# Deep Hashing Network for Unsupervised Domain Adaptation

Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, Sethuraman Panchanathan
Center for Cognitive Ubiquitous Computing, Arizona State University, Tempe, AZ, USA
{hemanthv, jeusebio, shayok.chakraborty, panch}@asu.edu

## Abstract

*In recent years, deep neural networks have emerged as a dominant machine learning tool for a wide variety of application domains. However, training a deep neural network requires a large amount of labeled data, which is an expensive process in terms of time, labor and human expertise. Domain adaptation or transfer learning algorithms address this challenge by leveraging labeled data in a different, but related source domain, to develop a model for the target domain. Further, the explosive growth of digital data has posed a fundamental challenge concerning its storage and retrieval. Due to its storage and retrieval efficiency, recent years have witnessed a wide application of hashing in a variety of computer vision applications. In this paper, we first introduce a new dataset, Office-Home, to evaluate domain adaptation algorithms. The dataset contains images of a variety of everyday objects from multiple domains. We then propose a novel deep learning framework that can exploit labeled source data and unlabeled target data to learn informative hash codes, to accurately classify unseen target data. To the best of our knowledge, this is the first research effort to exploit the feature learning capabilities of deep neural networks to learn representative hash codes to address the domain adaptation problem. Our extensive empirical studies on multiple transfer tasks corroborate the usefulness of the framework in learning efficient hash codes which outperform existing competitive baselines for unsupervised domain adaptation.*

## 1. Introduction

Deep learning algorithms automatically learn a discriminating set of features and have depicted commendable performance in a variety of computer vision applications. Unfortunately, training a deep model necessitates a large volume of labeled data, which can be time consuming and expensive to acquire. However, labeled data from a different, but related domain is often available, which has motivated the development of algorithms which can leverage

labeled data in a source domain to develop a machine learning model for the target domain. Learning a discriminative model in the presence of the shift between training and test distributions is known as transfer learning or domain adaptation [17]. Unsupervised domain adaptation is a challenging setting, where labeled data is available only in the source domain; no labeled data is available in the target domain. Conventional shallow transfer learning methods develop their models in two stages, feature extraction followed by domain adaptation. The features are fixed and then a model is trained to align the source and target domains [16, 20, 33, 38, 42, 43, 44]. On the other hand, deep transfer learning procedures exploit the feature learning capabilities of deep networks to learn transferable feature representations for domain adaptation and have demonstrated impressive empirical performance [17, 18, 31, 34, 46].

The explosive growth of digital data in the modern era has posed fundamental challenges regarding their storage, retrieval and computational requirements. Against this backdrop, hashing has emerged as one of the most popular and effective techniques due to its fast query speed and low memory cost [48]. Hashing techniques transform high dimensional data into compact binary codes and generate similar binary codes for similar data items. Motivated by this fact, we propose to train a deep neural network to output binary hash codes (instead of probability values), which can be used for classification. We see two advantages to estimating a hash value instead of a standard probability vector in the final layer of the network: (i) the hash values are used to develop a unique loss function for target data in the absence of labels and (ii) during prediction, the hash value of a test sample can be compared against the hash values of the training samples to arrive at a more robust category prediction.

In this paper, we first introduce a new dataset, *Office-Home*, which we use to evaluate our algorithm. The *Office-Home* dataset is an object recognition dataset which contains images from 4 domains. It has around $15,500$ images organized into $65$ categories. We further propose a novel deep learning framework called Domain Adaptive Hash-

ing (DAH) to learn informative hash codes to address the problem of unsupervised domain adaptation. We propose a unique loss function to train the deep network with the following components: (i) supervised hash loss for labeled source data, which ensures that source samples belonging to the same class have similar hash codes; (ii) unsupervised entropy loss for unlabeled target data, which imposes each target sample to align closely with exactly one of the source categories and be distinct from the other categories and (iii) a loss based on multi-kernel Maximum Mean Discrepancy (MK-MMD), which seeks to learn transferable features within the layers of the network to minimize the distribution difference between the source and target domains. Figure 1 illustrates the different layers of the DAH and the components of the loss function.

