

This CVPR workshop paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

Visual Domain Bridge: A source-free domain adaptation for cross-domain few-shot learning

Moslem Yazdanpanah Parham Moradi University of Kurdistan

{M.Yadanpanah, P.Moradi}@UOK.ac.ir

Abstract

Due to the covariate shift, deep neural networks performance always degrades when applied to novel domains. In order to mitigate this problem, domain adaptation techniques require samples from target data during the feature extraction training, which is not always applicable in realworld scenarios. Batch Normalization is a known component of computer vision models, aiming at reducing the training-time covariate shift. However, facing distribution shift results in an internal state mismatch inside the Batch-Norm layers during the inference time. In favor of alleviating the induced mismatch, this paper proposes a sourcefree, lightweight and straightforward approach by introducing the "Visual Domain Bridge" concept reducing the BatchNorm's internal mismatch in the cross-domain settings. Compared to the other BatchNorm-based sourcefree domain adaptation techniques such as AdaBN and Prediction-BN, our method formed a new state-of-the-art cross-domain few-shot fine-tuning method neglecting extra augmentations; while improving the performance in neardomain settings too. The proposed method can integrate with other domain adaptation methods and enhance their performance requiring just a few lines of modification in the BatchNorm's implementation. Implementations are available in https://github.com/MosyMosy/VDB

1. Introduction

Although Convolutional Neural Network (CNN) [10] achieved outstanding results in different deep learning tasks, when facing domain shift, their performance falls drastically. The problem of applying a learned representation from a source domain to a target domain under the distribution shift gets even worse when the target domain is rarely annotated [4]. In order to mitigate this shortcoming, there is a large body of research as well as challenges, competitions, and benchmarks, conducted toward transfer learning across distant domains [31].



Figure 1. Illustration of the Visual Domain Bridge (VDB). Inside a CNN's Batch Normalization (BN) layers, the intermediate features get transferred to the source domain distribution when facing a target domain with domain shift.

Domain adaption leveraging the statistical distribution of data is a well-studied technique in the literature [11]. However, most of these efforts focus on computationally expensive approaches benefiting target samples to make a robust representation under domain shift [20].

Batch Normalization (BN) [8] is a well-known component of deep CNN which speeds up model convergence during the training. Adding BN layers to deep learning models stabilizes the distribution of layer input features by mod-



Figure 2. The output of a pre-trained autoencoder when fed with different datasets in inference time. The top row shows the output when all normalization layers are the Batch Normalization (BN). The bottom row depicts the output for the same inputs but, the BN layers are replaced with the VDB. Best viewed in color.

ulating their mean and variance. Inspired from "label information is usually stored in network's weight matrix and the statistics of the batch norm layer represents the domainrelated knowledge." [13], there are attempts to utilize the BN layers toward filling the domain gap between the source and target domains. Utilizing the BN layers by the mean of domain adaptation as a tempting lightweight approach attracted a portion of the literature; Studies such as the AdaBN [13], Prediction-BN [17], TENT [24] and Core [30]. Despite the effectiveness of these proposed methods, they suffer from induced internal state mismatch. To alleviate this problem, they need steps of feature extractor update during the inference time. Both of these shortcomings have turned these methods into less real-world practical approaches.

Few-shot learning (FSL) as meta-testings followed by meta-learning steps refer to the problem of learning a task from just a few samples of the task's distribution. Despite its effectiveness and recent advances, the FSL underperform when there exists a significant shift between base and novel class domains [4]. This work proposes a source-free fewshot domain adaptation method by adapting the BN statistic as a lightweight interpretable approach, as depicted in Figure 1. Based on our knowledge, the Prediction-BN is the current state-of-the-art (SOTA) in source-free BN based methods, which does not utilize additional augmentations or backbone updating; the proposed method formed a new SOTA in cross-domain few-shot learning (CDFSL), outperforming the Prediction-BN on the ImageNet source CDFSL benchmark score, by more than two percents.

While most of the Domain Adaptation approaches aim at adapting a pre-trained model towards a target domain (with a possible distribution shift in input features); Inspired from [7], we propose a contradictory viewpoint by transferring the input's statistical elements to the training times accumulated statistics, inside BatchNorm Layers. This process accrues in test time without modifying the training procedure or the pre-trained backbone. From this, we attempt to decrease the internal state mismatch of the BNs, rooted in the shift between the source and target data (Figure 1).

The rest of the paper is structured as follows: First, the needed notations are represented, followed by the formal definition of the proposed method. Then we recall a brief history of the current related methods. After that, experiments setup and evaluation results are presented, followed by the analysis of the result. In the end, the paper is summed up with a conclusion.

2. Methodology

2.1. Preliminaries

For consistency with the literature, we adapt the notations from [25].

2.1.1 Domain

A domain \mathcal{D} is composed of a feature space \mathcal{X} and a marginal probability distribution P(X), where $X = \{x_1, ..., x_n\} \in \mathcal{X}$. Having a specific domain $\mathcal{D} = \{\mathcal{X}, P(X)\}$, a task \mathcal{T} consists of a feature space \mathcal{Y} and a conditional probability distribution P(Y|X); learnable in a supervised manner from the labeled data $\{x_i, y_i\}$, where $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$. Given two domains: a training dataset with adequate labeled data as the source domain $\mathcal{D}^s = \{\mathcal{X}^s, P(X)^s\}$, and a test dataset with insufficient labeled data as the target domain $\mathcal{D}^t = \{\mathcal{X}^t, P(X)^t\}$; with partially labeled part, \mathcal{D}^{tl} , and the unlabeled parts, \mathcal{D}^{tu} . Form the entire target domain, that is, $\mathcal{D}^t = \mathcal{D}^{tl} \cup$ \mathcal{D}^{tu} . Each domain comes with its task: The source task $\mathcal{T}^s = \{\mathcal{Y}^s, P(Y^s|X^s)\}$, and the target one is $\mathcal{T}^t = \{\mathcal{Y}^t, P(Y^t|X^t)\}$. **Distribution Shift (Domain shift)** refers to situations where feature distribution differs from D^s to D^t : $P(X)^s \neq P(X)^t$

