

A Multi Camera Unsupervised Domain Adaptation Pipeline for Object Detection in Cultural Sites through Adversarial Learning and Self-Training

Giovanni Pasqualino^a, Antonino Furnari^{a,*}, Giovanni Maria Farinella^{a,b,c}

^a*Department of Mathematics and Computer Science, University of Catania, Viale Andrea Doria 6, Catania, 95125, Italy*

^b*CUTGAN, University of Catania, Via Santa Sofia 98, Catania, 95123, Italy*

^c*ICAR-CNR, National Research Council, Via Ugo la Malfa 153, Palermo, 90146, Italy*

Abstract

Object detection algorithms allow to enable many interesting applications which can be implemented in different devices, such as smartphones and wearable devices. In the context of a cultural site, implementing these algorithms in a wearable device, such as a pair of smart glasses, allow to enable the use of augmented reality (AR) to show extra information about the artworks and enrich the visitors' experience during their tour. However, object detection algorithms require to be trained on many well annotated examples to achieve reasonable results. This brings a major limitation since the annotation process requires human supervision which makes it expensive in terms of time and costs. A possible solution to reduce these costs consist in exploiting tools to automatically generate synthetic labeled images from a 3D model of the site. However, models trained with synthetic data do not generalize on real images acquired in the target scenario in which they are supposed to be used. Furthermore, object detectors should be able to work with different wearable devices or different mobile devices, which makes generalization even harder. In this paper, we present a new dataset collected in a cultural site to study the problem of domain adaptation for object detection in the presence of multiple unlabeled target domains corresponding to different cameras and a labeled source domain obtained considering synthetic images for training purposes. We present a new domain adaptation method which outperforms current state-of-the-art approaches combining the benefits of aligning the domains at the feature and pixel level with a self-training process. We release the dataset at the following link <https://iplab.dmi.unict.it/OBJ-MDA/> and the code of the proposed architecture at <https://github.com/fpv-iplab/STMDA-RetinaNet>.

Keywords: Object Detection, Cultural Sites, First Person Vision, Unsupervised Domain Adaptation

*Corresponding author:

Email address: furnari@dmf.unict.it (Antonino Furnari)

1. Introduction

In recent years, wearable and mobile devices have increasingly attracted the interest of the scientific community because of their ability to capture human-centric data which reflects the intent and interests of the users (e.g., a picture shot with a mobile phone or a video acquired with an action camera). Given their ever-increasing computing capabilities, different computer vision algorithms have been integrated into these devices allowing the development of new applications in different scenarios. In particular, previous works, have shown that human-centric devices such as mixed reality glasses can be exploited in a cultural site to improve the fruition of artworks [1, 2] by showing extra information to visitors using augmented reality (AR), or to track users' behavior [3]. While all these applications require the ability to detect objects in the scene, training object detectors is still costly and time consuming because it requires images to be labeled by human annotators. To reduce the data annotation costs, the authors of [4] proposed to create large datasets of images of a cultural sites by generating synthetic labeled images in a simulated environment. Despite this approach speeds up data collection and reduces the annotation costs, object detectors trained with synthetic images achieve poor performance when tested with real images due to the domain shift between the data used for training (source domain) and the data used at test time (target domain) [5]. Furthermore in real-workly scenarios, object detection algorithms often need to be deployed to different devices, which are generally equipped with different cameras. This constraint further reduces the generalization ability of the object detection methods in real scenarios. Figure 1 reports some qualitative results of a standard object detector trained and tested on different domains of images of a cultural sites: a set of synthetic images, real images acquired with an HoloLens device, and a set of images collected with a GoPro. As can be noted, the detection of the artworks works perfectly only if the the training and test set belong to the same data distribution.

Domain adaptation techniques [6] can be used to reduce the domain difference between source and target sets. However, in a real scenario, the algorithm should also generalize to images collected using multiple cameras as in the example in Figure 1, which may present subtle characteristics capable of affecting model performance. We propose to tackle this problem as a multi-target unsupervised domain adaptation task in which there is a labeled source domain (the synthetic data) and more than one unlabeled target domains (the target images acquired using different cameras). We note that, since target unlabeled images can be acquired with a little effort, the task setup involves a small additional overhead as compared to single-target domain adaptation. We hence investigate whether the presence of more than one target domains can assist the domain adaptation process in the considered settings. To analyze the problem, we introduce a new dataset of both synthetic and real images collected in a cultural site and suitable to study unsupervised multi-camera domain adaptation. We perform

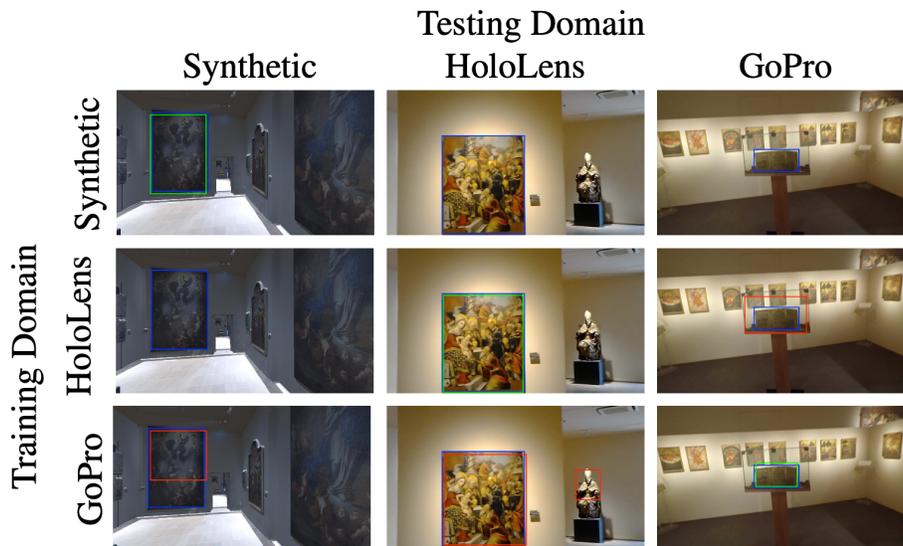


Figure 1: Qualitative results of a standard object detector trained and tested on different domains. Blue bounding boxes represent the ground truth, green boxes represent correct detections, whereas red boxes indicate wrong detections (either object localization or classification). The model was trained using the domain indicated in the rows and tested using the domain reported in the columns.

experiments to assess the ability of current domain adaptation approaches to generalize across multiple cameras. We hence investigate a generalization of current state-of-the-art methods which is shown to outperform current methods. The contributions of this paper are as follows. 1) We present the first dataset to study multi-camera domain adaptation in cultural sites. The dataset has been acquired by real visitors in a cultural site using two different wearable cameras. 2) We propose a domain adaptation approach which takes advantage of multiple unlabeled camera domains and a self-training procedure to improve cross-domain generalization. The proposed method outperforms the results of current state-of-the-art methods by up to +23% mAP. We discuss the limit of the proposed technique and present possible future research directions. The remainder of this paper is organized as follows. In Section 2 we discuss related work. Section 3 presents the proposed dataset. Section 4 discusses the proposed pipeline and method. Section 5 reports the experimental settings and discusses results. Section 6 concludes the paper and summarises the main findings of our study.

