

Attending to Discriminative Certainty for Domain Adaptation

Vinod Kumar Kurmi*, Shanu Kumar*, Vinay P. Namboodiri
 Indian Institute of Technology Kanpur
 India

{vinodkk, sshanu, vinaypn}@iitk.ac.in

Abstract

In this paper, we aim to solve for unsupervised domain adaptation of classifiers where we have access to label information for the source domain while these are not available for a target domain. While various methods have been proposed for solving these including adversarial discriminator based methods, most approaches have focused on the entire image based domain adaptation. In an image, there would be regions that can be adapted better, for instance, the foreground object may be similar in nature. To obtain such regions, we propose methods that consider the probabilistic certainty estimate of various regions and specify focus on these during classification for adaptation. We observe that just by incorporating the probabilistic certainty of the discriminator while training the classifier, we are able to obtain state of the art results on various datasets as compared against all the recent methods. We provide a thorough empirical analysis of the method by providing ablation analysis, statistical significance test, and visualization of the attention maps and t-SNE embeddings. These evaluations convincingly demonstrate the effectiveness of the proposed approach.

1. Introduction

With the advent of deep learning, there has been substantial progress for solving image classification tasks with state of the art methods obtaining lesser than 3% error (on top five results) on the imagenet dataset. However, it was observed that these results do not transfer to other datasets [25] due to the effect of dataset bias [44]. The classifiers trained on one dataset (termed source dataset) show a significant drop in accuracy when tested on another dataset (termed target dataset). To address this issue, some methods have been proposed for adapting domains [7, 50]. One of the more successful approaches towards addressing this domain shift has been based on the adversarial adaptation of features us-

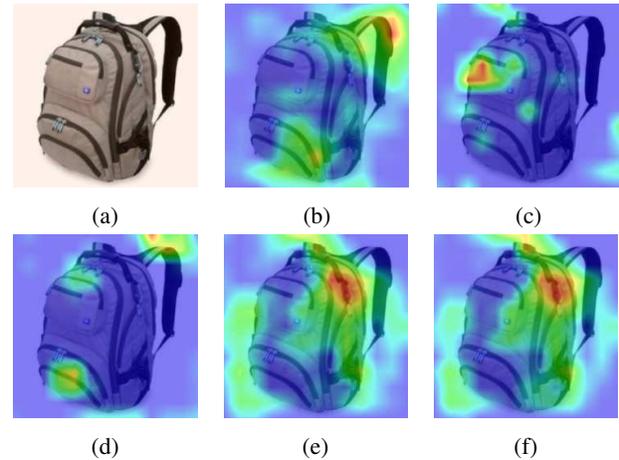


Figure 1: Visualization of the uncertainty and certainty maps of the discriminator during the midst of training is provided for the input image (a). The aleatoric and predictive uncertain regions of the discriminator are shown in image (b) and (d). While aleatoric and predictive certain regions of the discriminator are shown in (c) and (e). From the figure, it is clear that the certain regions of the discriminator during training mostly corresponds to (f) classifier’s activation map based on true label at the end of training.

ing an adversarial discriminator [12] that aims to distinguish between samples drawn from the source and target datasets. Due to the adversarial training, the feature representations are brought close such that the discriminator is not able to distinguish between samples from source and target dataset. However, all approaches that are based on this idea consider the whole image as being adapted. This usually is not the case as there are predominant regions in an image that may be better adapted and useful for improving classification on target dataset. We address this issue and propose a simple approach for solving this problem.

To specify regions that can be adapted we propose the use of certainty of a probabilistic discriminator. During training, we identify regions where the discriminator is certain, i.e., the probabilistic uncertainty for these regions is

*Equal contributions from both authors.

low. These regions can be adapted because there exists a clear distinction between the source and target regions. Figure 1 shows that using measures such as data uncertainty (known as aleatoric uncertainty) [19] and predictive uncertainty [24], we can obtain regions that can be adapted better. We also observe from the Figure 1 and 5 that for most of the duration during training, discriminator is certain on the foreground regions, as the foreground is hard to adapt. Hence, when the classifier is trained with the emphasis being placed on these regions, then we observe that the classifier focuses on these regions during prediction and therefore generalizes better on target dataset. Quantitatively we observe results that are up to 7.4% better than the current state of the art methods on Office-Home dataset [49].

To summarize, through this paper we make the following contributions:

- We propose a method to identify adaptable regions using the certainty estimate of the discriminator, and this is evaluated using various certainty estimates.
- We use these certainty estimate weights for improving the classifier performance on target dataset by focusing the training of the classifier on the adaptable regions of the source dataset.
- We provide a thorough evaluation of the method by considering detailed comparison on standard benchmark datasets against the state of the art methods and also provide an empirical analysis using statistical significance analysis, visualization and ablation analysis of the proposed method.
- An additional observation is that by using Bayesian classifiers we also improve the robustness of the classifier in addition to obtaining certainties of classification accuracy. This aids in better understanding of the proposed method.

2. Literature Survey

Some studies have examined different adaptation methods. One study by [47] examined domain adaptation by minimizing the maximum mean discrepancy distance. The maximum mean discrepancy based approaches were further extended to multi-kernel MMD in [25]. In adversarial learning framework [12] has proposed a method to minimize the source target discrepancy using a gradient reversal layer at discriminator. Recently many adversarial methods have been applied in the domain adaptation task to bring the source and target distribution closer. Adversarial discriminative domain adaptation [46] considers an inverted label GAN loss. Wasserstein distance based discriminator was used in [39] to bring the two distributions closer. Domain confusion network [45] was also used to solve the adaptation problem in two domains by minimizing the discriminate distance between two domains. The discriminative feedback of the discriminator also applied

