtained from the detection network for each respective RoI-pooled feature as illustrated in the image with Weighted Gradient Reversal Layer (WGRL) loss.

mation while highlighting the object regions. This image-level weighting through weak localization is termed in the paper as Image-level Categorical Regularization (ICR). Subsequently, the instance-level regularization is applied to the RoI-pooled features. Specifically, the instance-level domain classification loss is weighted using the difference between the RCN prediction probability and the multi-label classification probability corresponding to the category of the respective pooled feature. This regularization is known as the Categorical Consistency Regularization (CCR) loss. Xu *et al.* showed the effectiveness of this categorical regularization by modifying both [25] and [57] alignment strategy with the addition of both ICR and CCR loss functions.

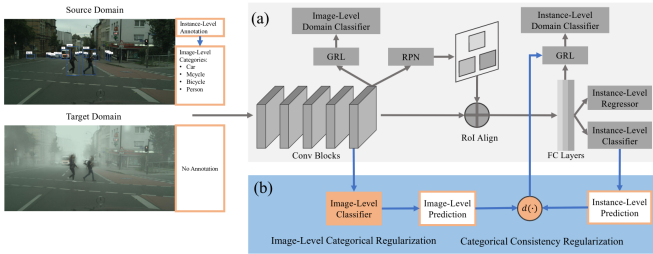


Fig. 12. The categorical regularization approach [69] explores the use of multi-label classification on source domain supervised data along with the adversarial feature learning. The multi-label classification can extract weakly localized object maps that can help guide the adversarial training. The method proposes two additional losses, namely Image-level Categorical Regularization (ICR) and Consistency Categorical Regularization (CCR), that can be applied to any adversarial feature learning pipeline to enhance the performance.

Zhao *et al.* [111] combined the multi-label classification with the weak global alignment [57]. Specifically, a multi-label classifier is additionally trained along with the detection network with the help of source labeled data. The probability scores predicted by the multi-label classifier are used to condition the domain discriminator at the final layer of the backbone network. The conditioning mechanism takes in both the feature map extracted from the backbone and the multi-label probability vector indicating the probability of all  $K$  objects being in the image, using a multi-linear mapping function. The weak global alignment is similar to the [57], i.e., employing focal loss to perform global

feature alignment. To regularize the feature alignment training further, the distance between renormalized prediction probability score vector from multi-label classifier and prediction probability score of RCN network is minimized via symmetric Kullback-Leibler divergence.

Sindagi *et al.* [72] considered adverse weather conditions as a special case of domain adaptation. They argue that in the case of adverse weather, well-defined models are available that mathematically formulate how the camera captures such conditions. Using these models, they extract weather-specific information termed as “prior”. These *priors* are then used to modify the conventional domain classification loss into a prior-prediction loss. The training follows a similar strategy of performing gradient reversal layer-based alignment using prior prediction loss. That is, the backbone network of the detection model tries to maximize the prior prediction loss, while the prior prediction network tries to minimize it. This strategy is termed as *prior adversarial training* in their paper and is shown to be effective, especially in the case of hazy and rainy weather conditions.

In contrast to existing methods, Hsu *et al.* [71] do not utilize Faster-RCNN detection framework and instead they follow FCOS [89] one-stage detection framework that has centerness loss in addition to object classification and bounding box regression losses for supervised detection. Hsu *et al.* [71] specifically exploit the FCOS framework to propose a center-aware feature alignment strategy. Since FCOS is trained with the centerness loss on labeled source data, it provides a coarse prediction of center-focused heatmap indicating the location of the objects. They utilize this center-information to create a class-agnostic map of objects and gate it with the final feature map of the backbone network. Both global and center-aware discriminators are applied on the respective feature maps to perform domain classification. The alignment is performed through gradient reversal layer-based adversarial training where the domain discriminators are responsible for center-aware feature alignment while rejecting any noisy information of the background with more focus on the prominent parts of the object instances. Similar to the other approaches, the feature alignment is performed at multiple levels of the backbone network.

VS *et al.* [73] pointed out that existing adversarial learning-based approaches are prone to the category-wise *negative transfer* problem. Training with gradient reversal and domain classifier only ensures that features become domain invariant. However, when both source and target domain contains multiple categories, which is often the case, there are chances of misalignment across different categories of target and source features. VS *et al.* [73] address this issue of negative transfer with by proposing memory-guided attention for category-aware domain adaptation (MeGA-CDA). Their key idea is to utilize  $K + 1$  category-specific domain discriminators for aligning each category separately. However, the label information is unavailable for the target dataset, which prohibits the usage of category-specific domain discriminators. To overcome this, the authors introduce a memory module in the feature map extracted by the backbone network and produces attention maps that attend to the respective category features while blocking information from the rest of the categories. Fig. 13 illustrates their overall approach of using memory modules for generating

category-specific information and its subsequent usage in class-specific domain alignment. The memory module is trained jointly with the other components of the network, and this end-to-end training ensures that the category-specific attention map is progressively improved over the training process. Note that they also use the category-agnostic domain alignment that is similar to the strong alignment strategy of [57]. This ensures that the overall feature maps are aligned between source and target domain, while the additional category-specific discriminators prevent the negative transfer by performing the class-wise alignment.

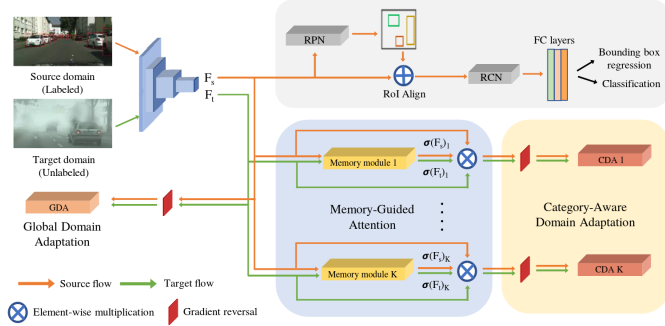


Fig. 13. Memory guided category wise adaptation [73] proposes a strategy to perform class-wise alignment between source and target domain. A memory module is utilized corresponding to each category that is iteratively updated to encode class-information. These memory modules are later used to create attention maps of respective categories. Through a global domain classifier and multiple category-specific domain classifier category-aware feature alignment is performed between source and target images.

Let  $F_b$  as the backbone network of the detection model,  $X_s \sim \mathcal{S}$  and  $X_t \sim \mathcal{T}$  be any arbitrary source and target domain image respectively, the typical weighted domain classification loss ( $\mathcal{L}_{adv}^{weighted}$ ) used for most of the methods discussed in this section can be formulated as:

$$\begin{aligned} \mathcal{L}_{adv}^{weighted}(X_s, X_t) = & \sum_{h=1}^H \sum_{w=1}^W y_d \cdot \sigma_{img_s}^{(h,w)} \cdot \left(1 - D\left(\sigma_{feat_s}^{(h,w)} \cdot F_b(X_s)^{(h,w)}\right)\right)^2 \\ & + (1 - y_d) \cdot \sigma_{img_t}^{(h,w)} \cdot \left(D\left(\sigma_{feat_t}^{(h,w)} \cdot F_b(X_t)^{(h,w)}\right)\right)^2, \end{aligned} \quad (23)$$

where  $D$  is the image-level (in some cases instance-level) domain classifier,  $y_d \in \{0, 1\}$  is domain label which is 1 for source and 0 for target domain images,  $\sigma_{img_s}, \sigma_{img_t} \in \mathbb{R}^{H \times W}$  are the image-level weights applied on the domain classification loss,  $\sigma_{feat_s}, \sigma_{feat_t} \in \mathbb{R}^{H \times W}$  are the feature-level weights applied to mask-out irrelevant information and highlight important features that aid source and target domain alignment. Depending on the method,  $\sigma_{img}$  or  $\sigma_{feat}$  are set to identity, and either image-level or feature level weights are applied for the loss computations. Methods such as [107], [60], [69], [72], [111] utilize only  $\sigma_{img}$  weights to appropriately balance the adversarial loss. Whereas other methods such as [71], [73] apply only feature-level weights  $\sigma_{feat}$ . The contribution of each method comes from the way either  $\sigma_{img}$  and  $\sigma_{feat}$  is modeled. For example, Xu *et al.* [69] models  $\sigma_{img}$  with the help of multi-label classification probability vector and VS *et al.* [73] models  $\sigma_{feat}$  with category-specific memory modules.

### 3.1.4 Additional methods

Chen *et al.* [117] provided a strategy to align source and target domain features for SSD detection framework, termed as I3Net. Their main strategy is similar to that of multi-level discriminators [25], [57], [72], [60], [107] that perform alignment at pixel-level and image-level as shown in Fig. 14. Furthermore, the domain loss used for training the image level discriminator, is weighted by the probability obtained from a multi-label classifier, similar to [69]. Additionally, I3Net enforces object pattern matching by exploiting the SSD architecture, which predicts probability at each feature map location extracted from the backbone network. These category-specific probability maps are then matched between source and target images for respective categories to enforce consistency between source and target activations. The category features are further improved by minimizing intra-class distance and maximizing inter-class distance with a margin between category-specific prototypes. The prototypes are calculated using category-wise probability patterns obtained by SSD detection framework and are updated through the exponential moving average.

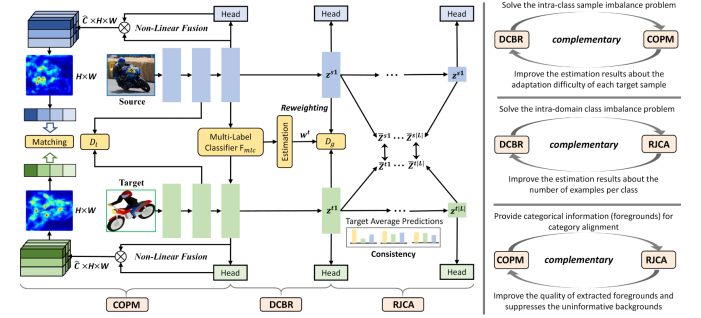


Fig. 14. Implicit instance variation network (I3Net) [117] involves four major components, 1) pixel-level and image-level feature alignment through gradient reversal layer, 2) category-aware object-pattern matching (COPM) through non-linear fusion, and 3) dynamic class-balanced re-weighting (DCBR) through multi-label classifier, 4) regularization through joint category alignment (RJCA). All these modules are combined together to adapt an SSD framework based one-stage detection model for the target domain data.

Wu *et al.* [120] points out that feature-level or pixel-level alignment strategies suffer from the risk of neglecting the instance-level object characteristics. To achieve this, the features learned through the training are required to be disentangled into domain-invariant and domain-specific parts. Wu *et al.* [120] proposed a progressive disentanglement strategy that performs stage-wise training of Faster-RCNN. As shown in Fig. 15, both image-level and instance-level domain classifiers are employed with gradient reversal layer similar to [25]. Mutual Information (MI) loss is applied to the separated features to disentangle domain-invariant from domain-specific features. The MI loss formulation is borrowed from the Mutual Information Neural Estimation (MINE) approach [127]. Specifically, MINE-based MI loss transforms the intractable mutual information maximization into a tractable binary classification objective that can be trained in an end-to-end manner. To further regularize the training, object-relational graph is created with the help of domain-invariant features for each input image, by considering the original and domain-invariant feature maps. Since both features are extracted from the same image, they



contain the same objects at the respective locations. Hence, a object Relational Consistency (RC) loss is enforced to maintain inter-class relationships across both feature maps.

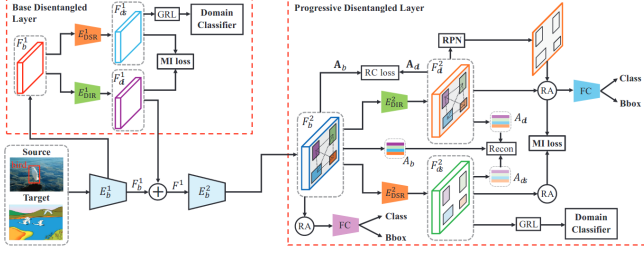


Fig. 15. Instance-invariant adaptation of object detector [120] introduces gradient reversal training at image-level and instance-level similar to [25]. Two disentanglement layers are proposed to achieve instance-invariance, which minimizes the Mutual Information (MI) loss between the domain-specific and domain-agnostic features. To further improve the disentanglement, a Relation Consistency (RC) loss is used by constructing an object relation graph. The overall model training is done in three stages resulting in an instance-invariant detection model that is better suited for the target domain data.

There are several subsequent works such as [66], [76], [63], [55], [110], [70], [114], [122] that utilize the adversarial feature learning strategies discussed in this section. Most of these approaches address cross-domain detection for the application of autonomous driving/surveillance. Some of the notable works that address the issue of domain shift in other applications include Yang *et al.* [113] which tackles the detector adaptation for medical tasks with the help of gradient reversal layer-based adversarial training for performing OCT lesion detection. In other tasks, Chen *et al.* [105] proposed adversarial feature training similar to [25] for adapting scene text detection models from synthetic data to outdoor settings, and Li *et al.* [108] utilized a multi-level feature alignment applied at multiple convolutional layers of the backbone network to train a detector for cross-domain detection on documents. Additionally, Li *et al.* [108] established a benchmark for the task of cross-domain detection in document space. Furthermore, there are many other interesting works available as pre-print [128], [129], [130], [75], [80], [131], [132]. Amongst different types of domain adaptive detection techniques, adversarial feature learning is the most popular approach with different variations available in the literature.

### 3.2 Pseudo-label based self-training

Self-training of object detectors on the target domain with pseudo label-based supervision is one of the simplest approaches to adapt the source-trained model to the target domain. Given a source pre-trained detection model,  $F^{src}$ , which is trained on the annotated source domain dataset  $\mathcal{S} = \{X_s^i, Y_s^i\}_{i=1}^{N_s}$ , it is used to generate pseudo-labels on the target domain data  $\mathcal{T} = \{X_t^i\}_{i=1}^{N_t}$ . The pseudo-labels,  $\hat{Y}_t^i = F^{src}(X_t^i)$ , along with the images from the target dataset forms a new dataset,  $\tilde{\mathcal{T}} = \{X_t^i, \hat{Y}_t^i\}_{i=1}^{N_t}$ . Here,  $\hat{Y}_t^i = F^{src}(X_t^i)$  denotes pseudo-labels obtained from a source-trained detector model  $F^{src}$  for  $i^{th}$  target domain image  $X_t^i$ . Self-training typically involves using these pseudo-labels to re-train the network on the target data. In reality, the pseudo-labels are potentially noisy and often incor-

rect. Hence, directly training the model with these pseudo-annotations can potentially lead to a situation where the errors keep getting reinforced into the network. A filtering strategy is employed to overcome this issue, which involves removing annotations that are not confident. Most of the existing works revolve around developing complex and accurate filtering techniques to deal with the noise present in the pseudo-labels. Note that this training strategy with pseudo-labels is simple yet very effective as it attempts to directly minimize the target domain prediction error, unlike adversarial feature learning where the strategy is to minimize an upper bound over the target prediction error.

