

**Disclaimer:**

This work has been accepted for publication in the IEEE Transactions on Pattern Analysis and Machine Intelligence:

doi: 10.1109/TPAMI.2019.2933829

link: <https://ieeexplore.ieee.org/document/8792192>

**Copyright:**

© 2019 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# Inferring Latent Domains for Unsupervised Deep Domain Adaptation

Massimiliano Mancini, Lorenzo Porzi, Samuel Rota Buló,  
Barbara Caputo and Elisa Ricci

**Abstract**—Unsupervised Domain Adaptation (UDA) refers to the problem of learning a model in a target domain where labeled data are not available by leveraging information from annotated data in a source domain. Most deep UDA approaches operate in a single-source, single-target scenario, *i.e.* they assume that the source and the target samples arise from a single distribution. However, in practice most datasets can be regarded as mixtures of multiple domains. In these cases, exploiting traditional single-source, single-target methods for learning classification models may lead to poor results. Furthermore, it is often difficult to provide the domain labels for all data points, *i.e.* latent domains should be automatically discovered. This paper introduces a novel deep architecture which addresses the problem of UDA by automatically discovering latent domains in visual datasets and exploiting this information to learn robust target classifiers. Specifically, our architecture is based on two main components, *i.e.* a side branch that automatically computes the assignment of each sample to its latent domain and novel layers that exploit domain membership information to appropriately align the distribution of the CNN internal feature representations to a reference distribution. We evaluate our approach on publicly available benchmarks, showing that it outperforms state-of-the-art domain adaptation methods.

**Index Terms**—Unsupervised Domain Adaptation, Batch Normalization, Domain Discovery, Object Recognition.

## 1 INTRODUCTION

IN the last few years, deep neural networks have enabled unprecedented progresses in many fields of artificial intelligence including computer vision, speech recognition, natural language processing and robotics. While deep models are extremely powerful by automatically learning discriminative representations of input data, one major limitation is that huge amounts of labeled data are typically required for training. However, compiling large scale annotated datasets is a tedious and extremely costly operation. To address this issue, over the years the research community has proposed different strategies, such as semi-supervised learning, transfer learning and domain adaptation.

Domain adaptation methods, in particular, are specifically designed in order to transfer knowledge from a *source* domain to a domain of interest, *i.e.* the *target* domain. In the specific case of Unsupervised Domain Adaptation (UDA), labeled data are only available in the source domain, while no annotations are provided for the target samples. Typically, knowledge transfer is achieved by learning domain-agnostic models or invariant feature representations. The

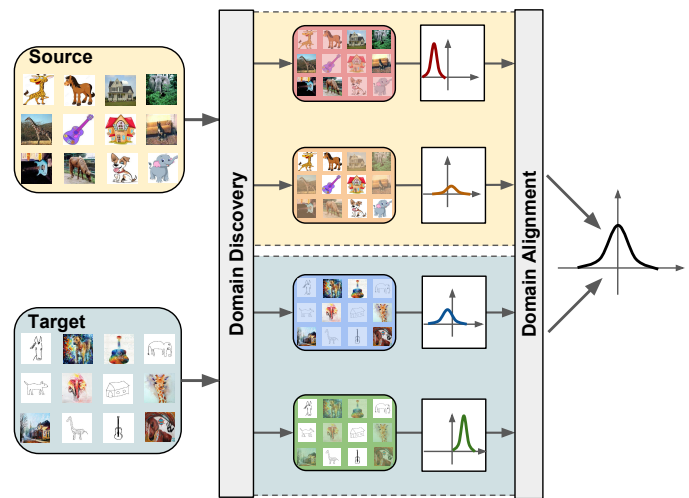


Fig. 1: The idea behind the proposed framework. We introduce a novel deep architecture which, given a set of images, automatically discover multiple latent domains and use this information to align the distributions of the internal CNN feature representations of sources and target domains for the purpose of domain adaptation. In this way, more accurate target classifiers can be learned. Image better seen at magnification.

- M. Mancini is with Department of Computer, Control, and Management Engineering, Sapienza University of Rome, Rome, Italy.(Email: mancini@diag.uniroma1.it)
- L. Porzi and S. Rota Buló are with Mapillary Research, Graz, Austria.(Email: {lorenzo,samuel}@mapillary.com)
- B. Caputo is with DAUIN Department of Control and Computer Engineering of Politecnico di Torino, Turin, Italy. (Email: barbara.caputo@polito.it)
- M. Mancini and B. Caputo are with Italian Institute of Technology, Turin, Italy.
- M. Mancini and E. Ricci are with Fondazione Bruno Kessler, Trento, Italy. (Email: eliricci@fbk.eu)
- E. Ricci is with Department of Information Engineering and Computer Science, University of Trento, Trento, Italy

problem of UDA has been widely studied and both theoretical results [4], [45] and several algorithms have been developed, both considering shallow models [13], [18], [20], [28], [38] and deep architectures [7], [8], [15], [17], [39], [40], [59]. While deep neural networks tend to produce more transferable and domain-invariant features with respect to shallow models, previous works have shown that the domain shift is only alleviated but not entirely removed [11].

Most previous works on UDA focus on a single-source and single-target scenario. However, in many computer vision applications labeled training data are often generated from multiple distributions, *i.e.* there are multiple source domains. Examples of multi-source DA problems arise when the source set corresponds to images taken with different cameras, collected from the web or associated to multiple points of views. In these cases, a naive application of single-source domain adaptation algorithms would not suffice, leading to poor results. Analogously, target samples may arise from more than a single distribution and learning multiple target-specific models may improve significantly the performance. Therefore, in the past several research efforts have been devoted to develop domain adaptation methods considering multiple source and target domains [12], [45], [58], [63]. However, these approaches assume that the multiple domains are known. A more challenging problem arises when training data correspond to latent domains, *i.e.* we can make a reasonable estimate on the number of source and target domains available, but we have no information, or only partial, about domain labels. To address this problem, known in the literature as *latent domain discovery*, previous works have proposed methods which simultaneously discover hidden source domains and use them to learn the target classification models [19], [27], [62].

This paper introduces the first approach based on deep neural networks able to automatically discover latent domains in multi-source, multi-target UDA setting. Our method is inspired by the recent works [8], [9], [43], which revisit Batch Normalization (BN) layers [29] for the purpose of domain adaptation and generalization, introducing specific Domain Alignment layers (DA-layers). The main idea behind DA-layers is to cope with domain shift by aligning representations of source and target distributions to a reference Gaussian distribution. Our approach develops from the same intuition. However, to address the additional challenges of discovering and handling multiple latent domains, we propose a novel architecture which is able to (i) learn a set of assignment variables which associate source and target samples to a latent domain and (ii) exploit this information for aligning the distributions of the internal CNN feature representations and learn robust target classifiers (Fig.1). Our experimental evaluation shows that the proposed approach alleviates the domain discrepancy and outperforms previous UDA techniques on popular benchmarks, such as Office-31 [52], PACS [35] and Office-Caltech [20].

To summarize, the contributions of this paper are three-fold. Firstly, we present a novel deep learning approach for unsupervised domain adaptation which operates in a multi-source, multi-target setting. Secondly, we propose a novel architecture which is not only able to handle multiple domains, but also permits to automatically discover them by grouping source and target samples. Thirdly, our experiments demonstrate that the proposed framework is superior to many state of the art single- and multi-source/target UDA methods.

This paper extends our earlier work [44] in many aspects. In particular, we modify the deep architecture described in [44] for handling not only different source distributions

but also multiple target domains. We also develop a novel formulation of the loss function in [44] which allows not only to further boost the performances, but to also stabilize the estimates of the domain discovery branches, being more robust to mode collapsing issues (*e.g.* all images assigned to a single domain). We also provide a more comprehensive review of related works. Finally, we significantly expand our experimental evaluation, considering more recent baseline methods and showing additional qualitative and quantitative results.

The remainder of this paper is organized as follows. We first introduce related works in Section 2 and then describe the proposed approach for latent domain discovery and unsupervised domain adaptation (Section 3). The results of our extensive experimental evaluation are provided in Section 4. We conclude the paper in Section 5.

## 2 RELATED WORK

In this section we review previous works on DA considering both methods based on hand-crafted features and more recent deep architectures, positioning our work in this context.

**DA methods with hand-crafted features.** Earlier DA approaches operate on hand-crafted features and attempt to reduce the discrepancy between the source and the target domains by adopting different strategies. For instance, instance-based methods [18], [28], [65] develop from the idea of learning classification/regression models by re-weighting source samples according to their similarity with the target data. A different strategy is exploited by feature-based methods, coping with domain shift by learning a common subspace for source and target data such as to obtain domain-invariant representations [13], [20], [38]. Parameter-based methods [66] address the domain shift problem by discovering a set of shared weights between the source and the target models. However, they usually require labeled target data which is not always available.

While most earlier DA approaches focus on a single-source and single-target setting, some works have considered the related problem of learning classification models when the training data spans multiple domains [12], [45], [58]. The common idea behind these methods is that when source data arises from multiple distributions, adopting a single source classifier is suboptimal and improved performance can be obtained by leveraging information about multiple domains. However, these methods assume that the domain labels for all source samples are known in advance. In practice, in many applications the information about domains is hidden and latent domains must be discovered into the large training set. Few works have considered this problem in the literature. Hoffman *et al.* [27] address this task by modeling domains as Gaussian distributions in the feature space and by estimating the membership of each training sample to a source domain using an iterative approach. Gong *et al.* [19] discover latent domains by devising a nonparametric approach which aims at simultaneously achieving maximum distinctiveness among domains and ensuring that strong discriminative models are learned for each latent domain. In [62] domains are modeled as manifolds and source images representations are learned decoupling information about semantic category

and domain. By exploiting these representations the domain assignment labels are inferred using a mutual information based clustering method.

**Deep Domain Adaptation.** Most recent works on DA consider deep architectures and robust domain-invariant features are learned using either supervised neural networks [7], [8], [15], [17], [39], [59], deep autoencoders [67] or generative adversarial networks [6], [56]. Research efforts can be grouped in terms of the number of source domains available at training time.

In the single-source DA setting, we can identify two main strategies. The first deals with *features* and aims at learning deep domain invariant representations. The idea is to introduce in the learning architecture different measures of domain distribution shift at a single or multiple levels [8], [9], [41], [57] and then train the network to minimize these measures while also reducing a task-specific loss, for instance for classification or detection. In this way the network produces features invariant to the domain shift, but still discriminative for the task at hand. Besides distribution evaluations, other domain shift measures used similarly are the error in the target sample reconstruction [17], or various coherence metrics on the pseudo-labels assigned by the source models to the target data [25], [53], [55]. Finally, a different group of feature-based methods rely on adversarial loss functions [16], [59]. The method proposed in [54], that push the network to be unable to discriminate whether a sample coming from the source or from the target, is an interesting variant of [16], where the domain difference is still measured at the feature level but passing through an image reconstruction step. Besides integrating the domain discrimination objective into end-to-end classification networks, it has also been shown that two-step networks may have practical advantages [2], [60]. The second popular deep adaptive strategy focuses on *images*. The described adversarial logic that demonstrated its effectiveness for feature-based methods, has also been extended to the goal of reducing the visual domain gap. Powerful GAN [21] methods have been exploited to generate new images or perturb existing ones to resemble the visual style of a certain domain, thus reducing the discrepancy at pixel level [5], [56]. Most of the works based on image adaptation aim at generating either target-like source images or source-like target images, but it has been recently shown that integrating both the transformation directions is highly beneficial [51].