in the paraphrase generation problem [33]. Another adversarial discriminator based model is [34], where multiple discriminators (MADA) have been used to solve the mode collapse problem in the domain adaptation. Some works closely related to MADA have been proposed in [4, 3]. The labeled discriminator [23] used to tackle the mode collapse problem in domain adaptation. The adversarial domain adaptation also explored in scene graph [22]. Other source and target discrepancy minimization based methods such as [37, 54] also address the domain adaptation problem. [32, 31] have proposed an exemplar based discrepancy minimization method. Recently [2, 5] have applied generative adversarial network [14] for the domain adaptation problems. Image generation methods are used to adapt the source and target domain [13, 38, 30]. Other work in [16] and [42] have used Cycle consistency [55] loss and deep coral loss [41, 40] for ensuring the closeness of source and target domain respectively. Deep Bayesian models have been used to play an important role in the estimation of the deep model uncertainty. The Bayesian formulation in domain adaptation natural language processing has been proposed in [9]. In [11], it has been justified that dropout can also work as an approximation of the deep Bayesian network. Works on the uncertainty estimation have been reported in [10, 19]. In [29] predictive uncertainty has been calculated over the prior networks. Uncertainty in the ensemble model along with adversarial training has been discussed in [24]. Another work on the Bayesian uncertainty estimation has been reported in the [43].

Attention-based networks have been widely applied in many computer vision applications such as image captioning [51, 53], visual question answer [52, 31, 18] and speech recognition [6]. The advantage of the attention model is that it helps to learn some set of weights over a set of representation input which has relatively more importance than others. Recently [48] showed that the attention mechanism can also be achieved by dispensing with recurrence and convolutions. A recent work [17] addresses the domain adaptation problem by obtaining the synthetic source and target images from CycleGAN [55], and then aligned the attention map of all the pairs.

3. Methodology

In the unsupervised domain adaptation problem, the source dataset $\mathcal{D}_s = (x_i^s, y_i^s)$ consists of data sample (x_i^s) with corresponding label (y_i^s) where $\mathcal{D}_s \in P_s$ and the target dataset $\mathcal{D}_t(x_i^t)$ consists of unlabeled data samples (x_i^t) where $\mathcal{D}_t \sim P_t$. P_s and P_t are the source and target distributions. We further assume that both the domains are complex and unknown. For solving this task, we are following the adversarial domain adaptation framework, where a discriminator is trained to learn domain invariant features domain invariant while a classifier is trained to learn class

discriminative features. In this paper, we are proposing a discriminator certainty based domain adaption model represented in the Figure 2, which consists of three major modules: Feature extractor, Bayesian Classifier, and Bayesian Discriminator. The feature extractor is pretrained on the Imagenet dataset, while both the classifier and discriminator are Bayesian neural networks (BNN). We have followed the approach defined in [11, 19, 20] for transforming deep neural networks into BNNs.

3.1. Bayesian Classifier

Bayesian framework is one of the efficient ways to predict uncertainty. Gal et.al [10] has shown that by applying dropout after every fully connected (fc) layer, we can perform probabilistic inference for deep neural networks. Hence we have followed a similar approach for defining the classifier. For estimating uncertainty, similar to [19], we trained the classifier to output class probabilities along with aleatoric uncertainty (data uncertainty). The predictive uncertainty includes both model uncertainty and data uncertainty, where model uncertainty results from uncertainty in model parameters. Estimation of aleatoric uncertainty for the classifier makes the features more robust for prediction, and estimation of predictive uncertainty provides a tool for visualizing model’s predictions.

For the input sample x_i , the feature extractor G_f outputs features \mathbf{f}_i , represented by $\mathbf{f}_i = G_f(x_i)$. The predicted class logits y_i^c and aleatoric uncertainty v_i^c are obtained as:

$$y_i^c = G_{cy}(G_c(\mathbf{f}_i)), \quad v_i^c = G_{cv}(G_c(\mathbf{f}_i)) \quad (1)$$

where G_{cy} and G_{cv} are the logits and aleatoric uncertainty prediction modules of the classifier G_c respectively. The classification loss for predicted logits is defined as:

$$\mathcal{L}_{cy} = \frac{1}{n_s} \sum_{x_i \in \mathcal{D}_s} \mathcal{L}(y_i^c, y_i) \quad (2)$$

where \mathcal{L} is the cross entropy loss function and y_i is the true class label for the input x_i . The total number of data samples in the source domain is denoted as n_s . The classifier aleatoric loss \mathcal{L}_{cv} for predicted uncertainty v_i^c is defined as:

$$\hat{y}_{i,t}^c = y_i^c + \sigma_i^c * \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, I) \\ \mathcal{L}_{cv} = -\frac{1}{n_s} \sum_{x_i \in \mathcal{D}_s} \log \frac{1}{T} \sum_t \mathcal{L}(\hat{y}_{i,t}^c, y_i) \quad (3)$$

where σ_i^c is the standard deviation, $v_i^c = (\sigma_i^c)^2$. The classifier is trained by jointly minimizing both the classification loss \mathcal{L}_{cy} and aleatoric loss \mathcal{L}_{cv} .

3.2. Bayesian Discriminator

In the proposed method, the discriminator is also modeled in the Bayesian framework similar to the Bayesian

classifier. The uncertainty in the discriminator network implies the region where it is uncertain about its prediction about the domain. The uncertainty estimation of the discriminator can guide the feature extractor more efficiently for domain adaptation. All real-world images contain some type of aleatoric uncertainty or noise. These regions which contain aleatoric uncertainty, are not adaptable. By aligning these uncertain regions, we are corrupting the feature representation, thus confusing the classifier during predictions for the target domain. So, by estimating aleatoric uncertainty, the discriminator is avoiding the learning of feature representations for these regions, which also reduces negative transfer [34]. The negative transfer introduces false alignment of the mode of two distributions across domains, which needs to be prevented during adaptation. Similarly, the predictive uncertainty tells us about the model’s incapability to classify the domains, as the discriminator is not sure about the domain. Predictive uncertainty occurs in the region where either it is already adapted, or there is noise which corresponds to aleatoric uncertainty. We obtain the discriminator predicated logits and variance using the following equations