Roychowdhry *et al.* [61] proposed a self-training based approach, where they compute pseudo-labels automatically using video data. Specifically, they exploit temporal consistency between adjacent frames to track detections by the source-trained model as illustrated in Fig. 16. This helps in mining more pseudo-labels which the detection model might have missed. Consequently, the pseudo-labels for the target dataset now contain predictions from the detection network and the tracker. It is trained on both labeled source dataset and pseudo-labeled target dataset for adapting the detector to target data. Irrespective of the way it was collected (i.e. from detector model or via tracking), the target domain pseudo-labels are assigned the same category label, i.e.,  $y_t^c = 1$  for positive class and  $y_t^c = 0$  for negative class. However, there are two types within the positive pseudo-labels: predictions extracted from the detection model and pseudo-labels extracted using the tracker. Both of them are assigned a score value to calculate an interpolated labels. The score assignment is given as:

$$s_t = \begin{cases} p_t, & \text{if pseudo-label from detector} \\ \theta, & \text{if pseudo-label from tracker} \end{cases} \quad (24)$$

Here,  $s_t$  is the score assigned to all pseudo-labels,  $p_t$  is probability obtained from detector model,  $\theta$  is confidence assigned to tracker pseudo-labels where regardless of their low confidence, the score is raised up to this constant value to emphasize their importance during training. Using the score value assigned to each pseudo-label, a soft-label is obtained by interpolating the label with score values. The interpolated labels are computed as follows:

$$\tilde{y}_t^c = \lambda s_t + (1 - \lambda) y_t^c, \quad (25)$$

where  $\lambda$  is interpolation parameter and  $y_t^c$  is hard-labels given in by the Eq. 24 for  $i^{th}$  bounding box. Finally, the detector model is trained on both source and target domain data with both ground-truth and pseudo-labels. For any image  $X^i$  the detection loss is given as:

$$\mathcal{L}_{det}^i = \begin{cases} \mathcal{L}_{det}^{frcnn}(Y_s^i, p_s^i), & \text{if } X^i \in \mathcal{S} \\ \mathcal{L}_{det}^{frcnn}(\tilde{Y}_t^i, \tilde{p}_t^i), & \text{if } X^i \in \mathcal{T} \end{cases}, \quad (26)$$

where  $\mathcal{L}_{det}^{frcnn}$  is Faster-RCNN detection loss,  $Y_s$  denotes annotation corresponding to respective source domain image that contains both bounding box and category label, and  $\tilde{Y}_t$  denotes pseudo-label corresponding to respective target domain image containing both bounding box information and interpolated category label  $\tilde{y}_t^c$ . As the training progresses, tracker-based pseudo-labels' emphasis helps in improving adaptation to target domain data.

Although mining pseudo-labels by exploiting the temporal consistency of video data is an effective strategy to obtain more annotations, it relies heavily on the availability



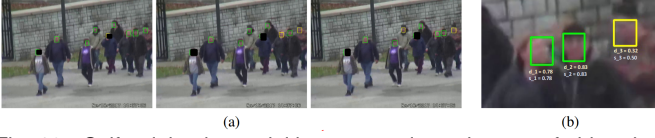


Fig. 16. Self-training by exploiting temporal consistency of video data [61] (a) Green boxes are the detection obtained by applying source-trained detector model. In contrast, yellow bounding boxes are obtained through tracking the bounding boxes across video, exploiting temporal consistency between adjacent frames. (b) Green boxes are high-confidence pseudo-labels and kept as is, whereas yellow boxes tracked by exploiting temporal frame consistency are assigned soft labels, i.e. instead of assigning hard category labels, they are softened through weighted addition with confidence.

of video data which is not always possible. Hence, it is important to develop single-image based pseudo-label training strategy. Khodabandeh *et al.* [62] proposed a method to specifically address the noise present in the pseudo-labels for the single-image target dataset. Their approach, termed as Robust Faster-RCNN, follows a three-phase training strategy. All phases are as illustrated with block diagrams corresponding to each phase in Fig. 17. In the first phase, the detection network is trained with a supervised detection loss  $\mathcal{L}_{det}^{frcnn}$  using a source domain dataset that has access to ground-truth annotations. In the second phase, the source-trained detector model is used to obtain pseudo-labels for the target domain dataset. Subsequently, the pseudo-labels are further refined using an additional classifier network that is pre-trained on a large-scale classification dataset. The refinement strategy utilizes both detector model prediction and classifier network predictions. With the help of refined labels, phase three involves training the detector model with a newly designed loss that accounts for potential noise in the pseudo-labels and helps the detector learn better. Let us consider,  $\mathbf{p}_c$  and  $\mathbf{p}_c^{img}$  be the probability vector of a prediction from detector model and pre-trained image-classification network, respectively. Furthermore, let  $\mathbf{s}_c$  and  $\mathbf{s}_c^{img}$  denote the logits (score vectors) of a prediction from detector model and pre-trained image-classification network, respectively. Also,  $\tilde{\mathbf{b}}^{pseudo}$  and  $\tilde{\mathbf{b}}^{current}$  denote the bounding box pseudo-label collected in the first phase and current bounding box prediction by the detector model, respectively. The Robust Faster-RCNN [62] training utilizes these predictions to refine the pseudo-annotations, which are then used as supervision to train the detection model on the target data. The following equations describe the refinement process:

$$\begin{aligned} \mathbf{p}_c^{refined} &= \text{softmax}((\mathbf{s}_c + \alpha \mathbf{s}_{img}) / (1 + \alpha)), \\ \mathbf{b}^{refined} &= ((\tilde{\mathbf{b}}^{current} + \alpha \tilde{\mathbf{b}}^{pseudo}) / (1 + \alpha)), \end{aligned} \quad (27)$$

where  $\alpha$  is a hyper-parameter that controls the trade-off between the two terms.  $\alpha$  starts with a large value and it is gradually decreased to smaller values as the third phase training progresses [62]. The theoretical formulation that lead to the refinement equations Eq. 27 can be found in [62]. Training with refined annotations help the detector model counter the noise in the pseudo-labels much better.

Previous methods to perform self-training were based on the Faster-RCNN framework and relied on subsequent assumptions to handle the noise in the obtained pseudo-labels. Kim *et al.* [97] proposed a strategy that focused on performing self-training of one-stage object detector models, specifically addressing it for the SSD-based model. The

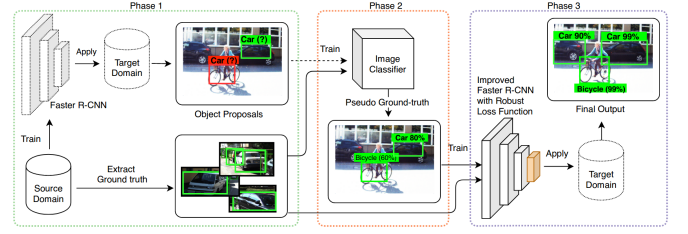


Fig. 17. Robust Faster-RCNN [62] utilizes the source-trained model to obtain noisy labels on the target domain dataset. A classification module is used to update the noisy labels and improve their quality. Finally, the detection model is trained on the target domain data with the help of a robust loss function that incorporates the uncertainty to account for possible mistakes present in the noisy pseudo-labels obtained through the source-trained model.

authors perform hard negative mining of pseudo-labels followed by a weak negative mining strategy as shown in Fig. 18. During the hard mining phase, a threshold  $\delta$  is utilized to filter out the false-negative bounding box proposals based on their IoU with final bounding box prediction. Weak mining phase further improves the pseudo-labels by reducing the risk of false positives with the help of assigning instance-level scores calculated using Support Region-based Reliable Score (SRRS). The score is calculated by a set of RoIs predicted by the detection model that satisfies the hard mining threshold  $\delta$ . If there are  $N_{support}$  RoIs that satisfy the  $\delta$ -IoU threshold corresponding to the any detector bounding box prediction  $\tilde{\mathbf{b}}$  having confidence score of  $\tilde{p}_b$ , the SRRS score corresponding to the respective bounding box prediction can be calculated as:

$$\text{SRRS}(\tilde{\mathbf{b}}) = \frac{1}{N_{support}} \sum_{i=1}^{N_{support}} \text{IoU}(\mathbf{b}_i, \tilde{\mathbf{b}}) \cdot p_{b_i}. \quad (28)$$

Based on the SRRS scores, the pseudo-labels are further filtered with another threshold  $v$  to reduce the false positives. They also employ an additional regularization, termed in the paper as Adversarial Background Score Regularization (ABSR). This regularization utilizes the gradient reversal layer discussed in Sec. 3.1.1; however, it is only applied on the target domain images. Without this regularization, the detector model risks incorrect prediction with high confidence. The addition of ABSR restricts the classifier sub-network of SSD to produce such overconfident predictions by avoiding alignment to non-transferable background regions. Together with pseudo-labels mined through hard and weak mining phase and the regularization loss, the detector model is trained in an end-to-end fashion. The self-training

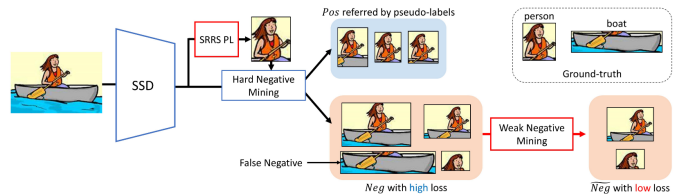


Fig. 18. Self-training one-stage detector [97] utilizes a pseudo-label based self-training to adapt SSD based detector model to target domain data. The pseudo-labels for training are generated following an eligibility criteria termed as, Supporting Region-based Reliable Score (SRRS) in the paper. Additionally, hard negative mining is used to get positive and negative examples which are later used to perform self-training of the SSD-based detector model.

is performed by minimizing  $\mathcal{L}_{det}^{ssd}$  as discussed in Sec. 2.1.2, where the loss is computed using the mined pseudo-labels.

Apart from these methods, there are many other interesting works available as pre-print [133], [134], [135] which utilize pseudo-label based self-training strategy for adapting detection model to target domain.

### 3.3 Image-to-image translation

As discussed earlier, the primary issue in unsupervised visual domain adaptation is that the target domain images are visually very distinct from the source domain. This causes a large gap in the feature space of the source-trained detector network, resulting in poor performance on the target domain. The method falling under adversarial feature learning (discussed in Sec. 3.1) and pseudo-label based self-training (discussed in Sec. 3.2) attempt to learn a feature representation that is more suited to perform detection on the target domain images. In contrast to feature alignment approaches, image-to-image translation-based methods to mitigate the domain gap at the input level. One of the most popular strategies used by image-to-image translation-based adaptation strategy is to use an unpaired image-to-image translation algorithm like Cycle-GAN [136], [137], UNpaired Image-to-image Translation [138], Multi-modal UNIT [139]. Methods that utilize the image-to-image translation-based approach often utilize additional strategies like self-training or adversarial feature learning to further boost the target detection performance. Image-to-image translation helps in bridging the gap at the input level while feature level adversarial alignment ensures a shared feature space for the target and source domain. Zhang *et al.* [99] specifically utilize Cycle-GAN to learn a mapping function between source and target domain images, in a method termed as Cycle-consistent Adaptive Faster-RCNN (CA-FRCNN) and illustrated in Fig. 19. Their approach extends the DA-Faster approach proposed in [25] (discussed in Sec. 3.1.2) by adding an image-to-image translation at the input level while adding other losses like gradient reversal at an instance and image-level intact. As shown in their experiments [99], the addition of an image-translation module further enhances the performance.

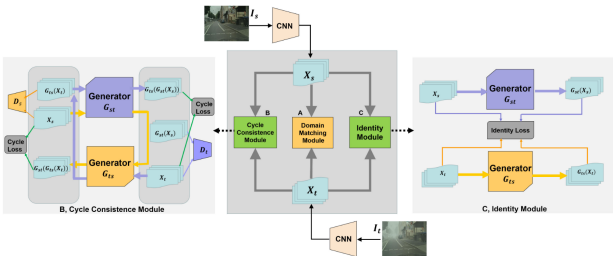


Fig. 19. Zhang *et al.* [99] extended the domain adaptive Faster-RCNN proposed in [25] with the addition of Cycle-GAN based image-to-image translation module at the input-level. The image-to-image translation module training is shown in the figure and domain matching module corresponds to the feature alignment strategy proposed in [25].

Hsu *et al.* [64] proposed a progressive adaptation strategy where they follow a similar strategy of using image-level translation. As shown in Fig. 20, a Cycle-GAN based image-to-image translation module is utilized to translate

target images into source-like images. Then with the help of the gradient reversal layer, the residual domain gap is further reduced in the feature space. Furthermore, with progressive adaptation [64], they showed that applying gradient reversal at only image-level is sufficient for adaptation rather than applying both image-level and instance-level losses. Arruda *et al.* [106] utilized the image translation module to specifically address the domain gap between daylight (source domain) and night-time (target domain) data, as shown in Fig. 21. Once the image translation module is ready, it is used to translate all daylight data to create a fake night-time dataset. Since ground-truth annotations are available for the daylight images, they can be used to supervise the detector model with fake night-time images as input. Such training will help learn a detector model that is better suited for the night-time scenario as compared to a model trained with fully supervised daylight images.

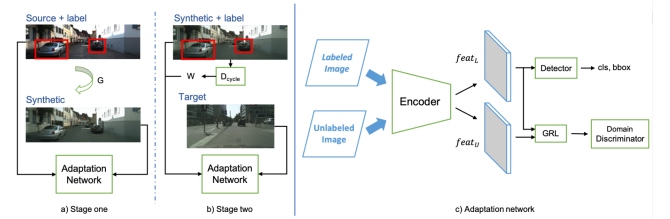


Fig. 20. Hsu *et al.* [64] proposed a progressive adaptation strategy that tries to reduce the gap between source and target domain at both input-level and feature-level. Specifically, before the detector training for the target domain, Hsu *et al.* [64] trains a Cycle-GAN model that can map source images to the target domain and vice versa. The Cycle-GAN model is then applied to target domain images to get translated images. Additionally, the gradient reversal layer-based learning strategy is used that further cuts the gap in the feature domain between the source domain and translated target domain, as shown in the figure.

Chen *et al.* [100] combined several of the previously proposed strategies along with image-to-image translation ideas in their approach termed as Harmonizing Transferability Calibration Network (HTCN). The proposed approach is illustrated in Fig. 22. First, a Cycle-GAN based model learns a mapping between the source and the target domain. HTCN utilizes this image-translation module to create source-like target images and target-like source images, also termed as “interpolated images”. The “interpolated images” helps HTCN fill in the distribution gap between domains and consequently reduce the source-bias of the decision boundaries. Additionally, HTCN employs local alignment loss similar to the one used in [57] for strong local alignment (see discussed in Sec. 3.1). Unlike other approaches based on image-to-image translation idea that directly utilize feature map outputs of a network, HTCN uses output of the local discrimination network to weigh the feature map of “interpolated images” forwarded to the subsequent layers. This strategy is termed as Importance Weighted Adversarial Training with Interpolated images (IWAT-I) and helps avoid negative transfer while promoting positive transfer by appropriately weighting the feature maps. The key motivation behind IWAT-I is that not all “interpolated images” have equal transferability and hence appropriate weights depending on the individual images are needed that can highlight the transferable regions while suppressing the noise in the feature map. This strategy closely follows the weighted

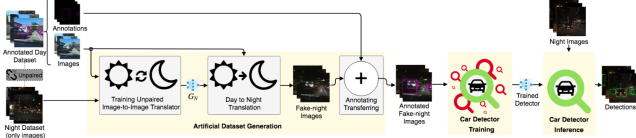


Fig. 21. Arruda *et al.* [106] utilized image-to-image translation based strategy to adapt detector model from daylight scenes to nighttime. As illustrated in the figure, a Cycle-GAN based module is used to translate daylight images to a fake nighttime one. The daylight images have ground-truth annotations and can train the detector with fake-nighttime in a supervised fashion. This results in a better adapted detector model for the nighttime scenario.

adversarial training discussed in Sec. 3.1.3. Furthermore, HCTN employs two additional discriminators to perform gradient reversal-based alignment of both masked features and global features at the output of the detection backbone network. The adversarial loss utilized in these two cases follows the formulation used in [57] for global alignment; however, HCTN uses binary cross-entropy loss as compared to focal loss used in [57]. Lastly, they perform context aware instance-level feature alignment that concatenates the features learned in the previous three discriminators with RoI-pooled features of detection network utilizing the formulation proposed in [25] for instance-level adaptation. This strategy is termed as Context-aware Instance-Level Alignment (CILA). By combining all of these losses, HCTN captures the transferable regions in the source, target and “interpolated image” domains to fully utilize the useful information for adapting images translated by Cycle-GAN. Shen *et al.* [121] utilize MUNIT [139] based Image-to-image translation to extend the method proposed in [58] (discussed in Sec. 3.1.2). Specifically, MUNIT-based translation module is learned by learning to map both whole images (image-level translation) and object regions (instance-level translation) between source and target domains. Furthermore, Shen *et al.* [121] also learns a style code bank to obtain an object, background and global style vectors containing rich spatial information to aid the translation network. In addition to the MUNIT-based image translation module, discriminator networks are applied at multiple layers to perform gradient reversal-based adversarial training. Further, another set of domain classification networks are trained without any gradient reversal layer. The features learned by these domain classifiers are concatenated with RoI-pooled features to improve the classification. Shen *et al.* [121] also noted that detaching these domain classification networks to restrict the gradient flow from the main detection pipeline further improves the detection performance. Other interesting works [65], [140], [141], [77] utilizing image-to-image translation strategy can be found as pre-print.