In practical applications one may be offered more than one source domain. This has triggered the study of multi-sources DA algorithms. The multi-source setting was initially studied from a theoretical point of view, focusing on theorems indicating how to optimally sub-select the data to be used in learning the source models [10], or proposing principled rules for combining the source-specific classifiers and obtain the ideal target class prediction [45]. Several other works followed this direction in the shallow learning framework. When dealing with shallow-methods the naïve model learned by collecting all the source data in single domain without any adaptation was usually showing low performance on the target. It has been noticed that this behavior changes when moving to deep learning, where the larger number of samples as well as their variability

supports generalization and usually provides good results on the target. Only very recently two methods presented multi-source deep learning approaches that improve over this reference. The approach proposed in [63] builds over [16] by replicating the adversarial domain discriminator branch for each available source. Moreover these discriminators are also used to get a perplexity score that indicates how the multiple sources should be combined at test time, according to the rule in [45]. A similar multi-way adversarial strategy is used also in [68], but this work comes with a theoretical support that frees it from the need of respecting a specific optimal source combination and thus from the need of learning the source weights.

While recent deep DA methods significantly outperform approaches based on hand-crafted features, the vast majority of them only consider single-source, single-target settings. Moreover, almost all work presented in the literature so far assume to have direct access to multiple source domains, where in many practical applications such knowledge might not be directly available, or costly to obtain in terms of time and human annotators. To our knowledge, this is the first work proposing a deep architecture for discovering latent source domains and exploiting them for improving classification performance on target data.

Finally, our work is also related to domain generalization [33], [47]. As in domain generalization, we assume the presence of multiple source domains. However, in domain generalization the domains are clearly separated and are leveraged to produce a model for any unseen target domain, *e.g.* learning domain invariant representations [17], [33], [34], [54], or combining domain specific information [35], [42], [43], [64]. In our scenario instead, the multiple source domains are mixed and our goal is to latently discover them while producing a model for the single or multiples (possibly mixed) target domains. We highlight that discovering latent domains can be useful even in the context of domain generalization, as shown in [64].

### 3 METHOD

#### 3.1 Problem Formulation and Notation

We assume to have data belonging to one of several domains. Specifically, we consider  $k_s$  *source* domains, characterized by unknown probability distributions  $p_{xy}^{s_1}, \dots, p_{xy}^{s_{k_s}}$  defined over  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X}$  is the input space (*e.g.* images) and  $\mathcal{Y}$  the output space (*e.g.* object or scene categories) and, similarly, we assume  $k_t$  *target* domains characterized by  $p_{xy}^{t_1}, \dots, p_{xy}^{t_{k_t}}$ . The numbers of source and target domains are not necessarily known a-priori, and are left as hyperparameters of our method.

During training we are given a set of labeled sample points from the source domains, and a set of unlabeled sample points from the target domains, while we can have partial or no information about the domain of the source sample points. We model the source data as a set  $\mathcal{S} = \{(x_1^s, y_1^s), \dots, (x_n^s, y_n^s)\}$  of i.i.d. observations from a mixture distribution  $p_{xy}^s = \sum_{i=1}^{k_s} \pi_{s_i} p_{xy}^{s_i}$ , where  $\pi_{s_i}$  is the unknown probability of sampling from a source domain  $s_i$ . Similarly, the target sample  $\mathcal{T} = \{x_1^t, \dots, x_m^t\}$  consists of i.i.d. observations from the marginal  $p_x^t$  of the mixture distribution over target domains. Furthermore, we denote

by  $x_S = \{x_1^s, \dots, x_n^s\}$  and  $y_S = \{y_1^s, \dots, y_n^s\}$ , the source data and label sets, respectively. We assume to know the domain label for a (possibly empty) sub-sample  $\hat{S} \subset S$  from the source domains and we denote by  $d_{\hat{S}}$  the domain labels in  $\{s_1, \dots, s_k\}$  of the sample points in  $x_{\hat{S}}$ . The set of source domains labels is given by  $\mathcal{D}_s = \{s_1, \dots, s_{k_s}\}$ , and similarly the set of target domains is denoted by  $\mathcal{D}_t$ .

Our goal is to learn a predictor that is able to classify data from the target domains. The major difficulties that this problem poses, and that we have to deal with, are: (i) the distributions of source and target domains can be drastically different, making it hard to apply a classifier learned on one domain to the others, (ii) we lack direct observation of target labels, and (iii) the assignment of each source and target sample point to its domain is unknown, or known for a very limited number of source sample points.

Several previous works [7], [8], [15], [17], [39], [59] have tackled the related problem of domain adaptation in the context of deep neural networks, dealing with (i) and (ii) in the single domain case for both source and target data (*i.e.*  $k_s = 1$  and  $k_t = 1$ ). In particular, some recent works have demonstrated a simple yet effective approach based on the replacement of standard BN layers with specific *Domain Alignment layers* [8], [9]. These layers reduce internal domain shift at different levels within the network by normalizing features in a domain-dependent way, matching their distributions to a pre-determined one. We revisit this idea in the context of multiple, unknown source and target domains and introduce a novel Multi-domain DA layer (mDA-layer) in Section 3.2, which is able to normalize the multi-modal feature distributions encountered in our setting. To do this, our mDA-layers exploit a side-output branch attached to the main network (see Section 3.3), which predicts domain assignment probabilities for each input sample. Finally, in Section 3.4 we show how the predicted domain probabilities can be exploited, together with the unlabeled target samples, to construct a prior distribution over the network’s parameters which is then used to define the training objective for our network.

### 3.2 Multi-domain DA-layers

DA-layers [8], [9], [36] are motivated by the observation that, in general, activations within a neural network follow domain-dependent distributions. As a way to reduce domain shift, the activations are thus normalized in a domain-specific way, shifting them according to a parameterized transformation in order to match their first and second order moments to those of a reference distribution, which is generally chosen to be normal with zero mean and unit standard deviation. While previous works only considered settings with two domains, *i.e.* source and target, the basic idea can in fact be applied to any number of domains, as long as the domain membership of each sample point is known. Specifically, denoting as  $q_x^d$  the distribution of activations for a given feature channel and domain  $d$ , an input  $x^d \sim q_x^d$  to the DA-layer can be normalized according to

$$\text{DA}(x^d; \mu_d, \sigma_d) = \frac{x^d - \mu_d}{\sqrt{\sigma_d^2 + \epsilon}},$$

where  $\mu_d = E_{x \sim q_x^d}[x]$ ,  $\sigma_d^2 = \text{Var}_{x \sim q_x^d}[x]$  are mean and variance of the input distribution, respectively, and  $\epsilon > 0$  is a small constant to avoid numerical issues. In practice, when the statistics  $\mu_d$  and  $\sigma_d^2$  are computed over the current mini-batch, we obtain the application of standard batch normalization separately to the sample points of each domain.

As mentioned above, this approach requires full domain knowledge, because for each domain  $d$ ,  $\mu_d$  and  $\sigma_d^2$  need to be calculated on a data sample belonging to the specific domain  $d$ . In our case, however, we do not know the domain of the source/target sample points, or we have only partial knowledge about that. To tackle this issue, we propose to model the layer’s input distribution as a mixture of Gaussians, with one component per domain. Specifically, we define a global input distribution  $q_x = \sum_d \pi_d q_x^d$ , where  $\pi_d$  is the probability of sampling from domain  $d$ , and  $q_x^d = \mathcal{N}(\mu_d, \sigma_d^2)$  is the domain-specific distribution for  $d$ , namely a normal distribution with mean  $\mu_d$  and variance  $\sigma_d^2$ . Given a mini-batch  $\mathcal{B} = \{x_i\}_{i=1}^b$ , a maximum likelihood estimate of the parameters  $\mu_d$  and  $\sigma_d^2$  is given by

$$\mu_d = \sum_{i=1}^b \alpha_{i,d} x_i, \quad \sigma_d^2 = \sum_{i=1}^b \alpha_{i,d} (x_i - \mu_d)^2, \quad (1)$$

where

$$\alpha_{i,d} = \frac{q_{d|x}(d | x_i)}{\sum_{j=1}^b q_{d|x}(d | x_j)}, \quad (2)$$

and  $q_{d|x}(d | x_i)$  is the conditional probability of  $x_i$  belonging to domain  $d$ , given  $x_i$ . Clearly, the value of  $q_{d|x}$  is known for all sample points for which we have domain information. In all other cases, the missing domain assignment probabilities are inferred from data, using the *domain prediction* network branch which will be detailed in Section 3.3. Thus, from the perspective of the alignment layer, these probabilities become an additional input, which we denote as  $w_{i,d}$  for the predicted probability of  $x_i$  belonging to  $d$ .

By substituting  $w_{i,d}$  for  $q_{d|x}(d | x_i)$  in (2), we obtain a new set of empirical estimates for the mixture parameters, which we denote as  $\hat{\mu}_d$  and  $\hat{\sigma}_d^2$ . These parameters are used to normalize the layer’s inputs according to

$$\text{mDA}(x_i, \mathbf{w}_i; \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}) = \sum_{d \in \mathcal{D}} w_{i,d} \frac{x_i - \hat{\mu}_d}{\sqrt{\hat{\sigma}_d^2 + \epsilon}}, \quad (3)$$

where  $\mathbf{w}_i = \{w_{i,d}\}_{d \in \mathcal{D}}$ ,  $\hat{\boldsymbol{\mu}} = \{\hat{\mu}_d\}_{d \in \mathcal{D}}$ ,  $\hat{\boldsymbol{\sigma}} = \{\hat{\sigma}_d^2\}_{d \in \mathcal{D}}$  and  $\mathcal{D}$  is the set of source/target latent domains. As in previous works [8], [9], [29], during back-propagation we calculate the derivatives through the statistics and weights, propagating the gradients to both the main input and the domain assignment probabilities.

### 3.3 Domain prediction

Our mDA-layers receive a set of domain assignment probabilities for each input sample point, which needs to be predicted, and different mDA-layers in the network, despite having different input distributions, share consistently the same domain assignment for the sample points. As a practical example, in the typical case in which mDA-layers are used in a CNN to normalize convolutional activations, the network would predict a single set of domain assignment

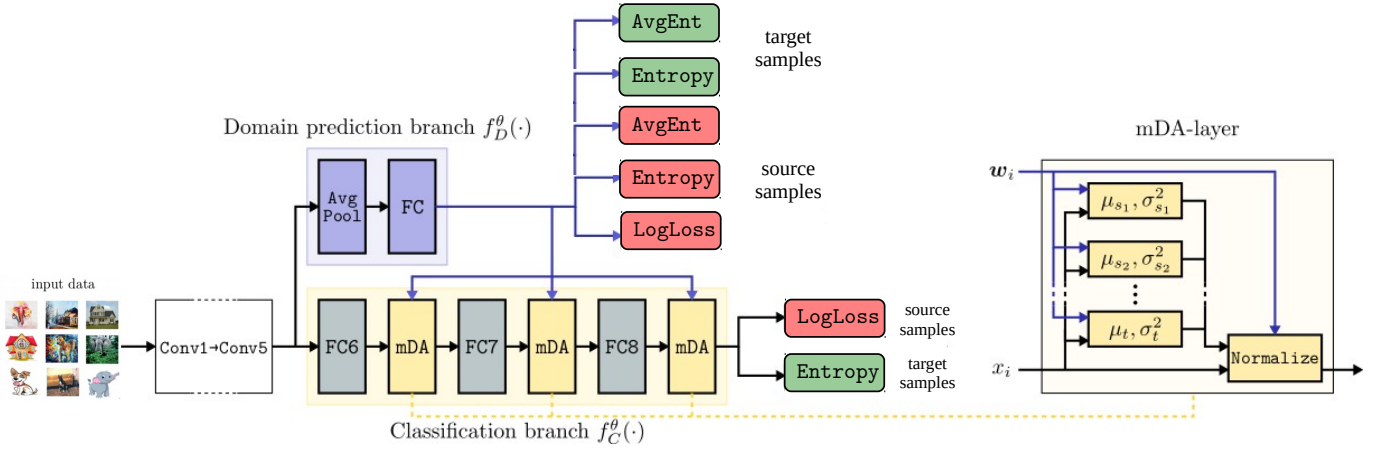


Fig. 2: Schematic representation of our method applied to the AlexNet architecture (left) and of an mDA-layer (right).

probabilities for each input image, which would then be fed to all mDA-layers and broadcasted across all spatial locations and feature channels corresponding to that image. We compute domain assignment probabilities using a distinct section of the network, which we call the *domain prediction* branch, while we refer to the main section of the network as the *classification* branch. The two branches share the bottom-most layers and parameters as depicted in Figure 2.