Near-Domain refers to situations where feature distribution are the same between D^s and D^t , but conditional probability distribution differs from \mathcal{T}^s to \mathcal{T}^t : $P(X)^s = P(X)^t$ and $P(Y^s|X^s) \neq P(Y^t|X^t)$

Cross-Domain refers to situations where feature distribution differs from D^s to D^t while conditional probability distribution differs from \mathcal{T}^s to \mathcal{T}^t too: $P(X)^s \neq P(X)^t$ and $P(Y^s|X^s) \neq P(Y^t|X^t)$

2.1.2 Batch Normalization

Having a batch of labeled examples $\{(x_i^s, y_i^s)\}_{i=1}^N$ of size N from a source domain \mathcal{D}^s where $x_i^s \in \mathcal{X}^s$ and $y_i^s \in \mathcal{Y}^s$, and Θ as a deep convolutional neural network consisting of L layers with weight matrices θ^l where l represents the layer index. If h represents the intermediate features of Θ for layer l, the Batch Normalization at layer l is computed for each channel and can be defined as:

$$BN(h_c^s) = \gamma^s \times \frac{h_c^s - \mu_c^s}{\sqrt{\sigma_c^{2^s} + \epsilon}} + \beta^s \tag{1}$$

here, subscript c represents the channel index, γ^s and β^s are learnable affine parameters and μ_c^s and σ_c^s are statistical mean and variance of h_c^s respectively defined as:

$$\mu_c^s = \frac{1}{NHW} \sum_{n,h,w} h_{nchw}^s \tag{2}$$

$$\sigma_c^s = \sqrt{\frac{1}{NHW} \sum_{n,h,w} \left(h_{nchw}^s - \mu_c^s\right)^2},\qquad(3)$$

where H and W are the spatial dimensions of h_c^s .

The superscript s stood for the relation of features, statistical elements, and learned parameters with the source domain \mathcal{D}^s .

2.1.3 Batch Normalization under distribution shift

During the inference time, features $x_i^t \in \mathcal{X}^t$ from a target domain \mathcal{D}^t with distribution shift to the source domain \mathcal{D}^s , are normalized using the training time domain's statistics and relevant fitted affine parameters as follow:

$$BN(h_c^t) = \gamma^s \times \frac{h_c^t - \mu_c^s}{\sqrt{\sigma_c^{2^s} + \epsilon}} + \beta^s$$
(4)

This results in an inconsistency between the input futures h_c^t statistical elements and source domain captured ones inside the BN layers referred to as Internal State Mismatch

[17]. We attempt to transfer the target features using the source and target domain's statistical elements to alleviate this internal state mismatch.

2.2. Proposed method

Assume there is \mathcal{D}^s and \mathcal{D}^t as cross-domains. While most of the Domain Adaptation approaches aim at adapting the $P(Y^s|X^s)$ toward the target domain \mathcal{D}^t , we propose a contradictory viewpoint by transferring the h_c^t 's statistical mean and variance to the source domain \mathcal{D}^s during the finetuning. From this, we attempt to decrease the internal state mismatch of the BNs, rooted in the shift between the $P(X)^s$ and $P(X)^t$ (Figure 1).

2.2.1 Visual Domain Bridge (VDB)

Having a target domain \mathcal{D}^t with distribution shift according to the training time source domain \mathcal{D}^s , before normalizing the intermediate feature h_c^t we transfer it using the source and target domain statistical elements as follows:

$$h_c^{trans} = \frac{h_c^t - \hat{\mu_c}^t}{\sqrt{\hat{\sigma_c^2}^t + \epsilon}} \times \sigma_c^{2^s} + \mu_c^s$$
(5)

in order to reduce the imposed mismatch of new features' statistics (as reported in [17]), we employ a weighted average of the target and source statistic. The $\hat{\mu_c}^t$ and $\hat{\sigma_c^2}^t$ are calculated as follow:

$$\hat{\mu_c}^t = \alpha \mu_c^t + (1 - \alpha) \mu_c^s \tag{6}$$

$$\hat{\sigma_c}^t = \alpha \sigma_c^t + (1 - \alpha) \sigma_c^s \tag{7}$$

where α equals to the momentum factor as the BN's in training time. $\hat{\mu_c}^t$ and $\hat{\sigma_c}^t$ are the same as $\hat{\mu_c}^s$ and $\hat{\sigma_c}^s$, but on target feature h_{nchw}^s .