### 3.4 Domain randomization

In the case of image-to-image translation based methods, the primary focus is to learn a mapping between source and target domain and reducing the domain gap at the input-level. In case a perfect mapping function is available that can translate target domain images to source and vice versa, the detector model can perfectly adapt to target domain without requiring any annotations. However, in practice, such mapping function are not necessarily accurate and

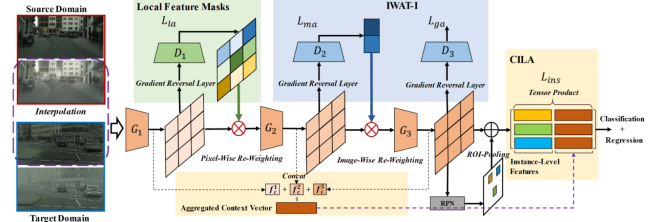


Fig. 22. Chen *et al.* [100] proposed a Cycle-GAN [136] based pre-processing step is introduced to cut the style gap between source and target domain images. Both source and target domain images are translated into “interpolated domain” to fill-in the input-level domain gap. Furthermore, they proposed harmonizing and transferability based adaptation of object detector that adopts multiple previously proposed strategies such as image-level and instance-level gradient reversal training similar to [25], context vectors from discriminators [57], and multiple image-level feature alignment. Here, IWAT-I denotes Importance Weighted Adversarial Training with Interpolated inputs, CILA denotes Context-aware Instance Alignment module,  $L_{la}$  denotes pixel-wise adaptation loss,  $L_{ma}$  and  $L_{ga}$  denote image-wise masked and global adaptation loss.

hence even after image translation, there might still exist a domain gap between original images and translated images, e.g., source images and translated target images or vice versa. Due to this, most image-to-image translation-based approaches utilize an additional gradient reversal based feature alignment or pseudo-label based self-training. To overcome this inability of learning accurate mapping function between source and target domain, domain randomization applies several transformations on the input image to derive multiple new domains. The goal for detection model is to then consider these new domains during training and learn feature representations that are invariant to all the domains including source and target. This strategy enforces strong constraints on feature representations of detection model as the network has to find features within images that remain invariant across a wide variety of image transformations. However, to achieve this, additional strategies such as adversarial feature learning or self-training are required.

In summary, image-to-image translation based methods aim to learn intermediate stages where input-level domain gap is reduced to improve model performance. Whereas, domain randomization synthetically generates domains that consist of multiple distinct styles (including source and target domain style) to ensure that the detection models learn features that are useful for detection regardless of style of the input images. Let us consider a source domain dataset  $\mathcal{S} = \{X_s^i, Y_s^i\}_{i=1}^{N_s}$  and target domain dataset  $\mathcal{T} = \{X_t^i\}_{i=1}^{N_t}$ . Domain randomization creates multiple new datasets derived from the given source domain dataset that consist of multiple distinct styles. Let us denote these source-derived stylized dataset as,  $\mathcal{S}^{style_m} = \{X_{s_m}^i, Y_s^i\}_{i=1}^{N_s}$ , where  $m \in \{1, \dots, M\}$  denotes  $m^{th}$  style and  $M$  denotes total number of unique styles. The stylization process ensures that the contents of the image are not changed and hence the object category and location information is preserved through the stylization process. All the stylized datasets are derived from source domain and hence, the  $i^{th}$  stylized image  $X_{s_m}^i$  can share the ground-truth annotations of corresponding  $i^{th}$  image from the original source dataset  $X_s^i$ , for any  $m \in \{1, \dots, M\}$  and  $i \in \{1, \dots, N_s\}$ .



Kim *et al.* [59] utilized domain randomization to adapt a Faster-RCNN based detection model. As illustrated in Fig. 23, the method has three components: domain diversification module, detection model, and multi-domain discriminator. The domain diversification module is based on generative adversarial networks [125] and is tasked to take in the source domain images and shift the domain to derive a diverse set of visually distinct domains. Further, they enforce certain constraints on the output of domain diversification networks such as reconstruction, color preservation and cycle consistency. This ensures that the synthetically generated domains do not destroy the contents of the image that might negatively impact the model performance. Subsequently, the detection model is trained on these synthetically generated domains data along with the source and target domain data. Supervised detection loss,  $\mathcal{L}_{det}^{frcnn}$ , is applied on the source and synthetic domains to train the detector model. To ensure that the base network of detection model learns domain invariant feature representations, a multi-domain discriminator network with gradient reversal layer is employed at the end of the base network. Typically, domain discriminators are tasked to perform binary classification to identify any image as either from source or target. However, the multi-domain discriminator used in this work is tasked to perform multi-class classification to identify whether the feature representations belong to source, target, or one of the synthetically generated domains. Hence, unlike binary cross-entropy loss that is used in [25], the multi-class cross-entropy loss is used for discriminator network. Here, the domain label for any  $i^{th}$  input image are given as  $y_d \in \{0, \dots, M+1\}$ , where, 0 and  $M+1$  denote source and target domain, and rest of the labels denote each synthetic domains generated by domain diversification module.

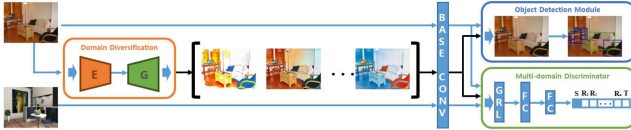


Fig. 23. Overview of the Diversify and Match methodology proposed by Kim *et al.* [59]. The domain diversification module creates stylized version of source images with multiple distinct styles which are used to train the detection model, in addition to passing the source and target domain images. The detection model is then trained to produce style invariant features with the help of a multi-style domain discriminator and gradient reversal based training. Unlike previous methods which use binary classification, the multi-style domain discriminator performs multi-class classification into source-style, target-style and pre-defined styles created by domain diversification module.

Rodriguez *et al.* [101] proposed a domain randomization approach for adapting SSD based one-stage object detectors. The key idea, illustrated in Fig. 24, is to utilize the style transfer network proposed by [142] to create a source-derived dataset with multiple distinct predefined styles. Multiple stylized version of the source-derived dataset is created with annotations borrowed from the source domain dataset. Furthermore, the source-trained detector model is evaluated on the target domain images to extract pseudo-labels, which are then used for self-training. In order to extract robust pseudo-labels for effective self-training, a positive and negative threshold is utilized to extract high-quality positive and negative examples. The base network

of the detector model, denoted as  $F_b$ , is encouraged to learn domain invariant features through feature consistency loss that minimizes the  $L2$ -norm between feature representations extracted from source data and stylized source data. This loss, termed as feature consistency loss, is defined as:

$$\mathcal{L}_{const}^{feat} = \frac{1}{N_s} \sum_{i=1}^{N_s} \sum_{m=1}^M \|F_b(X_s^i) - F_b(X_{s_m}^i)\|_2^2. \quad (29)$$

Let us denote source domain as  $\mathcal{S} = \{X_s^i, Y_s^i\}_{i=1}^{N_s}$ , source-derived  $m^{th}$ -style dataset as  $\mathcal{S}^{style_m} = \{X_{s_m}^i, Y_{s_m}^i\}_{i=1}^{N_s}$  with  $m \in \{1, \dots, M\}$ , and target domain dataset as  $\mathcal{T} = \{X_t^i, \tilde{Y}_t^i\}_{i=1}^{N_t}$ . Here,  $M$  is the pre-defined number of styles used to create distinct source-derived dataset,  $Y_s$  and  $\tilde{Y}_t$  denote ground-truth source annotations and pseudo-labels for an arbitrary source/ $m^{th}$ -style and target image, respectively. The final training loss for the object detector can be described as:

$$\begin{aligned} \mathcal{L}_{det}^{final} = & \lambda_s \sum_{i=1}^{N_s} \mathcal{L}_{det}^{ssd}(X_s^i, Y_s^i) \\ & + \sum_{m=1}^M \lambda_{s_m} \sum_{i=1}^{N_s} \mathcal{L}_{det}^{ssd}(X_{s_m}^i, Y_{s_m}^i) \\ & + \lambda_t \sum_{i=1}^{N_t} \mathcal{L}_{det}^{ssd}(X_t^i, \tilde{Y}_t^i) + \lambda_c \mathcal{L}_{const}^{feat}, \end{aligned} \quad (30)$$

where  $\lambda_s$ ,  $\{\lambda_{s_m}\}_{m=1}^M$ ,  $\lambda_t$ , and  $\lambda_c$  are trade-off parameters for source supervised loss, stylized-source supervised loss, pseudo-label training loss on target domain, and feature consistency loss, respectively.

Apart from these methods, there is an interesting work available as pre-print [143], which also utilizes domain randomization approach for adaptation of object detectors.

### 3.5 Mean teacher training

Knowledge distillation has been demonstrated to be effective for exploiting unlabeled data in transfer learning [144], [145], [146], semi-supervised learning [147], [148], [149], and domain adaptation [150], [151], [152], [153]. Most of the domain adaptation work that utilize student-teacher training strategy has considered only the task of image classification. Their success has inspired many works which employ student-teacher training framework to perform unsupervised domain adaptation of object detector models. A recent and notable work based on this strategy is that of

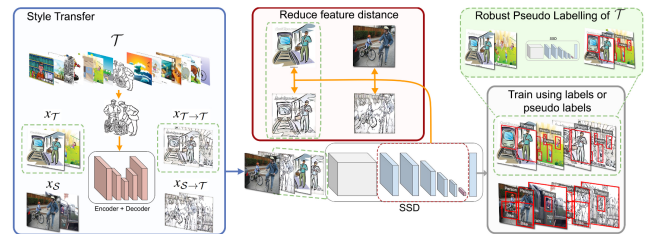


Fig. 24. Rodriguez *et al.* [101] proposed to adapt SSD based detector model by training with synthetically generated dataset using image style transfer. As shown in the figure, the stylization is applied to source domain images. Subsequently, supervised source, supervised stylized-source domain, and pseudo-label supervision from the target domain are used to train the detector model to learn style/domain invariant features for object detection in target domain data.



the unbiased mean-teacher strategy proposed by Deng *et al.* [68], that specifically utilizes mean-teacher framework [146] training for adapting object detector to the target domain. As it can be observed from Fig. 25, Deng *et al.* [68] combine multiple strategies like image-to-image translation (discussed in Sec. 3.3) and adversarial feature learning (discussed in Sec. 3.1) with mean-teacher framework. First, the method trains a Cycle-GAN module to learn image mapping between source and target domain images, which is then used to create a source-like target and target-like source images. The student model pipeline is trained with the original source domain, target domain and target-like source images. Since both source and target-like source images are fully labeled, they can be used for training using the supervised detection loss, shown in Fig. 25 as source detection loss and target-like detection loss. Training the student pipeline with source and target domain images helps mitigate bias in the student model. Target-like source image training encourages student models to be more favorable towards target domain images. The model predictions of source-like target images are matched with model predictions of original target domain images to perform knowledge distillation. Further, the teacher parameters are updated with Exponential Moving Average as:

$$\Theta_{\mathcal{F}^{tch}}^i = \alpha \Theta_{\mathcal{F}^{tch}}^{i-1} + (1 - \alpha) \Theta_{\mathcal{F}^{stu}}^i, \quad (31)$$

where  $\Theta_{\mathcal{F}^{tch}}^i$  and  $\Theta_{\mathcal{F}^{stu}}^i$  are network parameters for teacher and student model respectively,  $\alpha$  denotes smoothing coefficient hyper-parameter that can be used to control the teacher updates, the super-script  $i$  and  $i - 1$  denote the indices for the current and previous training iterations respectively. To further decrease the domain gap in the feature space of the student model, gradient reversal-based adversarial training involving strong local and weak global feature alignment [57]. Together with distillation, parameter updates supervised detection loss and adversarial feature training; the entire training is performed in an end-to-end fashion. Other

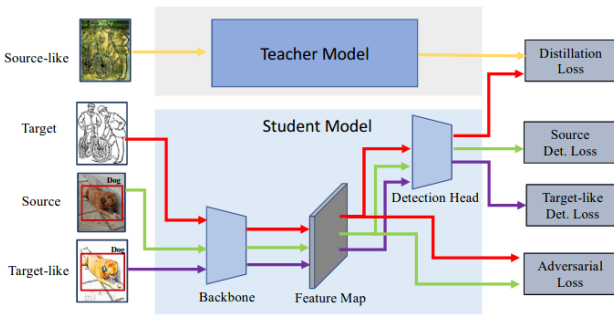


Fig. 25. Deng *et al.* [68] devised a strategy for cross-domain detection adaptation based on mean-teacher framework. As illustrated in the figure, it has three main components. First: image-to-image translation-based input processing on translating source into target-like images and vice versa. Second: distillation loss to transfer knowledge between teacher and student network training pipelines. Third: adversarial loss inspired from [57] to further reduce domain gap in feature space.

works that utilize the student-teacher framework include Liu *et al.* [149] which extended the mean-teacher framework to train an object detector in a semi-supervised fashion by exploiting large unlabeled data to improve the overall performance. Though Liu *et al.* focuses on semi-supervised detection with labeled and unlabeled data drawn from the

similar distributions, it can be easily adopted to perform unsupervised domain adaptive training of object detectors. Cai *et al.* [96] proposed cross-domain detection by utilizing a mean-teacher framework that relies on learning object relations through graph structures. Another interesting work [154], available as a pre-print, utilizes mean-teacher framework is also available as pre-print.