The domain prediction branch is implemented as a minimal set of layers followed by two softmax operations with  $k_s$  and  $k_t$  outputs for the source and target latent domains, respectively (more details follow in Section 4). The rationale of keeping the domain prediction separated between source and target derives from the knowledge that we have about the source/target membership of a sample point that we receive in input, while it remains unknown the specific source or target domain it belongs to. Furthermore, for each sample point  $x_i$  with known domain membership  $\hat{d}$ , we fix in each mDA-layer  $w_{i,d} = 1$  if  $d = \hat{d}$ , otherwise  $w_{i,d} = 0$ .

We split the network into a domain prediction branch and classification branch at some low level layer. This choice is motivated by the observation [1] that features tend to become increasingly more domain invariant going deeper into the network, meaning that it becomes increasingly harder to compute a domain membership as a function of deeper features. In fact, as pointed out in [8], this phenomenon is even more evident in networks that include DA-layers.

### 3.4 Training the network

In order to exploit unlabeled data within our discriminative setting, we follow the approach sketched in [8], where unlabeled data is used to define a regularizer over the network’s parameters. By doing so, we obtain a loss for  $\theta$  that takes the following form:

$$L(\theta) = L_{\text{cls}}(\theta) + L_{\text{dom}}(\theta), \quad (4)$$

where  $L_{\text{cls}}$  is a loss term that penalizes based on the final classification task, while  $L_{\text{dom}}$  accounts for the domain classification task.

**Classification loss  $L_{\text{cls}}$ .** The classification loss consists of two components, accounting for the supervised sample

from the source domain  $\mathcal{S}$  and the unlabeled target sample  $\mathcal{T}$ , respectively:

$$L_{\text{cls}}(\theta) = -\frac{1}{n} \sum_{i=1}^n \log f_C^\theta(y_i^s; x_i^s) + \frac{\lambda_C}{m} \sum_{i=1}^m H(f_C^\theta(\cdot; x_i^t)). \quad (5)$$

The first term on the right-hand-side is the average log-loss related to the supervised examples in  $\mathcal{S}$ , where  $f_C^\theta(y_i^s; x_i^s)$  denotes the output of the *classification branch* of the network for a source sample, *i.e.* the predicted probability of  $x_i^s$  having class  $y_i^s$ . The second term on the right-hand-side of (5) is the entropy  $H$  of the classification distribution  $f_C^\theta(\cdot; x_i^t)$ , averaged over all unlabeled target examples  $x_i^t$  in  $\mathcal{T}$ , scaled by a positive hyper-parameter  $\lambda_C$ .

**Domain loss  $L_{\text{dom}}$ .** Akin to the classification loss, the domain loss presents a component exploiting the supervision deriving from the known domain labels in  $\hat{\mathcal{S}}$  and a component exploiting the domain classification distribution on all sample points lacking supervision. However, the domain loss has in addition a term that tries to balance the distribution of sample points across domains, in order to avoid predictions to collapse into trivial solutions such as constant assignments to a single domain. Accordingly, the loss takes the following form:

$$L_{\text{dom}}(\theta) = -\frac{\lambda_D}{|\hat{\mathcal{S}}|} \sum_{x_i \in x_{\mathcal{S}}} \log f_{D_s}^\theta(d_i; x_i) - \lambda_B H(\bar{f}_{D_s}^\theta(\cdot)) + \frac{\lambda_E}{|\mathcal{S} \setminus \hat{\mathcal{S}}|} \sum_{x \in x_{\mathcal{S} \setminus \hat{\mathcal{S}}}} H(f_{D_s}^\theta(\cdot; x)) - \lambda_B H(\bar{f}_{D_t}^\theta(\cdot)) + \frac{\lambda_E}{m} \sum_{i=1}^m H(f_{D_t}^\theta(\cdot; x_i^t)). \quad (6)$$

Here,  $f_{D_s}^\theta$  and  $f_{D_t}^\theta$  denote the outputs of the *domain prediction branch* for data points from the source and target domains, respectively, while  $\bar{f}_{D_s}^\theta$  and  $\bar{f}_{D_t}^\theta$  denote the distributions of predicted domain classes across  $\mathcal{S}$  and  $\mathcal{T}$ , respectively, *i.e.*

$$\bar{f}_{D_s}^\theta(y) = \frac{1}{n} \sum_{i=1}^n f_{D_s}^\theta(y; x_i^s), \quad \bar{f}_{D_t}^\theta(y) = \frac{1}{m} \sum_{i=1}^m f_{D_t}^\theta(y; x_i^t).$$

The first term in (6) enforces the correct domain prediction on the sample points with known domain and it is scaled by a positive hyper-parameter  $\lambda_D$ . The terms scaled by the positive hyper-parameter  $\lambda_E$  enforce domain predictions with low uncertainty for the data points with unknown domain labels, by minimizing the entropy of the output distribution. Finally, the terms scaled by the positive hyper-parameter  $\lambda_B$  enforce balanced distributions of predicted domain classes across the source and target sample, by maximizing the entropy of the averaged distribution of domain predictions. Interestingly, since the classification branch has a dependence on the domain prediction branch via the mDA-layers, by optimizing the proposed loss, the network learns to predict domain assignment probabilities that result in a low classification loss. In other words, the network is free to predict domain memberships that do not necessarily reflect the real ones, as long as this helps improving its classification performance.

We optimize the loss in (4) with stochastic gradient descent. Hence, the samples  $\mathcal{S}$ ,  $\mathcal{T}$ ,  $\hat{\mathcal{S}}$  that are considered in the computation of the gradients are restricted to a random subsets contained in the mini-batch. In Section 4.1 we provide more details on how each mini-batch is sampled.

## 4 EXPERIMENTS

### 4.1 Experimental Setup

#### 4.1.1 Datasets

In our evaluation we consider several common DA benchmarks: the combination of USPS [14], MNIST [32] and MNIST-m [15]; the Digits-five benchmark in [63]; Office-31 [52]; Office-Caltech [20] and PACS [33].

**MNIST, MNIST-m and USPS** are three standard datasets for digits recognition. USPS [14] is a dataset of digits scanned from U.S. envelopes, MNIST [32] is a popular benchmark for digits recognition and MNIST-m [15] its counterpart obtained by blending the original images with colored patches extracted from BSD500 photos [3]. Due to their different representations (*e.g.* colored vs gray-scale), these datasets have been adopted as a DA benchmark by many previous works [6], [7], [15]. Here, we consider a multi source DA setting, using MNIST and MNIST-m as sources and USPS as target, training on the union of the training sets and testing on the test set of USPS.

**Digits-five** is an experimental setting proposed in [63] which considers 5 datasets of digits recognition. In addition to MNIST, MNIST-m and USPS, it includes SVHN [48] and Synthetic numbers datasets [16]. SVHN [48] contains pictures of real-world house numbers, collected from Google Street View. Synthetic numbers [16] is built from computer generated digits, including multiple sources of variations (*i.e.* position, orientation, background, color and amount of blur), for a total of 500 thousands images. We follow the experimental setting described in [63]: the train/test split comprises a subset of 25000 images for training and 9000 for testing for each of the domains, except for USPS for which the entire dataset is used. As in [63], we report the results when either SVHN or MNIST-m are used as targets and all the other domains are taken as sources.

**Office-31** is a standard DA benchmark which contains images of 31 object categories collected from 3 different sources: Webcam (W), DSLR camera (D) and the Amazon website (A). Following [62], we perform our tests in the multi-source setting, where each domain is in turn considered as target, while the others are used as source.

**Office-Caltech** [20] is obtained by selecting the subset of 10 common categories in the Office31 and the Caltech256 [24] datasets. It contains 2533 images, about half of which belong to Caltech256. The different domains are Amazon (A), DSLR (D), Webcam (W) and Caltech256 (C). In our experiments we consider the set of source/target combinations used in [19].

**PACS** [33] is a recently proposed DA benchmark which is especially interesting due to the significant domain shift between its domains. It contains images of 7 categories (*dog, elephant, giraffe, guitar, horse*) and 4 different visual styles: *i.e.* Photo (P), Art paintings (A), Cartoon (C) and Sketch (S). We employ the dataset in two different settings. First, following the experimental protocol in [33], we train our model considering 3 domains as sources and the remaining as target, using all the images of each domain. Differently from [33] we consider a DA setting (*i.e.* target data is available at training time) and we do not address the problem of domain generalization. Second, we use 2 domains as sources and the remaining 2 as targets, in a multi-source multi-target scenario. In this setting the results are reported as average accuracy between the 2 target domains.

In all experiments and settings, we assume to have no domain labels (*i.e.*  $\hat{\mathcal{S}} = \emptyset$ ), unless otherwise stated.

#### 4.1.2 Networks and training protocols

We apply our approach to four different CNN architectures: the MNIST and SVHN networks described in [15], [16], AlexNet [31] and ResNet [26]. We choose AlexNet due to its widespread use in many relevant DA works [8], [15], [39], [40], while ResNet is taken as an exemplar for modern state-of-the-art architectures employing batch-normalization layers. Both AlexNet and ResNet are first pre-trained on ImageNet and then fine-tuned on the datasets of interest. The MNIST and SVHN architectures are chosen for fair comparison with previous works considering digits datasets [16], [63]. Unless otherwise noted, we optimize our networks using Stochastic Gradient Descent with momentum 0.9 and weight decay  $5 \times 10^{-4}$ .

For the evaluation on MNIST, MNIST-m and USPS datasets, we employ the MNIST network described in [15], adding an mDA-layer after each convolutional and fully-connected layer. The domain prediction branch is attached to the output of `conv1`, and is composed of a convolution with the same meta-parameters as `conv2`, a global average pooling, a fully-connected layer with 100 output channels and finally a fully-connected classifier. Following the protocol described in [8], [15], we set the initial learning rate  $l_0$  to 0.01 and we anneal it through a schedule  $l_p$  defined by  $l_p = \frac{l_0}{(1+\gamma p)^\beta}$  where  $\beta = 0.75$ ,  $\gamma = 10$  and  $p$  is the training progress increasing linearly from 0 to 1. We rescale the input images to  $32 \times 32$  pixels, subtract the per-pixel image mean of the dataset and feed the networks with random crops of size  $28 \times 28$ . A batch-size of 128 images per domain is used.



For the Digits-five experiments we employ the SVHN architecture of [16], which is the same architecture adopted by [63], augmented with mDA-layers and a domain prediction branch in the same way as the MNIST network described in the previous paragraph. We train the architecture for 44000 iterations, with a batch-size of 32 images per domain, an initial learning rate of  $10^{-4}$  which is decayed by a factor of 10 after 80% of the training process. We use Adam as optimizer with a weight decay  $5 \times 10^{-5}$ , and pre-process the input images like in the MNIST, MNIST-m, USPS experiments.

For the experiments on Office-31 and Office-Caltech we employ the AlexNet architecture. We follow a setup similar to the one proposed in [8], [9], fixing the parameters of all convolutional layers and inserting mDA-layers after each fully-connected layer and before their corresponding activation functions. The domain prediction branch is attached to the last pooling layer `pool5`, and is composed of a global average pooling, followed by a fully connected classifier to produce the final domain probabilities. The training schedule and hyper-parameters are set following [8].