2. Related Work

In this section, we discuss the lines of research related to this work: use of wearable devices in cultural site, unsupervised domain adaptation techniques and unsupervised domain adaptation approaches for object detection.

2.1. Use of wearable device in cultural site

Wearable device applications allow to improve the fruition and the perception of the reality around us. In the context of cultural sites, the authors of [1, 7] focused on the creation of virtual guides to enrich the visitors' experience with the fruition of multimedia materials. The applications are based on the detection and recognition of objects to trigger the presentation of associated information. Training these kinds of algorithms requires many labeled images that must be acquired and manually annotated, thus increasing development times and costs. Due to the lack of data in this field, the authors of [8] proposed a dataset of first person videos acquired using Microsoft HoloLens to study different problems in the context of cultural sites. The authors of [4] presented a tool to generate synthetic labeled images from a 3D reconstruction of a cultural site. However, the generated images differ in color and shape with respect to the real counterparts. For this reason, object detection algorithms, which have to work at inference time on real images, produce poor result if trained only with synthetic data. In this paper, we study the problem of object detection in the presence of multiple target domains acquired using different wearable cameras. We present an approach to train an object detector to maximize its detection and recognition accuracy on both target camera domains using labeled synthetic images and unlabeled real images captured with the two cameras.

2.2. Unsupervised Domain Adaptation

Unsupervised domain adaptation (UDA) algorithms transfer the knowledge learned from a source labeled domain to a target unlabeled domain. Most state-of-art domain adaptation approaches take into account labeled source and an unlabeled target domains. Past works in this field focused on reducing some divergence statistics between representation of example belonging to the two domains. The authors of [9] proposed to use the MMD [10] metric to reduce the distributions difference between source and target feature distributions. The authors of [11] proposed the CORAL metric and integrated it inside a CNN to align the covariances of the source and target feature distributions. Other methods propose solutions based on adversarial training. The authors of [6] presented the Gradient Reversal Layer (GRL) which mimics the behavior of a generative adversarial network (GAN) [12] to encourage the extraction of indistinguishable features from images belonging to the source and target domains. The gradient reversal layer allows to train the architecture end-to-end instead of using the usual alternate optimization of the generator and the discriminator as in GANs [12]. The authors of [13] proposed a method based on two stages that combines discriminative modeling, untied weight sharing, and a adversarial loss. The authors of [14] presented a method based on image to image translation that, in the absence of paired example, learns a mapping between the two domains through the use of the adversarial loss. The authors of [15] proposed a clustering based method to generate pseudo labels for the target domain, than the method minimizes the discrepancy of the gradients generated by the source and target images. Other works studied the unsupervised domain adaptation

problem in the presence of many labeled source domains and only one unlabeled target domain. The authors of [16] proposed a method based on the gradient reversal layer that discriminates between all of $Target - Source_n$ pairs where $d = 0, \dots, D$ and D is the number of source domains. The authors of [17] presented a method that consists of three components: feature extractor, moment matching module and a final classifiers. These methods require the access to multiple labeled source dataset which in some cases maybe available or easy to produce. Another line of research focused on a more realistic scenario where there is only one source domain and multiple target domains. In this case, the presence of multiple target domains emulate a real scenario where, for example, different devices with different lenses and images generation pipelines (IGP) produce different target domain. The authors of [18] proposed a method based on an autoencoder which finds a latent space which can capture domain invariant and domain dependent features that can generalize over multiple target domains. The authors of [19] presented a method that extends the idea proposed by [13] replacing a binary discrimination with a multi-class discrimination. The authors of [20] proposed a method based on an iterative multi-teacher knowledge distillation from multiple teachers to a common student.

The presented methods which tackle the multi target domain adaptation problem are used to solve classification task and they are not directly applicable to the object detection problem. In this work, we study how to exploit the information coming from the different target domains for the object detection adapting an adversarial training scheme similar to the work of the authors of [6].

2.3. Unsupervised Domain Adaptation for Object Detection

The methods described in the previous section can be adapted to consider the object detection task. The authors of [5] presented DA-Faster RCNN which is a modified version of Faster RCNN [21] which aligns features at the image and instance levels exploiting gradient reversal layers [6]. The authors of [22] proposed to adapt the high-and low-level features. [23] proposed to extend the architecture presented by the authors of [5] adding more discriminators with gradient reversal layers to the Faster RCNN backbone. The authors of [24] presented an approach composed of two stages: 1) a domain diversification stage where the distribution of the labeled data is diversified by generating various distinctive domains shifted from the source domain using image to image techniques; 2) multi-domain-invariant representation learning, where adversarial learning is applied with a multi-domain discriminator to encourage feature to be indistinguishable across domains. The authors of [25] proposed to translate images from the source domain into the target domain using CycleGAN and trained an object detector using a self-training procedure to create pseudo label for the target dataset. The authors of [26] introduced in a SSD architecture a novel self-training method called weak self-training (WST) combined with the adversarial background score regularization (BSR) to prevent the degeneration of the performance due to incorrect pseudo label obtained using a naive approach reducing the amount of false negative and positive detections. The authors of [27] presented a framework which combines intermediate domains

to progressively adapt feature alignment for object detection and a weighted task loss which weights the samples in the intermediate domain. The authors of [28] presented a method based on SSD which is divided in three steps: in the first step the SSD detector is pretrained using the source images; in the second step the source images are converted to real with CycleGAN; in the third step SSD is trained using the converted source images, the target images and using the weak self-training method proposed in [26]. The authors of [29] proposed a Implicit Instance-Invariant Network (*I³Net*), a single stage object detector which adapt the source and the target domain considering: 1) a strategy to assign large weights to those sample-scarce categories and easy-to-adapt samples considering the intra-class and intra-domain variation, 2) a module to suppress uninformative background features boosting the foreground object matching, 3) a module that align the category at different domain specific layers and regularize the average prediction of different layer respect to the same category. The authors of [30] introduced a generic approach based on an attention mechanism which allows to detect the important regions of the feature map extracted from the backbone on which adaptation should focus. The authors of [31] proposed a method which works at image and instance level aligning the two distributions so that well-aligned and poor-aligned samples are adaptively weighted based on the uncertainty of each sample. The authors of [32] presented a feature alignment method based on Faster RCNN which consist of three modules: 1) a global discriminator which align the feature extracted from the backbone; 2) category wise discriminators which aligns the features of each class belonging to the source and the target domains; 3) a memory guided attention mechanism which aids the category-wise discriminators to align category specific features between the two domains.