## 2. Related Work

There have been many approaches to address the problem of domain-shift in unsupervised domain adaptation. One straightforward approach is, to modify a classifier trained for the source data by adapting it to classify target data [1, 4] or learn a transformation matrix to linearly transform the source data, so that it is aligned with the target [27, 42]. Some other procedures re-weight the data points in the source domain, to select source data that is similar to the target, when training a domain adaptive classifier, [9, 10, 19]. A standard procedure to reduce domain discrepancy is, to project the source and target data to a common subspace, thereby aligning their principal axes [16, 44]. Reducing domain disparity through nonlinear alignment of data has been possible with Maximum Mean Discrepancy (MMD) - a measure that provides the distribution difference between two datasets in a reproducing-kernel Hilbert space [13]. Kernel-PCA based methods apply the MMD to achieve nonlinear alignment of domains [32, 33, 38]. Manifold based approaches are also popular in domain adaptation for computer vision, where the subspace of a domain is treated as a point on the manifold and transformations are learned to align two domains [20, 23]. A survey of popular domain adaptation techniques for computer vision is provided in [41] and a more generic survey of transfer learning approaches can be found in [39].

All of the above techniques can be termed as shallow learning procedures, since the models are learned using predetermined features. In recent years deep learning has become very successful at learning highly discriminative features for computer vision applications [8]. Deep learning systems like deep CNNs learn representations of data that capture underlying factors of variation between different tasks in a multi-task transfer learning setting [3]. These representations also disentangle the factors of variation allowing for the transfer of knowledge between tasks [12, 18, 37]. Yosinski et al. [49] demonstrated how the lower layers of a network produce generic features and the upper layers output task specific features. Based on this, deep learning procedures for domain adaptation train networks to learn transferable features in the fully connected final layers of a network [31, 46]. In other approaches to deep domain adaptation, Ganin et al. [17] trained domain adversarial networks to learn features that make the source and target domain indistinguishable and Long et al. [34], trained a network to do both feature adaptation and classifier adaptation using residual transfer networks.

Unsupervised hashing techniques have been developed to extract unique hash codes for efficient storage and retrieval of data [22, 25]. Neural network based hashing has led the way in state-of-the-art unsupervised hashing techniques [7, 11, 14]. The closest work incorporating hashing and adaptation appears in cross-modal hashing, where deep hashing techniques embed multi-modal data and learn hash codes for two related domains, like text and images [5, 6, 29]. However, these algorithms are not unsupervised and they are mainly applied to extract common hash codes for multi-modal data for retrieval purposes. To the best of our knowledge, there has been no work in unsupervised domain adaptation using deep hashing networks. We now present the Domain Adaptive Hashing (DAH) network for unsupervised domain adaptation through deep hashing.

## 3. Domain Adaptive Hashing Networks

In unsupervised domain adaptation, we consider data from two domains; *source* and *target*. The source consists of labeled data, $\mathcal{D}_s = \{\boldsymbol{x}_i^s, y_i^s\}_{i=1}^{n_s}$ and the target has only unlabeled data $\mathcal{D}_t = \{\boldsymbol{x}_i^t\}_{i=1}^{n_t}$. The data points $\boldsymbol{x}_i^*$ belong to $X$, where $X$ is some input space. The corresponding labels are represented by $y_i^* \in Y \coloneqq \{1, \ldots, C\}$. The paradigm of domain adaptive learning attempts to address the problem of *domain-shift* in the data, where the data distributions of the source and target are different, i.e. $P_s(X, Y) \neq P_t(X, Y)$. The domain-shift notwithstanding, our goal is to train a deep neural network classifier $\psi(.)$, that can predict the labels $\{\hat{y}_i^t\}_{i=1}^{n_t}$, for the target data.