### 3.6 Graph-reasoning

Graph-reasoning is a widely popular technique in the computer vision community for visual recognition [155], [156], video understanding [157], [158] etc. It is also utilized extensively for unsupervised adaptation of classification networks [159], [160]. The ability of graphs to model inter-image and intra-image relationships of underlying object categories can be extremely useful in the case of adapting object detectors. Furthermore, modeling graph structures is mutually exclusive and hence can be combined with other existing adaptation strategies. Motivated by these benefits, Cai *et al.* [96] proposed a domain adaptive object detection approach, termed as Mean-Teacher Object Relations (MTOR), that utilizes graph-reasoning based strategy combined with student-teacher training. As illustrated in the Fig. 26, MTOR consists of four major components: 1) Regional-Level Consistency (RLC), 2) Inter-Graph Consistency (InterGC), 3) Intra-Graph Consistency (IntraGC), 4) Mean-Teacher Training (MTT). Here, the first three components are devoted to graph structures that model relationship between the objects within an image. Given a source domain dataset,  $\mathcal{S} = \{X_s^i, Y_s^i\}_{i=1}^{N_s}$ , target domain dataset,  $\mathcal{T} = \{X_t^i\}_{i=1}^{N_t}$ , detection model,  $F$ , the base network,  $F_{base}$ , teacher network,  $F^t$  and student network  $F^s$ , a graph,  $\mathcal{G}_{X_t} = \{\mathcal{V}_{X_t}, \mathcal{E}_{X_t}\}$  is constructed. Here  $\mathcal{V}_{X_t}$  are vertices denoting set of region proposal predictions by detection network and  $\mathcal{E}_{X_t}$  ( $|\mathcal{V}_{X_t}| \times |\mathcal{V}_{X_t}|$ ) is an affinity matrix of the corresponding the graph. Note that  $\mathcal{G}_{X_t}^t$  and  $\mathcal{G}_{X_t}^s$  correspond to the graphs constructed on target domain images using teacher and student models, respectively. Each vertex in the graph corresponds to a region proposal and is assigned a probability vector denoting the probability of that region belonging to one of the  $K_{X_t}$  categories. Before selecting the region proposal, they are filtered by eliminating all proposals having a maximum probability value falling below a threshold. For any arbitrary target image  $X_t$ , consider that there are  $N_b$  such region proposals and for  $i^{th}$  region proposal, the corresponding probability vector is  $\mathbf{p}_c^i$ , the region level consistency loss is defined as:

$$\mathcal{L}_t^{RCL} = \frac{1}{N_b} \sum_{i=1}^{N_b} \|\mathbf{p}_c^{t_i} - \mathbf{p}_c^{s_i}\|_2^2, \quad (32)$$

where  $\mathbf{p}_c^{t_i}$  and  $\mathbf{p}_c^{s_i}$  denote probability vectors extracted from teacher and student model and correspond to teacher graph  $\mathcal{G}_{X_t}^t$  and student graph  $\mathcal{G}_{X_t}^s$ , respectively.

To compute the inter-graph consistency, InterGC, the affinity matrix of both student and teacher graph is calculated using cosine similarity. If we denote teacher graph as  $\mathcal{G}_{X_t}^t = \{\mathcal{V}_{X_t}^t, \mathcal{E}_{X_t}^t\}$  for any arbitrary target image  $X_t$  having  $N_b$  region proposals, where any  $i^{th}$  region proposal consists of the corresponding RoI-pooled feature  $\mathbf{f}_t^i$  extracted from the teacher model. The affinity matrix is calculated using

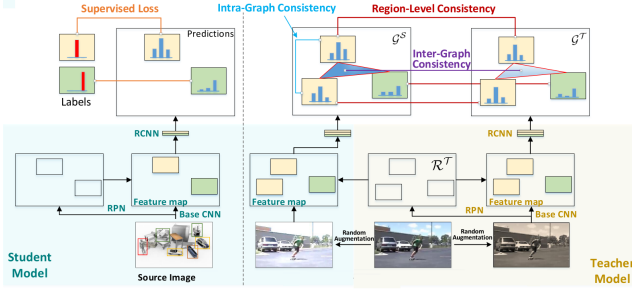


Fig. 26. Cai *et al.* [96] explores graph reasoning on Faster-RCNN based detector model to control the information flow between source and target domain pipelines in a mean teacher framework. The source domain pipeline is used as student and target domain pipeline is used as teacher. Specifically, graph reasoning enforces (1) Regional-level consistency to align regional level predictions between teacher and student, (2) Inter-graph consistency to match the graph structures between source and target pipelines, and (3) Intra-graph consistency that promotes the similarity between same categories within graph constructed from source domain model.

the feature similarity between any two proposal features in a given image and is defined as:

$$\mathcal{E}_{X_t}^t(i, j) = \frac{\mathbf{f}_t^i \cdot \mathbf{f}_t^j}{\|\mathbf{f}_t^i\|_2 \|\mathbf{f}_t^j\|_2}. \quad (33)$$

Similarly, Eq. 33 can be used to calculate affinity matrix for student model  $\mathcal{E}_{X_t}^s = \{\mathcal{V}_{X_t}^s, \mathcal{E}_{X_t}^s\}$  for the same target domain image. Once we have both student and teacher model affinity matrix for a given target domain image, the inter-graph consistency loss is calculated as follows:

$$\mathcal{L}_t^{InterGC} = \frac{1}{N_b} \cdot \|\mathcal{E}_{X_t}^t - \mathcal{E}_{X_t}^s\|_2^2. \quad (34)$$

This inter-graph consistency enforces alignment of the graph structure between student and teacher pipeline for a target image. Within each image, there can be multiple categories that are semantically related. Intra-graph consistency reinforces the similarity between proposals of the same category for student graphs with supervision from the teacher. Specifically, the teacher model is used to predict whether two region proposals contain an object of the same category. If  $i^{th}$  and  $j^{th}$  proposal contain same category, label  $y_{ij}^{graph}$  is assigned a value of 1 or 0 otherwise. Given affinity matrix of student model for an arbitrary target image,  $\mathcal{E}_{X_t}^s$ ,

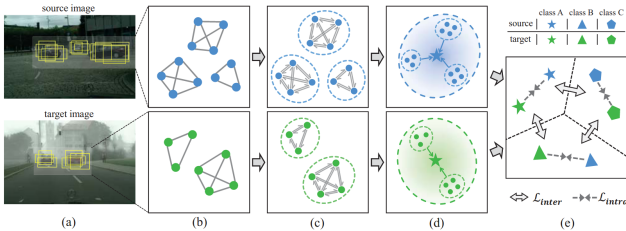


Fig. 27. The figure provides an overview of the method proposed by Xu *et al.* [95]. (a) Source and target domain images with respective object region proposals. (b) Utilizing the proposals to create a relation graph for both source and target image. (c) Information propagation between proposals of the same category to get accurate instance-level feature representations. (d) Merging category-wise information with confidence weighting to create a prototype for the respective object class. (e) Category-level adversarial feature alignment by enhancing intra-class compactness  $\mathcal{L}_{intra}$  and inter-class separability  $\mathcal{L}_{inter}$ .

the intra-graph consistency loss is defined as:

$$\mathcal{L}_t^{IntraGC} = \frac{\sum_{i,j} y_{ij}^{graph} \cdot (1 - \mathcal{E}_{X_t}^s(i, j))}{\max(1, \sum_{i,j} y_{ij}^{graph})}. \quad (35)$$

In addition to these graph-based consistency losses, the student network is trained with supervised detection loss. The teacher network parameter updates follow mean-teacher formulation discussed in Sec. 3.5.

Recently, Xu *et al.* [95] proposed a graph-based adaptation approach where they construct a relational graph between object proposals to compute category prototypes that help alignment between source and target domain. The training strategy involves graph construction and alignment as illustrated in Fig. 27. The proposed method is based on the Faster-RCNN detection framework and utilizes its two-stage detector training to compute relational graphs. The relational graphs are used to compute category-wise prototypes, which are then matched to enforce alignment between source and target domain at each stage of the detector model. The overall alignment strategy is termed as Graph-induced Prototype Alignment (GPA) in the paper [95]. Let us consider that, for an arbitrary image, the detection model produces  $N_b$  number of region proposals with varying probabilities of containing an object of any given category in the dataset. Based on these proposals, a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  is constructed with  $\mathcal{V}$  as nodes and  $\mathcal{E}$  denoting adjacency matrix/edge-values assigned to the graph. Each node represents a region proposal predicted by the detection model and the adjacency matrix values are computed by calculating Intersection over Union (IoU) between pairs of region proposals. Specifically, let's consider two region proposals for an arbitrary input image denoted as bounding boxes  $\mathbf{b}_i$  and  $\mathbf{b}_j$ . The edge value between them is computed using the following equation:

$$\mathcal{E}_{ij} = \text{IoU}(\mathbf{b}_i, \mathbf{b}_j) = \frac{\mathbf{b}_i \cap \mathbf{b}_j}{\mathbf{b}_i \cup \mathbf{b}_j}, \quad (36)$$

where  $\cap$  denotes intersection operation and  $\cup$  denotes union operation. This creates an adjacency matrix of size  $N_b \times N_b$ . Once the relation graph is constructed, the adjacency matrix is used to aggregate the RoI-pooled features and their respective class probabilities for each region proposals. Let us denote each RoI-pooled feature having  $d$  dimension as  $\mathbf{f} \in \mathbb{R}^{1 \times d}$  and respective probability vector as  $\mathbf{p} \in [0, 1]^{1 \times K}$ . We can stack these features and probability values to create feature matrix  $\mathbf{F} \in \mathbb{R}^{N_b \times d}$  and probability matrix  $\mathbf{P} \in [0, 1]^{N_b \times d}$ . The aggregated features  $\tilde{\mathbf{F}}$  and probabilities  $\tilde{\mathbf{P}}$  can be represented as:

$$\begin{aligned} \tilde{\mathbf{F}} &= \mathbf{D}^{-\frac{1}{2}} \mathcal{E} \mathbf{D}^{-\frac{1}{2}} \mathbf{F}, \\ \tilde{\mathbf{P}} &= \mathbf{D}^{-\frac{1}{2}} \mathcal{E} \mathbf{D}^{-\frac{1}{2}} \mathbf{P}, \end{aligned} \quad (37)$$

where matrix  $\mathbf{D}$  is a diagonal matrix with diagonal elements given as  $\mathbf{D}_{ii} = \sum_j \mathcal{E}_{ij}$ . This aggregation updates the feature and probability vectors to represent the instance-level information accurately. Furthermore, the updated features and probabilities are used to calculate category-wise prototypes. The prototype feature for  $k^{th}$  category can be computed as:

$$\mathbf{v}_k = \frac{\sum_{i=1}^{N_b} \tilde{\mathbf{P}}_{ik} \cdot \tilde{\mathbf{F}}_i^T}{\sum_{i=1}^{N_b} \tilde{\mathbf{P}}_{ik}}, \quad (38)$$

where  $\mathbf{v}_k \in \mathbb{R}^{1 \times d}$  denotes the prototype corresponding to  $k^{th}$  category and are calculated for both source and target domain. The alignment is performed by matching source and target prototypes by minimizing the distance between prototypes of the same category while maximizing the distance between different categories. Additionally, the same category prototypes are matched between source and target domain by minimizing respective category prototypes.

## 4 EVALUATION PROTOCOLS

In this section, we discuss the details of the evaluation protocol used in evaluating the performance of the various domain adaptive detection approaches discussed earlier. First, we provide details of various benchmark datasets that are used for the evaluation purpose, followed by a discussion of various adaptation settings such as synthetic-to-real, clean-to-adverse weather, etc. Next, we define the metrics used for evaluation and comparison. Finally, we provide a detailed discussion of the results of various methods.

### 4.1 Datasets

A variety of datasets have been used for evaluating the domain adaptive object detection approaches. Most of these datasets have been extensively used for evaluating object detection approaches, and hence, can naturally be used for the purpose of evaluating domain adaptive methods by simply defining different sets of data as source and target domain. These datasets can be classified into the following categories: (1) General objects [19], [94], (2) Self-driving [161], [162], [163], [164], (3) Face detection [165], [166], (4) Weather degradation [163], [164], [165], [167], and (5) Synthetic datasets [168]. Fig. 28 illustrates an overview of the different categories of datasets and Table 4 provides summary of various datasets train and test samples. In what follows, we discuss the various datasets in detail.

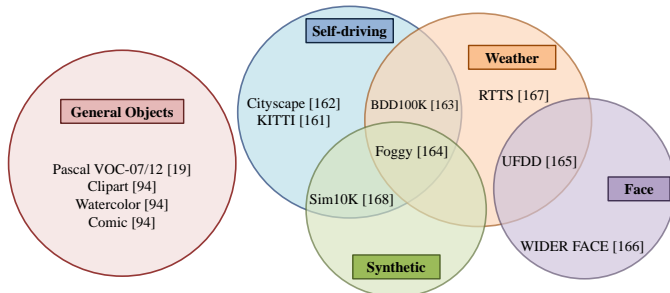


Fig. 28. In unsupervised domain adaptive object detection, the benchmark datasets can be clustered into general objects, self-driving, face, weather degradation and synthetic categories. The dataset can be from purely one category, or it can be a combination of multiple categories. From the figure, we can infer that the Cityscapes dataset is purely from the self-driving category. In contrast, the FoggyCityscapes dataset is a combination of weather degradation, synthetic and self-driving.

**Cityscapes.** Scene understanding of complex urban streets is an important problem statement for a wide range of applications. To address this issue, the Cityscapes dataset was released in 2016 by a group of researchers in Daimler and TU Dresden [162]. Cityscapes has 8 categories: car, truck, motorcycle/bike, train, bus, rider, and person.

The dataset is collected in 50 cities and covers a variety of seasons (spring, summer, fall). The Cityscapes dataset contains 2975 images for training and 500 images for testing.

**FoggyCityscapes.** Analyzing scenes of urban streets under adverse weather is a challenging problem. To consider such scenarios, the FoggyCityscapes was introduced [164], by applying a fog filter over the Cityscapes dataset. The FoggyCityscapes has the same 8 categories of Cityscapes dataset. The FoggyCityscapes dataset contains 2975 images for training and 500 images for testing.

**Sim10K.** Considering that the collection of real-world data is time-consuming and tedious, advancements in computer graphics have been exploited to generate photo-realistic data which can be easily rendered and annotated. This alternative to real-world data collection is inexpensive and efficient. In 2017, Sim10K [168], a synthetic dataset rendered by the gaming engine Grand Theft Auto was released. It has 10K images and only a car category having 58,701 car instances. The dataset provides annotations in the Pascal VOC format for ease of use.

**KITTI.** KITTI [161] is one of the most popular datasets for self-driving and mobile robotics. The dataset contains hours of videos of traffic scenarios recorded with various high-quality sensors like RGB, grayscale and depth sensors. The KITTI object detection dataset consists of 7,481 training images and 7,518 test images, comprising of 80,256 labeled objects which span over 8 categories: Car, Van, Truck, Pedestrian, Person sitting, Cyclist, Tram, Misc.

**Pascal VOC2007/2012.** PASCAL VOC [19] is a real-world dataset containing general objects. It was introduced for a variety of tasks such as large-scale classification, object detection, and segmentation. The PASCAL VOC image set has been updated multiple times between 2005 to 2012. Among these multiple versions, the most popular ones are 2007 and 2012 image sets with training and validation split containing a total of 15k images. The images in the dataset contain over 20 general object categories: airplane, bicycle, bird, boat, bottle, bus, car, cat, chair, cow, table, dog, horse, bike, person, plant, sheep, sofa, train, tv.

**Clipart, Watercolor, Comic.** Clipart [94], Watercolor [94] and Comic [94] datasets contain abstract, artistic and comical images, respectively. Understanding these abstract-type images will allow for the direct study of how a model infer high-level semantic information. Consequently, these datasets are extensively used in dissimilar domain adaptation problems. The Clipart contains a total of 1K images with 20 categories same as Pascal VOC. The watercolor and comic dataset contain 6 categories: bike, bird, car, cat, dog, person. Moreover, both datasets contain 1K training images and 1K testing images, respectively.

**BDD100K.** The BDD100K dataset [163] is an autonomous driving dataset contains large diversity in terms of scenes, geographical regions. It is the largest (till-date) driving video dataset with 100K videos containing scenes from New York, Berkeley, San Francisco and Bay Area. The



TABLE 2

Quantitative comparison (mAP) of existing domain adaptive object detection methods. CS: Cityscapes, F: FoggyCityscapes, S10K: Sim 10K, KI: KITTI, VOC: Pascal VOC, Clip: Clipart, WC: Watercolor, BDD: BDD100K, SS: Shorter side, LR: Learning Rate, Iter: Iteration, FRCNN: Faster-RCNN. Red and blue color indicate best and second-best methods in respective adaptation scenario in terms of mAP. † denotes that the corresponding method uses ResNet-50 backbone in the detection model, rest of the methods utilize VGG16 backbone.