For the experiments on the PACS dataset we consider the ResNet architecture in the 18-layers setup described in [26], denoted as ResNet18. This architecture comprises an initial  $7 \times 7$  convolution, denoted as `conv1`, followed by 4 main modules, denoted as `conv2 - conv5`, each containing two residual blocks. To apply our approach, we replace each Batch Normalization layer in the residual blocks of the network with an mDA-layer. The domain prediction branch is attached to `conv1`, after the pooling operation. The branch is composed of a residual block with the same structure as `conv2`, followed by global average pooling and a fully connected classifier. In the multi-target experiments we add a second, identical domain prediction branch to discriminate between target domains. We also add a standard BN layer after the final domain classifiers, which we found leads to a more stable training process in the multi-target case. In both cases, we adopt the same training meta-parameters as for AlexNet, with the exception of weight-decay which is set to  $10^{-6}$  and learning rate which is set to  $5 \cdot 10^{-4}$ . The network is trained for 600 iterations with a batch-size of 48, equally divided between the domains, and the learning rate is scaled by a factor 0.1 after 75% of the iterations.

Regarding the hyper-parameters of our method, we set the number of source domains  $k$  equal to  $Q - 1$ , where  $Q$  is the number of different datasets used in each single experiment. In the multi-source multi-target scenarios, since we always have the domains equally split between source and target, we consider  $k$  equal  $Q/2$  for both source and target. Following [8], in the experiments with AlexNet we fix  $\lambda_C = \lambda_E = 0.2$  with  $\lambda_B = 0.1$ . Similarly, for the experiments on digits classification, we set  $\lambda_C = \lambda_E = 0.1$  and  $\lambda_B = 0.05$  for MNIST, MNIST-m and USPS, and  $\lambda_C = 0.01$  and  $\lambda_E = \lambda_B = 0.05$  for Digits-five, with  $\lambda_E = 0.01$  if  $\lambda_B = 0$ , which we found leading to a more stable minimization of the loss of the domain branch. In the experiments involving ResNet18 we select the values  $\lambda_C = 0.1$  and  $\lambda_E = \lambda_B = 0.0001$  through cross-validation, following the procedure adopted in [8], [38]. Similarly, in the multi-target ResNet18 experiments we select  $\lambda_C = \lambda_E = \lambda_B = 0.1$ .

When domain labels are available for a subset of source samples, we fix  $\lambda_D = 0.5$ .

We implement<sup>1</sup> all the models with the Caffe [30] framework and our evaluation is performed using an NVIDIA GeForce 1070 GTX GPU. We initialize both AlexNet and ResNet18 from models pre-trained on ImageNet, taking AlexNet from the Caffe model zoo, and converting ResNet18 from the original Torch model<sup>2</sup>. For all the networks and experiments, we add mDA layers and their variants in place of standard BN layers.

## 4.2 Results

In this section, we first analyze the proposed approach, demonstrating the advantages of considering multiple sources/targets and discovering latent domains. We then compare the proposed method with state-of-the-art approaches. For all the experiments we report the results in terms of accuracy, repeating the experiments at least 5 times and averaging the results. In the multi-target experiments, the reported accuracy is the average of the accuracies over the target domains. As for standard deviations, since we do not tune the hyperparameters of our model and baselines by employing the accuracy on the target domain, their values can be high in some settings. For this reason, in order to provide a more appropriate analysis of the significance of our results, we propose to adopt the following approach. In particular, let us model the accuracy of an algorithm as a random variable  $X_a$  with unknown distribution. The accuracy of a single run of the algorithm is an observation from this distribution. Therefore, in order to compare two algorithms we consider the two sets of associated observations  $A = \{a_1, \dots, a_n\}$  and  $B = \{b_1, \dots, b_m\}$  and estimate the probability that one algorithm is better than the other as:

$$p(X_a > X_b) = \frac{\sum_{a \in A} \sum_{b \in B} \delta(a > b)}{|A| \cdot |B|}$$

where  $\delta$  is the Dirac function. In the following we use this metric to compare our approach with respect to a baseline where no latent domain discovery process is implemented (specifically, the method DIAL [9], see below) considering five runs for each experiment. For sake of clarity, we denote this probability estimate as  $p^*$ .

In the following we first analyze the performances of the proposed approach with  $\lambda_B = 0$  (denoted as Ours  $\lambda_B = 0$ ), *i.e.* the algorithm we presented in [44], and then we describe the impact of the loss term we introduce in this paper setting  $\lambda_B > 0$  (denoted simply as Ours).

### 4.2.1 Experiments on the Digits datasets

In a first series of experiments, reported in Table 1, we test the performance of our approach on the MNIST, MNIST-m to USPS benchmark (M-Mm to U). The comparison includes: (i) the baseline network trained on the union of all source domains (*Unified sources*); (ii) training separate networks for each source, and selecting the one the performs the best on the target (*Best single source*); (iii) DIAL [9], trained on the union of the sources (*DIAL [9] - Unified sources*); (iv)

1. Code available at: [https://github.com/mancinimassimiliano/latent\\_domains\\_DA.git](https://github.com/mancinimassimiliano/latent_domains_DA.git)

2. <https://github.com/HolmesShuan/ResNet-18-Caffemodel-on-ImageNet>



TABLE 1: Digits datasets: comparison of different models in the multi-source scenario. MNIST (M) and MNIST-m (Mm) are taken as source domains, USPS (U) as target.

Method	M-Mm to U
Unified sources	57.1
Best single source	59.8
DIAL [9] - Unified sources	81.7
DIAL [9] - Best single source	81.9
Ours $\lambda_B = 0$ $k = 2$	82.5
Ours $\lambda_B = 0$ $k = 3$	82.2
Ours $\lambda_B = 0$ $k = 4$	82.7
Ours $\lambda_B = 0$ $k = 5$	82.4
Ours ( $k = 2$ )	82.4
Multi-source DA	84.2

DIAL, trained separately on each source and selecting the best performing model on the target (*DIAL [9] - Best single source*). We also report the results of our approach in the ideal case where the multiple source domains are known and we do not need to discover them (*Multi-source DA*). For our approach with  $\lambda_B = 0$ , we consider several different values of  $k$ , *i.e.* the number of discovered source domains.

By looking at the table several observations can be made. First, there is a large performance gap between models trained only on source data and DA methods, confirming that deep architectures by themselves are not enough to solve the domain shift problem [11]. Second, in analogy with previous works on DA [12], [45], [58], we found that considering multiple sources is beneficial for reducing the domain shift with respect to learning a model on the unified source set. Finally, and more importantly, when the domain labels are not available, our approach is successful in discovering latent domains and in exploiting this information for improving accuracy on target data, partially filling the performance gap between the single source models and *Multi-source DA*. Interestingly, the performance of our algorithm changes only slightly for different values of  $k$ , motivating our choice to always fix  $k$  to the known number of domains in the next experiments. Importantly, comparing our approach with DIAL we achieve higher accuracy in most of the runs, *i.e.*  $p^* = 0.65$ . In this experiment, the introduction of the loss term forcing a uniform assignment among clusters (denoted as Ours) leads to comparable performances to our method with  $\lambda_B = 0$ . This behaviour can be ascribed to the fact that the separation among different domains is quite clear in this case and adding constraints to the domain discovery process is not required. In the following, we show that the proposed loss is beneficial in more challenging datasets.

In a second set of experiments (Table 2), we compare our approach with previous and recently proposed single and multi-source unsupervised DA approaches. Following [63], we perform experiments on the Digits-five dataset, considering two settings with SVHN and MNIST-m as targets. As in the previous case, we evaluate the performance of the baseline network (with and without BN layers) and of DIAL when trained on the union of the sources, and, as an upper bound, our Multi-source DA with perfect domain knowledge. Moreover, we consider the Deep Cocktail Network (DCTN) [63] multi-source DA model, as well as the “source only” baseline and the single source DA models reported in [63]: Reverse gradient (RevGrad) [15] and Domain

TABLE 2: Digits-five [63] setting, comparison of different single source and multi-source DA models. The first row indicates the target domain with the others used as sources.

	Method	SVHN	MNIST-m	Mean
Unified sources	Source only (ours)	74.1	64.4	69.3
	Source only+BN (ours)	77.7	59.4	68.6
	Source only from [63]	72.2	64.1	68.2
	RevGrad [15]	68.9	71.6	70.3
	DAN [39]	71.0	66.6	68.8
	DIAL [9]	82.2	68.8	75.5
	Ours $\lambda_B = 0$	82.4	69.1	75.8
	Ours	<b>82.6</b>	<b>70.1</b>	<b>76.4</b>
Multi-source	Only Source [63]	64.6	60.7	62.7
	RevGRAD [15]	61.4	71.1	66.3
	DAN [39]	62.9	62.6	62.8
	DCTN [63]	77.5	70.9	74.2
	Multi-source DA	84.1	69.4	76.8

Adaptation Networks (DAN) [39]. For all single source DA models we consider two settings: “Unified Sources”, where all source domains are merged, and “Multi-Source”, where a separate model is trained for each source domain, and the final prediction is computed as an ensemble. As we can see, the Unified Sources DIAL already achieves remarkable results in this setting, outperforming DCTN, and Multi-source DA only provides a modest performance increase. As expected, the performance of our approach lies between these two ( $p^*$  equal to 0.56 and 0.64 for SVHN and MNIST-m respectively, with  $\lambda_B = 0$ ).

#### 4.2.2 Experiments on PACS

**Comparison with state of the art.** In our main PACS experiments we compare the proposed approach with the baseline ResNet18 network, and with ResNet18 + DIAL [9], both trained on the union of source sets. As in the digits experiments, we also report the performance of our method when perfect domain knowledge is available (Multi-source DA). Table 3 shows our results. In general, DA models are especially beneficial when considering the PACS dataset, and multi-source DA networks significantly outperform the single source one. Remarkably, our model is able to infer domain information automatically without supervision. In fact, its accuracy is either comparable with Multi-source DA (Photo, Art and Cartoon) or in between DIAL and Multi-source DA (Sketch). The average  $p^*$  is 0.67. Looking at the partial results, it is interesting to note that the improvements of our approach and Multi-source DA w.r.t. DIAL are more significant when either the Sketch or the Cartoon domains are employed as target set (average  $p^* = 0.81$ ). Since these domains are less represented in the ImageNet database, we believe that the corresponding features derived from the pre-trained model are less discriminative, and DA methods based on multiple sources become more effective. Setting  $\lambda_B > 0$ , allows to obtain a further boost of performances in the Sketch scenario, where the source domains are closer in appearances. In the other settings, the domain shift is mostly among the Sketch domain and all the others and it can be easily captured by our original formulation in [44].

To analyze the performances of our approach in a multi-source multi-target scenario, we perform a second set of experiments on the PACS dataset considering 2 domains as sources and the other 2 as targets. The results, shown in Table 4, comprise the same baselines as in Table 3. Note

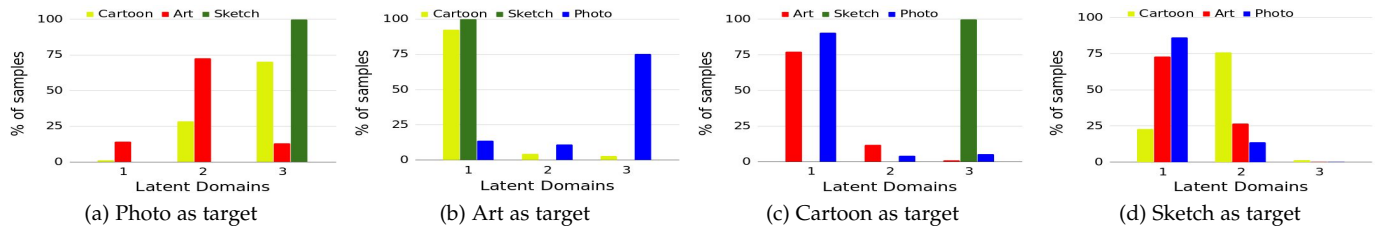


Fig. 3: Distribution of the assignments produced by the domain prediction branch for each latent domain in all possible settings of the PACS dataset. Different colors denote different source domains.