We implement the neural network as a deep CNN which consists of 5 convolution layers *conv*1 - *conv*5 and 3 fully connected layers *fc*6 - *fc*8 followed by a loss layer. In our model, we introduce a hashing layer *hash-fc*8 in place of the standard *fc*8 layer to learn a binary code $\boldsymbol{h}_i$, for every data point $\boldsymbol{x}_i$, where $\boldsymbol{h}_i \in \{-1, +1\}^d$. The *hash-fc*8 layer is driven by two loss functions, (i) *supervised hash loss* for the source data, (ii) *unsupervised entropy loss* for the target data. The supervised hash loss ensures hash values that are distinct and discriminatory, i.e. if $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ belong to the same category, their hash values $\boldsymbol{h}_i$ and $\boldsymbol{h}_j$ are similar and different otherwise. The unsupervised entropy loss aligns the target hash values with source hash values based on the similarity of their feature representations. The output of the
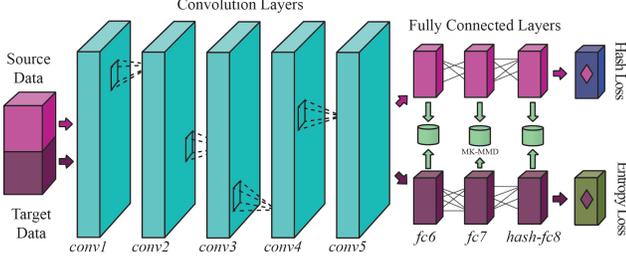
**Figure 1:** The Domain Adaptive Hash (DAH) network that outputs hash codes for the source and the target. The network is trained with a batch of source and target data. The convolution layers *conv*1 - *conv*5 and the fully connected layers *fc*6 and *fc*7 are fine tuned from the VGG-F network. The MK-MMD loss trains the DAH to learn feature representations which align the source and the target. The *hash-fc*8 layer is trained to output vectors of $d$ dimensions. The supervised hash loss drives the DAH to estimate a unique hash value for each object category. The unsupervised entropy loss aligns the target hash values to their corresponding source categories. Best viewed in color.

network is represented as $\psi(\boldsymbol{x})$, where $\psi(\boldsymbol{x}) \in \mathbb{R}^d$, which we convert to a hash code $\boldsymbol{h} = \text{sgn}(\psi(\boldsymbol{x}))$, where sgn(.) is the sign function. Once the network has been trained, the probability of $\boldsymbol{x}$ being assigned a label $y$ is given by $f(\boldsymbol{x}) = p(y|\boldsymbol{h})$. We train the network using $\mathcal{D}_s$ and $\mathcal{D}_t$ and predict the target data labels $\hat{y}_*^t$ using $f(.)$.

In order to address the issue of domain-shift, we need to align the feature representations of the target and the source. We do that by reducing the domain discrepancy between the source and target feature representations at multiple layers of the network. In the following subsections, we discuss the design of the domain adaptive hash (DAH) network in detail.

## 3.1. Reducing Domain Disparity

Deep learning methods have been very successful in domain adaptation with state-of-the-art algorithms [17, 31, 34, 46] in recent years. The feature representations transition from generic to task-specific as one goes up the layers of a deep CNN [49]. The convolution layers *conv*1 to *conv*5 have been shown to be generic and so, readily transferable, whereas the fully connected layers are more task-specific and need to be adapted before they can be transferred. In the DAH algorithm, we attempt to minimize the MK-MMD loss to reduce the domain difference between the source and target feature representations for fully connected layers, $\mathcal{F} = \{fc6, fc7, fc8\}$. Such a loss function has been used in previous research [31, 34]. The multi-layer MK-MMD loss is given by,

$$\mathcal{M}(u_s, u_t) = \sum_{l \in \mathcal{F}} d_k^2(u_s^l, u_t^l), \qquad (1)$$

where, $u_s^l = \{\boldsymbol{u}_i^{s,l}\}_{i=1}^{n_s}$ and $u_t^l = \{\boldsymbol{u}_i^{t,l}\}_{i=1}^{n_t}$ are the set of output representations for the source and target data at

layer $l$, where $\boldsymbol{u}_i^{*,l}$ is the output representation of $\boldsymbol{x}_i^*$ for the $l^{th}$ layer. The final layer outputs are denoted as $u_s$ and $u_t$. The MK-MMD measure $d_k^2(.)$ is the multi-kernel maximum mean discrepancy between the source and target representations, [24]. For a nonlinear mapping $\phi(.)$ associated with a reproducing kernel Hilbert space $\mathscr{H}_k$ and kernel $k(.)$, where $k(\boldsymbol{x}, \boldsymbol{y}) = \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{y}) \rangle$, the MMD is defined as,

$$d_k^2(u_s^l, u_t^l) = \left\| \mathbb{E}[\phi(\boldsymbol{u}^{s,l})] - \mathbb{E}[\phi(\boldsymbol{u}^{t,l})] \right\|_{\mathscr{H}_k}^2. \qquad (2)$$

The characteristic kernel $k(.)$, is determined as a convex combination of $\kappa$ PSD kernels, $\{k_m\}_{m=1}^{\kappa}$, $\mathcal{K} := \big\{k : k = \sum_{m=1}^{\kappa} \beta_m k_m, \sum_{m=1}^{\kappa} \beta_m = 1, \beta_m \geq 0, \forall m\big\}$. We set $\beta_m = 1/\kappa$ according to [34] and it works well in practice.