$$y_i^d = G_{dy}(G_d(\mathbf{f}_i)), \quad v_i^d = G_{dv}(G_d(\mathbf{f}_i)) \quad (4)$$

where G_{dy} and G_{dv} predict domain class logits y_i^d and domain aleatoric uncertainty v_i^d respectively using features from G_d . The domain classification loss \mathcal{L}_{dy} is defined as:

$$\mathcal{L}_{dy} = \frac{1}{n_s + n_t} \sum_{x_i \in \mathcal{D}_s \cup \mathcal{D}_t} \mathcal{L}(y_i^d, d_i) \quad (5)$$

where \mathcal{L} is the cross entropy loss function, d_i is the true domain of the image, and n_s and n_t are the number of source and target samples. The domain label d_i is defined to be 0 if $x_i \in \mathcal{D}_s$ and 1 if $x_i \in \mathcal{D}_t$. The discriminator aleatoric loss \mathcal{L}_{dv} is defined as:

$$\hat{y}_{i,t}^d = y_i^d + \sigma_i^d * \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, I) \\ \mathcal{L}_{dv} = -\frac{1}{n_s + n_t} \sum_{x_i \in \mathcal{D}_s \cup \mathcal{D}_t} \log \frac{1}{T} \sum_t \mathcal{L}(\hat{y}_{i,t}^d, d_i) \quad (6)$$

where $v_i^d = (\sigma_i^d)^2$. Discriminator is trained by jointly minimizing both the domain classification loss \mathcal{L}_{dy} and discriminator aleatoric loss \mathcal{L}_{dv} .

3.3. Certainty Based Attention

Uncertainty estimation of the discriminator can help in identifying those regions which can be adapted, cannot be adapted, or already adapted. The regions which are already aligned will confuse the discriminator for predicting the domain. Hence discriminator will be highly uncertain on these regions. The discriminator will also be highly uncertain on

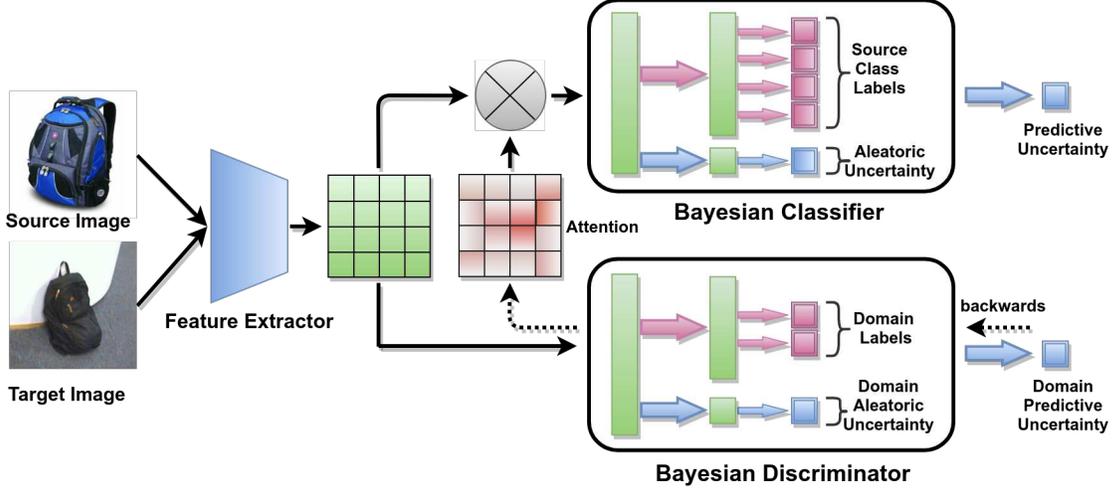


Figure 2: The architecture of Certainty based Attention for Domain Adaptation (CADA), consists of a shared feature extractor, Bayesian classifier and Bayesian discriminator where both the classifier and discriminator predict the variance value along with the prediction score. Discriminator’s predictive or aleatoric uncertainty is used to highlight the regions where the discriminator is certain about its predictions.

the regions containing aleatoric uncertainty, and these regions can’t be adapted. Uncertainty estimation can also help to identify regions where discriminator is certain or the regions which can be further aligned.

In most of the datasets, the discriminator can easily discriminate between the source and target images by only attending on the background during the initial phase of the training. Hence, the discriminator will be more certain in these regions, which results in easier adaptation of background regions after some adversarial training. But the foreground regions are difficult to adapt, as foreground varies a lot across all the classes and images. Therefore for most of the span during training, the discriminator will be certain on the transferable regions of the foreground. Thus, if the classifier attends to the certain regions of the discriminator, it will focus more on the transferable regions of the foreground during training.

The Certainty Attention based Domain Adaption (CADA) model is shown in Figure 2. In the proposed work, we propose the two variants of CADA: aleatoric certainty based attention (CADA-A) and predictive certainty based attention (CADA-P).