Method	CS → F	S10K → CS	KI → CS	CS → KI	VOC → Clip	VOC → WC	VOC → Comic	CS → BDD	Framework	LR	Additional Comments
Faster-RCNN in the wild [25]	27.6	38.9	38.5	64.1	-	-	-	-	FRCNN	0.001	SS: 500
Diversify and match [59]	34.6	-	-	-	41.8	52.0	34.5	-	FRCNN	0.001	SS: 600
Multi-adversarial adaptation [60]	34.0	41.1	41.0	72.1	-	-	-	-	FRCNN	0.001	-
Progressive adaptation [64]	36.9	-	43.9	-	-	-	-	24.3	FRCNN	0.001	-
Strong weak distribution alignment [57]	34.3	40.7	-	-	38.1	53.3	-	-	FRCNN	0.001	SS: 600
Mean teacher with object relations [96]	35.1	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Large-scale instance-level [121]	40.3	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Curriculum self-paced learning† [123]	-	47.6	43.8	-	37.8	-	-	-	FRCNN	0.001	SS: 600
Image-instance full alignment network [66]	36.2	47.1	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Forward-Backward cyclic adaptation [113]	36.7	42.7	-	-	38.5	53.6	-	-	FRCNN	0.001	SS: 600, Epochs: 10
Categorical regularization [69]	37.4	-	-	-	38.3	-	-	26.9	FRCNN	0.001	SS: 600
Collaborative training [98]	35.9	44.5	43.6	-	-	-	-	-	FRCNN	0.001	SS: 600
Robust learning from noisy labels [62]	36.4	42.5	42.9	77.6	-	-	-	-	FRCNN	0.001	SS: 600
Unbiased mean-teacher [68]	41.4	43.1	-	-	42.7	56.9	-	-	FRCNN	0.001	SS: 600
Prior-based DA Object Detection [72]	39.3	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Coarse-to-fine adaptation [107]	38.6	43.8	-	41.0	-	-	-	-	FRCNN	0.001	SS: 600
Dual Multi-Label Prediction [111]	38.8	-	-	-	-	56.0	33.5	-	FRCNN	0.001	SS: 600, Epochs: 10
Instance-invariant progressive disentanglement [120]	36.4	-	-	-	42.1	56.9	-	-	FRCNN	0.001	SS: 600
Graph-induced prototype alignment† [95]	39.5	47.6	47.9	-	-	-	-	-	FRCNN	0.001	SS: 600, Epochs: 20
Harmonizing transferability [100]	39.8	42.5	-	-	40.3	-	-	-	FRCNN	0.001	SS: 600
Asymmetric Tri-way Faster-RCNN [67]	38.7	42.8	42.1	73.5	42.1	54.9	-	-	FRCNN	0.001	SS: 600
Selective cross-domain alignment [58]	33.8	43.0	42.5	-	-	-	-	-	FRCNN	0.0001	SS: 512, Epochs: 25
Spatial attention pyramid network [110]	40.9	44.9	43.4	75.2	42.2	55.2	-	-	FRCNN	0.00001	SS: 600, Iter: 90K
Cycle-consistent adaptation [99]	33.2	41.5	41.7	-	-	-	-	-	FRCNN	0.002	SS: 600
Multi-level adaptation [63]	36.0	42.8	-	-	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Multi-scale robust discrimination [70]	37.0	43.4	-	-	-	-	-	-	FRCNN	0.002	SS: 600, Iter: 90K
Domain invariant region proposal [114]	39.2	45.5	44.0	-	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Memory-guided category-wise adaptation [73]	41.8	44.8	43.0	75.5	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Conditional normalization network [109]	36.6	49.3	44.9	-	-	-	-	-	FRCNN	0.00625	SS: 512, Epochs: 12
Self-training for one-stage detector [97]	-	-	-	-	35.7	49.9	26.8	-	SSD	0.001	SS: 300, Iter: 120K
Implicit invariant one-stage network [117]	-	-	-	-	37.8	51.5	30.1	-	SSD	0.001	SS: 300
Object detection via style consistency [101]	34.3	-	-	-	44.8	57.3	39.4	-	SSD	0.00001	SS: 300, Epochs: 20
Every pixel matters [71]	36.0	49.0	43.2	-	-	-	-	-	FCOS	0.005	SS: 800, Epochs: 10
Adaptive transformer-based detector [169]	43.5	55.3	-	-	-	-	-	-	DETR	0.0002	Epochs: 50

TABLE 3

Quantitative comparison ( $\Delta$ mAP) of existing domain adaptive object detection methods. CS: Cityscapes, F: FoggyCityscapes, S10K: Sim 10K, KI: KITTI, VOC: Pascal VOC, Clip: Clipart, WC: Watercolor, BDD: BDD100K, SS: Shorter side, LR: Learning Rate, Iter: Iteration, FRCNN: Faster-RCNN. Red and blue color indicate best and second-best methods in respective adaptation scenario in terms of  $\Delta$ mAP. † denotes that the corresponding method uses ResNet-50 backbone in the detection model, rest of the methods utilize VGG16 backbone.

Method	CS → F	S10K → CS	KI → CS	CS → KI	VOC → Clip	VOC → WC	VOC → Comic	CS → BDD	Framework	LR	Additional Comments
Faster-RCNN in the wild [25]	8.80	8.85	8.30	10.6	-	-	-	-	FRCNN	0.001	SS: 500
Diversify and match [59]	16.7	-	-	-	16.9	12.2	13.1	-	FRCNN	0.001	SS: 600
Multi-adversarial adaptation [60]	15.2	11.0	10.8	18.6	-	-	-	-	FRCNN	0.001	-
Progressive adaptation [64]	17.3	-	15.1	-	-	-	-	3.50	FRCNN	0.001	-
Strong weak distribution alignment [57]	14.0	6.10	-	-	10.3	8.7	-	-	FRCNN	0.001	SS: 600
Mean teacher with object relations [96]	8.20	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Large-scale instance-level [121]	20.0	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Curriculum self-paced learning† [123]	-	16.9	12.3	-	11.7	-	-	-	FRCNN	0.001	SS: 600
Image-instance full alignment network [66]	15.2	12.0	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Forward-Backward cyclic adaptation [113]	17.9	11.5	-	-	10.7	9.00	-	-	FRCNN	0.001	SS: 600, Epochs: 10
Categorical regularization [69]	15.4	-	-	-	11.3	-	-	3.50	FRCNN	0.001	SS: 600
Collaborative training [98]	9.70	10.0	8.7	-	-	-	-	-	FRCNN	0.001	SS: 600
Robust learning from noisy labels [62]	4.5	11.5	11.8	21.4	-	-	-	-	FRCNN	0.001	SS: 600
Unbiased mean-teacher [68]	19.6	8.80	-	-	13.6	8.00	-	-	FRCNN	0.001	SS: 600
Prior-based DA Object Detection [72]	14.9	-	-	-	-	-	-	-	FRCNN	0.001	SS: 600
Coarse-to-fine adaptation [107]	17.8	8.80	-	7.60	-	-	-	-	FRCNN	0.001	SS: 600
Dual Multi-Label Prediction [111]	15.4	-	-	-	-	11.4	13.8	-	FRCNN	0.001	SS: 600, Epochs: 10
Instance-invariant progressive disentanglement [120]	13.7	-	-	-	14.3	12.3	-	-	FRCNN	0.001	SS: 600
Graph-induced prototype alignment† [95]	12.6	13.0	10.3	-	-	-	-	-	FRCNN	0.001	SS: 600, Epochs: 20
Harmonizing transferability [100]	19.5	7.90	-	-	12.5	-	-	-	FRCNN	0.001	SS: 600
Asymmetric Tri-way Faster-RCNN [67]	18.4	8.20	11.9	20.0	14.3	10.3	-	-	FRCNN	0.001	SS: 600
Selective cross-domain alignment [58]	7.60	9.10	5.10	-	-	-	-	-	FRCNN	0.0001	SS: 512, Epochs: 25
Spatial attention pyramid network [110]	17.6	10.3	13.2	21.7	14.4	10.6	-	-	FRCNN	0.00001	SS: 600, Iter: 90K
Cycle-consistent adaptation [99]	8.78	6.70	4.00	-	-	-	-	-	FRCNN	0.002	SS: 600
Multi-level adaptation [63]	13.2	8.50	-	-	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Multi-scale robust discrimination [70]	17.1	9.90	-	-	-	-	-	-	FRCNN	0.002	SS: 600, Iter: 90K
Domain invariant region proposal [114]	16.6	11.3	9.30	-	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Memory-guided category-wise adaptation [73]	17.4	10.5	12.8	22.0	-	-	-	-	FRCNN	0.002	SS: 600, Epochs: 10
Conditional normalization network [109]	10.5	15.0	7.80	-	-	-	-	-	FRCNN	0.00625	SS: 512, Epochs: 12
Self-training for one-stage detector [97]	-	-	-	-	9.00	2.80	4.90	-	SSD	0.001	SS: 300, Iter: 120K
Implicit invariant one-stage network [117]	-	-	-	-	11.1	4.4	8.20	-	SSD	0.001	SS: 300
Object detection via style consistency [101]	8.70	-	-	-	14.5	7.70	14.5	-	SSD	0.00001	SS: 300, Epochs: 20
Every pixel matters [71]	17.2	18.9	13.0	-	-	-	-	-	FCOS	0.005	SS: 800, Epochs: 10
Adaptive transformer-based detector [169]	9.5	4.80	-	-	-	-	-	-	DETR	0.0002	Epochs: 50



TABLE 4

Summary of various datasets used in domain adaptive object detection experiments. All the datasets except RTTS contain labeled training images and test images. The RTTS dataset contains unlabeled training images and labeled test images.

Dataset	Train		Test	
	Images	Catagories	Images	Catagories
Cityscapes	2975	8	500	8
VOC-2007	5011	20	5011	20
VOC-2012	11540	20	11540	20
SIM10K	10000	1	10000	1
KITTI	7481	1	7481	1
FoggyCityscapes	2975	8	500	8
Clipart	1000	20	1000	20
Watercolor	1000	6	1000	6
BDD100K	36728	10	5258	10
Comic	1000	6	1000	6
WIDER FACE	32000	1	32000	1
UFDD	442	1	442	1
RTTS	4807	5	4322	5

dataset contains 70k training images and 10k validation images. The authors ensure that the images encompass different types of weather, six different scenes, three separate times of the day, and 10 object categories and bounding box annotations. For the domain adaptive detection task specifically, a subset of the BDD100k dataset that contains images labeled as daytime is used. This subset contains 36,728 training and 5,258 validation images.

**WIDER FACE.** Face detection is an important application in modern computer vision. However, due to lack of variance in the face detection dataset, in 2015, Shuo [166] released the WIDER FACE, which is the largest benchmark dataset consisting of 32,000 images and 199K labeled faces. The authors ensure high degree of variability in scale, expression, illumination, pose, makeup and occlusion.

**UFDD.** Face detection under adverse conditions is very important, especially for video surveillance. Considering this, the UFDD dataset [165] consisting of images collected under a variety of weather conditions was released. The authors attempt to specifically capture the following conditions: rain, snow, haze, lens impediments, blur, illumination. The dataset contains a total of 6424 images.

**RTTS.** Realistic Single Image DEhazing (RESIDE) [167] is a large-scale dataset consisting of both real and synthetic hazy images. RTTS is an object detection dataset and is a subset of the RESIDE. It contains 4,807 un-annotated and 4,322 annotated real-world hazy images covering mostly traffic and driving scenarios. RTTS has total 5 categories, namely motorcycle/bike, person, bicycle, bus and car.

## 4.2 Adaptation scenarios

In this section, we describe the various adaptation scenarios and protocols followed by the existing approaches.

### 4.2.1 Adverse weather conditions

Stable detection performance in different weather conditions is important for safety-critical applications like self-driving cars. Weather conditions introduce image artifacts which can negatively impact the detection performance. FoggyCityscapes and Cityscapes can be utilized as target and source domains to evaluate the effectiveness of adaptation methods in adverse weather. Moreover, under the real hazy condition, one can extend the domain adaptation setting for Cityscapes to the RTTS dataset, used as the source and target domain, respectively. Similarly, for face detection, such weather conditions prove challenging when the face detector is trained with clean weather data. This adaptation scenario can be explored with the source domain as WIDER-Face and target domain as UFDD-Haze.

### 4.2.2 Synthetic data adaptation

Synthetic data offers an inexpensive alternative to real data collection as it is easier to collect, and with appropriate engineering, annotations can be made readily available for synthetic data with very little labor cost. In spite of the advancements in computer graphics, photo-realistic synthetic data generated using state-of-the-art rendering engines suffer from subtle image artifacts, which can result in sub-optimal performance on real-world data. The adaptation from Sim10K (source domain) to Cityscapes (target domain) is used in the literature to analyze this setting.

### 4.2.3 Cross-camera adaptation

Differences in the intrinsic and extrinsic camera properties like resolution, distortion, orientation, location result in images that capture the objects differently from each other in terms of quality, scale, and viewing angle. While the collected data can be real, these domain differences will potentially result in severe performance degradation. To evaluate the domain adaptation methods under cross-camera adaptation settings, most literature considers KITTI to Cityscapes and KITTI to Cityscapes.

### 4.2.4 Adaptation between dissimilar domains

In the shift from real-world to artistic images, underlying features, the texture of real-world images completely differ in artistic images. Hence, a lot of methods consider adaptation of dissimilar domains with three adaptation settings: PASCAL-VOC to Clipart, PASCAL-VOC to Watercolor and PASCAL-VOC to Comic. All these domain shifts are from real-world images to synthetic as well as artistic images.

### 4.2.5 Adaptation to Large-scale Dataset

The increasing availability of cheap and good quality cameras has made collecting large-scale datasets much easier. However, annotating them for the detection task is labor-intensive. Therefore, exploring the possibility of adopting from a smaller dataset to a larger dataset has significant real-world implications. To evaluate this adaptation setting, methods in most domain adaptation literature consider Cityscapes to BDD100k. In the BDD100K dataset, only the daylight subset is considered as the target domain and the Cityscapes dataset as the source domain.

### 4.3 Evaluation Metric

Typically, object detectors are evaluated using the average precision (AP) that was introduced in VOC2007 [19]. The average precision is computed for each category by calculating the area under the precision-recall curve. Precision and recall are defined as follows:

$$\begin{aligned} \text{Precision} &= \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \\ \text{Recall} &= \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \end{aligned} \quad (39)$$

Precision is a measure of how accurate are the predictions of an object detector. Recall identifies how many of positive predictions of the object detector are correctly predicted. In order to determine if the prediction of a detector matches with the ground truth bounding box, the intersection over union (IoU) measure is used (see Fig. 29 for details). This measure is defined as follows:

$$\text{IoU} = \frac{\text{Area of overlap}}{\text{Area of union}}. \quad (40)$$

The IoU measure is then used to compute if a particular

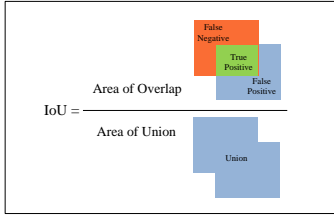


Fig. 29. Illustration of intersection over union (IoU).

prediction is a true positives or false positive given a pre-defined threshold, using the following equation:

$$\begin{aligned} TP_{ij} &= 1, \text{ if } \text{IoU}(t)_{ij} > \tau \\ FP_{ij} &= 1, \text{ if } \text{IoU}(t)_{ij} < \tau. \end{aligned} \quad (41)$$

Here,  $i$  is the index of the  $i$ -th image while  $j$  is the index of the  $j$ -th object and  $\tau$  is the threshold value. Hence, by applying a threshold over IoU, one can measure the precision and recall by identifying all the object proposals as a true positive or false positive. Furthermore, to compare the overall performance among all categories, the mean Average Precision (mAP) metric is used, which involves calculating the average of all categorywise AP.

$$\text{mAP} = \frac{\sum_{k=1}^K \text{AP}_k}{K}, \quad (42)$$

where  $K$  is the total number of category.