Our work investigates whether the presence of multiple unlabeled target domains can improve the generalization of current methods. We further present an architecture based on feature alignment, image to image translation and self-training to tackle multi-camera unsupervised domain adaptation for object detection in cultural sites.

3. Dataset

To study the problem, we created a dataset¹ that contains images of 16 artworks included in the cultural site “Galleria Regionale di Palazzo Bellomo²”. The collection covers different types of artworks, as well as books, sculptures and paintings. We considered three domains: i) synthetic images generated from a 3D model of the cultural site and automatically labeled during the generation process, ii) real images collected by 10 visitors with a HoloLens device and manually labeled, iii) real images collected by the same visitors with a GoPro and manually labeled. Figure 2 shows some examples of images belonging to the

¹The dataset is available at <https://iplab.dmi.unict.it/OBJ-MDA>

²<http://www.regione.sicilia.it/beniculturali/palazzobellomo/>

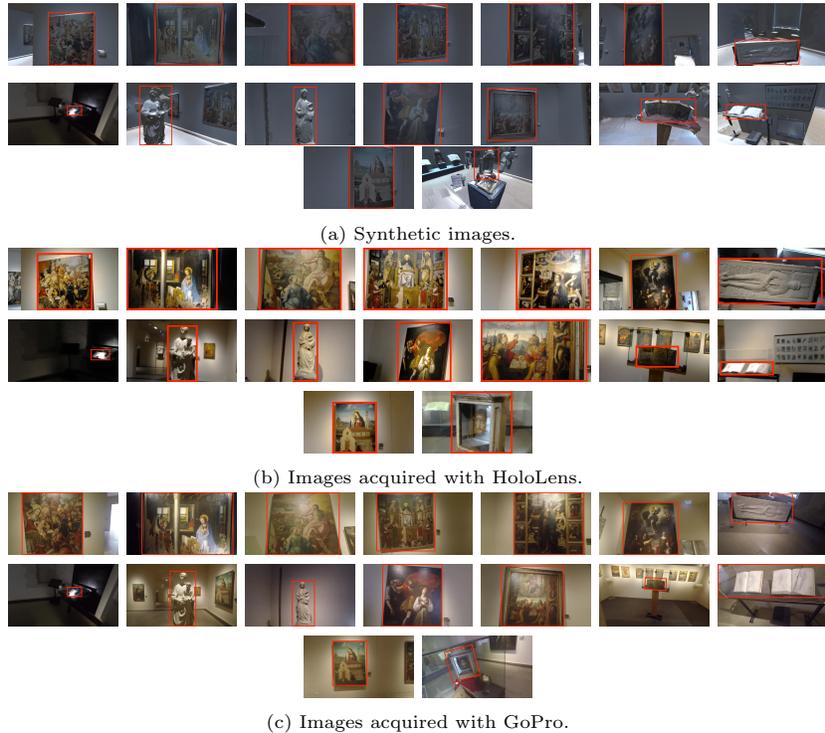


Figure 2: Example of labeled images of the 16 artworks with respect to the three considered domains: (a) Synthetic images, (b) images acquired with HoloLens, (c) images acquired with GoPro.

three domains. As can be noted, synthetic images differ from the real images in style, shapes of the 3D objects (e.g., observe the statues of Figure 2) and field of view. Similarly, real images acquired with two the different devices differ only in style and field of view. The three sets of images have been collected as detailed in the following:

- Synthetic labeled images (used as source domain): these images have been generated using the tool proposed by [4]. The tool allows to annotate in 3D the position of artworks in the 3D model of a cultural site and simulates an agent navigating the environment while acquiring egocentric images of the observed artworks. The acquired images are automatically labeled by projecting the 3D bounding boxes of the objects onto the generated 2D images. This set contains 75244 images divided in 51284 training images and 23960 test images.
- Target images acquired using a HoloLens: this set of data has been sampled from the work of [8] where data has been manually annotated drawing a bounding box around each of the 16 object to match the same artworks present in the synthetic set. This set contains 2190 images divided in 1502 for the training and 688 for the test;

Table 1: Statistics of the proposed dataset for unsupervised multi-target domain adaptation. The average occupied area (last column) is the average percentage of the image occupied by the bounding boxes of the considered object class.

Object Instances	Synthetic Domain (Source)		HoloLens Domain (Target)		GoPro Domain (Target)		Total object instances for each class	Average occupied area
	Training	Test	Training	Test	Training	Test		
Annunciazione	1301	605	191	69	211	74	2451	42.87%
Libro d'Ore miniato	1628	722	105	30	146	42	2673	8.02%
Lastra tombale di Giovanni Cabastida	2313	1181	200	100	247	114	4155	24.58%
Madonna del Cardillo	2345	1264	106	40	166	66	3987	9.74%
Disputa di San Tommaso	2202	965	100	46	155	67	3535	28.17%
Traslazione della Santa Casa	1904	964	161	46	225	71	3371	22.24%
Madonna col Bambino	2135	1044	119	47	161	46	3552	21.93%
L'immacolata Concezione e Dio Padre in Gloria	2557	1139	77	39	100	54	3966	35.70%
Adorazione dei Magi	1517	478	64	36	69	39	2203	30.35%
Sant'Elena e Costantino e Madonna con Bambino in gloria fra angeli	3285	1031	94	44	153	61	4668	33.72%
Taccuini di disegni	1617	513	59	33	75	39	2336	22.34%
Martirio di S. Lucia	3567	2353	106	36	184	45	6291	22.55%
Volto di Cristo	990	519	25	26	50	36	1646	11.74%
Dipinti di Sant'Orsola	2721	1897	83	69	125	86	4981	30.56%
Immacolata e i santi Chiara, Francesco, Antonio, Abate, Barbara e Maria Maddalena	3824	2424	104	69	187	89	6697	32.36%
Storia della Genesi	927	375	55	14	57	15	1443	22.79%
Total object instances for each split	34833	17474	1649	744	2311	944		

- Target images acquired using a Gopro: the dataset was created similarly to the previous one HoloLens. The images have been collected by the same visitor which have visited the site wearing both HoloLens and GoPro wearable cameras. This set contains a total of 2707 images splitted into 1911 for the training and 796 for the test.