3.3.1 Aleatoric Certainty based Attention

The aleatoric uncertainty (v_i^d) of the domain discriminator occurs due to corruption or noise in the regions. These regions are not suited for the adaptation as well as for object classification and should be less focused as compared to the certain regions for the classification task. For identifying the aleatoric uncertain regions, we compute the gradient of aleatoric uncertainty with respect to the features f_i . These gradients ($\frac{\partial v_i^d}{\partial f_i}$) will flow back through the gradient rever-

sal layer, and will correspond to the gradients of aleatoric certainty, i.e., $-\frac{\partial v_i^d}{\partial f_i}$.

$$p_i = f_i * -\frac{\partial v_i^d}{\partial f_i} \quad (7)$$

Therefore the positive sign of the product of features and gradients of the aleatoric certainty denotes the positive influence of aleatoric certainty on these features. i.e the discriminator is certain on these regions.

$$a_i = \text{ReLU}(p_i) + c * \text{ReLU}(-p_i) \quad (8)$$

For obtaining the regions where the discriminator is certain, the product p_i is passed through a ReLU function. But for ignoring the negative values that corresponds to uncertain regions, $-p_i$ is again passed through a ReLU function. This is then multiplied with a large number c , such that after applying softmax over a_i all negative values of p_i become zero and all the positive values are normalized.

$$w_i = (1 - v_i^d) * \text{Softmax}(a_i) \quad (9)$$

To focus on more attentive (certain) regions, we follow the residual setting [27] for obtaining the effective weighted features h_i . Images with high aleatoric uncertainty (lower certainty) should be less attentive, and it is obtain by multiplying the normalized softmax attention value to its certainty value ($1 - v_i^d$) using the Eq. 9. The weighted feature h_i is generated as follows:

$$h_i = f_i * (1 + w_i) \quad (10)$$

3.3.2 Predictive Certainty based Attention

The predictive uncertainty measures the model’s capability for prediction. It occurs in the regions which are either

already domain invariant or which contain noise. The regions corresponding to discriminator’s predictive certainty are domain transferable regions and should be attended during classification. We follow the approach proposed in [19] for computing the predictive uncertainty of discriminator. It is obtained by the entropy of the average class probabilities of the Monte Carlo (MC) samples. We sample weights from \mathbf{G}_d , and perform MC simulations over domain probabilities $p(y_{i,t}^d)$ and aleatoric uncertainty $v_{i,t}^d$ for estimating predictive uncertainty $\text{Var}^d(f_i)$.

$$g_{i,t} = G_d^t(\mathbf{f}_i), \quad G_d^t \sim \mathbf{G}_d \quad (11)$$

Using the sampled model, we calculate domain probabilities and aleatoric uncertainty.

$$p(y_{i,t}^d) = \text{Softmax}(G_{dy}(g_{i,t})), \quad v_{i,t}^d = G_{dv}(g_{i,t}) \quad (12)$$

$$H(y_{i,t}^d) = - \sum_{c=1}^2 p(y_{i,t}^d = c) * \log(p(y_{i,t}^d = c)) \quad (13)$$

$$\text{Var}^d(f_i) \approx \frac{1}{T} \sum_{t=1}^T (v_{i,t}^d + H(y_{i,t}^d)) \quad (14)$$

where $H(y_{i,t}^d)$ is the Entropy of the $p(y_{i,t}^d)$. For identifying the predictive uncertain regions, we compute gradients of the predictive uncertainty $\text{Var}^d(f_i)$ of the discriminator with respect to the features f_i i.e. $\frac{\partial \text{Var}^d(f_i)}{\partial f_i}$, and negative of these gradients which returns through the gradient reversal layer will correspond to the gradients of the predictive certainty, i.e., $-\frac{\partial \text{Var}^d(f_i)}{\partial f_i}$.

$$p_i = f_i * -\frac{\partial \text{Var}^d(f_i)}{\partial f_i} \quad (15)$$

$$a_i = \text{ReLU}(p_i) + c * \text{ReLU}(-p_i) \quad (16)$$

$$w_i = (1 - \text{Var}^d(f_i)) * \text{Softmax}(a_i) \quad (17)$$

Similar to aleatoric certainty based attention, for obtaining the predictive certain regions we apply ReLU function to p_i and ignoring its negative values (corresponds to uncertain regions), a ReLU function is applied to negative of p_i , and multiply it with a large number c using Eq. 16. After applying the Softmax, the features that are activated by the predictive uncertainty will have zero weight and for features that are highly activated by the predictive certainty will get more weight. The residual weighted features are obtain by following equations

$$h_i = f_i * (1 + w_i) \quad (18)$$

Therefore the images with high predictive uncertainty have lower value of w_i , have less attention, while images have high predictive certainty have high value of w_i , produce high attentive features. This will ensure that already adapted regions or non- adaptive region (both cases have high uncertainty) have the lower attention value.

3.4. Training Algorithm

We employ Certainty Attention based Domain Adaption (CADA) models for solving the task of unsupervised domain adaptation. Both the CADA-P and CADA-A models jointly learn domain invariant and class invariant features, by focusing the classifier’s attention on the discriminator’s certain regions. So, the class probabilities y_i^c and aleatoric uncertainty v_i^c for the classifier will be estimated using weighted feature h_i .

$$y_i^c = G_{cy}(G_c(\mathbf{h}_i)), \quad v_i^c = G_{cv}(G_c(\mathbf{h}_i)) \quad (19)$$

The final objective function J for optimizing both the models is defined as:

$$J = \mathcal{L}_{cy} + \mathcal{L}_{cv} - \lambda * (\mathcal{L}_{dy} + \mathcal{L}_{dv}) \quad (20)$$

where λ is a trade-off parameter between classifier and discriminator. The optimization problem is to find the parameters $\theta_f, \theta_c, \theta_{cy}, \theta_{cv}, \theta_d, \theta_{dy}, \theta_{dv}$ that jointly satisfy:

$$(\hat{\theta}_f, \hat{\theta}_c, \hat{\theta}_{cy}, \hat{\theta}_{cv}) = \arg \min_{\substack{\theta_f, \theta_c, \\ \theta_{cy}, \theta_{cv}}} J(\theta_f, \theta_c, \theta_{cy}, \theta_{cv}, \theta_d, \theta_{dy}, \theta_{dv})$$

$$(\hat{\theta}_d, \hat{\theta}_{dy}, \hat{\theta}_{dv}) = \arg \max_{\theta_d, \theta_{dy}, \theta_{dv}} J(\theta_f, \theta_c, \theta_{cy}, \theta_{cv}, \theta_d, \theta_{dy}, \theta_{dv})$$

The implementation details are provided in supplementary material and other details are provided in the project page ¹

4. Experiments and Results

4.1. Datasets

Office-31 Dataset The Office-31 [36] consists of three domains: Amazon, Webcam and DSLR with 31 classes in each domain. There are 2817 images in Amazon (A) domain, 795 images in Webcam (W) and 498 images are in DSLR (D) domain makes total 4,110 images. To enable unbiased evaluation, we evaluate all methods on 6 transfer tasks $A \rightarrow W, D \rightarrow A, W \rightarrow A, A \rightarrow D, D \rightarrow W$ and $W \rightarrow D$.