The batch normalization to the transferred intermediate feature h_c^{trans} is applied using the new emerging transferred domain \mathcal{D}^{trans} 's statistical elements:

$$BN_{transfer}(h_c^t) = \gamma^s \times \frac{h_c^{trans} - \hat{\mu}_c^{trans}}{\sqrt{\hat{\sigma}_c^{2trans} + \epsilon}} + \beta^s \qquad (8)$$

similarly, the $\hat{\mu_c}^{trans}$ and $\hat{\sigma_c^2}^{trans}$ are weighted average of statistic elements of the training-time source domain and the transferred features h_{nchw}^{trans} :

$$\hat{\mu_c}^{trans} = \alpha \mu_c^{trans} + (1 - \alpha) \, \mu_c^s \tag{9}$$

$$\hat{\sigma_c}^{trans} = \alpha \sigma_c^{trans} + (1 - \alpha) \sigma_c^s \tag{10}$$

Where

$$\mu_c^{trans} = \frac{1}{NHW} \sum_{n,h,w} h_{nchw}^{trans} \tag{11}$$

$$\sigma_c^{trans} = \sqrt{\frac{1}{NHW} \sum_{n,h,w} \left(h_{nchw}^{trans} - \mu_c^{trans} \right)^2} \quad (12)$$



Figure 3. t-SNE plot of output logits of a ResNet18 pre-traine on ImageNet depicted as Kernel Density Estimations. Blue: with BN fed with the Imagenet, Orange: with BN fed with the ISIC, Green: with VDB fed with the ISIC. Hperparameters: learning-rate=200, iteration=1000, and perplexity=30. The + marks represent the density centers.

3. Related work

Despite the fame of meta-learning domain adaptation methods, Guo et al. [4] revealed that the pure fine-tuning exceeds current SOTA meta-learning performance when fronting a significant distribution shift. Problems with a significant domain gap between the source and target data are the main focus of recent research in FSL [4]. Also, progress in semi-supervised and self-supervised learnings methods resulted in advances for the CDFSL problems; They always need significant computational resources and multiple steps of adjustment. STARTUP [18] is a well-known technique in distant domain problems that utilizes a mixture of self-supervised and self-training elements for CDFSL. In contrast to the STARTUP, the proposed method is simpler and more lightweight by a large margin. Furthermore, a branch of domain adaptation methods is specially focused on employing batch normalization to bridge between domains [13]. AutoDIAL [3], and TransNorm [26] offered techniques to train the feature extractor parameters using a mixture of source and known target domains. However, the target domain is not always known during the training time. These methods are constrained to re-train the feature extractor, which is much heavier than just statistically transferring the input features. A recent work [6] offered a twostep proposal for transferring the input feature vector into a more Gaussian-shaped distribution followed by a transitive approach to adapting the feature vector to the test-time data. Despite its effectiveness, its transductive nature limits its application to new targets without re-adaptation, which is a limited approach compared to the proposed method. As a generalized source-free domain adaptation technique, [28] proposes a method to adapt to new domains while keeping the performance on the source domain. Source-free refers to methods independent of the training-time data during the adaptation phase. Focusing on the source-free methods as more practical domain adaptation approaches, the AdaBN [13] proposes adapting the BatchNorm's moments at test time based on the target's domain statistics, demanding access to the whole target samples. Although the AdaBN is not explicitly designed for the few-shot setting, it is known as an influential paper in the field of BatchNorm-based domain adaptation. Despite its effectiveness and simplicity, its performance decay when facing a distribution shift in test-time data. The AdaIN method [7] exploited the InstanceNorm [21] layers statistics by mean of style transfer. Although the style transfer is rarely related to the fewshot adaptation, the core idea is similar to the AdaBN and is a source of inspiration in this work. The Prediction-BN [17] proposes a similar method to AdaBN by replacing the batch-wise target statistics instead of the whole domain's statistics. The main motivation for this refinement is that the AdaBN imposes a large discrepancy between the Batch-Norm's statistics and weight layers, which is alleviated by batch-wise statistics updates. Evaluations reveal that the Prediction-BN performance decrease through time and finally converges to the AdaBN's level. In order to adapt the backbone's wights to the imposed discrepancy from the BN's statistics replacement and Inspired by Prediction-BN, [24] added the BN's affine parameters adapting by entropy minimization named as TENT, which requires the backbone updating during the inference time. However, entropy minimization is not a good measure to keep the feature extractor's discrimination power. As TENT does not preserve the discriminative ability of the feature extractor, and the Prediction-BN causes a mismatch with the network state, [30] proposes a test-time calibration of the BN's statistics followed by an affine parameter adaption named as Core. Despite the effectiveness of the TENT and Core methods, both require backpropagation across the feature extractor during the inference time, which is not always applicable. This study investigated an alternative BN's statistic transfer in the CDFSL setting, exposing a slight internal state mismatch to the backbone but independent of updating the feature extractor. Recently the FeaturNorm [29] was proposed as an alternative for the BatchNorm layers during the training time, resulting in a more generalizable representation under distribution-shift. As this method neglects the learnable affine parameters of the BN layers, it needs the



Figure 4. Output distribution of selected Normalization layers from a ResNet18, pre-trained on ImageNet (Best viewed in color). Red: ImageNet on BN layers, Blue: EuroSAT on BN Layers, and Green: EuroSat on VDB layers. Left: layer 0, Center: layer 7, and Right: layer 18 which are selected based on a uniform probability. When facing distribution shift, the model with VDB results in more picked and centered distributions in per-layer aggregated feature values. Also it as apparent that VDB maintains the distribution shape when fed with training time domain. We witnessed the same pattern across all layers.

backbone to be trained from scratch and is not comparable with our source-free and simple method. However, we have evaluated the combination of the FeatureNorm with our proposed method to demonstrate the ability of our method to combine with other straightforward approaches. Compared to the AdaBN, Prediction-BN, TENT, and Core, our method transfers the statistical elements of the input feature vector. It does not exert any modification or update to the backbone, making it a suitable lifelong and source-free domain adaptation method.

4. Evaluation

4.1. Why VDB work?

Visualization is a convenient and interpretable approach to understanding method effects. Here, the VDB is visualized to investigate the effect of the proposed method on the feature representation of an encoder component of a convolutional autoencoder architecture. Then we present the representation and normalization layers' output distributions of a pre-trained model when facing the training-time source domain and an arbitrary target domain both for BN and VDB as normalization layers.