Table 1 shows the distribution of the object instances in the proposed dataset. As can be noted, the HoloLens and GoPro domains have a number of object instances less than ten times smaller than the synthetic domain. The table also highlights that the proposed dataset is challenging for domain adaption for object detection due to the average size of each object. Indeed, the biggest object present in the dataset occupies only the 42.87% of the images' area while the smallest occupies 8.02% of the frame.

4. Methods

In this section, we first give a formal definition of the considered problem. We then discuss the compared methods and present the proposed one.

4.1. Problem Definition

Let be $S = \{(x_s^n, y_s^n)\}_{n=1}^{N_s}$ the set of N_s labeled images related to the source domain where x_s^n indicates the n^{th} source image and y_s^n the corresponding annotation. Let $T = \{T_1, T_2, \dots, T_D\}$ be the set of targets domains where $T_i = \{x_{T_i}^n\}_{n=1}^{D_{T_i}}$ corresponds to the T_i^{th} target domain. We set $D_{T_i} = 2$ in our experiments. The goal of unsupervised multi-camera domain adaptation for object detection is to maximize the mAP of the object detector across all the target

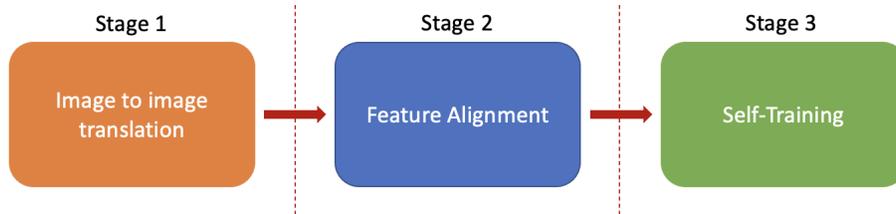


Figure 3: Pipeline of the proposed method. In the Stage 1 synthetic domains is translated to the real domains to reduce the gap at pixel level. In the Stage 2 the gap at feature level is reduced using feature alignment method. In Stage 3 an iterative self-training procedure is used to produce pseudo labels for the target domains.

domains training only on labeled images of the source domain and unlabeled images of the target domains.

4.2. Baselines without domain adaptation

We analyze the behaviour of two state-of-the-art object detectors: RetinaNet [33] and Faster RCNN [21]. We train and test both detectors on the target domains to produce “Oracle results” and assess the performance drop observed when the algorithms are trained on synthetic images and tested on the real images of the target domains.

4.3. Domain adaptation based on feature alignment

State-of-the-art domain adaptation methods for object detectors commonly consider only one source domain and one target domain. To study whether these state-of-the-art methods can be used to tackle multi-camera domain adaptation, we consider a naive approach which merges the two target domains into a single one. In particular, we considered the following unsupervised domain adaptation methods for object detection: DA-Faster RCNN [5], Strong Weak [22], DA-RetinaNet [34] and CDSSL [25].

4.4. Domain adaptation through feature alignment and image to image translation

Feature alignment methods aim to reduce the difference between source and target domains at the feature level without taking into account the difference at pixel level (like style, color, shape etc.) which are present between the source and targets domains. For this reason, we combine feature alignment methods with image to image translation methods to reduce the gap also at the pixel level. For the image to image translation task we used the CycleGAN algorithm [14] to translate synthetic images to real.



Figure 4: Qualitative results obtained using CycleGAN as image to image translation method (Stage 1). Synthetic images (left) are translated to the merged target domains (center). Real images similar to the translated ones are also reported for reference (right).

4.5. Proposed Method

The training of the proposed method comprises three stages that will be discussed in order of execution in the following sections. Each of them contributes to improving the performance of the object detector and works to adapt the two distribution at different levels. Figure 3 shows an overview of the general pipeline of the proposed method.

4.5.1. Image to Image Translation

Synthetic images generated from a model acquired using a 3D scanner, such as Matterport³, differ in general from real images in the style and shape of the object which can affect object detection performance. To reduce this diversity, the first step of our method consists in mitigating the style and shape differences using an image to image translation method. In particular, we used CycleGAN to transform training synthetic images into the real. In the later stages of our pipeline, the object detection model will be trained on the transformed images and tested directly on the real images. This step is optional in our pipeline for two reasons: 1) it can be computationally expensive when the datasets are large; 2) when the target and the source domain are not too similar, this transformation can be not sufficiently accurate. Figure 4 shows some qualitative results of this translation. As can be noted, the transformed images look more similar to the real counterpart after the transformation.

4.5.2. Feature Alignment

Although image to image translation can be used to reduce the differences in terms of style and shape, the features extracted from the two domains can

³<https://matterport.com/>

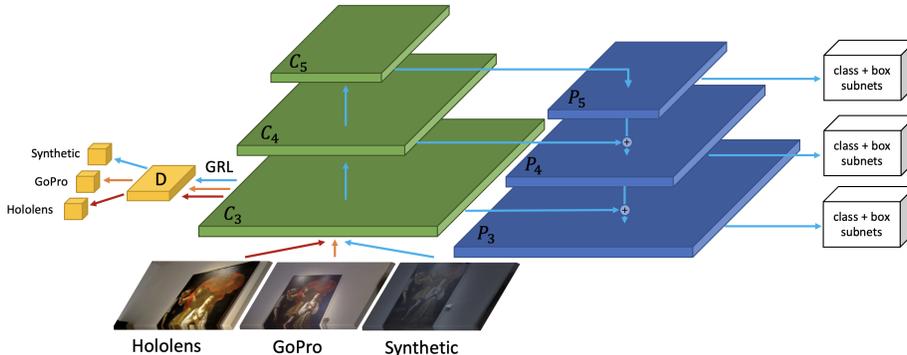


Figure 5: Architecture of the proposed MDA-RetinaNet model.