Office-Home Dataset We also evaluated our model on the Office-Home dataset [49] for unsupervised domain adaptation. This dataset consists of four domains: Art (Ar), Clip-art (Cl), Product (Pr) and Real-World (Rw). Each domain has common 65 categories and total 15,500 images. We evaluated our model by considering all the 12 adaptation tasks. The performance is reported in the Table 3.

ImageCLEF Dataset ImageCLEF-2014 dataset consists of three datasets: Caltech-256 (C), ILSVRC 2012 (I), and Pascal VOC 2012 (P). There are 12 common classes and total 600 images in each domain. We evaluate our model on all the 6 transfer tasks and results are reported in Table 4.

¹<https://delta-lab-iitk.github.io/CADA/>

Table 1: Classification accuracy (%) on *Office-31* dataset for unsupervised domain adaptation (AlexNet[21])

Method	A → W	D → W	W → D	A → D	D → A	W → A	Average
Alexnet[21]	60.6 ± 0.4	95.0 ± 0.2	99.5 ± 0.1	64.2 ± 0.3	45.5 ± 0.5	48.3 ± 0.5	68.8
MMD[47]	61.0 ± 0.5	95.0 ± 0.3	98.5 ± 0.3	64.9 ± 0.4	47.2 ± 0.5	49.4 ± 0.4	69.3
RTN[27]	73.3 ± 0.3	96.8 ± 0.2	99.6 ± 0.1	71.0 ± 0.2	50.5 ± 0.3	51.0 ± 0.1	74.1
DAN[25]	68.5 ± 0.4	96.0 ± 0.3	99.0 ± 0.2	66.8 ± 0.2	50.0 ± 0.4	49.8 ± 0.3	71.7
GRL [12]	73.0 ± 0.5	96.4 ± 0.3	99.2 ± 0.3	72.3 ± 0.3	52.4 ± 0.4	50.4 ± 0.5	74.1
JAN [28]	75.2 ± 0.4	96.6 ± 0.2	99.6 ± 0.1	72.8 ± 0.3	57.5 ± 0.2	56.3 ± 0.2	76.3
CDAN[26]	77.9 ± 0.3	96.9 ± 0.2	100.0 ± 0	74.6 ± 0.2	55.1 ± 0.3	57.5 ± 0.4	77.0
MADA[34]	78.5 ± 0.2	99.8 ± 0.1	100.0 ± 0	74.1 ± 0.1	56.0 ± 0.2	54.5 ± 0.3	77.1
CADA-W	82.3 ± 0.3	99.2 ± 0.1	99.6 ± 0.1	75.9 ± 0.2	57.7 ± 0.1	53.3 ± 0.2	78.0
CADA-A	84.1 ± 0.2	99.2 ± 0.2	99.8 ± 0.2	77.3 ± 0.1	61.3 ± 0.2	54.1 ± 0.3	79.3
CADA-P	83.4 ± 0.2	99.8 ± 0.1	100.0 ± 0	80.1 ± 0.1	59.8 ± 0.2	59.5 ± 0.3	80.4

Table 2: Classification accuracy (%) on *Office-31* dataset for unsupervised domain adaptation (ResNet-50 [15])

Method	A → W	D → W	W → D	A → D	D → A	W → A	Average
ResNet-50[15]	68.4 ± 0.2	96.7 ± 0.1	99.3 ± 0.1	68.9 ± 0.2	62.5 ± 0.3	60.7 ± 0.3	76.1
DAN[25]	80.5 ± 0.4	97.1 ± 0.2	99.6 ± 0.1	78.6 ± 0.2	63.6 ± 0.3	62.8 ± 0.2	80.4
RTN[27]	84.5 ± 0.2	96.8 ± 0.1	99.4 ± 0.1	77.5 ± 0.3	66.2 ± 0.2	64.8 ± 0.3	81.6
DANN[12]	82.0 ± 0.4	96.9 ± 0.2	99.1 ± 0.1	79.7 ± 0.4	68.2 ± 0.4	67.4 ± 0.5	82.2
ADDA [46]	86.2 ± 0.5	96.2 ± 0.3	98.4 ± 0.3	77.8 ± 0.3	69.5 ± 0.4	68.9 ± 0.5	82.9
JAN[28]	85.4 ± 0.3	97.4 ± 0.2	99.8 ± 0.2	84.7 ± 0.3	68.6 ± 0.3	70.0 ± 0.4	84.3
MADA[34]	90.0 ± 0.1	97.4 ± 0.1	99.6 ± 0.1	87.8 ± 0.2	70.3 ± 0.3	66.4 ± 0.3	85.2
SimNet[35]	88.6 ± 0.5	98.2 ± 0.2	99.7 ± 0.2	85.2 ± 0.3	73.4 ± 0.8	71.6 ± 0.6	86.2
GTA[38]	89.5 ± 0.5	97.9 ± 0.3	99.8 ± 0.4	87.7 ± 0.5	72.8 ± 0.3	71.4 ± 0.4	86.5
DAAA [17]	86.8 ± 0.2	99.3 ± 0.1	100.0 ± 0.0	88.8 ± 0.4	74.3 ± 0.2	73.9 ± 0.2	87.2
CDAN[26]	93.1 ± 0.1	98.6 ± 0.1	100.0 ± 0.0	93.4 ± 0.2	71.0 ± 0.3	70.3 ± 0.3	87.7
CADA-W	93.9 ± 0.1	99.1 ± 0.2	99.6 ± 0.2	93.2 ± 0.3	68.9 ± 0.1	68.3 ± 0.2	87.2
CADA-A	96.8 ± 0.2	99.0 ± 0.1	99.8 ± 0.1	93.4 ± 0.1	71.7 ± 0.2	70.5 ± 0.3	88.5
CADA-P	97.0 ± 0.2	99.3 ± 0.1	100.0 ± 0	95.6 ± 0.1	71.5 ± 0.2	73.1 ± 0.3	89.5