Experiment Setup: We utilized a SegNet [1] based convolutional autoencoder from [14] and replaced the encoder part's BN layers with the VDB, just before loading the pre-trained state dictionary. As the state dictionary is pre-trained on the ImageNet like datasets, a visual transfer of any input to the real-world day-time color schema

is expected. We feed batches of samples from challenging datasets (EuroSAT [5], Under Water [12], Diabetic Retinopathy [27], HRF [2], ExDark [15], PCam [22], FairFace [9], and the ImageNet as blank domain) first with BN for all of the normalization layers and then with VDB replaced in encoder part (momentum set to 1). The result is illustrated in Figure 2. An empirical evaluation of the BN and VDB output's distribution is conducted using a publicly available pre-trained ResNet18 fed with batches of samples from the ImageNet as the source domain and the EuroSAT as the target domain. The output values of normalization layers are channel-wise aggregated as presented in Figure 4. In Figure 3, the representation gap induced from the distribution gap is considered. We plot the calculated t-SNE¹ from logits of a pre-trained ResNet18, equipped with the BN layers fed with samples from the source domain and the ISIC. Then, the Kernel Density Estimation of the t-SNE plot is compared with the same model where the BN layers are replaced with the VDB and fed with the ISIC domain.

Results: From Figure 2, it is evident that the autoencoder with VDB visually transfers any arbitrary target domains to the source ImageNet like domain. From the right side of the figure, the sample from ImageNet remains unchanged, whereas the samples from other domains visually get transferred toward the source domain, resulting in a lower distribution gap in CNN's layers inputs. The effect of VDB is demonstrated as more picked distributions at the normaliza-