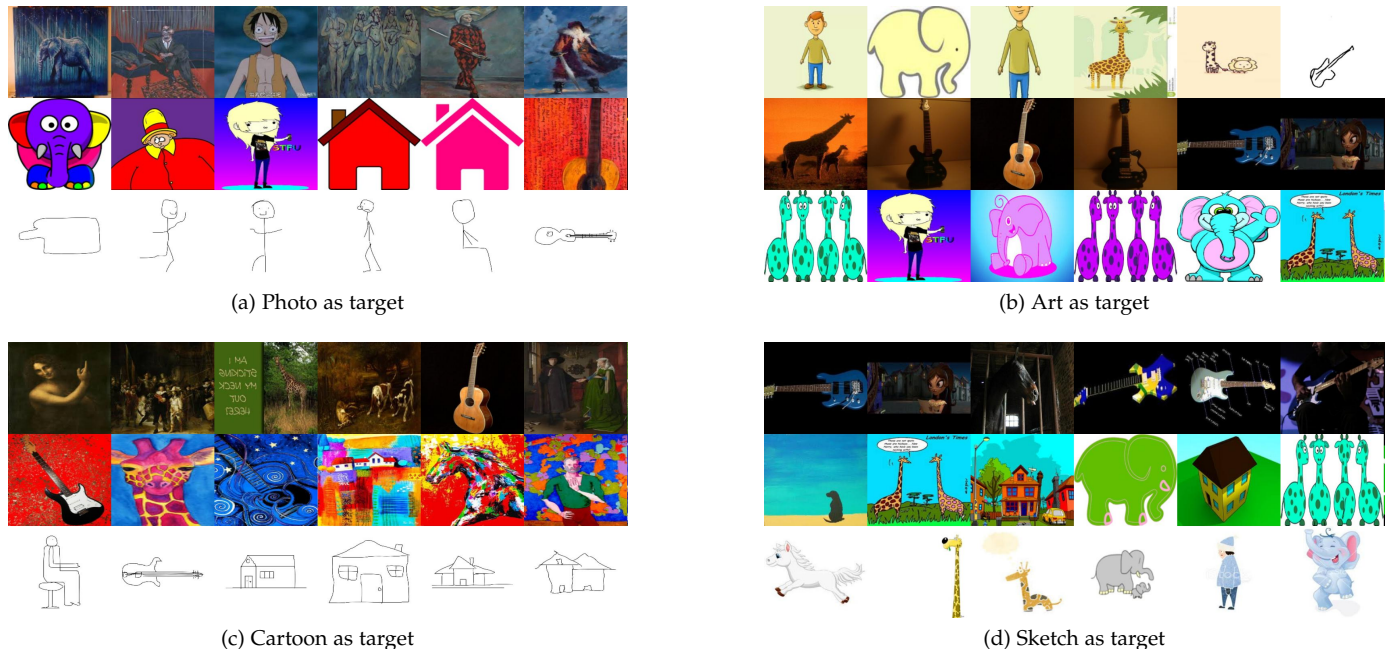


Fig. 4: Top-6 images associated to each latent domain for the different sources/target combinations. Each row corresponds to a different latent domain.

TABLE 3: PACS dataset: comparison of different methods using the ResNet architecture. The first row indicates the target domain, while all the others are considered as sources.

Method	Sketch	Photo	Art	Cartoon	Mean
ResNet [26]	60.1	92.9	74.7	72.4	75.0
DIAL [9]	66.8	<b>97.0</b>	87.3	85.5	84.2
Ours $\lambda_B = 0$	69.6	<b>97.0</b>	<b>87.7</b>	<b>86.9</b>	85.3
Ours	<b>70.7</b>	<b>97.0</b>	87.4	86.3	<b>85.4</b>
Multi-source DA	71.6	96.6	87.5	87.0	85.7

that, apart from the difficulty of providing useful domain assignments both in the source and target sets during training, the domain prediction step is required even at test time, thus having a larger impact on the final performances of the model. The performance gap between DIAL and our approach increases in this setting compared to Table 3. Our hypothesis is that not accounting for multiple domains has a larger impact on the *unlabeled* target than on the *labeled* source. Looking at the partial results, when Photo is considered as one of the target domains there are no particular differences in the final performances of the various DA models: this may be caused by the bias of the pretrained

network towards this domain. However, when the other domains are considered as targets, the gain in performances produced by our model are remarkable. When Sketch is one of the target domains, our model completely fills the gap between the unified source/target DA method and the multi-source multi-target upper bound with a gain of more than 7% when Art and Cartoon considered as other target. Setting  $\lambda_B > 0$  in this setting allows to obtain a further boost of performances. This is evident in the scenario where Photo and Art are both the source or target domains, with Cartoon-Sketch correspond to the other pair. In this scenario the source/target pairs are quite close and enforcing a uniform assignment among the latent domains allows to obtain a better estimate of the latent domain.

**Ablation study.** We exploit the challenging multisource-multitarget scenario of Table 4 in order to assess the impact of the various components of our algorithm. In particular we show how the performance are affected if (i) a random domain is assigned to each sample; (ii) no loss is applied to the domain prediction branch; (iii) no entropy loss is applied to the classification of unlabeled target samples. From Table 4 we can easily notice that if we drop either

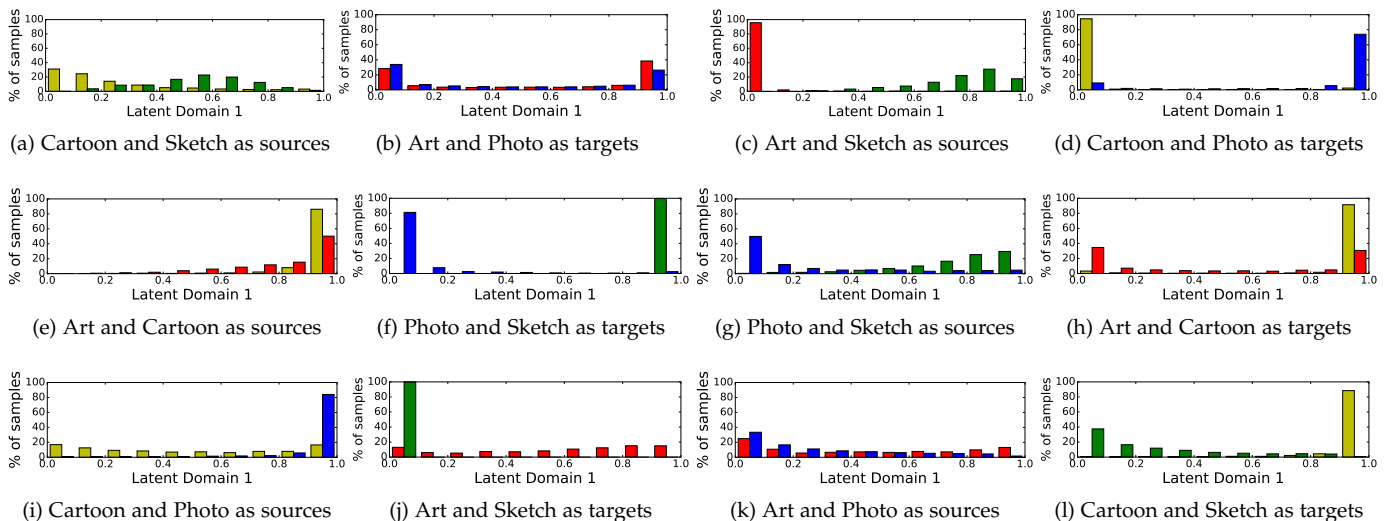


Fig. 5: Distribution of the assignments produced by the domain prediction branch in all possible multi-target settings of the PACS dataset. Different colors denote different source domains (red: Art, yellow: Cartoon, blue: Photo, green: Sketch).

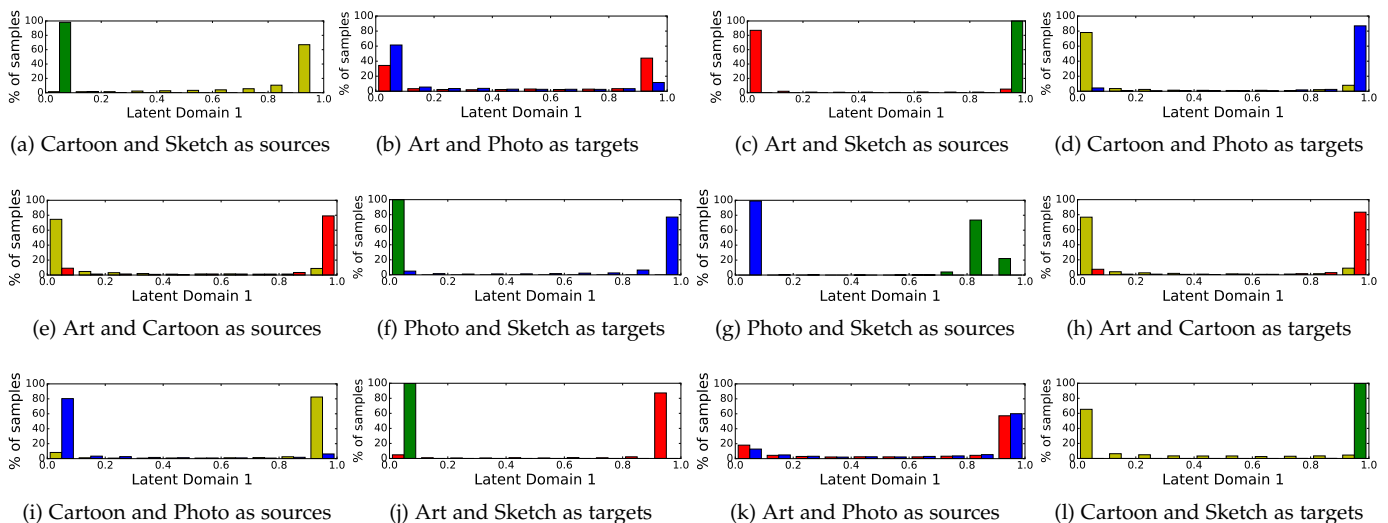


Fig. 6: Distribution of the assignments produced by the domain prediction branch trained with the additional constraint on the entropy loss in all possible multi-target settings of the PACS dataset. Different colors denote different source domains (red: Art, yellow: Cartoon, blue: Photo, green: Sketch).

the domain prediction branch (random assignment) or the losses on top of it ( $\lambda_E = \lambda_B = 0$ ), the performances of the model become comparable to the ones obtain by the DIAL baseline. This shows not only the importance of discovering latent domains, but also that both the domain branch and our losses allow to extract meaningful subsets from the data. Moreover, this demonstrates the fact that our improvements are not only due to the introduction of multiple normalization layers, but also to the latent domain discovering procedure. For what concerns the classification branch, without the entropy component on unlabelled target samples ( $\lambda_C = 0$ ), the performance of the model significantly decreases (*i.e.* from 82.6 to 76.4 in average). This confirms the findings of previous works [8], [9] about the impact that this loss for normalization based DA approaches. In particular, assuming that source and target samples of different domains are independently normalized, the entropy

loss generates a gradient flow through unlabeled samples based in the direction of its most confident prediction. This is particularly important to learn useful features even for the target domain/s, for which no supervision is available.