Table 3: Classification accuracy (%) on *Office-Home* dataset for unsupervised domain adaptation (ResNet-50 [15])

Method	Ar → Cl	Ar → Pr	Ar → Rw	Cl → Ar	Cl → Pr	Cl → Rw	Pr → Ar	Pr → Cl	Pr → Rw	Rw → Ar	Rw → Cl	Rw → Pr	Avg
ResNet-50[15]	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
DAN[25]	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
DANN[12]	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
JAN[28]	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3
CDAN[26]	50.6	65.9	73.4	55.7	62.7	64.2	51.8	49.1	74.5	68.2	56.9	80.7	62.8
CADA-A	56.9	75.4	80.2	61.7	74.6	74.9	62.9	54.4	80.9	74.3	61.1	84.4	70.1
CADA-P	56.9	76.4	80.7	61.3	75.2	75.2	63.2	54.5	80.7	73.9	61.5	84.1	70.2

4.2. Results

Following the common setting in unsupervised domain adaptation, we used the pre-trained Alexnet [21] and pre-trained ResNet [15] architecture as our base model. For Office-31 dataset, results are reported in the Table 1 and

Table 2. From the table, it is clear that the proposed CADA outperforms the other methods on most transfer tasks, where CADA-P is the top-performing variant for both Alexnet and Resnet model. On average, we obtain improvements of 3.3% and 1.8% over the state of the art methods

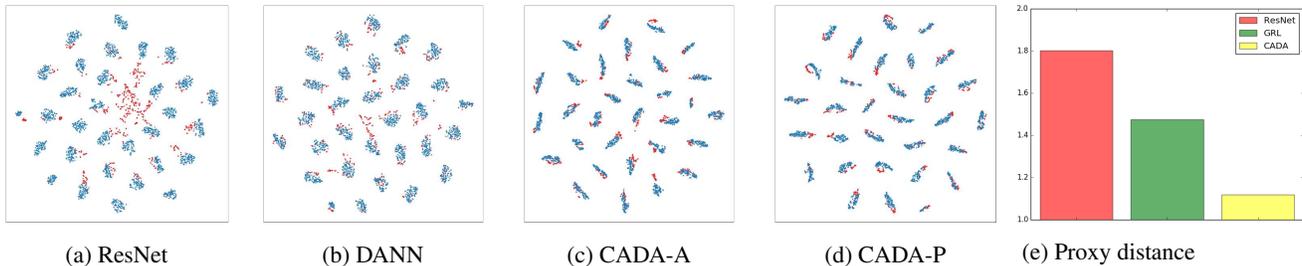


Figure 3: The t-SNE visualization of representations learned by (a) ResNet, (b) DANN, (c) CADA-A, and (d) CADA-P (blue: A; red: W), (e) shows Proxy distance for A \rightarrow W task for method ResNet [15], GRL [12] and the proposed model CADA-P

Table 4: Classification accuracy (%) on *ImageCLEF* dataset for unsupervised domain adaptation (ResNet-50 [15])

Method	I \rightarrow P	P \rightarrow I	II \rightarrow C	C \rightarrow I	I \rightarrow C	P \rightarrow C	Avg
ResNet [15]	74.8	83.9	91.5	78.0	65.5	91.2	80.7
DAN[25]	75.0	86.2	93.3	84.1	69.8	91.3	83.3
RTN[27]	75.6	86.8	95.3	86.9	72.7	92.2	84.9
GRL [12]	75.0	86.0	96.2	87.0	74.3	91.5	85.0
JAN[28]	76.8	88.0	94.7	89.5	74.2	91.7	85.8
MADA [34]	75.0	87.9	96.0	88.8	75.2	92.2	85.8
CDAN[26]	77.2	88.3	98.3	90.7	76.7	94.0	87.5
CADA-A	78.0	91.5	96.3	91.0	77.1	95.3	88.2
CADA-P	78.0	90.5	96.7	92.0	77.2	95.5	88.3

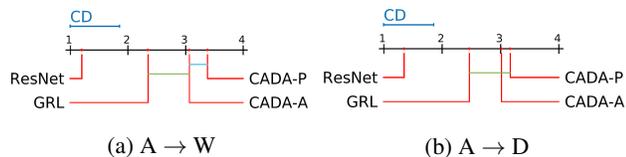


Figure 4: Analysis of statistically significant difference for A \rightarrow W and A \rightarrow D in ResNet [15], GRL [12], CADA-A, and CADA-P methods, with a significance level of 0.05.

as it can be seen that the difference between other methods are usually less than 1% and therefore this amount of improvement is fairly significant. In some cases such as, in Amazon-Webcam (A-W) we obtain almost 4% improvement over the state of the art method. Note that for DSLR-Amazon (D-A) and Webcam-Amazon (W-A), we do not obtain a state of the art. A very recent work [17] obtain state of the art results for these two cases. The difference between the domains is significant in these cases, and our method was not trained to optimally for these cases. The proposed method has obtained better results in all other cases, and even in these two cases, our results are competitive.