¹t-distributed stochastic neighbor embedding

	5 Ways # Shots	CropDisease	EuroSAT	ISIC	ChestX
pure tuning AdaBN Prediction-BN VDB (our)	1	$\begin{array}{c} 68.46 \pm 0.87 \\ 68.19 \pm 0.85 \\ 70.57 \pm 0.84 \\ 71.98 \pm 0.82 \end{array}$	$\begin{array}{c} 59.18 \pm 0.85 \\ 57.46 \pm 0.80 \\ 61.60 \pm 0.83 \\ 63.60 \pm 0.87 \end{array}$	$\begin{array}{c} 33.11 \pm 0.60 \\ 34.72 \pm 0.65 \\ 35.52 \pm 0.65 \\ 35.32 \pm 0.65 \end{array}$	$\begin{array}{c} 22.54 \pm 0.42 \\ 22.32 \pm 0.41 \\ 22.60 \pm 0.42 \\ 22.99 \pm 0.44 \end{array}$
STARTUP STARTUP + VDB		$\begin{array}{c} 74.56 \pm 0.85 \\ 77.43 \pm 0.81 \end{array}$	$\begin{array}{c} 64.00 \pm 0.88 \\ 63.61 \pm 0.87 \end{array}$	$\begin{array}{c} 35.12 \pm 0.64 \\ 35.54 \pm 0.64 \end{array}$	$\begin{array}{c} 22.93 \pm 0.43 \\ 23.20 \pm 0.43 \end{array}$
FN FN + VDB		$\begin{array}{c} 70.93 \pm 0.87 \\ 74.75 \pm 0.79 \end{array}$	$\begin{array}{c} 62.75 \pm 0.88 \\ 65.60 \pm 0.89 \end{array}$	$\begin{array}{c} 32.77 \pm 0.60 \\ 34.94 \pm 0.63 \end{array}$	$\begin{array}{c} 22.37 \pm 0.41 \\ 22.54 \pm 0.43 \end{array}$
MAML* ProtoNet* pure tuning AdaBN Prediction-BN VDB (our) STARTUP STARTUP + VDB	5	$\begin{array}{c} 78.05 \pm 0.68 \\ 79.72 \pm 0.67 \\ 89.86 \pm 0.50 \\ 90.12 \pm 0.50 \\ 89.80 \pm 0.51 \\ 90.77 \pm 0.49 \\ 92.86 \pm 0.43 \\ 93.13 \pm 0.45 \end{array}$	$\begin{array}{c} 71.70 \pm 0.72 \\ 73.29 \pm 0.71 \\ 79.99 \pm 0.64 \\ 80.21 \pm 0.60 \\ 81.71 \pm 0.57 \\ 82.06 \pm 0.63 \\ 82.51 \pm 0.62 \\ 82.11 \pm 0.64 \end{array}$	$\begin{array}{c} 40.13 \pm 0.58 \\ 39.57 \pm 0.57 \\ 45.53 \pm 0.59 \\ 48.88 \pm 0.62 \\ 50.42 \pm 0.65 \\ 48.72 \pm 0.65 \\ 48.54 \pm 0.63 \\ 49.72 \pm 0.65 \end{array}$	$\begin{array}{c} 23.48 \pm 0.96 \\ 24.05 \pm 1.01 \\ 26.66 \pm 0.42 \\ 25.66 \pm 0.42 \\ 25.61 \pm 0.43 \\ 26.62 \pm 0.45 \\ 27.17 \pm 0.44 \\ 27.48 \pm 0.45 \end{array}$
FN FN + VDB		$\begin{array}{c} 91.32 \pm 0.46 \\ 92.20 \pm 0.45 \end{array}$	$\begin{array}{c} 80.75 \pm 0.63 \\ 82.68 \pm 0.61 \end{array}$	$\begin{array}{c} 44.80 \pm 0.57 \\ 47.79 \pm 0.60 \end{array}$	$\begin{array}{c} 26.32 \pm 0.43 \\ 26.56 \pm 0.46 \end{array}$
MAML* ProtoNet* pure tuning AdaBN Prediction-BN VDB (our) STARTUP STARTUP + VDB FN	20	$\begin{array}{c} 89.75 \pm 0.42 \\ 88.15 \pm 0.51 \\ 96.01 \pm 0.28 \\ 96.26 \pm 0.28 \\ 95.92 \pm 0.28 \\ 96.36 \pm 0.27 \\ 97.43 \pm 0.23 \\ 97.42 \pm 0.22 \\ 96.68 \pm 0.25 \end{array}$	$\begin{array}{c} 81.95 \pm 0.55 \\ 82.27 \pm 0.57 \\ 88.02 \pm 0.47 \\ 88.94 \pm 0.44 \\ 89.82 \pm 0.41 \\ 89.42 \pm 0.45 \\ 89.63 \pm 0.43 \\ 89.65 \pm 0.45 \\ 87.88 \pm 0.45 \end{array}$	$\begin{array}{c} 52.36 \pm 0.57 \\ 49.50 \pm 0.55 \\ 55.64 \pm 0.57 \\ 58.98 \pm 0.58 \\ 60.43 \pm 0.61 \\ 59.09 \pm 0.59 \\ 59.98 \pm 0.59 \\ 60.37 \pm 0.57 \\ 56.43 \pm 0.56 \end{array}$	$\begin{array}{c} 27.53 \pm 0.43 \\ 28.21 \pm 1.15 \\ 31.98 \pm 0.44 \\ 31.11 \pm 0.46 \\ 30.91 \pm 0.45 \\ 31.87 \pm 0.44 \\ 33.54 \pm 0.46 \\ 33.49 \pm 0.49 \\ 32.36 \pm 0.46 \end{array}$
FN + VDB ProtoNet* pure tuning AdaBN Prediction-BN VDB (our) STARTUP STARTUP + VDB FN	50	$\begin{array}{c} 97.09 \pm 0.25\\ \hline 90.81 \pm 0.43\\ 97.59 \pm 0.22\\ 97.90 \pm 0.19\\ 97.61 \pm 0.20\\ 97.89 \pm 0.19\\ 98.53 \pm 0.16\\ 98.48 \pm 0.17\\ 98.09 \pm 0.19\\ \end{array}$	$\begin{array}{c} 89.70 \pm 0.43 \\ \hline 80.48 \pm 0.57 \\ 91.04 \pm 0.37 \\ 92.08 \pm 0.36 \\ 92.62 \pm 0.33 \\ 92.24 \pm 0.35 \\ 92.59 \pm 0.33 \\ 92.46 \pm 0.35 \\ 91.01 \pm 0.38 \end{array}$	59.51 ± 0.57 51.99 ± 0.52 61.38 ± 0.57 63.61 ± 0.59 64.34 ± 0.55 64.02 ± 0.58 65.90 ± 0.56 65.21 ± 0.55 62.64 ± 0.57	$\begin{array}{c} 32.31 \pm 0.46 \\ \hline 29.32 \pm 1.12 \\ 35.28 \pm 0.48 \\ 34.21 \pm 0.45 \\ 34.14 \pm 0.46 \\ 35.55 \pm 0.45 \\ 37.67 \pm 0.47 \\ 38.17 \pm 0.46 \\ 36.32 \pm 0.47 \end{array}$
FN + VDB pure tuning AdaBN Prediction-BN VDB (our)	avg	$\begin{array}{c} 98.35 \pm 0.17 \\ 87.98 \pm 0.39 \\ 88.12 \pm 0.38 \\ 88.47 \pm 0.38 \\ \textbf{89.25} \pm \textbf{0.38} \end{array}$	$\begin{array}{c} 92.29 \pm 0.34 \\ \hline 79.56 \pm 0.43 \\ 79.67 \pm 0.42 \\ 81.44 \pm 0.41 \\ \textbf{81.83} \pm \textbf{0.43} \end{array}$	65.81 ± 0.56 48.91 ± 0.43 51.55 ± 0.44 52.68 ± 0.44 51.79 ± 0.44	$\begin{array}{c} 36.90 \pm 0.48 \\ \hline 29.11 \pm 0.37 \\ 28.33 \pm 0.37 \\ 28.32 \pm 0.38 \\ \textbf{29.26} \pm \textbf{0.38} \end{array}$
STARTUP STARTUP + VDB FN FN + VDB		$\begin{array}{l} 90.85 \pm 0.37 \\ \textbf{91.62} \pm \textbf{0.36} \\ 89.26 \pm 0.38 \\ \textbf{90.60} \pm \textbf{0.36} \end{array}$	$\begin{array}{l} \textbf{82.18} \pm \textbf{0.43} \\ \textbf{81.96} \pm \textbf{0.43} \\ \textbf{80.60} \pm \textbf{0.43} \\ \textbf{82.57} \pm \textbf{0.43} \end{array}$	$52.39 \pm 0.44 \\ 52.71 \pm 0.44 \\ 49.16 \pm 0.43 \\ 52.01 \pm 0.44$	$\begin{array}{l} 30.33 \pm 0.38 \\ \textbf{30.58} \pm \textbf{0.38} \\ 29.35 \pm 0.38 \\ \textbf{29.58} \pm \textbf{0.38} \end{array}$

Table 1. The CDFSL classification accuracy results of a ResNet10 backbone pre-trained on miniImageNet for 5 ways and # number of shots. FN refers to the FeatureNorm method. avg: Average over all shots. * refers to results from [4].