### 4.4 Discussion of results

In this section, we provide detailed comparison and discussion of the performance of various domain adaptive object detection methods. For comparison, we only consider adaptation scenarios used by two or more works in the literature. Most of the approaches conduct adaptation experiments that consist of various scenarios created by carefully selecting a source and target domain from the datasets discussed in the previous section. To indicate an adaptation scenario, we follow dataset1→dataset2 format, where dataset1 is used as a labeled source domain and dataset2 is used as an unlabeled target domain.

The performance comparison reported in Table 2 uses the absolute mAP performance as reported in the respective

papers. It is important to note that this comparison is not necessarily fair for the following reasons: Each method uses different version of the baseline detector. This is true even for the methods that use Faster-RCNN as the baseline detector. These methods use different hyperparameters for training the baseline. For example, Chen *et al.* [25] use Faster-RCNN with ROI pooling and images with a shorter side of 500 pixels, whereas Cai *et al.* [96] uses Faster-RCNN with ROI-Align and images with a shorter side of 600 pixels. Additionally, there are several other differences like the number of epochs, learning rate, etc across different approaches. Hence, to provide a fair, comprehensive comparison, we consider relative mAP improvement, denoted as  $\Delta\text{mAP}$ , which is the difference between the performance of a particular approach and that of the corresponding source-only baseline as reported in the respective paper. Though  $\Delta\text{mAP}$  would still not be a perfect comparison metric, we argue that it is better than comparing absolute mAP of methods that are trained with different hyperparameters. Since all methods in the literature train both source-only baseline and final method with same hyperparameters,  $\Delta\text{mAP}$  would obviate most of the influences on the performance caused by hyperparameters of the respective methods. The  $\Delta\text{mAP}$  metric for all methods is reported in Table 3.

Comparing these two tables provide some interesting insights into the performance trends over the years for any particular adaptation scenario. For example, consider the case of Cityscapes→BDD100K, where, from Table. 2, it seems that the forward-backward cyclic adaptation [113] approach improves over the progressive adaptation [64] strategy. However, considering the performance in terms of  $\Delta\text{mAP}$ , both these methods obtain similar performance improvements. This suggests that the Cityscapes→BDD100K scenario needs to be explored further by future works. The case of adverse weather adaptation case (Cityscapes→FoggyCityscapes), and has observed  $\sim 20$  mAP improvement since the first work [25]. From the absolute mAP comparison in Table. 2, the transformer-based adaptation [78] and the memory-guided adaptation [73] provide the best and second-best performance. Whereas from the  $\Delta\text{mAP}$  comparison, it can be observed that the unbiased mean-teacher [68] and harmonizing transferability [100] obtains the most improvement over their respective source-only baseline. A similar trend is observed in other adaptation scenarios as well, e.g., based on absolute mAP for Sim10K→Cityscapes, the transformer-based approach [78] performs best. In contrast, the approach based on every-pixel matters [71] performs best when  $\Delta\text{mAP}$  is considered. Note that for Sim10K→Cityscapes, conditional normalization [109] has second-best performance for both absolute mAP and  $\Delta\text{mAP}$ . Similar results can be observed for the PascalVOC→Comic adaptation scenario where style-consistency [101] is the best performing method across the absolute mAP and the  $\Delta\text{mAP}$  measures.

#### 4.4.1 Is one category of adaptation methods better than the other category of approaches?

After discussing all the class of methods (i.e. adversarial feature learning, pseudo-label based self-training etc) and comparing their performances, it is natural to wonder if one class of methods is better than the rest.

From Table. 2 and Table. 3, we can see that there is no clear winner among different class of methods. For example, robust learning [62] is a pseudo-label based self-training approach that outperforms adversarial feature learning-based approach adaptive faster-rcnn [25] on many adaptation experiments. However, many adversarial feature learning-based methods [71], [72], [73], [109] outperform robust learning [62] on a variety of adaptation scenarios. Each class of adaptation strategy has its benefit and drawbacks. For example, adversarial feature learning (discussed in Sec. 3.1) based training is known to be unstable in practice and requires careful parameter tuning and regularization to stabilize the feature adaptation process. However, under appropriate hyperparameters and with correct regularization, adversarial feature learning can provide significant performance improvements. Most likely for this reason, the most popular strategy used in the literature is adversarial feature learning-based methods [25], [57], [69], [71], [73]. On the contrary, self-training with pseudo-labels (discussed in Sec. 3.2) has a stable training curve as supervised detection loss can be used to train the network with supervision from pseudo-labels obtained through the source-trained model. However, self-training based methods need to be careful of the noise in the pseudo-labels as such noise can potentially get re-enforced into the model resulting in sub-optimal performance. For this reason, most self-training based approaches utilize only confident predictions [62], [97].

Image-to-image translation based methods (discussed in Sec. 3.3) try to avoid such training related issues by addressing the domain gap in the input-space. Typically, a image-translation module is added that transforms source samples to target-like samples and vice versa. Since target-like samples will have supervision readily available, the model can be trained better to perform detection on original samples from the target domain. However, this strategy heavily relies on the efficacy of the image-translation module to learn a perfect mapping between source and target domain. For this reason, most image-to-image translation methods additionally utilize adversarial feature learning and/or self-training strategy that can compensate some of the mistakes made by the image-translation module [64], [99], [100].

Domain randomization (discussed in Sec. 3.4) based approaches on the other hand, create multiple stylized version of the source domain in order to ensure that the detection model is devoid of style-bias from source domain and focus on other important discriminative features to detect the objects. However, to completely de-bias the detector from any specific domain, many random domains need to be included during training, which might not be feasible. To address this, domain randomization methods also additionally apply adversarial feature learning between source, stylized-source, and target domain.

Mean-teacher training (discussed in 3.5) is another popular strategy used to adapt detectors to a target domain through mean-teacher updates [146]. Mean-teacher [146] has been successfully used for transfer learning and was initially proposed to utilize unlabeled data of similar domains. These approaches utilize additional regularization or feature matching strategy to extend the mean-teacher for unsupervised domain adaptation [68], [96].

Similarly, graph-reasoning (discussed in 3.6) based meth-

ods aim to understand object relations within any given sample of the source domain and try to leverage that information to make sure the target domain also follow similar relations. This helps the knowledge transfer between source and target domain by enforcing certain constraints on the target features. Most of these different categories of approaches are mutually exclusive and can be combined in different ways to leverage the benefits of each strategy while mitigating the drawbacks of other strategies. For example, unbiased mean-teacher [68] combines mean-teacher training, adversarial feature learning and image-to-image translation and demonstrated that it could outperform most other works in the literature.

## 5 RESEARCH DIRECTIONS

As discussed in the earlier sections, various methods have been proposed to address the challenging problem of adapting deep object detectors to different scenarios. In this section, we explore outstanding issues along with potential solutions and future research directions. In the following subsections, we discuss two broad categories of issues and directions: (i) Comprehensive evaluation and (ii) Improving generalization with real-world constraints.

### 5.1 Comprehensive evaluation

We identify some of the issues concerning the evaluation protocol. We believe that addressing these issues would result in a more robust and rigorous evaluation process.

- **Generalizing to other detection frameworks:** Most of the existing domain adaptation approaches for object detection are evaluated only on the Faster-RCNN detection framework [15]. Further, some of these approaches, such as [25] are designed with the assumption of a two-stage detection process. However, considering the popularity of single shot detection approaches such as YOLO [16], SSD [17], and more recent ones such as FCOS [89] and DETR [90], it is important to evaluate the performance of the adaptation methods using other frameworks as well. Evaluating on this additional category of detection approaches will verify if the adaptation approaches are generalizable to different detectors.
- **Complex real-world datasets:** Initial approaches such as DA-Faster [25], SWDA [57] for domain adaptive object detection designed the evaluation protocol based on existing datasets such as FoggyCityscapes, Sim10K, Pascal VOC, *etc* and subsequent approaches followed similar protocols. It is important to note that some of these datasets and protocols were not specifically constructed to reflect the real-world distribution gap. For example, consider the case of building a self-driving system where it is important to construct source and target datasets such that the distribution gap arises from a variety of factors such as differences in weather conditions, geographical locations, types of vehicles, types of backgrounds, densities of vehicles, *etc*. Considering this, it is crucial to construct complex datasets reflecting real-world scenarios in order to

enable a more rigorous and robust evaluation process.

- **Considering other applications:** Existing methods in the literature mainly focus on applications such as autonomous navigation [25], [59], [71], [73], crowd surveillance [72], and commonly found general objects [57], [67], [69]. Whereas multiple other real world applications are rarely explored in the context of domain adaptive detection, such as medical imaging [113], scene text [105], document objects [108]. For example, consider a lesion detection model trained on Optical Coherence Tomography (OCT) images captured from a particular device. As OCT image distribution shifts when the respective image is captured from another device. It becomes crucial to perform an unsupervised adaptation of the lesion detection model to improve cross-device performance in such cases. Hence, it becomes important to develop adaptation strategies for these rarely explored applications, as they have significant real-world implications.
- **Consistent training and inference strategy:** Although most of the domain adaptation methods for detection demonstrated their effectiveness on the Faster-RCNN detection framework, we observed that the training and inference strategies are not uniform. As it can be observed from Table 2, hyperparameters like the learning rate and training duration varies from method to method. While one may argue that this results from hyperparameter tuning, it is important to note that these hyperparameters are more associated with the backbone network, and changes in them can potentially be the cause of the performance changes. Hence, it is important to maintain consistent hyperparameters related to the backbone network in order to perform a fair comparison of the methods. Similarly, different methods use a different resizing ratio of the input images, which directly affects the detection performance. Considering these observations, it is important to establish consistent and uniform training/inference strategies that will enable the reader to obtain a fair perspective of the evaluation results.

## 5.2 Improving generalization with real-world constraints

In this section, we focus on potential future research directions that consider some real-world settings typically ignored in the existing approaches. Most of the existing methods are based on assumptions that need not necessarily be true in the real-world. For example, it is often assumed that (i) the number of samples across classes in both source and target domain dataset is assumed to be balanced, (ii) all the classes existing in the source domain dataset are present in the target domain dataset and vice versa, (iii) source domain dataset is always available during target domain adaptation, (iv) all target domain samples are unlabeled, and (v) large number of unlabeled target domain samples are available at the training. We discuss various problem settings where such assumptions are relaxed, the resulting challenges and

possible strategies to overcome them. Furthermore, we also discuss additional strategies such as multi-source domain adaptation, test-time adaptation, and continuous adaptation necessary for real-world practice settings.

- **Weak/semi-supervised domain adaptation:** Unsupervised domain adaptation is specifically useful because annotating the target domain samples is labor-intensive and costly for the object detection task. However, in some scenarios, it might be possible to obtain bounding-box annotations for a subset of target samples or provide weak image-level annotations indicating the presence/absence of the categories. In such cases, this additional knowledge can be leveraged to improve the performance on target domain further. For weakly supervised domain adaptation, each image is annotated at image-level to indicate which categories are present/absence in the image and no bounding box annotations are provided. Inoue *et al.* [94] is the only work in the literature addressing this issue with the help of image-to-image translation and pseudo-label training guided by weak annotations. A subset of the target dataset is labeled fully with bounding-box and respective category annotation for semi-supervised adaptation. Liu *et al.* [149] proposed the use of annotated subset and a relatively large unlabeled subset to enhance the detection performance. Though Liu *et al.* [149] does not extensively evaluate the adaptive detection benchmarks, it provides a potentially useful strategy that can be easily adopted for the domain adaptation case. Hence, the use of both weak-supervision and semi-supervision is important to explore further to bridge the performance gap between fully supervised and adaptive training.
- **One/few-shot domain adaptation:** Conventional unsupervised domain adaptation setting assumes availability of a large number of unlabeled samples from the target domain. However, this might not hold true in many realistic conditions where data is an incoming stream of images [112] resulting in either one or few unlabeled samples available from the target domain during adaptation training. To address this sample scarcity in the target domain-specific adaptation strategy is needed to get the best out of the available dataset. D’Innocente *et al.* [112] is the only approach to explore the one-shot adaptation problem by adding an auxiliary self-supervision loss and few-shot adaptation is not yet explored in the domain adaptive detection literature. A straightforward strategy to address one/few-shot adaptation for object detectors would be to apply image stylizing to increase the pool of unlabeled target domain samples and apply a conventional feature-matching strategy by making detector domain invariant [170]. Another strategy would be to label the one or few images available in the target domain and perform a supervised domain adaptation with annotated sample-rich source domain and sample-scarce target domain [171], [172], [172].
- **Imbalanced classes:** Object frequency in the real



world often follows a power law, and this imbalance class problem is typically ignored by existing domain adaptive detection methods which assume aligned class spaces. This assumption limits the methods' performances on imbalanced tasks encountered frequently in the real world, especially in the object detection task. Hence, it is important to avoid the dominance of any particular class while performing the domain alignment. An obvious solution to this problem is to re-weight the classes based on the frequency [46], [173]. However, this strategy is infeasible in the case of the target dataset where samples are unlabeled. Hence, it is important to develop methods that can estimate the class distribution in the target dataset, which can then be used for re-weighting the classes.

- **Partial domain adaptation:** A related problem to the class imbalance issue is that of partial domain adaptation. A typical assumption in the domain adaptation literature is that the label space between the source and target dataset is same. That is, the source and target dataset have the same  $K$  classes. However, a real-world setting might not completely reflect this scenario. It might often be the case that the target dataset has much fewer classes than the ones in the source dataset. Aligning the source domain completely with the target domain might result in negative transfer in such cases. Hence, it is important to mitigate the problem of negative transfer by down-weighting the data samples belonging to the outlier source classes, thereby promoting feature alignment only in the shared label spaces (positive transfer) [174], [175], [176].
- **Open-set domain adaptation:** Similar to the issue of partial domain adaptation, open-set DA deals with the issue where the label space is not completely shared between the source and the target domain. However, in this particular case, the target domain contains unknown classes that are not observed in the source domain. This issue also results in the negative transfer; however, identify the unknown class samples in the target dataset are challenging and we need to borrow principles from the open set recognition task in order to automatically down-weight such samples during the domain alignment [177], [178].
- **Source-free domain adaptation:** A common assumption in the existing literature is that the samples from the source domain are available during domain adaptation training. However, in real-world scenarios, gaining access to source data might not be practical due to privacy concerns, legal issues, and inefficient data transmission. To this end, we must tackle the problem of source-free domain adaptive object detection, where, there is no access to the source data but only the source trained model. An obvious approach would be based on the pseudo-label based self-training, which involves first generating pseudo-labels on the target and re-using them to supervise the network on the target data. However, it is important to generate highly reliable pseudo-labels by filtering out

incorrect labels [74], [179], [180].

- **Multi-source domain adaptation:** All the existing approaches for domain adaptive object detection assume that the labeled training data are sampled from a single domain. This neglects a more practical scenario where training data are collected from multiple sources. For example, consider the case of a self-driving system where source data is available from multiple cities, and the overall goal is to adapt to a new city. Since the source dataset comes from multiple cities, these data points might potentially correspond to multiple sub-distributions. Aligning the target data with all these sub-distributions might not necessarily result in optimal performance. Hence, it would be essential to perform a selective adaptation where the only relevant source samples are considered [181], [182].
- **Continuous and test time adaptation:** Existing domain adaptive detection approaches consider the world to be separated into stationary domains. However, this may not be the case in many real-world applications, where samples arise from a continuously evolving underlying process. Examples include videos with gradually changing lighting/weather conditions. Hence, a more practical setting for domain adaptation is to consider the lifelong learning problem of adapting a pre-trained model to dynamically changing environmental conditions. This setting reflects a real-world scenario where we encounter unlabeled images from new target environments that are not observed during training. Hence, it is important to develop strategies to adapt at test time where we have access only to a single or a few instances of target domain data [183], [184], [185]. Note that these strategies should consider additional constraints such as computationally inexpensive adaptation, restricted access to source domain data, and the issue of catastrophic forgetting, which is a prominent issue that results from continuous learning [160], [186], [187], [188]. Such constraints are required to adapt to dynamically changing environments.