## 3.2. Supervised Hashing for Source Data

The Hamming distance for a pair of hash values $\boldsymbol{h}_i$ and $\boldsymbol{h}_j$ has a unique relationship with the dot product $\langle \boldsymbol{h}_i, \boldsymbol{h}_j \rangle$, given by: $\text{dist}_H(\boldsymbol{h}_i, \boldsymbol{h}_j) = \frac{1}{2}(d - \boldsymbol{h}_i^\top \boldsymbol{h}_j)$, where $d$ is the hash length. The dot product $\langle \boldsymbol{h}_i, \boldsymbol{h}_j \rangle$ can be treated as a similarity measure for the hash codes. Larger the value of the dot product (high similarity), smaller is the distance $\text{dist}_H$ and smaller the dot product (low similarity), larger is the distance $\text{dist}_H$. Let $s_{ij} \in \{0, 1\}$ be the similarity between $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$. If $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ belong to the same category, $s_{ij} = 1$ and 0, otherwise. The probability of similarity between $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ given the corresponding hash values $\boldsymbol{h}_i$ and $\boldsymbol{h}_j$, can be expressed as a likelihood function, given by,

$$p(s_{ij}|\boldsymbol{h}_i, \boldsymbol{h}_j) = \begin{cases} \sigma(\boldsymbol{h}_i^\top \boldsymbol{h}_j), & s_{ij} = 1 \\ 1 - \sigma(\boldsymbol{h}_i^\top \boldsymbol{h}_j), & s_{ij} = 0, \end{cases} \qquad (3)$$

where, $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function. As the dot product $\langle \boldsymbol{h}_i, \boldsymbol{h}_j \rangle$ increases, the probability of $p(s_{ij} = 1|\boldsymbol{h}_i, \boldsymbol{h}_j)$ also increases, i.e., $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ belong to the same category. As the dot product decreases, the probability $p(s_{ij} = 1|\boldsymbol{h}_i, \boldsymbol{h}_j)$ also decreases, i.e., $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ belong to different categories. We construct the $(n_s \times n_s)$ similarity matrix $\mathcal{S} = \{s_{ij}\}$, for the source data with the provided labels, where $s_{ij} = 1$ if $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ belong to the same category and 0, otherwise. Let $\mathbf{H} = \{\boldsymbol{h}_i\}_{i=1}^{n_s}$ be the set of source data hash values. If the elements of $\mathbf{H}$ are assumed to be i.i.d., the negative log likelihood of the similarity matrix $\mathcal{S}$ given $\mathbf{H}$ can be written as,

$$\min_{\mathbf{H}} \mathcal{L}(\mathbf{H}) = -\log p(\mathcal{S}|\mathbf{H})$$
$$= - \sum_{s_{ij} \in \mathcal{S}} \Big( s_{ij} \boldsymbol{h}_i^\top \boldsymbol{h}_j - \log\big(1 + \exp(\boldsymbol{h}_i^\top \boldsymbol{h}_j)\big) \Big). \qquad (4)$$

By minimizing Equation (4), we can determine hash values $\mathbf{H}$ for the source data which are consistent with the

similarity matrix $\mathcal{S}$. The hash loss has been used in previous research for supervised hashing [30, 50]. Equation (4) is a discrete optimization problem that is challenging to solve. We introduce a relaxation on the discrete constraint $h_i \in \{-1, +1\}^d$ by instead solving for $u_i \in \mathbb{R}^d$, where $\mathcal{U}_s = \{u_i\}_{i=1}^{n_s}$ is the output of the network and $u_i = \psi(x_i)$ (the superscript denoting the domain has been dropped for ease of representation). However, the continuous relaxation gives rise to (i) approximation error, when $\langle h_i, h_j \rangle$ is substituted with $\langle u_i, u_j \rangle$ and, (ii) quantization error, when the resulting real codes $u_i$ are binarized [50]. We account for the approximation error by having a tanh(.) as the final activation layer of the neural network, so that the components of $u_i$ are bounded between $-1$ and $+1$. In addition, we also introduce a quantization loss $||u_i - \text{sgn}(u_i)||_2^2$ along the lines of [22], where sgn(.) is the sign function. The continuous optimization problem for supervised hashing can now be outlined;