For the Office-Home dataset, the results are reported in Table 3. For this more challenging dataset we have achieved state-of-art performance. It is note worthy that the proposed model provides the classification accuracy that are substantially better on this Office-Home dataset which is harder

dataset for domain adaptation problem obtaining on average an improvement of 7.4% and 7.3% over the state of the art methods using CADA-P and CADA-A respectively.

The results on the ImageCLEF are reported in Table 4. Both CADA-P and CADA-A outperform the other state of the art models for all the transfer tasks except I \rightarrow C, with 0.8% and 0.7% improvement on average over the state of the art methods respectively. The room for improvement is smaller than the Office-Home Dataset, as ImageCLEF only have 12 classes and datasets in each domain and class category is equal, making it much easier for domain adaptation.

5. Ablation Study

We investigate the Bayesian model with and without attention for both Alexnet and ResNet model on the office-31 dataset. From Table 1 and 2, it is clear that the Bayesian model without attention (CADA-W) performs significantly better than the most of the other previous models, as predicting uncertainty for discriminator reduces negative transfer, by neglecting the regions which contain data uncertainty. Table 1 and 2 demonstrates that CADA-P (predictive certainty) performs better than CADA-A (aleatoric certainty), as predictive uncertainty includes both model and aleatoric uncertainty, providing a better estimate of certain regions for the discriminator.

6. Empirical Analysis

We further provide empirical analysis in terms of qualitative analysis of attention maps, feature visualization, discrepancy distance and statistical significance for additional insights about the performance of our method.

6.1. Qualitative Analysis of Attention Maps

To provide the effectiveness of proposed certainty based adaptation, we provide the certainty map of the discriminator at different training stages (chosen randomly) in the Figure 5. In the figure, we see that at the initial phase of the training (after 4 epochs), the discriminator discriminates the source and target domains just by some random location. As the training progress, the discriminator learns

the domain by attending the background of the images (In $A \rightarrow W$, domains are mostly dissimilar in the background). After some more training, the background is adapted (after 125 epochs), and now the discriminator attends to foreground part of the image to differentiate the domains (after 535 epochs). However, the foreground varies a lot across all the images. Hence discriminator is highly certain in the class object regions. Now with further training of the model, these class object regions are also adapted (after 1300 epochs). The remaining regions of the image cannot be further adapted because there is data uncertainty. At the end of the training, the discriminator will be highly uncertain regarding the domain, and the attention weight on the regions which are not adaptable will have low weights as we are using the certainty of the discriminator as the measure for the weights. Note that at the time of inference, we do not use the attention weights obtained by certainty for aiding classifier. These are used only at the time of training. The results show that these attention weights based training aids the classifier to better generalize to the target domain. We have provided more visualization examples in the supplementary material for further justification.

6.2. Feature Visualization

Adaptability of target to source features can be visualized using the t-SNE embeddings of images feature. We follow similar setting in [47, 12, 34] and plot t-SNE embeddings of the dataset in the Figure 3. We can observe that the proposed model correctly aligns the source and target domain images with exactly 31 clusters which is equal to the number of class labels with clear boundaries.

6.3. Discrepancy Distance

\mathcal{A} -distance is a measure of cross domain discrepancy[1], which, together with the source risk, will bound the target risk. The proxy \mathcal{A} -distance is defined as $d_{\mathcal{A}} = 2(1 - 2\epsilon)$, where ϵ is the generalization error of a classifier(e.g. kernel SVM) trained on the binary task of discriminating source and target. Figure 3 (e) shows $d_{\mathcal{A}}$ on tasks $A \rightarrow W$ with features of ResNet[15], GRL[12], and our model. We observe that $d_{\mathcal{A}}$ using our model features is much smaller than $d_{\mathcal{A}}$ using ResNet and GRL features, which suggests that our features can reduce the cross-domain gap more effectively.

6.4. Statistical Significance Test

We analyzed statistical significance [8] for variants of proposed method against GRL [12]. The Critical Difference (CD) for Nemenyi test depends upon the given confidence level (which is 0.05 in our case) for average ranks and number of test datasets. If the difference in the rank of the two methods lies within CD, then they are not significantly different. Figure 4 visualizes the post hoc analysis using the CD diagram for $A \rightarrow W$ dataset. From the figures, it is clear that our models are significantly different from GRL[12].

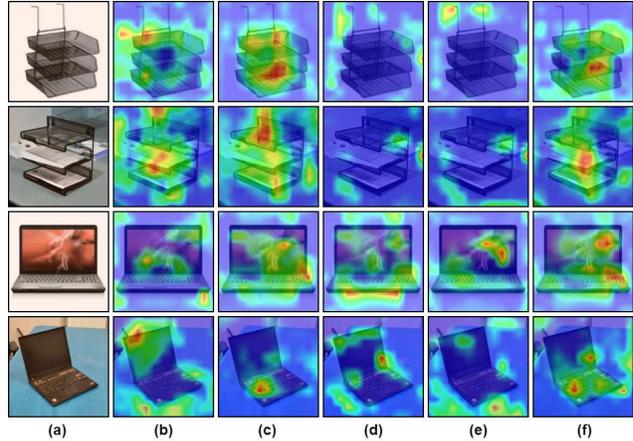


Figure 5: Attention visualization of the last convolutions layer of the proposed model CADA-P. The first and the third row shows the image from source domain (A) whereas the second and the fourth row shows the image from target domain (W). In each row, the leftmost image (a) represents the original image and the rightmost image (f) represents the classifier’s activation maps for ground truth class label at the end of the training. From left to right, the attention map of discriminator’s predictive certainty is illustrated at different training stages: (b) 4 epochs, (c) 125 epochs, (d) 535 epochs, and (e) 1300 epochs. We can see as the training progress, the discriminator’s certainty activation map changes from the background to the foreground, and then to the regions which can not be adapted further.