still be different. For this reason, in this second stage we propose an object detection architecture which jointly adapts the features during the training. State-of-the-art domain adaptation methods for object detection do not consider the existence of multiple target domains. To take advantage of multiple unlabeled target domains during the training, we propose a model, that we called MDA-RetinaNet, to address the problem of unsupervised domain adaptation for object detection based on adversarial learning [6]. Figure 5 shows the architecture of the proposed method which builds on RetinaNet [33]. To reduce the domain gap present at the feature level, we attach a domain discriminator with a gradient reversal layer to the feature map C_3 obtained from the ResNet backbone [35]. In particular, to adapt multiple domains (in this case 1 source and 2 targets in our experiments) we consider a multi-class classifier D which discriminates among all of them. The discriminator has 3 convolutional layers with kernel size equal to 1, followed by a ReLU activation function. Following [6], we place a gradient reversal layer at the input of the discriminator and train the model by minimizing the following loss function:

$$L = L_{class} + L_{box} - \lambda(L_D)$$

where L_{class} and L_{box} are the regression losses of RetinaNet, L_D is the loss of the discriminator module and λ is an hyper-parameter that balances the object detection and domain adaptation losses. This approach differs from standard methods that use a binary classifier used to only discriminate features belonging to the source and target domains, hence ignoring the presence of multiple targets. We hypothesize that, providing a multi-classes discriminator, the model will learn to extract features which are not only indistinguishable across synthetic and real domains, but also indistinguishable across the different real cameras. It is important to highlight that this type of adaptation allow to learn a combination of weights that extract feature maps using the backbone that generalize to the different domains. No adaptation is directly enforced for the layers involved in the classification and regression of the bounding boxes (Figure 5 white modules).

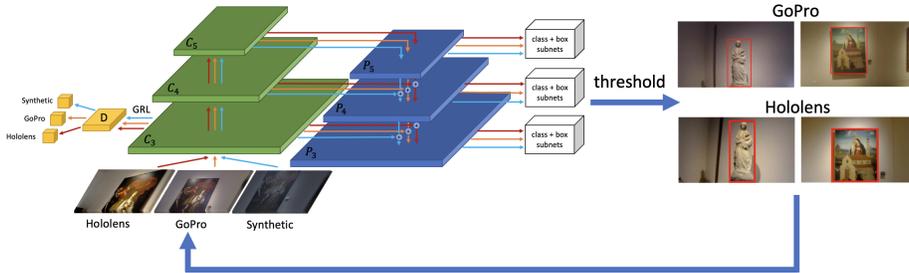


Figure 6: Self-training module for MDA-RetinaNet.

Algorithm 1: Proposed multi-target domain adaptation for object detection algorithm.

Input: $S = \{x_s^n, y_s^n\}$ the source domain, $T = \{T_1, T_2, \dots, T_D\}$ target domains, $converge = false$;

Step 1: transform the set S into T using CycleGAN;

Step 2: train MDA-RetinaNet using S' and T to adapt the features;

Step 3: set the threshold $t = 0.75$ and produce the pseudo labels y_d^n for each images in each target T_1, T_2, \dots, T_D using MDA-RetinaNet;

Step 4:

```

while !converge do
    train MDA-RetinaNet using  $y_d^n$ ;
    produce the new pseudo label  $y_d^n$ ;
    if  $t < 0.9$  then
        |  $t = t + 0.05$ ;
    else
        |  $converge = true$ ;
    end
end

```

end

Output: B - the set of predicted bounding boxes.

4.5.3. Self-Training

As noted in the previous section 4.5.2, the adaptation provided in the second stage is at the level of the features extracted by the backbone, whereas the classification and regression layers which detect the objects are trained only with the synthetic images due to the absence of labels for the real images. To tackle this limitation, we generate pseudo labels for the real images by retraining the model with the predictions on real images above a given confidence threshold produced by the model trained at the second stage of our pipeline (see Figure 3). This allows to train MDA-RetinaNet in a supervised way as illustrated in Figure 6 by exploiting the obtained pseudo labels. This latest module is trained in an iterative way, gradually increasing the threshold used to generate the pseudo labels to reduce the error of potentially wrong predicted labels. Algorithm 1 reports the complete procedure of the proposed method.

4.6. Experimental Settings

All the compared models were trained for 60K iterations using weights pre-trained on ImageNet [36]. We set the learning rate to 0.0002 for the first 30K iterations, then we multiply it by 0.1 for the remains 30K iterations. DA-Faster RCNN, Strong Weak and CDSSL were trained with the same parameters proposed by the authors in their respective papers [5, 22, 25]. The batch size was set to 4 for RetinaNet, 6 for DA-RetinaNet (4 source and 2 target images, 1 HoloLens and 1 GoPro) and 8 for MDA-RetinaNet⁴ (4 source and 4 target images divided in 2 HoloLens and 2 GoPro images). MDA-RetinaNet was implemented using Detectron2 [37]. To reduce the noise of the initial training of the Discriminator D, we adapt the λ hyperparameter following the update rule proposed in [6]. The second stage (Section 4.5.2) is performed only one time to produce initial pseudo labels. The Self-Training stage (Section 4.5.3) is executed 4 times gradually increasing the threshold used to generate the new pseudo labels which will be used in the next iteration. CycleGAN was trained for 60 epochs using the default parameters.

5. Results

This Section reports and analyzes the results of the experimental analysis.

5.1. Feature Alignment Results

Table 2 reports the results of the feature alignment based models. The first two rows show the results of the baseline Faster RCNN and RetinaNet modules trained with synthetic images and tested on HoloLens and GoPro without any domain adaptation technique. It is worth noting that RetinaNet is less sensitive to the domain gap, obtaining an mAP $\sim 7\%$ higher than Faster RCNN (14.10% vs 7.61% on HoloLens and 30.39% vs 37.13% on GoPro). For this reason, we focused our further experiments considering RetinaNet as backbone for the object detector in our proposed methods. The second group of rows (rows 3-6 of Table 2) report the results of state-of-the-art methods adapted for this specific task. In particular, due to the fact that these methods are able to work only with a single target, we merged the HoloLens and GoPro datasets into one. The proposed MDA-RetinaNet performs better than the other models and outperforms the best state-of-the-art method, DA-RetinaNet, by $\sim 3\%$ for HoloLens (34.97% vs 31.63%) and $\sim 2\%$ for GoPro (50.81% vs 48.37%). The last row shows the results of MDA-RetinaNet combined with the self-training procedure. As the results highlight, this combination allows to increase the performances of $\sim 23\%$ if compared with DA-RetinaNet (54.36% vs 31.63%) for HoloLens, $\sim 11\%$ (59.51% vs 48.37%) for GoPro and $\sim 20\%$ if compared with MDA-RetinaNet without self-training (54.36% vs 34.94%) for HoloLens and $\sim 9\%$ (59.51% vs 50.81%) for GoPro. Furthermore, the performance gap between HoloLens and GoPro with this last model it is almost negligible.