	5 Ways # Shots	CropDisease	EuroSAT	ISIC	ChestX
pure tuning AdaBN Prediction-BN VDB (our)	1	$\begin{array}{c} 71.71 \pm 0.87 \\ 67.00 \pm 0.84 \\ 69.84 \pm 0.77 \\ 75.46 \pm 0.76 \end{array}$	$\begin{array}{c} 64.58 \pm 0.82 \\ 55.58 \pm 0.86 \\ 59.72 \pm 0.86 \\ 67.76 \pm 0.83 \end{array}$	$\begin{array}{c} 30.30 \pm 0.52 \\ 29.43 \pm 0.51 \\ 31.65 \pm 0.54 \\ 33.22 \pm 0.58 \end{array}$	$\begin{array}{c} 22.16 \pm 0.40 \\ 21.75 \pm 0.38 \\ 21.55 \pm 0.39 \\ 22.28 \pm 0.41 \end{array}$
FN FN + VDB		$\begin{array}{c} 76.98 \pm 0.83 \\ 79.68 \pm 0.74 \end{array}$	$\begin{array}{c} 66.45 \pm 0.78 \\ 69.67 \pm 0.80 \end{array}$	$\begin{array}{c} 30.09 \pm 0.57 \\ 32.96 \pm 0.57 \end{array}$	$\begin{array}{c} 23.14 \pm 0.42 \\ 22.64 \pm 0.41 \end{array}$
pure tuning AdaBN Prediction-BN VDB (our) FN FN + VDB	5	$\begin{array}{c} 91.20 \pm 0.46 \\ 91.09 \pm 0.47 \\ 91.30 \pm 0.43 \\ 93.11 \pm 0.42 \\ 93.54 \pm 0.42 \\ 94.63 \pm 0.37 \end{array}$	$\begin{array}{c} 84.88 \pm 0.57 \\ 77.73 \pm 0.67 \\ 80.10 \pm 0.63 \\ 85.29 \pm 0.52 \\ 86.10 \pm 0.51 \\ 87.31 \pm 0.50 \end{array}$	$\begin{array}{c} 44.01\pm0.57\\ 43.26\pm0.55\\ 44.11\pm0.52\\ 47.48\pm0.61\\ 44.48\pm0.61\\ 47.48\pm0.59\end{array}$	$\begin{array}{c} 25.23 \pm 0.43 \\ 23.94 \pm 0.40 \\ 24.00 \pm 0.41 \\ 25.25 \pm 0.42 \\ 26.15 \pm 0.41 \\ 25.55 \pm 0.43 \end{array}$
pure tuning AdaBN Prediction-BN VDB (our)	20	$96.52 \pm 0.27 96.84 \pm 0.25 96.91 \pm 0.24 97.61 \pm 0.21 07.62 \pm 0.22 07.61 \pm 0.22 07.63 \pm 0.22 07.64 \pm 0.22 07.65 \pm 0.22 \\0.25 \pm 0.27 \\0.25 \pm $	$91.51 \pm 0.37 \\ 86.69 \pm 0.53 \\ 87.61 \pm 0.48 \\ 91.93 \pm 0.37 \\ 02.28 \pm 0.22 \\ 0.21 \\ 0.21 \\ 0.21 \\ 0.22 \\ 0.21 \\ 0.21 \\ 0.22 \\ 0.21 \\ 0.22 \\ 0.$	$55.66 \pm 0.57 55.33 \pm 0.53 55.35 \pm 0.52 58.89 \pm 0.59 55.61 \pm 0.50 $	$29.19 \pm 0.41 27.95 \pm 0.41 27.43 \pm 0.42 29.49 \pm 0.42 20.14 \pm 0.44 $
FIN FN + VDB		97.62 ± 0.23 98.20 ± 0.19	92.28 ± 0.33 92.85 ± 0.34	56.61 ± 0.39 59.65 ± 0.58	30.14 ± 0.44 30.28 ± 0.45
pure tuning AdaBN Prediction-BN VDB (our) FN	50	$98.05 \pm 0.20 97.97 \pm 0.19 97.89 \pm 0.22 98.40 \pm 0.16 98.56 \pm 0.17 08.82 \pm 0.12 $	$93.34 \pm 0.28 \\ 89.63 \pm 0.43 \\ 90.55 \pm 0.39 \\ 93.95 \pm 0.30 \\ 94.30 \pm 0.27 \\ 04.70 \pm 0.27 \\ 04.70 \pm 0.27 \\ 0.27 \\ 0.100 \\ 0.10$	$61.42 \pm 0.56 59.83 \pm 0.55 60.25 \pm 0.53 64.23 \pm 0.58 63.72 \pm 0.57 65.27 \pm 0.56 $	$\begin{array}{c} 31.80 \pm 0.47 \\ 29.75 \pm 0.45 \\ 29.12 \pm 0.45 \\ 32.37 \pm 0.47 \\ 33.60 \pm 0.46 \\ 22.24 \pm 0.45 \end{array}$
pure tuning AdaBN Prediction-BN VDB (our) FN	Average	$\begin{array}{c} 89.37 \pm 0.13 \\ 89.37 \pm 0.38 \\ 88.23 \pm 0.37 \\ 88.99 \pm 0.36 \\ \textbf{91.15} \pm \textbf{0.35} \\ 91.68 \pm 0.36 \end{array}$	83.58 ± 0.40 77.40 ± 0.45 79.49 ± 0.43 84.73 ± 0.40 84.78 ± 0.39	$\begin{array}{c} 55.27 \pm 0.30 \\ 47.85 \pm 0.42 \\ 46.96 \pm 0.42 \\ 47.84 \pm 0.41 \\ \textbf{50.96} \pm \textbf{0.43} \\ 48.73 \pm 0.43 \end{array}$	$\begin{array}{c} 53.24 \pm 0.43 \\ \hline 27.10 \pm 0.37 \\ 25.85 \pm 0.36 \\ \hline 25.53 \pm 0.36 \\ \hline 27.35 \pm 0.37 \\ \hline 28.26 \pm 0.37 \end{array}$
FN + VDB		$\textbf{92.83} \pm \textbf{0.34}$	$\textbf{86.15} \pm \textbf{0.39}$	$\textbf{51.34} \pm \textbf{0.43}$	27.93 ± 0.37