7. Conclusion

In this paper, we propose the use of certainty estimates of the discriminator to aid the generalization of the classifier by increasing the attention of the classifier on these regions. As it can be observed through our results, the attention maps obtained through certainty agree well with the classifier certainty for true labels and this aids in the generalization of the classifier for the target domain as well. The proposed method is thoroughly evaluated by comparison with state of the art methods and shows improved performance over all the other methods. Further, the analysis is provided in terms of statistical significance tests, discrepancy distance, and visualizations for better insight about the proposed method. The proposed method shows a new direction of using probabilistic measures for domain adaptation, and in the future, we aim to further explore this approach.

Acknowledgment: We acknowledge travel support from Microsoft Research India and Google Research India. We also acknowledge resource support from Delta Lab, IIT Kanpur. Vinod Kurmi acknowledges support from TCS Research Scholarship Program.

References

- [1] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010. 8
- [2] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 7, 2017. 2
- [3] Zhangjie Cao, Mingsheng Long, Jianmin Wang, and Michael I. Jordan. Partial transfer learning with selective adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [4] Zhangjie Cao, Lijia Ma, Mingsheng Long, and Jianmin Wang. Partial adversarial domain adaptation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 135–150, 2018. 2
- [5] Yunjey Choi, Minje Choi, and Munyoung Kim. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. 2
- [6] Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. In *Advances in neural information processing systems (NIPS)*, pages 577–585, 2015. 2
- [7] Gabriela Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint arXiv:1702.05374*, 2017. 1
- [8] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30, 2006. 8
- [9] Jenny Rose Finkel and Christopher D Manning. Hierarchical bayesian domain adaptation. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 602–610. Association for Computational Linguistics, 2009. 2
- [10] Yarin Gal. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016. 2, 3
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016. 2, 3
- [12] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning*, pages 1180–1189, 2015. 1, 2, 6, 7, 8
- [13] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *Proceedings of the IEEE international conference on computer vision*, pages 2551–2559, 2015. 2
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6, 7, 8
- [16] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning*, pages 1994–2003, 2018. 2
- [17] Guoliang Kang, Liang Zheng, Yan Yan, and Yi Yang. Deep adversarial attention alignment for unsupervised domain adaptation: the benefit of target expectation maximization. In *The European Conference on Computer Vision (ECCV)*, September 2018. 2, 6, 7
- [18] Vahid Kazemi and Ali Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. *arXiv preprint arXiv:1704.03162*, 2017. 2
- [19] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017. 2, 3, 5
- [20] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7482–7491, 2018. 3
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 2012. 6
- [22] Shanu Kumar, Shubham Atreja, Anjali Singh, and Mohit Jain. Adversarial adaptation of scene graph models for understanding civic issues. *arXiv preprint arXiv:1901.10124*, 2019. 2
- [23] Vinod Kumar Kurmi and Vinay P. Namboodiri. Looking back at labels: A class based domain adaptation technique, 2019. 2
- [24] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pages 6402–6413, 2017. 2
- [25] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International Conference on Machine Learning*, pages 97–105, 2015. 1, 2, 6, 7
- [26] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pages 1640–1650, 2018. 6, 7
- [27] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation with residual transfer networks. In *Advances in Neural Information Processing Systems*, pages 136–144, 2016. 4, 6, 7
- [28] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In *International Conference on Machine Learning*, pages 2208–2217, 2017. 6, 7

- [29] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems*, pages 7047–7058, 2018. 2
- [30] Zak Murez, Soheil Kolouri, David Kriegman, Ravi Ramamoorthi, and Kyungnam Kim. Image to image translation for domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4500–4509, 2018. 2
- [31] Badri Patro and Vinay P. Nambodiri. Differential attention for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [32] Badri Narayana Patro, Sandeep Kumar, Vinod Kumar Kurmi, and Vinay Nambodiri. Multimodal differential network for visual question generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4002–4012, 2018. 2
- [33] Badri Narayana Patro, Vinod Kumar Kurmi, Sandeep Kumar, and Vinay Nambodiri. Learning semantic sentence embeddings using sequential pair-wise discriminator. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2715–2729, 2018. 2
- [34] Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *Advances in Artificial Intelligence (AAAI)*, 2018. 2, 3, 6, 7, 8
- [35] Pedro O Pinheiro and AI Element. Unsupervised domain adaptation with similarity learning. 6
- [36] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010. 5
- [37] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018. 2
- [38] Swami Sankaranarayanan, Yogesh Balaji, Arpit Jain, Ser Nam Lim, and Rama Chellappa. Learning from synthetic data: Addressing domain shift for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3752–3761, 2018. 2, 6
- [39] Jian Shen, Yanru Qu, Weinan Zhang, and Yong Yu. Wasserstein distance guided representation learning for domain adaptation. In *AAAI*, 2018. 2
- [40] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *AAAI*, volume 6, page 8, 2016. 2
- [41] Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain adaptation. *Domain Adaptation in Computer Vision Applications*, page 153, 2017. 2
- [42] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *European Conference on Computer Vision*, pages 443–450. Springer, 2016. 2
- [43] Mattias Teye, Hossein Azizpour, and Kevin Smith. Bayesian uncertainty estimation for batch normalized deep networks. In *International Conference on Machine Learning*, pages 4914–4923, 2018. 2
- [44] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE, 2011. 1
- [45] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *Computer Vision (ICCV), 2015 IEEE International Conference on*, pages 4068–4076. IEEE, 2015. 2
- [46] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 4, 2017. 2, 6
- [47] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 2, 6, 8
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017. 2
- [49] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proc. CVPR*, pages 5018–5027, 2017. 2, 5
- [50] Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 2018. 1
- [51] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015. 2
- [52] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016. 2
- [53] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4651–4659, 2016. 2
- [54] Zhen Zhang, Mianzhi Wang, Yan Huang, and Arye Nehorai. Aligning infinite-dimensional covariance matrices in reproducing kernel hilbert spaces for domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3437–3445, 2018. 2
- [55] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2223–2232, 2017. 2