⁴code available at <https://github.com/fpv-iplab/STMDA-RetinaNet>

Table 2: Results of baseline and feature alignment methods. S refers to Synthetic, H refers to HoloLens and G to GoPro. ST indicates the self-training procedure.

Model	Source	Target	Test H	Test G
Faster RCNN [21]	S	-	7.61%	30.39%
RetinaNet [33]	S	-	14.10%	37.13%
DA-Faster RCNN [5]	S	H+G	10.53%	48.23%
Strong Weak [22]	S	H+G	26.68%	48.55%
CDSSL [25]	S	H+G	28.66%	45.33%
DA-RetinaNet [34]	S	H+G	31.63%	48.37%
MDA-RetinaNet	S	H, G	34.97%	50.81%
MDA-RetinaNet + ST	S	H, G	54.36%	59.51%
Faster RCNN [21] (Oracle)	H	-	91.97%	76.88%
Faster RCNN [21] (Oracle)	G	-	68.65%	89.21%
RetinaNet [33] (Oracle)	H	-	92.44%	77.96%
RetinaNet [33] (Oracle)	G	-	69.70%	89.69%

5.2. Feature Alignment and Image to Image translation Results

Table 3 shows the results obtained combining the baseline and feature alignment methods with CycleGAN. The first two rows report the results of Faster RCNN and RetinaNet when trained on synthetic images transformed to the merged HoloLens and GoPro domain. As can be noted, pixel level domain adaptation allows to significantly increase the performance of Faster RCNN and RetinaNet respectively by about 8% (7.61% vs 15.34%) and 16% (14.10% vs 31.43%) on HoloLens and by about 33% (30.39% vs 63.60%) and 32% (37.13% vs 69.59%) on GoPro, reducing the gap between synthetic and real images. The middle part of the table shows the results of the methods based on feature alignment. Also in this case, MDA-RetinaNet achieves a higher mAP with respect to the best state-of-the-art method, CDSSL, (53.06% vs 58.11% for HoloLens and 71.17% vs 71.39% for GoPro) which further improves if we introduce the self-training procedure (58.11% vs 66.64% for HoloLens and 71.39% vs 72.22% for GoPro). It is worth noting that, with self-training the gap in performances between HoloLens and GoPro is reduced from $\sim 13\%$ to $\sim 6\%$ which suggest that the model acquires knowledge from the GoPro images that is useful to detect object in the HoloLens domain. Furthermore, the performance of MDA-RetinaNet with self-training is really close to the performance of the RetinaNet oracles when trained with the labeled HoloLens domain and tested on GoPro and vice versa (66.64% vs 69.70% for HoloLens and 72.22% vs 77.96% for GoPro). However, there is still space of improvement if we consider the performances of the oracles trained and tested in their respective domains, which makes proposed dataset still challenging (66.64% vs 92.44% for HoloLens and 72.22% vs 89.69% for GoPro).

5.3. Ablation Study

Table 4 reports the ablation study of the proposed MDA-RetinaNet model and compares the results with respect to the DA-RetinaNet architecture. We evaluated the models on HoloLens domain, which is more challenging if compared to GoPro, analyzing the impact of the placement of the discriminator at different levels of the feature map extracted from the RetinaNet backbone (see Figure 5 on the paper). As can be noted, each single discriminator increases the performances of the standard RetinaNet architecture and obtain better performances than DA-RetinaNet. The discriminator attached to the first feature map C_3 , allows to achieve better results than the other two discriminators attached to the C_4 and C_5 feature maps. Moreover, considering more than one discriminator to align the feature at different levels does lead to obtain better results in our experiments as in the case of the single domain DA-RetinaNet but only decreases the performance. The best combination and optimal number of discriminators was found empirically and, as shown in Table 4, it is achieved using only one discriminator at the C_3 level. We hypothesize that considering more discriminators at the same time could unbalance the models training, obtaining features that are aligned but less effective for the main object detection task. In Table 4 we also report an ablation study of the impact of each discriminator attached at P_i or at C_i levels. As can be noted, in each case, the performances achieved by the models that use the discriminator at P_i levels are lower than their counterparts which use discriminators at the C_i levels. Table 5 shows the results obtained with different linear schedules of the values of the threshold. We noted that, due to the domain gap between source and target domains, it is convenient to use a low threshold in the first iterations of self-training, where a set of initial pseudo-labels is needed, and increasing this threshold to an higher value as training proceeds. Indeed, we achieve best results for using a threshold value starting at 0.75 and ending at 0.9. Table 6 reports the results of adapting DA-Faster RCNN [5] and Strong Weak [22] to multiple target domains using the same methodology proposed for MDA-RetinaNet. Specifically, instead of merging the to dataset into one and use the binary discriminator proposed by the authors in their papers, we replaced it with our multi classes discriminator and considered the target domains individually instead of merging them. As can be noted, the performances of the other two methods improves by 3-4% if compared with the results of Table 2. Nevertheless, the best results are still obtained by the proposed MDA-RetinaNet architecture. These results suggest that using a multi class discriminator instead of a binary discriminator allows to consistently improve performances with different architectures.

5.4. Comparison between MDA-RetinaNet and DA-RetinaNet

Table 7 compares the results of the proposed MDA-RetinaNet with DA-RetinaNet. It is worth noting that training the model using only one target domain at a time results in worse performance in both domains despite they are very similar. This happens because the model overfits with respect to the considered target domain used for training. Using both domains during training,

Table 3: Results of feature alignment methods combined with CycleGAN. H refers to HoloLens while G to GoPro. “{G, H}” refers to synthetic images translated to the merged HoloLens and GoPro domains. ST indicates self-training procedure.