**In-depth analysis.** The ability of our approach to discover latent domains is further investigated on PACS. First, in Figure 3, we show how our approach assigns source samples to different latent domains in the single target setting. The four plots correspond to a single run of the experiments of Table 3. Interestingly, when either Cartoon (Figure 3c) or Sketch (Figure 3d) is the target, samples from Photo and Art tend to be associated to the same latent domain and, similarly, when either Photo (Figure 3a) or Art (Figure 3b) is the target, samples from Cartoon and Sketch are mostly grouped together. These results confirms the ability of our approach to automatically assign images of similar visual appearance to the same latent distribution. In Figure 4, we

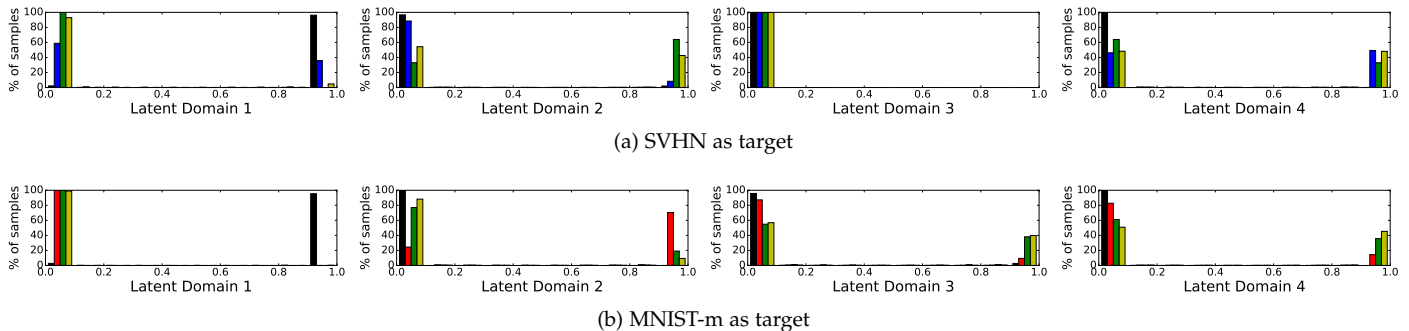


Fig. 7: Distribution of the assignments produced by the domain prediction branch for each latent domain in all target settings of the Digits-five dataset. Different colors denote different source domains (black: MNIST, blue: MNIST-m, green: USPS, red: SVHN, yellow: Synthetic numbers).

show the top-6 images associated to each latent domain for each sources/target setting. In most cases, images associated to the same latent domain have similar appearance, while there is high dissimilarity between images associated to different latent domains. Moreover, images assigned to the same latent domain tend to be associated with one of the original domains. For instance, the first row of Figure 4a contains only images from Art, while the third contains only images from Sketch. Note that no explicit domain supervision is ever given to our method in this setting.

In Figure 5, we show the histograms of the domain assignment probabilities predicted by our model with  $\lambda_B = 0$  in the various multi-source, multi-target settings of Table 3. As the figures shows, in most cases the various pairs of target domains tend to be very well separated: this justifies the large gain of performances produced by our model in this scenario. The only cases where the separation is less marked is when Art and Photo, which have very similar visual appearance, are considered as targets. On the other hand, source domains are not always as clearly separated as the targets. In particular the pairs Photo-Cartoon, Art-Photo and Art-Cartoon, tend to receive similar assignments when they are considered as source. A possible explanation is that the supervised source loss could have a stronger influence on the domain assignment than the unsupervised target one. In any case, note that these results do not detract from the validity of our approach. In fact, our main objective is to obtain a good classification model for the target set, independently from the actual domain assignments we learn.

In Figure 6, the same analysis is performed on our method with the additional constraint of having a uniform assignment distribution among domains. As the figure shows, this constraint allows to obtain a clearer domain separation in most of the cases, overcoming the difficulties that the domain prediction branch experienced in separating domain pairs such as Photo-Cartoon and Photo-Art.

We perform a similar analysis in another dataset, Digits-five. The results are reported in Figure 7. As the figure shows, when SVHN is the target domain, one of the latent domains (latent domain 1) receives very confident assignments for the samples of the MNIST dataset. The samples of the other source datasets receive assignment spread through all the latent domains, with the exceptions of USPS which receives the most confident predictions for the second latent domain and MNIST-m, which partially

influences the first latent domain, the one with confidence assignments to MNIST. One latent domain does not receive assignments from any of the sources (latent domain three): this might happen if the entropy term overcomes the uniform assignment constraints in the early stages of training. Similarly, when MNIST-m is the target domain, the first two latent domains receive confident assignments for samples belonging to MNIST and SVHN datasets respectively, while the third and the fourth receive higher assignments for samples of the remaining source domains.

#### 4.2.3 Experiments on Office-31

In our Office-31 experiments we consider the following baselines, trained on the union of the source sets: (i) a plain AlexNet network; (ii) AlexNet with BN inserted after each fully-connected layer; and (iii) AlexNet + DIAL [9]. Additionally, we consider single source domain adaptation approaches, using the results reported in [63]. The methods are Transfer Component Analysis (TCA) [50], Geodesic Flow Kernel (GFK) [20], Deep Domain Confusion (DDC) [59], Deep Reconstruction Classification Networks (DRCN) [17] and Residual Transfer Network (RTN) [40], as well as the Reversed Gradient (RevGrad) [15] and Domain Adaptation Network (DAN) [39] algorithms considered in the digits experiments. For these algorithms we report the performances obtained in the “Best single source” and “Unified sources” settings, as available from [63]. As in the previous experiments, Multi-source DA with perfect domain knowledge can be regarded as a performance upper bound for our method. Finally, we include results reported in [63] for different multi-source DA models: Deep Cocktail Network (DCTN) [63], the two shallow methods in [61] (sFRAME) and [22] (SGF), and an ensemble of baseline networks trained on each source domain separately (Source only). These results are summarized in Table 5.

We note that, in this dataset, the improvements obtained by adopting a multi-source model instead of a single-source one are small. This is in accordance with findings in [33], where it is shown that the domain shift in Office-31, when considering deep features, is indeed quite limited if compared to PACS, and it is mostly linked to changes in the background (Webcam-Amazon, DSLR-Amazon) or acquisition camera (DSLR-Webcam). This is further supported by the smaller gap between DIAL and our method in this case compared to the previous experiments (average  $p^*$  of 0.54).



TABLE 4: PACS dataset: comparison of different methods using the ResNet architecture on the multi-source multi-target setting. The first row indicates the two target domains.

Method	Photo-Art	Photo-Cartoon	Photo-Sketch	Art-Cartoon	Art-Sketch	Cartoon-Sketch	Mean
ResNet [26]	71.4	84.2	81.4	62.2	70.3	54.2	70.6
DIAL [9]	86.7	86.5	86.8	77.1	72.1	67.7	79.5
Random assignment	86.6	86.7	85.9	76.2	69.1	69.4	79.1
Ours $\lambda_E = \lambda_B = 0$	86.8	86.5	86.7	78.6	73.8	68.7	80.2
Ours $\lambda_B = \lambda_C = 0$	82.4	85.0	83.7	71.7	74.0	68.8	76.4
Ours $\lambda_B = 0$	86.1	87.9	87.9	<b>79.3</b>	79.9	74.9	82.6
Ours	<b>87.2</b>	<b>88.1</b>	<b>88.7</b>	<b>77.7</b>	<b>81.3</b>	<b>77.0</b>	<b>83.3</b>
Multi-source/target DA	87.7	88.9	86.8	79.0	79.8	75.6	83.0

TABLE 5: Office-31 dataset: comparison of different methods using AlexNet. In the first row we indicate the source (top) and the target domains (bottom).

	Method	Source Target	A-W D	A-D W	W-D A	Mean
Best single source [63]	TCA [50]		95.2	93.2	51.6	68.8
	GFK [20]		95.0	95.6	52.4	68.7
	DDC [59]		98.5	95.0	52.2	70.7
	DRCN [17]		99.0	96.4	56.0	73.6
	RevGrad [15]		99.2	96.4	53.4	74.3
	DAN [39]		99.0	96.0	54.0	72.9
	RTN [40]		99.6	96.8	51.0	73.7
	Unified sources	Source only from [63]		98.1	93.2	50.2
Source only (ours)			94.6	89.1	49.1	77.6
Source only+BN (ours)			91.9	92.7	46.5	77.0
RevGrad [63]			<b>98.8</b>	<b>96.2</b>	54.6	83.2
DAN [63]			<b>98.8</b>	95.2	53.4	82.5
Single BN			92.9	95.2	60.1	82.7
DIAL [9]			93.8	94.3	62.5	83.5
Ours $\lambda_B = 0$			93.7	94.6	<b>62.6</b>	83.6
Ours			93.6	93.6	62.4	83.2
Multi-source		Source only [63]		98.2	92.7	51.6
	sFRAME [61]		54.5	52.2	32.1	46.3
	SGF [22]		39.0	52.0	28.0	39.7
	DCTN [63]		99.6	96.9	54.9	83.8
	Multi-source DA		94.8	95.8	62.9	84.5

In this setting, introducing our uniform loss term does not provide boost in performances. We ascribe this behaviour to the fact that in this scenario, each batch is built with a non-uniform number of samples per domain (following [8]) while our current objective assumes a balanced sampling among domains.

In a final Office-31 experiment, we consider a setting where the true domain of a subset of the source samples is known at training time. Figure 8 shows the average accuracy obtained when a different amount of domain labels are available. Interestingly, by increasing the level of domain supervision the accuracy quickly saturates towards the value of Multi-source DA, completely filling the gap with as few as 5% of the source samples.

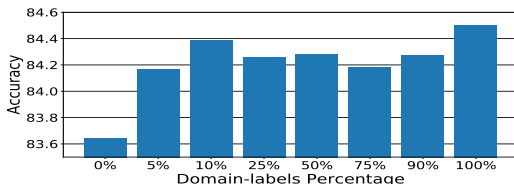


Fig. 8: Office31 dataset. Performance at varying number of domain labels (%) for source samples.

#### 4.2.4 Comparison with S.o.t.A. on inferring latent domains

In this section we compare the performance of our approach with previous works on DA which also consider the problem of inferring latent domains [19], [27], [62]. As stated in Section 2, there are no previous works adopting deep

TABLE 6: Office-31: comparison with state-of-the-art algorithms. In the first row we indicate the source (top) and the target domains (bottom).

Method	Sources Target	A-D W	A-W D	W-D A	Mean
Hoffman <i>et al.</i> [27]		24.8	42.7	12.8	26.8
Xiong <i>et al.</i> [62]		29.3	43.6	13.3	28.7
Gong <i>et al.</i> (AlexNet) [19]		91.8	94.6	48.9	78.4
Ours $\lambda_B = 0$		93.1	94.3	64.2	83.9
Ours		<b>94.5</b>	<b>94.9</b>	<b>64.9</b>	<b>84.8</b>
Gopalan <i>et al.</i> [23]		51.3	36.1	35.8	41.1
Nguyen <i>et al.</i> [49]		64.5	68.6	41.8	58.3
Lin <i>et al.</i> [37]		73.2	81.3	41.1	65.2

learning models (i) in a multi-source setting and (ii) discovering hidden domains. Therefore, the methods we compare to all employ handcrafted features. For these approaches we report results taken from the original papers. Furthermore, we evaluate the method of Gong *et al.* [19] using features from the last layer of the AlexNet architecture. For a fair comparison, when applying our method we freeze AlexNet up to `fc7`, and apply mDA layers only after `fc7` and the classifier.

We first consider the Office-31 dataset, as this benchmark has been used in [27], [62], showing the results in Table 6. Our model outperforms all the baselines, with a clear margin in terms of accuracy. Importantly, even when the method in [19] is applied to features derived from AlexNet, still our approach leads to higher accuracy. For the sake of completeness, in the same table we also report results from previous multi-source DA methods [23], [37], [49] based on shallow models. While these approaches significantly outperform [27] and [62], still their accuracy is much lower than ours. Moreover, introducing our novel loss term provides higher performances with respect to the our approach with  $\lambda_B = 0$ .