$$\min_{\mathcal{U}_s} \mathcal{L}(\mathcal{U}_s) = - \sum_{s_{ij} \in \mathcal{S}} \left( s_{ij} u_i^\top u_j - \log\left(1 + \exp(u_i^\top u_j)\right) \right)$$
$$+ \sum_{i=1}^{n_s} ||u_i - \text{sgn}(u_i)||_2^2. \qquad (5)$$

### 3.3. Unsupervised Hashing for Target Data

In the absence of target data labels, we use the similarity measure $\langle u_i, u_j \rangle$, to guide the network to learn discriminative hash values for the target data. An ideal target output $u_i^t$, needs to be similar to many of the source outputs from the $j^{th}$ category $\left(\{u_k^{s_j}\}_{k=1}^K\right)$. We assume without loss of generality, $K$ source data points for every category $j$ where, $j \in \{1, \ldots, C\}$ and $u_k^{s_j}$ is the $k^{th}$ source output from category $j$. In addition, $u_i^t$ must be dissimilar to most other source outputs $u_k^{s_l}$ belonging to a different category ($j \neq l$). Enforcing similarity with all the $K$ data points makes for a more robust target data category assignment. We outline a probability measure to capture this intuition. Let $p_{ij}$ be the probability that input target data point $x_i$ is assigned to category $j$ where,

$$p_{ij} = \frac{\sum_{k=1}^K \exp(u_i^{t\top} u_k^{s_j})}{\sum_{l=1}^C \sum_{k=1}^K \exp(u_i^{t\top} u_k^{s_l})} \qquad (6)$$

The exp(.) has been introduced for ease of differentiability and the denominator ensures $\sum_j p_{ij} = 1$. When the target data point output is similar to one category only and dissimilar to all the other categories, the probability vector $p_i = [p_{i1}, \ldots, p_{iC}]^T$ tends to be a one-hot vector. A one-hot vector can be viewed as a low entropy realization of $p_i$. We can therefore envisage all the $p_i$ to be one-hot vectors (low entropy probability vectors), where the target data point outputs are similar to source data point outputs in one and only one category. To this end we introduce a loss

to capture the entropy of the target probability vectors. The entropy loss for the network outputs is given by,

$$\mathcal{H}(\mathcal{U}_s, \mathcal{U}_t) = -\frac{1}{n_t} \sum_{i=1}^{n_t} \sum_{j=1}^C p_{ij} \log(p_{ij}) \qquad (7)$$

Minimizing the entropy loss gives us probability vectors $p_i$ that tend to be one-hot vectors, i.e., the target data point outputs are similar to source data point outputs from any one category only. Enforcing similarity with $K$ source data points from a category, guarantees that the hash values are determined based on a common similarity between multiple source category data points and the target data point.

### 3.4. Domain Adaptive Hash Network

We propose a model for deep unsupervised domain adaptation based on hashing (DAH) that incorporates unsupervised domain adaptation between the source and the target (1), supervised hashing for the source (5) and unsupervised hashing for the target (7) in a deep convolutional neural network. The DAH network is trained to minimize

$$\min_{\mathcal{U}} \mathcal{J} = \mathcal{L}(\mathcal{U}_s) + \gamma \mathcal{M}(\mathcal{U}_s, \mathcal{U}_t) + \eta \mathcal{H}(\mathcal{U}_s, \mathcal{U}_t), \qquad (8)$$

where, $\mathcal{U} := \{\mathcal{U}_s \cup \mathcal{U}_t\}$ and $(\gamma, \eta)$ control the importance of domain adaptation (1) and target entropy loss (7) respectively. The hash values $\mathbf{H}$ are obtained from the output of the network using $\mathbf{H} = \text{sgn}(\mathcal{U})$. The loss terms (5) and (7) are determined in the final layer of the network with the network output $\mathcal{U}$. The MK-MMD loss (1) is determined between layer outputs $\{\mathcal{U}_s^l, \mathcal{U}_t^l\}$ at each of the fully connected layers $\mathcal{F} = \{fc6, fc7, fc8\}$, where we adopt the linear time estimate for the unbiased MK-MMD as described in [24] and [31]. The DAH is trained using standard backpropagation. The detailed derivation of the derivative of (8) w.r.t. $\mathcal{U}$ is provided in the supplementary material.