Model	Source	Target	Test H	Test G
Faster RCNN [21]	{G, H}	-	15.34%	63.60%
RetinaNet [33]	{G, H}	-	31.43%	69.59%
DA-Faster RCNN [5]	{G, H}	H+G	32.13%	65.19%
Strong Weak [22]	{G, H}	H+G	41.11%	66.45%
DA-RetinaNet [34]	{G, H}	H+G	52.07%	71.14%
CDSSL [25]	{G, H}	H+G	53.06%	71.17%
MDA-RetinaNet	{G, H}	H, G	58.11%	71.39%
MDA-RetinaNet + ST	{G, H}	H, G	66.64%	72.22%
Faster RCNN [21] (Oracle)	H	-	91.97%	76.88%
Faster RCNN [21] (Oracle)	G	-	68.65%	89.21%
RetinaNet [33] (Oracle)	H	-	92.44%	77.96%
RetinaNet [33] (Oracle)	G	-	69.70%	89.69%

Table 4: Ablation study about the impact of each discriminator D_i and comparison between each discriminator D_i placed at C_i and P_i level.

Model	C_3	P_3	C_4	P_4	C_5	P_5	mAP
RetinaNet							14.10%
DA-RetinaNet					✓		15.84%
MDA-RetinaNet					✓		19.54%
MDA-RetinaNet						✓	16.29%
DA-RetinaNet			✓				16.38%
MDA-RetinaNet			✓				19.88%
MDA-RetinaNet				✓			17.01%
DA-RetinaNet	✓						28.61%
MDA-RetinaNet	✓						34.97%
MDA-RetinaNet		✓					31.44%
DA-RetinaNet	✓		✓				30.52%
MDA-RetinaNet	✓		✓				34.09%
MDA-RetinaNet		✓		✓			30.85%
DA-RetinaNet	✓		✓		✓		31.04%
MDA-RetinaNet	✓		✓		✓		32.11%
MDA-RetinaNet		✓		✓		✓	30.18%

as the proposed MDA-RetinaNet model does, allows to generalize over both target domains with a single model, which also results in improved performance.

Table 5: Comparison performance considering different threshold.

Model	Threshold	Test H	Test G
MDA-RetinaNet + ST	0.90	47.48%	52.25%
MDA-RetinaNet + ST	0.85 to 90	49.21%	54.90%
MDA-RetinaNet + ST	0.80 to 0.90	52.49%	57.67%
MDA-RetinaNet + ST	0.75 to 0.90	54.36%	59.51%

Table 6: Comparison between DA-Faster RCNN, Strong Weak and MDA-RetinaNet when modified using multiclass discriminators. S refers to Synthetic, H refers to Hololens and G to GoPro.

Model	Source	Target	Test H	Test G
DA-Faster RCNN [5]	S	H, G	13.79%	48.35%
Strong Weak [22]	S	H, G	29.52%	49.06%
MDA-RetinaNet	S	H, G	34.97%	50.81%

5.5. Qualitative Results

Figure 7 compares some qualitative detection results obtained by the proposed MDA-RetinaNet with and without Self-Training with respect to RetinaNet baseline (the ground truth is the blue bounding box). RetinaNet fails the detection in many cases. Indeed, it does not detect any artwork or produce a wrong classification and/or regression. MDA-RetinaNet well recognize small and large artworks but fails in the last two rows. MDA-RetinaNet with Self-Training improve the performance of the standard RetinaNet and MDA-RetinaNet with a more accurate detection of the artworks.

6. Conclusion

We studied the problem of unsupervised multi-camera domain adaptation for object detection in cultural sites. To perform the study, we have collected and publicly released a new challenging dataset with the aim to encourage the community to continue researching on the problem. We proposed a new method which combines feature alignment, pixel level and self-training methods that outperforms current state-of-the-art methods.

Acknowledgments

This research has been supported by the project VALUE (N. 08CT6209090207 - CUP G69J18001060007) - PO FESR 2014/2020 - Azione 1.1.5., by the project MEGABIT - Research Program Pia.ce.ri. 2020/2022 Linea 2 - University of Catania, and by Project HERO - Huge and Easy Reproduction of Objects.

Table 7: Comparison between DA-RetinaNet trained using one target set at a time and MDA-RetinaNet. S refers to Synthetic, H refers to Hololens and G to GoPro.

Model	Source	Target	Test H	Test G
RetinaNet [33]	S	-	14.10%	37.13%
DA-RetinaNet [34]	S	H	31.01%	36.60%
DA-RetinaNet [34]	S	G	21.63%	45.86%
MDA-RetinaNet	S	H, G	34.97%	50.81%

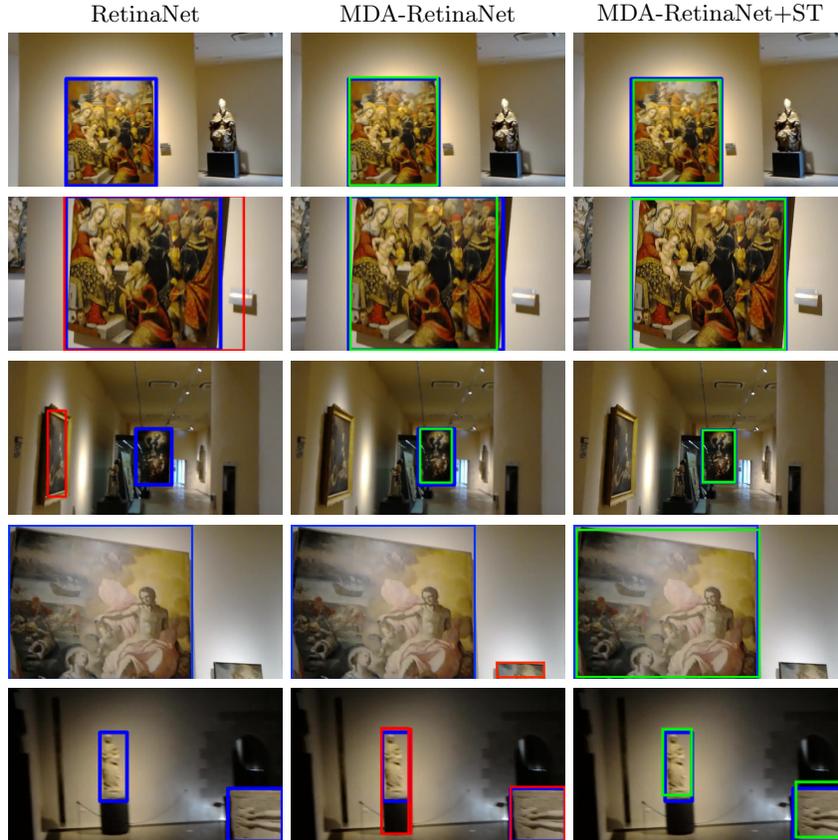


Figure 7: Qualitative results of RetinaNet, MDA-RetinaNet and MDA-RetinaNet with self-training (ST). The blue box represents ground truth, the red box indicates a wrong detection (object localization or classification), the green box represents correct detections.