To provide a comparison in a multi-target scenario, we also consider the Office-Caltech dataset, comparing our model with [19], [27]. Following [19], we test both single target (Amazon) and multi-target (Amazon-Caltech and Webcam-DSLR) scenarios. As for the PACS multi-source/multi-target case, the assignment of each sample to the source or target set is assumed to be known, while the assignment to the specific domain is unknown. We again want to remark that, since we do not assume to know the target domain to which a sample belongs, the task is even harder since we require a domain prediction step also at test time. As in the Office-31 experiments, our approach outperforms all baselines, including the method in [19] applied to AlexNet features. In this scenario, introducing our uniform loss provides a boost in performances in the multi-target setting, where the two source/target pairs have

TABLE 7: Office-Caltech dataset: comparison with state-of-the-art algorithms. In the first row we indicate the source (top) and the target domains (bottom).

Method	Source Target	A-C W-D	W-D A-C	C-W-D A	Mean
Gong <i>et al.</i> [19] - original		41.7	35.8	41.0	39.5
Hoffman <i>et al.</i> [27] - ensemble		31.7	34.4	38.9	35.0
Hoffman <i>et al.</i> [27] - matching		39.6	34.0	34.6	36.1
Gong <i>et al.</i> [19] - ensemble		38.7	35.8	42.8	39.1
Gong <i>et al.</i> [19] - matching		42.6	35.5	44.6	40.9
Gong <i>et al.</i> (AlexNet) [19] - ensemble		87.8	87.9	93.6	89.8
Ours $\lambda_B = 0$		93.5	88.2	93.7	91.8
Ours		<b>95.0</b>	<b>88.7</b>	<b>93.9</b>	<b>92.5</b>

similar appearance. This is inline to what reported for the multitarget experiments on PACS (Table 4).

## 5 CONCLUSIONS

In this work we presented a novel deep DA model for automatically discovering latent domains within visual datasets. The proposed deep architecture is based on a side-branch which computes the assignment of source and target samples to their associated latent domain. These assignments are then exploited within the main network by novel domain alignment layers which reduce the domain shift by aligning the feature distributions of the discovered sources and the target domains. Our experimental results demonstrate the ability of our model to efficiently exploit the discovered latent domains for addressing challenging domain adaptation tasks. Future works will investigate other architectural design choices for the domain prediction branch, as well as the possibility to integrate it into other CNN models for unsupervised domain adaptation [15].

## ACKNOWLEDGMENTS

We acknowledge financial support from ERC grant 637076 - RoboExNovo and project DIGIMAP, grant 860375, funded by the Austrian Research Promotion Agency (FFG). This work was carried out under the "Vision and Learning joint Laboratory" between FBK and UNITN.

## REFERENCES

- [1] R. Aljundi and T. Tuytelaars. Lightweight unsupervised domain adaptation by convolutional filter reconstruction. In *ECCV TASK-CV Workshops*, 2016.
- [2] G. Angeletti, B. Caputo, and T. Tommasi. Adaptive deep learning through visual domain localization. In *ICRA*, 2018.
- [3] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE T-PAMI*, 33(5):898–916, 2011.
- [4] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- [5] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with gans. In *CVPR*, 2017.
- [6] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *CVPR*, 2017.
- [7] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *NIPS*, 2016.
- [8] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. Rota Bulò. Autodial: Automatic domain alignment layers. *ICCV*, 2017.
- [9] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. Rota Bulò. Just dial: Domain alignment layers for unsupervised domain adaptation. In *ICIAP*, 2017.
- [10] K. Crammer, M. Kearns, and J. Wortman. Learning from multiple sources. *JMLR*, 9:1757–1774, 2008.
- [11] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*, 2014.
- [12] L. Duan, I. W. Tsang, D. Xu, and T.-S. Chua. Domain adaptation from multiple sources via auxiliary classifiers. In *ICML*, 2009.
- [13] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *ICCV*, 2013.
- [14] J. Friedman, T. Hastie, and R. Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [15] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. *ICML*, 2015.
- [16] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(59):1–35, 2016.
- [17] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li. Deep reconstruction-classification networks for unsupervised domain adaptation. In *ECCV*, 2016.
- [18] B. Gong, K. Grauman, and F. Sha. Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation. In *ICML*, 2013.
- [19] B. Gong, K. Grauman, and F. Sha. Reshaping visual datasets for domain adaptation. In *NIPS*, 2013.
- [20] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, 2012.
- [21] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [22] R. Gopalan, R. Li, and R. Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *ICCV*, 2011.
- [23] R. Gopalan, R. Li, and R. Chellappa. Unsupervised adaptation across domain shifts by generating intermediate data representations. *IEEE T-PAMI*, 36(11):2288–2302, 2014.
- [24] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007.
- [25] P. Haeusser, T. Frerix, A. Mordvintsev, and D. Cremers. Associative domain adaptation. In *ICCV*, 2017.
- [26] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [27] J. Hoffman, B. Kulis, T. Darrell, and K. Saenko. Discovering latent domains for multisource domain adaptation. In *ECCV*, 2012.
- [28] J. Huang, A. Gretton, K. M. Borgwardt, B. Schölkopf, and A. J. Smola. Correcting sample selection bias by unlabeled data. In *NIPS*, 2006.
- [29] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [30] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACM-Multimedia*, 2014.
- [31] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [32] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [33] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. *ICCV*, 2017.
- [34] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, 2018.
- [35] W. Li, Z. Xu, D. Xu, D. Dai, and L. Van Gool. Domain generalization and adaptation using low rank exemplar svms. *IEEE T-PAMI*, 2017.
- [36] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou. Revisiting batch normalization for practical domain adaptation. *ICLR-WS*, 2017.
- [37] Y. Lin, J. Chen, Y. Cao, Y. Zhou, L. Zhang, Y. Y. Tang, and S. Wang. Cross-domain recognition by identifying joint subspaces of source domain and target domain. *IEEE Transactions on Cybernetics*, 47(4):1090–1101, 2017.
- [38] M. Long, G. Ding, J. Wang, J. Sun, Y. Guo, and P. S. Yu. Transfer sparse coding for robust image representation. In *CVPR*, 2013.
- [39] M. Long and J. Wang. Learning transferable features with deep adaptation networks. In *ICML*, 2015.
- [40] M. Long, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. *NIPS*, 2016.
- [41] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Deep transfer learning with joint adaptation networks. In *ICML*, 2017.
- [42] M. Mancini, S. R. Bulò, B. Caputo, and E. Ricci. Best sources forward: domain generalization through source-specific nets. In *ICIP*, 2018.
- [43] M. Mancini, S. R. Bulò, B. Caputo, and E. Ricci. Robust place

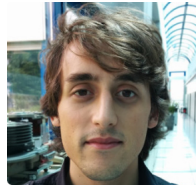
categorization with deep domain generalization. *IEEE Robotics and Automation Letters*, 3(3):2093–2100, 2018.

- [44] M. Mancini, L. Porzi, S. R. Bulò, B. Caputo, and E. Ricci. Boosting domain adaptation by discovering latent domains. In *CVPR*, 2018.
- [45] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain adaptation with multiple sources. In *NIPS*, 2009.
- [46] P. Morerio, J. Cavazza, and V. Murino. Minimal-entropy correlation alignment for unsupervised deep domain adaptation. *ICLR*, 2018.
- [47] K. Muandet, D. Balduzzi, and B. Schölkopf. Domain generalization via invariant feature representation. In *ICML*, 2013.
- [48] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS-WS on deep learning and unsupervised feature learning*, 2011.
- [49] H. V. Nguyen, H. T. Ho, V. M. Patel, and R. Chellappa. Dash-n: Joint hierarchical domain adaptation and feature learning. *IEEE T-IP*, 24(12):5479–5491, 2015.
- [50] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2011.
- [51] P. Russo, F. M. Carlucci, T. Tommasi, and B. Caputo. From source to target and back: symmetric bi-directional adaptive gan. In *CVPR*, 2018.
- [52] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. *ECCV*, 2010.
- [53] K. Saito, Y. Ushiku, and T. Harada. Asymmetric tri-training for unsupervised domain adaptation. *arXiv preprint arXiv:1702.08400*, 2017.
- [54] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chellappa. Generate to adapt: Aligning domains using generative adversarial networks. In *CVPR*, 2018.
- [55] O. Sener, H. O. Song, A. Saxena, and S. Savarese. Learning transferable representations for unsupervised domain adaptation. In *NIPS*. 2016.
- [56] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, 2017.
- [57] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. In *AAAI*, 2016.
- [58] Q. Sun, R. Chattopadhyay, S. Panchanathan, and J. Ye. A two-stage weighting framework for multi-source domain adaptation. In *NIPS*, 2011.
- [59] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *ICCV*, 2015.
- [60] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Adversarial discriminative domain adaptation. In *CVPR*, 2017.
- [61] J. Xie, W. Hu, S.-C. Zhu, and Y. N. Wu. Learning sparse frame models for natural image patterns. *IJCV*, 114(2-3):91–112, 2015.
- [62] C. Xiong, S. McCloskey, S.-H. Hsieh, and J. J. Corso. Latent domains modeling for visual domain adaptation. In *AAAI*, 2014.
- [63] R. Xu, Z. Chen, W. Zuo, J. Yan, and L. Lin. Deep cocktail network: Multi-source unsupervised domain adaptation with category shift. In *CVPR*, 2018.
- [64] Z. Xu, W. Li, L. Niu, and D. Xu. Exploiting low-rank structure from latent domains for domain generalization. In *ECCV*, 2014.
- [65] M. Yamada, L. Sigal, and M. Raptis. No bias left behind: Covariate shift adaptation for discriminative 3d pose estimation. In *ECCV*, 2012.
- [66] J. Yang, R. Yan, and A. G. Hauptmann. Adapting svm classifiers to data with shifted distributions. In *ICDM-WS 2007*, 2007.
- [67] X. Zeng, W. Ouyang, M. Wang, and X. Wang. Deep learning of scene-specific classifier for pedestrian detection. In *ECCV*, 2014.
- [68] H. Zhao, S. Zhang, G. Wu, J. ao P. Costeira, J. M. F. Moura, and G. J. Gordon. Multiple source domain adaptation with adversarial learning. In *ICLR-WS*, 2018.



**Massimiliano Mancini** Massimiliano Mancini is a Ph. D. student with Sapienza University of Rome, in collaboration with Fondazione Bruno Kessler (FBK). He received his master degree in Artificial Intelligence and Robotics from Sapienza University of Rome in 2016. Currently, he is a member of the Visual Learning and Multimodal Applications Laboratory, led by Prof. Barbara Caputo, and the Technologies of Vision Laboratory in FBK. He was also visiting Ph.D. student in the Robotics, Perception and Learning

Laboratory at KTH Royal Institute of Technology in Stockholm.



**Lorenzo Porzi** Lorenzo Porzi is a researcher with Mapillary Research, focusing on deep learning applied to semantic and 3D vision. He received his PhD in computer science from the University of Perugia and Fondazione Bruno Kessler (FBK) in 2016. Before joining Mapillary Research, Lorenzo held a post-doctoral research position at IRI-UPC in Barcelona, working with the group of Dr. Francesc Moreno-Noguer. While in FBK, he participated in the VENTURI and REPLICATE EU projects.



**Samuel Rota Buló** Samuel Rota Buló received the PhD in computer science at University of Venice in 2009. He worked there as PostDoc and held teaching positions until 2013. He then became a researcher for FBK in computer vision and machine learning. Since 2017, he is a senior researcher with Mapillary Research. He was awarded the prestigious Marr Prize in 2015. He serves on the editorial board for "Pattern Recognition" and "International Journal of Machine Learning and Cybernetics" and is regularly on the program committee of international conferences of his field. He participated to several EU projects (SIMBAD, VENTURI, REPLICATE).