**Network Architecture**: Owing to the paucity of images in a domain adaptation setting, we circumvent the need to train a deep CNN with millions of images by adapting the pre-trained VGG-F [8] network to the DAH. The VGG-F has been trained on the ImageNet 2012 dataset and it consists of 5 convolution layers (*conv*1 - *conv*5) and 3 fully connected layers (*fc*6, *fc*7, *fc*8). We introduce the hashing layer *hash-fc*8 that outputs vectors in $\mathbb{R}^d$ in the place of *fc*8. To account for the hashing approximation, we introduced a tanh() layer. However, we encounter the issue of vanishing gradients [26] when using tanh() as it saturates with large inputs. We therefore preface the tanh() with a batch normalization layer which prevents the tanh() from saturating. In effect, *hash-fc*8 := {*fc*8 → *batch-norm* → *tanh*()}. The *hash-fc*8 provides greater stability when fine-tuning the learning rates than the deep hashing networks [30, 50]. Figure 1 illustrates the proposed DAH network.

**Table 1:** Statistics for the *Office-Home* dataset. **Min: #** is the minimum number of images amongst all the categories, **Min: Size** and **Max: Size** are the minimum and maximum image sizes across all categories and **Acc.** is the classification accuracy.

| Domain. | Min: # | Min: Size | Max: Size | Acc |
|---|---|---|---|---|
| Art | 15 | 117×85 pix. | 4384×2686 pix. | 44.99±1.85 |
| Clipart | 39 | 18×18 pix. | 2400×2400 pix. | 53.95±1.45 |
| Product | 38 | 75×63 pix. | 2560×2560 pix. | 66.41±1.18 |
| Real-World | 23 | 88×80 pix. | 6500×4900 pix. | 59.70±1.04 |

## 4. The *Office-Home* Dataset

Supervised deep learning models require a large volume of labeled training data. Unfortunately, existing datasets for vision-based domain adaptation are limited in their size and are not suitable for validating deep learning algorithms. The standard datasets for vision based domain adaptation are, facial expression datasets *CKPlus* [35] and *MMI* [40], digit datasets *SVHN* [36], *USPS* and *MNIST*[28], head pose recognition datasets *PIE* [33], object recognition datasets *COIL*[33], *Office* [42] and *Office-Caltech* [20]. These datasets were created before deep-learning became popular and are insufficient for training and evaluating deep learning based domain adaptation approaches. For instance, the object-recognition dataset *Office* has 4110 images across 31 categories and *Office-Caltech* has 2533 images across 10 categories.

We release the ***Office-Home*** dataset for domain adaptation based object recognition, that can be used to evaluate deep learning algorithms for domain adaptation. The *Office-Home* dataset consists of 4 domains, with each domain containing images from 65 categories of everyday objects and a total of around 15,500 images. The domains include, Art: artistic depictions of objects in the form of sketches, paintings, ornamentation, etc.; Clipart: collection of clipart images; Product: images of objects without a background, akin to the Amazon category in *Office* dataset; Real-World: images of objects captured with a regular camera.

Public domain images were downloaded from websites like www.deviantart.com and www.flickr.com to create the Art and Real-World domains. Clipart images were gathered from multiple clipart websites. The Product domain images were exclusively collected from www.amazon.com using web-crawlers. The collected images were manually filtered on the basis of quality, size and content. The dataset has an average of around 70 images per category and a maximum of 99 images in a category. The primary challenge in creating this dataset was acquiring sufficient number of public domain images across all the 4 domains. Figure 2 depicts a sampling of 16 categories from the *Office-Home* dataset and Table 1 outlines some meta data for the dataset. The **Acc.** column in the Table 1 refers to classification accuracies using the LIBLINEAR SVM [15] classifier (5-fold cross validation) with deep features extracted using the VGG-F network. The dataset is publicly available for research [1].