## 6 CONCLUSION

In this paper, we considered the unsupervised domain adaptation of deep object detectors and presented an extensive survey of existing approaches for this task. We have reviewed several approaches that were published in the last few years. We have provided a comprehensive taxonomy of the existing approaches, followed by a detailed analysis of the various methods along with their merits and demerits. We have also provided a detailed discussion on the existing datasets, evaluation protocols and comprehensive performance comparison. Finally, we have identified some of the outstanding issues and promising directions to drive future research.

## REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need,"

- in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 4171–4186.
  - [3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *arXiv preprint arXiv:2005.14165*, 2020.
  - [4] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
  - [5] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel *et al.*, “Mastering atari, go, chess and shogi by planning with a learned model,” *arXiv preprint arXiv:1911.08265*, 2019.
  - [6] T. Xiao, E. Jang, D. Kalashnikov, S. Levine, J. Ibarz, K. Hausman, and A. Herzog, “Thinking while moving: Deep reinforcement learning with concurrent control,” in *International Conference on Learning Representations*, 2019.
  - [7] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
  - [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
  - [9] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
  - [10] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
  - [11] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
  - [12] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
  - [13] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.
  - [14] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
  - [15] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
  - [16] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
  - [17] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *European conference on computer vision*. Springer, 2016, pp. 21–37.
  - [18] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
  - [19] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
  - [20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
  - [21] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, “A theory of learning from different domains,” *Machine learning*, vol. 79, no. 1, pp. 151–175, 2010.
  - [22] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa, “Visual domain adaptation: A survey of recent advances,” *IEEE Signal Processing Magazine*, vol. 32, no. 3, pp. 53–69, 2015.
  - [23] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, “Domain-adversarial training of neural networks,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
  - [24] J. Hoffman, D. Wang, F. Yu, and T. Darrell, “Fcns in the wild: Pixel-level adversarial and constraint-based adaptation,” *arXiv preprint arXiv:1612.02649*, 2016.
  - [25] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, “Domain adaptive faster r-cnn for object detection in the wild,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3339–3348.
  - [26] A. Shrivastava, S. Shekhar, and V. M. Patel, “Unsupervised domain adaptation using parallel transport on grassmann manifold,” in *IEEE winter conference on applications of computer vision*. IEEE, 2014, pp. 277–284.
  - [27] B. Sun, J. Feng, and K. Saenko, “Return of frustratingly easy domain adaptation,” in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
  - [28] K. Saito, Y. Ushiku, and T. Harada, “Asymmetric tri-training for unsupervised domain adaptation,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2988–2997.
  - [29] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, “Cycada: Cycle-consistent adversarial domain adaptation,” in *International Conference on Machine Learning*, 2018, pp. 1989–1998.
  - [30] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, “Adversarial discriminative domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7167–7176.
  - [31] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, “Maximum classifier discrepancy for unsupervised domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3723–3732.
  - [32] C.-Y. Lee, T. Batra, M. H. Baig, and D. Ulbricht, “Sliced wasserstein discrepancy for unsupervised domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 285–10 295.
  - [33] W.-G. Chang, T. You, S. Seo, S. Kwak, and B. Han, “Domain-specific batch normalization for unsupervised domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7354–7362.
  - [34] S. Roy, A. Siarohin, E. Sangineto, S. R. Buló, N. Sebe, and E. Ricci, “Unsupervised domain adaptation using feature-whitening and consensus loss,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9471–9480.
  - [35] M. Long, H. Zhu, J. Wang, and M. I. Jordan, “Unsupervised domain adaptation with residual transfer networks,” in *NIPS*, 2016.
  - [36] K. Saito, Y. Ushiku, T. Harada, and K. Saenko, “Adversarial dropout regularization,” in *International Conference on Learning Representations*, 2018.
  - [37] S. Lee, D. Kim, N. Kim, and S.-G. Jeong, “Drop to adapt: Learning discriminative features for unsupervised domain adaptation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 91–100.
  - [38] Z. Murez, S. Kolouri, D. Kriegman, R. Ramamoorthi, and K. Kim, “Image to image translation for domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4500–4509.
  - [39] Y. Sun, E. Tzeng, T. Darrell, and A. A. Efros, “Unsupervised domain adaptation through self-supervision,” 2019.
  - [40] R. Volpi, P. Morerio, S. Savarese, and V. Murino, “Adversarial feature augmentation for unsupervised domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5495–5504.
  - [41] L. Hu, M. Kan, S. Shan, and X. Chen, “Duplex generative adversarial network for unsupervised domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1498–1507.
  - [42] V. K. Kurmi, S. Kumar, and V. P. Namboodiri, “Attending to discriminative certainty for domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 491–500.

- [43] Y.-H. Tsai, W.-C. Hung, S. Schuster, K. Sohn, M.-H. Yang, and M. Chandraker, "Learning to adapt structured output space for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7472–7481.
- [44] Y.-H. Chen, W.-Y. Chen, Y.-T. Chen, B.-C. Tsai, Y.-C. Frank Wang, and M. Sun, "No more discrimination: Cross city adaptation of road scene segmenters," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1992–2001.
- [45] S. Sankaranarayanan, Y. Balaji, A. Jain, S. Nam Lim, and R. Chellappa, "Learning from synthetic data: Addressing domain shift for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3752–3761.
- [46] Y. Zou, Z. Yu, B. Kumar, and J. Wang, "Unsupervised domain adaptation for semantic segmentation via class-balanced self-training," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 289–305.
- [47] Y. Chen, W. Li, and L. Van Gool, "Road: Reality oriented adaptation for semantic segmentation of urban scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7892–7901.
- [48] Y. Luo, L. Zheng, T. Guan, J. Yu, and Y. Yang, "Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2507–2516.
- [49] Y.-H. Tsai, K. Sohn, S. Schuster, and M. Chandraker, "Domain adaptation for structured output via discriminative patch representations," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1456–1465.
- [50] G. Li, G. Kang, W. Liu, Y. Wei, and Y. Yang, "Content-consistent matching for domain adaptive semantic segmentation," 2020.
- [51] L. Du, J. Tan, H. Yang, J. Feng, X. Xue, Q. Zheng, X. Ye, and X. Zhang, "Ssf-dan: Separated semantic feature based domain adaptation network for semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 982–991.
- [52] W.-L. Chang, H.-P. Wang, W.-H. Peng, and W.-C. Chiu, "All about structure: Adapting structural information across domains for boosting semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1900–1909.
- [53] Y. Li, L. Yuan, and N. Vasconcelos, "Bidirectional learning for domain adaptation of semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6936–6945.
- [54] S. Paul, Y.-H. Tsai, S. Schuster, A. K. Roy-Chowdhury, and M. Chandraker, "Domain adaptive semantic segmentation using weak labels," *arXiv preprint arXiv:2007.15176*, 2020.
- [55] W. Hong, Z. Wang, M. Yang, and J. Yuan, "Conditional generative adversarial network for structured domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1335–1344.
- [56] P. Zhang, B. Zhang, P. Zhang, D. Chen, Y. Wang, and F. Wen, "Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation," *arXiv preprint arXiv:2101.10979*, 2021.
- [57] K. Saito, Y. Ushiku, T. Harada, and K. Saenko, "Strong-weak distribution alignment for adaptive object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6956–6965.
- [58] X. Zhu, J. Pang, C. Yang, J. Shi, and D. Lin, "Adapting object detectors via selective cross-domain alignment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 687–696.
- [59] T. Kim, M. Jeong, S. Kim, S. Choi, and C. Kim, "Diversify and match: A domain adaptive representation learning paradigm for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 456–12 465.
- [60] Z. He and L. Zhang, "Multi-adversarial faster-rcnn for unrestricted object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6668–6677.
- [61] A. RoyChowdhury, P. Chakrabarty, A. Singh, S. Jin, H. Jiang, L. Cao, and E. Learned-Miller, "Automatic adaptation of object detectors to new domains using self-training," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 780–790.
- [62] M. Khodabandeh, A. Vahdat, M. Ranjbar, and W. G. Macready, "A robust learning approach to domain adaptive object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 480–490.
- [63] R. Xie, F. Yu, J. Wang, Y. Wang, and L. Zhang, "Multi-level domain adaptive learning for cross-domain detection," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [64] H.-K. Hsu, C.-H. Yao, Y.-H. Tsai, W.-C. Hung, H.-Y. Tseng, M. Singh, and M.-H. Yang, "Progressive domain adaptation for object detection," in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 749–757.
- [65] E. Tzeng, K. Burns, K. Saenko, and T. Darrell, "Splat: semantic pixel-level adaptation transforms for detection," *arXiv preprint arXiv:1812.00929*, 2018.
- [66] C. Zhuang, X. Han, W. Huang, and M. Scott, "ifan: Image-instance full alignment networks for adaptive object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 122–13 129.
- [67] Z. He and L. Zhang, "Domain adaptive object detection via asymmetric tri-way faster-rcnn," in *European Conference on Computer Vision*. Springer, 2020, pp. 481–497.
- [68] J. Deng, W. Li, Y. Chen, and L. Duan, "Unbiased mean teacher for cross domain object detection," 2021.
- [69] C.-D. Xu, X.-R. Zhao, X. Jin, and X.-S. Wei, "Exploring categorical regularization for domain adaptive object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 724–11 733.
- [70] Y. Pan, A. J. Ma, Y. Gao, J. Wang, and Y. Lin, "Multi-scale adversarial cross-domain detection with robust discriminative learning," in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 1324–1332.
- [71] C.-C. Hsu, Y.-H. Tsai, Y.-Y. Lin, and M.-H. Yang, "Every pixel matters: Center-aware feature alignment for domain adaptive object detector," in *European Conference on Computer Vision*. Springer, 2020, pp. 733–748.
- [72] V. A. Sindagi, P. Oza, R. Yasarla, and V. M. Patel, "Prior-based domain adaptive object detection for hazy and rainy conditions," in *European Conference on Computer Vision*. Springer, 2020, pp. 763–780.
- [73] V. VS, P. Oza, V. A. Sindagi, V. Gupta, and V. M. Patel, "Megacda: Memory guided attention for category-aware unsupervised domain adaptive object detection," 2021.
- [74] X. Li, W. Chen, D. Xie, S. Yang, P. Yuan, S. Pu, and Y. Zhuang, "A free lunch for unsupervised domain adaptive object detection without source data," *arXiv preprint arXiv:2012.05400*, 2020.
- [75] C. Liang, Z. Zhao, J. Liu, and J. Zhang, "Domain adaptive object detection via feature separation and alignment," *arXiv preprint arXiv:2012.08689*, 2020.
- [76] D. Guan, J. Huang, A. Xiao, S. Lu, and Y. Cao, "Uncertainty-aware unsupervised domain adaptation in object detection," 2021.
- [77] T. Sun, J. Chen, and F. Ng, "Multi-target domain adaptation via unsupervised domain classification for weather invariant object detection," *arXiv preprint arXiv:2103.13970*, 2021.
- [78] J. Zhang, J. Huang, Z. Luo, G. Zhang, and S. Lu, "Da-detr: Domain adaptive detection transformer by hybrid attention," *arXiv preprint arXiv:2103.17084*, 2021.
- [79] T. Scheck, A. P. Grassi, and G. Hirtz, "Unsupervised domain adaptation from synthetic to real images for anchorless object detection," *arXiv preprint arXiv:2012.08205*, 2020.
- [80] H. Yang, S. Jiang, X. Zhu, M. Huang, Z. Shen, C. Liu, and J. Shi, "Channel-wise alignment for adaptive object detection," *arXiv preprint arXiv:2009.02862*, 2020.
- [81] M. Wang and W. Deng, "Deep visual domain adaptation: A survey," *Neurocomputing*, vol. 312, pp. 135–153, 2018.
- [82] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A comprehensive survey on transfer learning," *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2020.
- [83] M. Toldo, A. Maracani, U. Michieli, and P. Zanuttigh, "Unsupervised domain adaptation in semantic segmentation: a review," *Technologies*, vol. 8, no. 2, p. 35, 2020.
- [84] S. Zhao, X. Yue, S. Zhang, B. Li, H. Zhao, B. Wu, R. Krishna, J. E. Gonzalez, A. L. Sangiovanni-Vincentelli, S. A. Seshia et al., "A review of single-source deep unsupervised visual domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [85] W. Li, F. Li, Y. Luo, and P. Wang, "Deep domain adaptive object detection: a survey," *arXiv preprint arXiv:2002.06797*, 2020.