References

- [1] L. Seidenari, C. Baccchi, T. Uricchio, A. Ferracani, M. Bertini, A. D. Bimbo, Deep artwork detection and retrieval for automatic context-aware audio guides, *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 13 (3s) (2017) 1–21.

- [2] R. Cucchiara, A. Del Bimbo, Visions for augmented cultural heritage experience, *IEEE MultiMedia* 21 (1) (2014) 74–82.
- [3] G. M. Farinella, G. Signorello, S. Battiato, A. Furnari, F. Ragusa, R. Leonardi, E. Ragusa, E. Scuderi, A. Lopes, L. Santo, M. Samarotto, VEDI: Vision exploitation for data interpretation, in: *International Conference on Image Analysis and Processing (ICIAP)*, 2019.
URL http://iplab.dmi.unict.it/VEDI_project/
- [4] S. A. Orlando, A. Furnari, G. M. Farinella, Egocentric visitor localization and artwork detection in cultural sites using synthetic data, *Pattern Recognition Letters* 133 (2020) 17–24.
- [5] Y. Chen, W. Li, C. Sakaridis, D. Dai, 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.
- [6] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by backpropagation, in: *International conference on machine learning*, 2015, pp. 1180–1189.
- [7] M. Portaz, M. Kohl, G. Quénot, J.-P. Chevallet, Fully convolutional network and region proposal for instance identification with egocentric vision, in: *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 2383–2391.
- [8] F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella, Ego-ch: Dataset and fundamental tasks for visitors behavioral understanding using egocentric vision, *Pattern Recognition Letters* 131 (2020) 150–157.
- [9] A. Rozantsev, M. Salzmann, P. Fua, Beyond sharing weights for deep domain adaptation, *IEEE transactions on pattern analysis and machine intelligence* 41 (4) (2018) 801–814.
- [10] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, A. J. Smola, A kernel method for the two-sample-problem, in: *Advances in neural information processing systems*, 2007, pp. 513–520.
- [11] B. Sun, J. Feng, K. Saenko, Correlation alignment for unsupervised domain adaptation, in: *Domain Adaptation in Computer Vision Applications*, Springer, 2017, pp. 153–171.
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [13] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7167–7176.

- [14] J.-Y. Zhu, T. Park, P. Isola, 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.
- [15] Z. Du, J. Li, H. Su, L. Zhu, K. Lu, Cross-domain gradient discrepancy minimization for unsupervised domain adaptation, in: CVPR, 2021.
- [16] H. Zhao, S. Zhang, G. Wu, J. M. Moura, J. P. Costeira, G. J. Gordon, Adversarial multiple source domain adaptation, *Advances in neural information processing systems* 31 (2018) 8559–8570.
- [17] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, B. Wang, Moment matching for multi-source domain adaptation, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 1406–1415.
- [18] B. Gholami, P. Sahu, O. Rudovic, K. Bousmalis, V. Pavlovic, Unsupervised multi-target domain adaptation: An information theoretic approach, *IEEE Transactions on Image Processing* 29 (2020) 3993–4002.
- [19] M. Ragab, Z. Chen, M. Wu, H. Li, C.-K. Kwok, R. Yan, X. Li, Adversarial multiple-target domain adaptation for fault classification, *IEEE Transactions on Instrumentation and Measurement* 70 (2020) 1–11.
- [20] L. T. Nguyen-Meidine, A. Belal, M. Kiran, J. Dolz, L.-A. Blais-Morin, E. Granger, Unsupervised multi-target domain adaptation through knowledge distillation, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2021, pp. 1339–1347.
- [21] S. Ren, K. He, R. Girshick, 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.
- [22] K. Saito, Y. Ushiku, T. Harada, 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.
- [23] R. Xie, F. Yu, J. Wang, Y. Wang, 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.
- [24] T. Kim, M. Jeong, S. Kim, S. Choi, 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. 12456–12465.
- [25] F. Yu, D. Wang, Y. Chen, N. Karianakis, P. Yu, D. Lymberopoulos, X. Chen, Unsupervised domain adaptation for object detection via cross-domain semi-supervised learning, *ArXiv abs/1911.07158*.

- [26] S. Kim, J. Choi, T. Kim, C. Kim, Self-training and adversarial background regularization for unsupervised domain adaptive one-stage object detection, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 6091–6100. doi:10.1109/ICCV.2019.00619.
- [27] H.-K. Hsu, C.-H. Yao, Y.-H. Tsai, W.-C. Hung, H.-Y. Tseng, M. Singh, M.-H. Yang, Progressive domain adaptation for object detection, in: 2020 IEEE Winter Conference on Applications of Computer Vision (WACV), 2020, pp. 738–746. doi:10.1109/WACV45572.2020.9093358.
- [28] K. Fujii, K. Kawamoto, Generative and self-supervised domain adaptation for one-stage object detection, Array 11 (2021) 100071. doi:https://doi.org/10.1016/j.array.2021.100071.
URL <https://www.sciencedirect.com/science/article/pii/S2590005621000199>
- [29] C. Chen, Z. Zheng, Y. Huang, X. Ding, Y. Yu, I3net: Implicit instance-invariant network for adapting one-stage object detectors, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- [30] Vidit, M. Salzmann, Attention-based domain adaptation for single stage detectors, ArXiv abs/2106.07283.
- [31] D. Guan, J. Huang, A. Xiao, S. Lu, Y. Cao, Uncertainty-aware unsupervised domain adaptation in object detection, ArXiv abs/2103.00236.
- [32] V. Vibashan, V. Gupta, P. Oza, V. A. Sindagi, V. Patel, Mega-cda: Memory guided attention for category-aware unsupervised domain adaptive object detection, in: CVPR, 2021.
- [33] T.-Y. Lin, P. Goyal, R. Girshick, K. He, P. Dollár, Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2980–2988.
- [34] G. Pasqualino, A. Furnari, G. Signorello, G. M. Farinella, An unsupervised domain adaptation scheme for single-stage artwork recognition in cultural sites, Image and Vision Computing 107 (2021) 104098. doi:https://doi.org/10.1016/j.imavis.2021.104098.
URL <https://www.sciencedirect.com/science/article/pii/S0262885621000032>
- [35] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [36] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, 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.
- [37] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, R. Girshick, Detectron2 (2019).
URL <https://github.com/facebookresearch/detectron2>