Table 2. The CDFSL results for a ResNet18 backbone pre-trained on ImageNet, evaluated over 600 episodes. pure tuning, linear classifier fine-tuned on labeled target samples. AdaBN is fine-tuned the same as the pure tuning except adapt the BN layers before fine-tuning, Prediction-BN and VDB similar to the AdaBN except for the BN adaptation occurs during the fine-tuning. Results are for 5-Way and # Shots classification accuracy. avg: Average over all shots.

	1-shot	5-shot	20-shot
pure tuning	54.56 ± 0.84	76.18 ± 0.69	84.53 ± 0.52
VDB (our)	$\textbf{56.08} \pm \textbf{0.86}$	$\textbf{76.61} \pm \textbf{0.72}$	84.59 ± 0.55

Table 3. Near-domain few-shot evaluation on pure tuning and VDB. Models are pre-trained on miniImageNet and evaluated on novel classes of ImageNet over 600 episodes.

tion layers outputs as illustrated in Figure 4. From this figure, we can see VDB maintains the distribution shape when fed with training time domain while the model with VDB results in more picked and centered distributions in perlayer aggregated feature value when facing distribution shift (we witnessed the same pattern for all the layers). From Figure 3 it is observant that the model with VDB reduces the representation gap on data with distribution-shift. Compared to the model with BN (orange), the VDB equipped green shape is expanded toward the source's shape (blue), and the center of the density is moved to the source domain's center of density. As previously explored in [19] distribution shift will appear as distance in representation, and reducing the representation distance will directly result in a performance upgrade. We repeated the plot with different hyperparameters resulted in similar patterns.

4.2. VDB in CDFSL

Experiment Setup: The FSL evaluation on the CDFSL benchmark is established as introduced in [4]. The mini-ImageNet [23] is employed as the base training-time source dataset as well as the more extensive ImageNet dataset. The target datasets are comprised of four datasets, each from a shifted distribution relative to the source domain (miniImageNet and ImageNet). The target datasets are composed of EuroSAT (satellite imagery to determine land usage), CropDiseases (plant images to identify botanical diseases), ChestX (chest X-rays to detect pathology), and ISIC2018 (images of skin abrasions to detect melanoma) as described in [4]. Methods such as STARTUP and AdaBN benefit randomly sampled 20% of unlabelled images from the target datasets to use in the representation learning stage, consistent with the setup of [18], remaining samples are used during the fine-tuning. Identical to [4], we conduct experiments in an FSL classification setting while the support set is comprised of 5 classes with k samples per class (5-way k-shot), where $k \in \{1, 5, 20, 50\}$ and the overall CDFSL score is an average of accuracies across all target datasets for $k \in \{5, 20, 50\}$. Evaluation of pre-trained models brings off over 600 episodes. 95% confidence intervals with mean accuracy are reported. Pure tuning stands for our baseline, is the method without any adaptation rather than the fine-tuning, which utilizes a pre-trained model regularly trained on the source domain. The CDFSL evaluations for ResNet10 pre-trained on the miniImageNet and ResNet18 pre-trained on ImageNet, are reported in Tables 1 and 2, respectively.

Results: When facing Cross-Domain targets, a ResNet10 architecture pre-trained on miniImageNet is evaluated based on the CDFSL benchmark's setting as presented in Table 1. From the reported values, we can see the VDB outperformed the last SOTA method, the Prediction-BN with a subtle margin. Also the VDB is able to improve the STARTUP method as a fancy domain adaptation method with heavy training process utilizing target samples and benefits from semi-supervised, and self-supervised losses beside the supervised one. Surprisingly, it is observant that the VDB rather than enhancing the FeaturNorm method, FeaturNorm + VDB perform near the STARTUP, which is a much heavier method. The CDFSL scores for Table 1

methods is: pure tuning (66.58 \pm 0.67), AdaBN (67.33 \pm 0.66), Prediction-BN (67.78 \pm 0.66), VDB (67.88 \pm 0.82) , STARTUP (68.86 \pm 0.65), STARTUP + VDB (68.97 \pm **0.81**), FN (67.05 \pm 0.66), and FN + VDB (68.43 \pm 0.81). When it comes to the common and more practical source domain in real-world scenarios, the ImageNet, we utilized a publicly available pre-trained ResNet18 and evaluated it in the CDFSL settings and on the target datasets. As it is evident from Table 2, Although the AdaBN and Prediction-BN, both underachieved the pure tuning, The VDB improved the pure tuning's performance by more than one percent in the CDFSL score. A gain roughly in free, capable of integrating with other domain adaptation methods. The CDFSL scores for Table 2 methods is: pure tuning (66.90 \pm 0.64), AdaBN (65.00 \pm 0.66), Prediction-BN (65.39 \pm 0.65), **VDB** (68.17 \pm 0.78), FN (68.09 \pm 0.63), and FN + VDB (69.01 \pm 0.77).