**Barbara Caputo** Barbara Caputo is Full Professor at the DAUIN Department of Control and Computer Engineering of Politecnico di Torino and Principal Investigator at the Italian Institute of Technology (IIT), where she leads the Visual Learning and Multimodal Applications Laboratory (VANDAL). Her main research interest is to develop algorithms for learning, recognition and categorization of visual and multimodal patterns for artificial autonomous systems. These features are crucial to enable robots to represent

and understand their surroundings, to learn and reason about it, and ultimately to equip them with cognitive capabilities. Her research is sponsored by the Swiss National Science Foundation (SNSF), the Italian Ministry for Education, University and Research (MIUR), the European Commission (EC) and the European Research Council (ERC).



**Elisa Ricci** Elisa Ricci is an associate professor at University of Trento and a researcher at Fondazione Bruno Kessler. She received her PhD from the University of Perugia in 2008. She has since been a post-doctoral researcher at Idiap Research Institute and an assistant professor at University of Perugia. She was also a visiting researcher at University of Bristol. Her research interests are mainly in the areas of computer vision and machine learning.



## APPENDIX

From the main text, we have the output of our mDA layer denoted by

$$y_i = \text{mDA}(x_i, \mathbf{w}_i; \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}) = \sum_{d \in \mathcal{D}} w_{i,d} \hat{x}_{i,d}, \quad (7)$$

where, for simplicity:

$$\hat{x}_{i,d} = \frac{x_i - \hat{\mu}_d}{\sqrt{\hat{\sigma}_d^2 + \epsilon}}, \quad (8)$$

and the statistics are given by

$$\begin{aligned} \hat{\mu}_d &= \sum_{i=1}^b \hat{w}_{i,d} x_i, \\ \hat{\sigma}_d^2 &= \sum_{i=1}^b \hat{w}_{i,d} (x_i - \hat{\mu}_d)^2, \end{aligned} \quad (9)$$

where  $\hat{w}_{i,d} = w_{i,d} / \sum_{j=1}^b w_{j,d}$ .

From the previous equations we can derive the partial derivative of the loss function with respect to both the input  $x_i$  and the domain assignment probabilities  $w_{i,d}$ . Let us denote  $\frac{\partial L}{\partial y_i}$  the partial derivative of the loss function  $L$  with respect to the output  $y_i$  of the mDA layer. We have:

$$\begin{aligned} \frac{\partial \hat{x}_{i,d}}{\partial \hat{\sigma}_{d^*}^2} &= -\mathbb{1}_{d=d^*} \frac{1}{2} (x_i - \hat{\mu}_{d^*}) \cdot (\hat{\sigma}_{d^*}^2 + \epsilon)^{-\frac{3}{2}}, \\ \frac{\partial \hat{x}_{i,d}}{\partial \hat{\mu}_{d^*}} &= -\mathbb{1}_{d=d^*} (\hat{\sigma}_{d^*}^2 + \epsilon)^{-\frac{1}{2}}, \end{aligned} \quad (10)$$

and

$$\frac{\partial \hat{\sigma}_d^2}{\partial x_i} = 2 \hat{w}_{i,d} \cdot (x_i - \hat{\mu}_d), \quad \frac{\partial \hat{\mu}_d}{\partial x_i} = \hat{w}_{i,d}. \quad (11)$$

Thus, the partial derivative of  $L$  w.r.t. the input  $x_i$  is:

$$\frac{\partial L}{\partial x_{i^*}} = \sum_{d \in \mathcal{D}} \frac{w_{i^*,d}}{\sqrt{\hat{\sigma}_d^2 + \epsilon}} \left[ \frac{\partial L}{\partial y_{i^*}} - A_d - \hat{x}_{i^*,d} B_d \right], \quad (12)$$

where:

$$\begin{aligned} A_d &= \sum_{i=1}^b \hat{w}_{i,d} \frac{\partial L}{\partial y_i}, \\ B_d &= \sum_{i=1}^b \hat{w}_{i,d} \hat{x}_{i,d} \frac{\partial L}{\partial y_i}. \end{aligned} \quad (13)$$

For the domain assignment probabilities  $w_{i,d}$  we have:

$$\frac{\partial \hat{\mu}_d}{\partial w_{i^*,d^*}} = \mathbb{1}_{d=d^*} x_{i^*}, \quad (14)$$

$$\frac{\partial \hat{\sigma}_d^2}{\partial w_{i^*,d^*}} = \mathbb{1}_{d=d^*} (x_{i^*} - \hat{\mu}_d)^2, \quad (15)$$

$$\frac{\partial \hat{w}_{i,d}}{\partial w_{i^*,d^*}} = \mathbb{1}_{d=d^*} \frac{\mathbb{1}_{i=i^*} \sum_{j=1}^b w_{j,d} - w_{i^*,d}}{(\sum_{j=1}^b w_{j,d})^2}. \quad (16)$$

Thus, the partial derivative of  $L$  w.r.t.  $w_{i,d}$  is:

$$\frac{\partial L}{\partial w_{i^*,d}} = \hat{x}_{i^*,d} \left( \frac{\partial L}{\partial y_{i^*}} - A_d \right) - \frac{1}{2} \left( \hat{x}_{i^*,d}^2 - \frac{\hat{\sigma}_d^2}{\hat{\sigma}_d^2 + \epsilon} \right) B_d, \quad (17)$$

where  $A_d$  and  $B_d$  are defined as in (13).

In this section, we plot the losses as the training progresses for the Digits-five experiments. The plots are shown

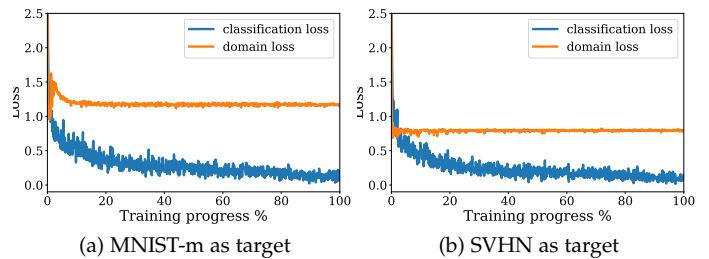


Fig. 9: Digits-five: plots of the domain (orange) and classification (blue) losses during the training phase.

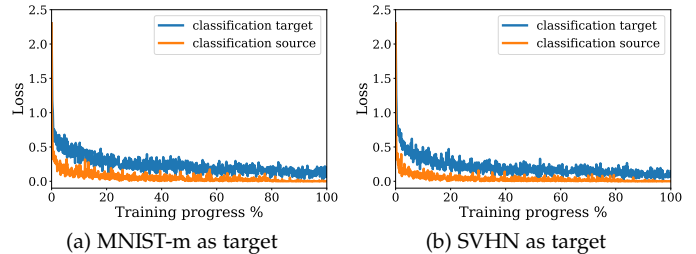


Fig. 10: Digits-five: plots of the cross-entropy loss on source samples (orange) and entropy loss on target sample (blue) for the semantic classifier during the training phase.

in Figure 9. For both MNIST-m and SVHN, the classification loss smoothly decreases, while the domain loss first decreases and then stabilizes around a fixed value. This is a consequence of the introduced balancing term on the domain assignments, which enforces the entropy to be low for the assignment of a single sample, but high for the assignments averaged across the entire batch. In Figures 10 and 11 we plot the single components of the classification and domain loss respectively. For the semantic part (Figure 10), both the entropy loss on target sample and the cross-entropy loss on source samples decrease smoothly. For the domain assignment part (Figure 11), we can see how the entropy loss on single samples rapidly decreases, while the average batch assignment keeps an high entropy, as expected. We highlight that when SVHN is used as target, the source domains are a bit closer to each other in appearance, thus the average batch entropy has a slightly lower value (*i.e.* the assignments are less balanced) with respect to the MNIST-m as target case.

Finally, it is worth noticing that the domain loss reaches a stable value earlier than the classification components. This is a design choice, since we want to learn a semantic predictor on stable and confident domain assignments.

A crucial problem in domain adaptation rarely addressed in the literature is how to tune model hyper-parameters. In fact, setting the hyper-parameters values based on the performance on the source domain is sub-optimal, due to the domain shift. Furthermore, assuming the presence of a validation set for the target domain is not realistic in practice [46]: in unsupervised domain adaptation we only assume the presence of a set of unlabelled target data. Despite recent research in this direction [46], there is no clear solution to this problem in the literature. This problem

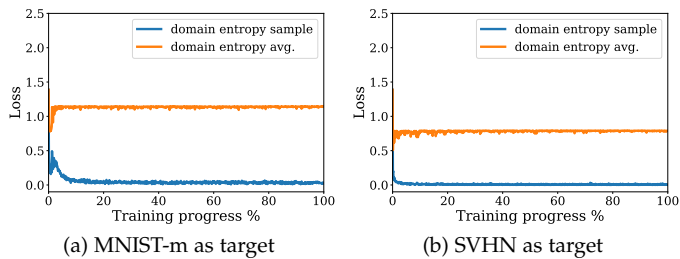


Fig. 11: Digits-five: plots of the entropy loss on single sample (blue) and on the average batch assignments (orange) for the domain classifier during the training phase.

TABLE 8: PACS dataset: comparison of different methods using the ResNet architecture. The first row indicates the target domain, while all the others are considered as sources. The numbers in parenthesis indicate the results using a target validation set for model selection.

Method	Sketch	Photo	Art	Cartoon	Mean
ResNet [26]	60.1	92.9	74.7	72.4	75.0
DIAL [9]	66.8 (71.3)	97.0 (97.4)	87.3 (87.5)	85.5 (87.0)	84.2 (85.8)
Ours	<b>70.7 (75.2)</b>	<b>97.0 (97.3)</b>	<b>87.4 (87.7)</b>	<b>86.3 (87.2)</b>	<b>85.4 (86.9)</b>
Multi-source DA	71.6 (78.1)	96.6 (97.2)	87.5 (88.7)	87.0 (87.4)	85.7 (87.9)

is more severe in our case, since it is not trivial to define a validation set for the latent domain discovery problem, due to the assumption that multiple source and target domains are mixed.

Nonetheless, for the sake of completeness, we analyze the performances of our model and the baselines if we assume the presence of a target validation set to perform model selection. We consider the PACS dataset, in both the single and multi-target scenarios. The results are reported in parenthesis in Table 8 and in Table 9. While our model and the baselines obviously benefit from the validation set, the overall trends remain the same, with our model achieving higher performances with respect to the baseline and close to the multi-source upper bound. Notice that a validation set is especially beneficial in the case of consistent domain shift: for instance, all the methods increase their results by almost 5% in Table 8 when Sketch is the target domain.

As a final note, we underline that the use of a validation set on the target domain for unsupervised domain adaptation is not a common practice in the community, thus these results can be regarded as an upper bound with respect to our model.

TABLE 9: PACS dataset: comparison of different methods using the ResNet architecture on the multi-source multi-target setting. The first row indicates the two target domains. The numbers in parenthesis indicate the results using a target validation set for model selection.

Method	Photo-Art	Photo-Cartoon	Photo-Sketch	Art-Cartoon	Art-Sketch	Cartoon-Sketch	Mean
ResNet [26]	71.4	84.2	81.4	62.2	70.3	54.2	70.6
DIAL [9]	86.7 (87.5)	86.5 (87.1)	86.8 (88.2)	77.1 (78.7)	72.1 (74.2)	67.7 (70.4)	79.5 (81.0)
Ours	<b>87.2 (87.7)</b>	<b>88.1 (88.5)</b>	<b>88.7 (89.7)</b>	<b>77.7 (79.6)</b>	<b>81.3 (82.2)</b>	<b>77.0 (79.3)</b>	<b>83.3 (84.5)</b>
Multi-source/target DA	87.7 (88.8)	88.9 (89.8)	86.8 (88.3)	79.0 (79.5)	79.8 (82.2)	75.6 (79.1)	83.0 (84.6)