- [86] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *International journal of computer vision*, vol. 128, no. 2, pp. 261–318, 2020.
- [87] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [88] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *arXiv preprint arXiv:1905.05055*, 2019.
- [89] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9627–9636.
- [90] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.
- [91] H. Daumé III, A. Kumar, and A. Saha, "Frustratingly easy semi-supervised domain adaptation," in *Proceedings of the 2010 Workshop on Domain Adaptation for Natural Language Processing*, 2010, pp. 53–59.
- [92] K. Saito, D. Kim, S. Sclaroff, T. Darrell, and K. Saenko, "Semi-supervised domain adaptation via minimax entropy," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8050–8058.
- [93] X. Xu, X. Zhou, R. Venkatesan, G. Swaminathan, and O. Majumder, "d-sne: Domain adaptation using stochastic neighborhood embedding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2497–2506.
- [94] N. Inoue, R. Furuta, T. Yamasaki, and K. Aizawa, "Cross-domain weakly-supervised object detection through progressive domain adaptation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5001–5009.
- [95] M. Xu, H. Wang, B. Ni, Q. Tian, and W. Zhang, "Cross-domain detection via graph-induced prototype alignment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 355–12 364.
- [96] Q. Cai, Y. Pan, C.-W. Ngo, X. Tian, L. Duan, and T. Yao, "Exploring object relation in mean teacher for cross-domain detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 457–11 466.
- [97] S. Kim, J. Choi, T. Kim, and C. Kim, "Self-training and adversarial background regularization for unsupervised domain adaptive one-stage object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6092–6101.
- [98] G. Zhao, G. Li, R. Xu, and L. Lin, "Collaborative training between region proposal localization and classification for domain adaptive object detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 86–102.
- [99] D. Zhang, J. Li, L. Xiong, L. Lin, M. Ye, and S. Yang, "Cycle-consistent domain adaptive faster rcnn," *IEEE Access*, vol. 7, pp. 123 903–123 911, 2019.
- [100] C. Chen, Z. Zheng, X. Ding, Y. Huang, and Q. Dou, "Harmonizing transferability and discriminability for adapting object detectors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8869–8878.
- [101] A. L. Rodriguez and K. Mikolajczyk, "Domain adaptation for object detection via style consistency," *British Machine Vision Conference*, 2019.
- [102] Y. Shan, W. F. Lu, and C. M. Chew, "Pixel and feature level based domain adaptation for object detection in autonomous driving," *Neurocomputing*, vol. 367, pp. 31–38, 2019.
- [103] G. Xu, Q. Zhang, D. Liu, G. Lin, J. Wang, and Y. Zhang, "Adversarial adaptation from synthesis to reality in fast detector for smoke detection," *IEEE Access*, vol. 7, pp. 29 471–29 483, 2019.
- [104] Y. Yu, X. Xu, X. Hu, and P.-A. Heng, "Dalocnet: Improving localization accuracy for domain adaptive object detection," *IEEE Access*, vol. 7, pp. 63 155–63 163, 2019.
- [105] D. Chen, L. Lu, Y. Lu, R. Yu, S. Wang, L. Zhang, and T. Liu, "Cross-domain scene text detection via pixel and image-level adaptation," in *International Conference on Neural Information Processing*. Springer, 2019, pp. 135–143.
- [106] V. F. Arruda, T. M. Paixão, R. F. Berriel, A. F. De Souza, C. Badue, N. Sebe, and T. Oliveira-Santos, "Cross-domain car detection using unsupervised image-to-image translation: From day to night," in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
- [107] Y. Zheng, D. Huang, S. Liu, and Y. Wang, "Cross-domain object detection through coarse-to-fine feature adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 766–13 775.
- [108] K. Li, C. Wington, C. Tensmeyer, H. Zhao, N. Barmpalios, V. I. Morariu, V. Manjunatha, T. Sun, and Y. Fu, "Cross-domain document object detection: Benchmark suite and method," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 915–12 924.
- [109] P. Su, K. Wang, X. Zeng, S. Tang, D. Chen, D. Qiu, and X. Wang, "Adapting object detectors with conditional domain normalization," in *European Conference on Computer Vision (ECCV)*, 2020.
- [110] C. Li, D. Du, L. Zhang, L. Wen, T. Luo, Y. Wu, and P. Zhu, "Spatial attention pyramid network for unsupervised domain adaptation," in *European Conference on Computer Vision*. Springer, 2020.
- [111] Z. Zhao, Y. Guo, H. Shen, and J. Ye, "Adaptive object detection with dual multi-label prediction," in *European Conference on Computer Vision*. Springer, 2020, pp. 54–69.
- [112] A. D'Innocente, F. C. Borlino, S. Bucci, B. Caputo, and T. Tommasi, "One-shot unsupervised cross-domain detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 732–748.
- [113] S. Yang, L. Wu, A. Wiliem, and B. C. Lovell, "Unsupervised domain adaptive object detection using forward-backward cyclic adaptation," in *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [114] X. Yang, S. Wan, and P. Jin, "Domain-invariant region proposal network for cross-domain detection," in *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2020, pp. 1–6.
- [115] D.-K. Nguyen, W.-L. Tseng, and H.-H. Shuai, "Domain-adaptive object detection via uncertainty-aware distribution alignment," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2499–2507.
- [116] H. Alqasir, D. Muselet, and C. Ducottet, "Region proposal oriented approach for domain adaptive object detection," in *International Conference on Advanced Concepts for Intelligent Vision Systems*. Springer, 2020, pp. 38–50.
- [117] C. Chen, Z. Zheng, Y. Huang, X. Ding, and Y. Yu, "T 3net: Implicit instance-invariant network for adapting one-stage object detectors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2021.
- [118] Y. Wang, R. Zhang, S. Zhang, M. Li, Y. Xia, X. Zhang, and S. Liu, "Domain-specific suppression for adaptive object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9603–9612.
- [119] H. Wang, S. Liao, and L. Shao, "Afan: Augmented feature alignment network for cross-domain object detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 4046–4056, 2021.
- [120] A. Wu, Y. Han, L. Zhu, and Y. Yang, "Instance-invariant adaptive object detection via progressive disentanglement," *arXiv preprint arXiv:1911.08712*, 2019.
- [121] Z. Shen, M. Huang, J. Shi, Z. Liu, H. Maheshwari, Y. Zheng, X. Xue, M. Savvides, and T. S. Huang, "Ctdt: A large-scale cross-domain benchmark for instance-level image-to-image translation and domain adaptive object detection," *International Journal of Computer Vision*, vol. 129, no. 3, pp. 761–780, 2021.
- [122] Y. Chen, H. Wang, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Scale-aware domain adaptive faster r-cnn," *International Journal of Computer Vision*, pp. 1–21, 2021.
- [123] P. Soviany, R. T. Ionescu, P. Rota, and N. Sebe, "Curriculum self-paced learning for cross-domain object detection," *Computer Vision and Image Understanding*, vol. 204, p. 103166, 2021.
- [124] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [125] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, 2014, pp. 2672–2680.
- [126] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chellappa, "Generate to adapt: Aligning domains using generative adversarial networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8503–8512.
- [127] M. I. Belghazi, S. Rajeswar, A. Baratin, D. Hjelm, and A. Courville, "Mine: Mutual information neural estimation," 2018.



- [128] F. Liu, X. Zhang, F. Wan, X. Ji, and Q. Ye, "Domain contrast for domain adaptive object detection," *arXiv preprint arXiv:2006.14863*, 2020.
- [129] M. Fu, Z. Xie, W. Li, and L. Duan, "Deeply aligned adaptation for cross-domain object detection," *arXiv preprint arXiv:2004.02093*, 2020.
- [130] H. Liu, P. Song, and R. Ding, "Wqt and dg-yolo: towards domain generalization in underwater object detection," *arXiv preprint arXiv:2004.06333*, 2020.
- [131] M. Salzmann *et al.*, "Attention-based domain adaptation for single stage detectors," *arXiv preprint arXiv:2106.07283*, 2021.
- [132] L. T. Nguyen-Meidine, M. Kiran, M. Pedersoli, J. Dolz, L.-A. Blais-Morin, and E. Granger, "Incremental multi-target domain adaptation for object detection with efficient domain transfer," *arXiv preprint arXiv:2104.06476*, 2021.
- [133] F. Yu, D. Wang, Y. Chen, N. Karianakis, P. Yu, D. Lymberopoulos, and X. Chen, "Unsupervised domain adaptation for object detection via cross-domain semi-supervised learning," *arXiv preprint arXiv:1911.07158*, 2019.
- [134] F. Munir, S. Azam, and M. Jeon, "Sstn: Self-supervised domain adaptation thermal object detection for autonomous driving," *arXiv preprint arXiv:2103.03150*, 2021.
- [135] X. Wang, T. E. Huang, B. Liu, F. Yu, X. Wang, J. E. Gonzalez, and T. Darrell, "Robust object detection via instance-level temporal cycle confusion," *arXiv preprint arXiv:2104.08381*, 2021.
- [136] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [137] Z. Yi, H. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2849–2857.
- [138] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 700–708.
- [139] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 172–189.
- [140] W. Li, F. Li, Y. Luo, and P. Wang, "Unsupervised image-generation enhanced adaptation for object detection in thermal images," *arXiv preprint arXiv:2002.06770*, 2020.
- [141] A. Abramov, C. Bayer, and C. Heller, "Keep it simple: Image statistics matching for domain adaptation," *arXiv preprint arXiv:2005.12551*, 2020.
- [142] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1501–1510.
- [143] F. Munir, S. Azam, M. A. Rafique, A. M. Sheri, and M. Jeon, "Thermal object detection using domain adaptation through style consistency," *arXiv preprint arXiv:2006.00821*, 2020.
- [144] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [145] S. Gupta, J. Hoffman, and J. Malik, "Cross modal distillation for supervision transfer," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2827–2836.
- [146] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1195–1204.
- [147] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," *arXiv preprint arXiv:1905.02249*, 2019.
- [148] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li, "Fixmatch: Simplifying semi-supervised learning with consistency and confidence," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [149] Y.-C. Liu, C.-Y. Ma, Z. He, C.-W. Kuo, K. Chen, P. Zhang, B. Wu, Z. Kira, and P. Vajda, "Unbiased teacher for semi-supervised object detection," *International Conference on Learning Representations*, 2021.
- [150] G. French, M. Mackiewicz, and M. Fisher, "Self-ensembling for visual domain adaptation," in *International Conference on Learning Representations*, no. 6, 2018.
- [151] Z. Deng, Y. Luo, and J. Zhu, "Cluster alignment with a teacher for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9944–9953.
- [152] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," 2016.
- [153] R. Shu, H. H. Bui, H. Narui, and S. Ermon, "A dirt-t approach to unsupervised domain adaptation," in *Proc. 6th International Conference on Learning Representations*, 2018.
- [154] S. Tang, Z. Cheng, S. Pu, D. Guo, Y. Niu, and F. Wu, "Learning a domain classifier bank for unsupervised adaptive object detection," *arXiv preprint arXiv:2007.02595*, 2020.
- [155] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," 2016.
- [156] Y. Yan, Q. Zhang, B. Ni, W. Zhang, M. Xu, and X. Yang, "Learning context graph for person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2158–2167.
- [157] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [158] Y. Yan, N. Zhuang, J. Zhang, M. Xu, Q. Zhang, Z. ZHENG, S. Cheng, Q. Tian, X. Yang, W. Zhang *et al.*, "Fine-grained video captioning via graph-based multi-granularity interaction learning," *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [159] X. Ma, T. Zhang, and C. Xu, "Gcan: Graph convolutional adversarial network for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8266–8276.
- [160] M. Mancini, S. R. Bulò, B. Caputo, and E. Ricci, "Adagraph: Unifying predictive and continuous domain adaptation through graphs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6568–6577.
- [161] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [162] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.
- [163] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2636–2645.
- [164] C. Sakaridis, D. Dai, and L. V. Gool, "Semantic foggy scene understanding with synthetic data," *International Journal of Computer Vision*, vol. 126, pp. 973–992, 2018.
- [165] H. Nada, V. A. Sindagi, H. Zhang, and V. M. Patel, "Pushing the limits of unconstrained face detection: a challenge dataset and baseline results," in *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2018, pp. 1–10.
- [166] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "Wider face: A face detection benchmark," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5525–5533.
- [167] B. Li, W. Ren, D. Fu, D. Tao, D. Feng, W. Zeng, and Z. Wang, "Benchmarking single-image dehazing and beyond," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 492–505, 2019.
- [168] M. Johnson-Roberson, C. Barto, R. Mehta, S. N. Sridhar, K. Rosaen, and R. Vasudevan, "Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks?" *arXiv preprint arXiv:1610.01983*, 2016.
- [169] X. Jiang, F. R. Yu, T. Song, Z. Ma, Y. Song, and D. Zhu, "Blockchain-enabled cross-domain object detection for autonomous driving: A model sharing approach," *IEEE Internet of Things Journal*, vol. 7, no. 5, pp. 3681–3692, 2020.
- [170] Y. Luo, P. Liu, T. Guan, J. Yu, and Y. Yang, "Adversarial style mining for one-shot unsupervised domain adaptation," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [171] S. Motiian, Q. Jones, S. M. Iranmanesh, and G. Doretto, "Few-shot adversarial domain adaptation," in *NIPS*, 2017.
- [172] J. Zhang, Z. Chen, J. Huang, L. Lin, and D. Zhang, "Few-shot structured domain adaptation for virtual-to-real scene parsing,"

- in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [173] M. A. Jamal, M. Brown, M.-H. Yang, L. Wang, and B. Gong, “Rethinking class-balanced methods for long-tailed visual recognition from a domain adaptation perspective,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7610–7619.
  - [174] Z. Cao, L. Ma, M. Long, and J. Wang, “Partial adversarial domain adaptation,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 135–150.
  - [175] Y. Kim, S. Hong, S. Yang, S. Kang, Y. Jeon, and J. Kim, “Associative partial domain adaptation,” *arXiv preprint arXiv:2008.03111*, 2020.
  - [176] Z. Cao, K. You, M. Long, J. Wang, and Q. Yang, “Learning to transfer examples for partial domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2985–2994.
  - [177] P. Panareda Busto and J. Gall, “Open set domain adaptation,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 754–763.
  - [178] H. Liu, Z. Cao, M. Long, J. Wang, and Q. Yang, “Separate to adapt: Open set domain adaptation via progressive separation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2927–2936.
  - [179] J. N. Kundu, N. Venkat, R. V. Babu *et al.*, “Universal source-free domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4544–4553.
  - [180] A. R. Nelakurthi, R. Maciejewski, and J. He, “Source free domain adaptation using an off-the-shelf classifier,” in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 140–145.
  - [181] H. Zhao, S. Zhang, G. Wu, J. M. Moura, J. P. Costeira, and G. J. Gordon, “Adversarial multiple source domain adaptation,” *Advances in neural information processing systems*, vol. 31, pp. 8559–8570, 2018.
  - [182] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, “Moment matching for multi-source domain adaptation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1406–1415.
  - [183] A. Royer and C. H. Lampert, “Classifier adaptation at prediction time,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1401–1409.
  - [184] E. M. Fredericks, B. DeVries, and B. H. Cheng, “Towards run-time adaptation of test cases for self-adaptive systems in the face of uncertainty,” in *Proceedings of the 9th International Symposium on Software Engineering for Adaptive and Self-Managing Systems*, 2014, pp. 17–26.
  - [185] D. Wang, E. Shelhamer, S. Liu, B. Olshausen, T. Darrell, U. Berkeley, and A. Research, “Tent: Fully test-time adaptation by entropy minimization,” in *International Conference on Learning Representations*, vol. 4, 2021, p. 6.
  - [186] J. Hoffman, T. Darrell, and K. Saenko, “Continuous manifold based adaptation for evolving visual domains,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 867–874.
  - [187] M. Wulfmeier, A. Bewley, and I. Posner, “Incremental adversarial domain adaptation for continually changing environments,” in *2018 IEEE International conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 4489–4495.
  - [188] Z. Wu, X. Wang, J. E. Gonzalez, T. Goldstein, and L. S. Davis, “Ace: Adapting to changing environments for semantic segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2121–2130.

**Poojan Oza** is a PhD student in the Department Of Electrical & Computer Engineering at The Johns Hopkins University. He graduated from IIT-Delhi with a Master's degree in Electronics and Computer Engineering. His research interests include deep learning based one-class methods, anomaly/novelty detection, open-set recognition, domain adaptation and object detection.

**Vishwanath A. Sindagi** is a PhD student in the Department Of Electrical & Computer Engineering at The Johns Hopkins University. Prior to joining Johns Hopkins University, he worked for Samsung R&D Institute-Bangalore. He graduated from IIIT-Bangalore with a Master's degree in Information Technology. His research interests include deep learning based crowd analytics, object detection, applications of generative modeling, domain adaptation and low-level vision.

**Vibashan VS** is a PhD student in the Department Of Electrical & Computer Engineering at The Johns Hopkins University. He graduated from NIT-Tiruchirappalli with a Bachelor's degree in Instrumentation and Control. His research interests include deep learning based object detection, domain adaptation.

**Vishal M. Patel** is an Assistant Professor in the Department of Electrical and Computer Engineering (ECE) at Johns Hopkins University. Prior to joining Hopkins, he was an A. Walter Tyson Assistant Professor in the Department of ECE at Rutgers University and a member of the research faculty at the University of Maryland Institute for Advanced Computer Studies (UMIACS). His current research interests include signal processing, computer vision, and pattern recognition with applications in bio-metrics and imaging. He has received a number of awards including the 2016 ONR Young Investigator Award, the 2016 Jimmy Lin Award for Invention, A. Walter Tyson Assistant Professorship Award, Best Paper Award at IEEE AVSS 2017, Best Paper Award at IEEE BTAS2015, Honorable Mention Paper Award at IAPR ICB 2018, two Best Student Paper Awards at IAPR ICPR 2018, and Best Poster Awards at BTAS 2015 and 2016. He is an Associate Editor of the IEEE Signal Processing Magazine, IEEE Biometrics Compendium, Pattern Recognition Journal, and serves on the Machine Learning for Signal Processing Technical Committee of the IEEE Signal Processing Society. He is serving as the Vice President (Conferences) of the IEEE Biometrics Council. He is a member of Eta Kappa Nu, Pi Mu Epsilon, and Phi Beta Kappa.