4.3. VDB in near-domain

Catastrophic forgetting [16] as a caveat of domain adaptation techniques in a way that a comprehension model tends to perform poorly in the source domain is a wellknown phenomenon among domain adaptation methods. In order to examine the effect of VDB on forgetting the source domain, we evaluated a pre-trained ResNet10 on miniImageNet in Near-Domain setting in which $P(X)^s = P(X)^t$ and $P(Y^s|X^s) \neq P(Y^t|X^t)$, with the same setting as the CDFSL benchmark but on novel classes of ImageNet datasets as target domain.

From Table 3, although it was not expected, the VDB maintains the in-distribution performance even with a slight superiority compared to the same pre-trained model without any modifications.

5. Conclusion

This work proposes a source-free domain-adaptation method, the VDB, in the cross-domain few-shot learning setting. A modification to the BN layer to be substituted with it across a CNN in fine-tuning time. Empirical evaluations reveal the effectiveness of the VDB against the domain gap through the reduction in the representation gap. Also, it performed as a new SOTA approach among the BN-based source-free domain-adaptation methods on a sever few-shot benchmark, the CDFSL. The VDB improved the performance on miniImageNet as well as ImageNet as a more general source domain wherein previous SOTA methods underperform the pure fine-tuning.

Acknowledgments

All the experiments recorded in Tables 3, 1, and 2 are conducted using Google Colab, the free version.

References

- Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 5
- [2] Attila Budai, Rüdiger Bock, Andreas Maier, Joachim Hornegger, and Georg Michelson. Robust vessel segmentation in fundus images. *International journal of biomedical imaging*, 2013, 2013. 5
- [3] Fabio Maria Carlucci, Lorenzo Porzi, Barbara Caputo, Elisa Ricci, and Samuel Rota Bulo. Autodial: Automatic domain alignment layers. In 2017 IEEE international conference on computer vision (ICCV), pages 5077–5085. IEEE, 2017. 4
- [4] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *European Conference on Computer Vision*, pages 124–141. Springer, 2020. 1, 2, 4, 6, 8
- [5] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019. 5
- [6] Yuqing Hu, Vincent Gripon, and Stéphane Pateux. Leveraging the feature distribution in transfer-based few-shot learning. In *International Conference on Artificial Neural Networks*, pages 487–499. Springer, 2021. 4
- [7] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 2, 4
- [8] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. 1
- [9] Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1548–1558, 2021. 5
- [10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25:1097–1105, 2012.
- [11] Chuanzi Li, Jining Feng, Li Hu, Junhong Li, and Haibin Ma. Review of image classification method based on deep transfer learning. In 2020 16th International Conference on Computational Intelligence and Security (CIS), pages 104–108. IEEE, 2020. 1
- [12] Chongyi Li, Chunle Guo, Wenqi Ren, Runmin Cong, Junhui Hou, Sam Kwong, and Dacheng Tao. An underwater image enhancement benchmark dataset and beyond. *IEEE Transactions on Image Processing*, 29:4376–4389, 2020. 5
- [13] Yanghao Li, Naiyan Wang, Jianping Shi, Xiaodi Hou, and Jiaying Liu. Adaptive batch normalization for practical domain adaptation. *Pattern Recognition*, 80:109–117, 2018. 2, 4

- [14] Yang Liu. Convolutional autoencoder with setnet in pytorch. https://github.com/foamliu/Autoencoder, 2018. 5
- [15] Yuen Peng Loh and Chee Seng Chan. Getting to know lowlight images with the exclusively dark dataset. *Computer Vision and Image Understanding*, 178:30–42, 2019. 5
- [16] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. 8
- [17] Zachary Nado, Shreyas Padhy, D Sculley, Alexander D'Amour, Balaji Lakshminarayanan, and Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate shift. arXiv preprint arXiv:2006.10963, 2020. 2, 3, 4
- [18] Cheng Perng Phoo and Bharath Hariharan. Self-training for few-shot transfer across extreme task differences, 2021. 4, 8
- [19] Karin Stacke, Gabriel Eilertsen, Jonas Unger, and Claes Lundström. Measuring domain shift for deep learning in histopathology. *IEEE journal of biomedical and health informatics*, 25(2):325–336, 2020. 8
- [20] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. arXiv preprint arXiv:2007.00644, 2020. 1
- [21] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022, 2016. 4
- [22] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant CNNs for digital pathology. June 2018. 5
- [23] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. Advances in neural information processing systems, 29:3630– 3638, 2016. 8
- [24] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*, 2021. 2, 4
- [25] Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 312:135–153, 2018. 2
- [26] Ximei Wang, Ying Jin, Mingsheng Long, Jianmin Wang, and Michael Jordan. Transferable normalization: Towards improving transferability of deep neural networks. 2019. 4
- [27] Zhiguang Wang and Jianbo Yang. Diabetic retinopathy detection via deep convolutional networks for discriminative localization and visual explanation. In Workshops at the thirty-second AAAI conference on artificial intelligence, 2018. 5
- [28] Shiqi Yang, Yaxing Wang, Joost van de Weijer, Luis Herranz, and Shangling Jui. Generalized source-free domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8978–8987, 2021. 4
- [29] Moslem Yazdanpanah, Aamer Abdul Rahman, Muawiz Chaudhary, Christian Desrosiers, Mohammad Havaei, Eugene Belilovsky, and Samira Ebrahimi Kahou. Revisiting learnable affines for batch norm in few-shot transfer learning.

In Proceedings of the IEEE/CVF International Conference on Computer Vision [Manuscript submitted for publication], 2022. 4

- [30] Fuming You, Jingjing Li, and Zhou Zhao. Test-time batch statistics calibration for covariate shift. *arXiv preprint arXiv:2110.04065*, 2021. 2, 4
- [31] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruyssen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. The visual task adaptation benchmark. 2019. 1