## 5. Experiments

In this section we conduct extensive experiments to evaluate the DAH algorithm. Since we propose a domain adaptation technique based on hashing, we evaluate objection recognition accuracies for unsupervised domain adaptation and also study the discriminatory capability of the learned hash codes for unsupervised domain adaptive hashing. The implementation details are available at https://github.com/hemanthdv/da-hash

### 5.1. Datasets

***Office*** [42]: This is currently the most popular benchmark dataset for object recognition in the domain adaptation computer vision community. The dataset consists of images of everyday objects in an office environment. It has 3 domains; Amazon (**A**), Dslr (**D**) and Webcam (**W**). The dataset has around 4,100 images with a majority of the images (2816 images) in the Amazon domain. We adopt the common evaluation protocol of different pairs of transfer tasks for this dataset [31, 34]. We consider 6 transfer tasks for all combinations of source and target pairs for the 3 domains. ***Office-Home***: We introduce this new dataset and evaluate it in a similar manner to the *Office* dataset. We consider 12 transfer tasks for the Art (**Ar**), Clipart (**Cl**), Product (**Pr**) and Real-World (**Rw**) domains for all combinations of source and target for the 4 domains. Considering all the different pairs of transfer enables us to evaluate the inherent bias between the domains in a comprehensive manner [45].

### 5.2. Implementation Details

We implement the DAH using the MatConvnet framework [47]. Since we train a pre-trained VGG-F, we fine-tune the weights of *conv*1-*conv*5, *fc*6 and *fc*7. We set their learning rates to $1/10^{th}$ the learning rate of *hash-fc*8. We vary the learning rate between $10^{-4}$ to $10^{-5}$ over 300 epochs with a momentum 0.9 and weight decay $5 \times 10^{-4}$. We set $K = 5$ (number of samples from a category). Since we have 31 categories in the *Office* dataset, we get a source batch size of $31 \times 5 = 155$. For the target batch, we randomly select 155 samples. The total batch size turns out to be 310. For the *Office-Home* dataset, with $K = 5$ and 65 categories, we get a batch size of 650. We set $d = 64$ (hash code length) for all our experiments. Since there is imbalance in the number of like and unlike pairs in $\mathcal{S}$, we set the values in similarity matrix $\mathcal{S}_{i,j} \in \{0, 10\}$. Increasing the similarity weight of like-pairs improves the performance of DAH. For the entropy loss, we set $\eta = 1$. For the MK-MMD loss, we follow the heuristics mentioned in [24], to

---

[1]https://hemanthdv.github.io/officehome-dataset/

**Figure 2:** Sample images from the ***Office-Home*** dataset. The dataset consists of images of everyday objects organized into 4 domains; `Art`: paintings, sketches and/or artistic depictions, `Clipart`: clipart images, `Product`: images without background and `Real-World`: regular images captured with a camera. The figure displays examples from 16 of the 65 categories.

determine the parameters. We estimate $\gamma$, by validating a binary domain classifier to distinguish between source and target data points and select $\gamma$ which gives largest error on a validation set. For MMD, we use a Gaussian kernel with a bandwidth $\sigma$ given by the median of the pairwise distances in the training data. To incorporate the multi-kernel, we vary the bandwidth $\sigma_m \in [2^{-8}\sigma, 2^8\sigma]$ with a multiplicative factor of 2. We define the target classifier $f(\boldsymbol{x}_i^t) = p(y|\boldsymbol{h}_i^t)$ in terms of 6. The target data point is assigned to the class with the largest probability, with $\hat{y}_i = \max_j(p_{ij})$ using the hash codes for the source and the target.

## 5.3. Unsupervised Domain Adaptation

In this section, we study the performance of the DAH for unsupervised domain adaptation, where labeled data is available only in the source domain and no labeled data is available in the target domain. We compare the DAH with state-of-the-art domain adaptation methods: (i) Geodesic Flow Kernel (**GFK**) [20], (ii) Transfer Component Analysis (**TCA**) [38], (iii) Correlation Alignment (**CORAL**) [44] and (iv) Joint Distribution Adaptation (**JDA**) [33]. We also compare the DAH with state-of-the-art deep learning methods for domain adaptation: (v) Deep Adaptation Network (**DAN**) [31] and (vi) Domain Adversarial Neural Network (**DANN**) [17]. For all of the shallow learning methods, we extract and use deep features from the *fc*7 layer of the VGG-F network that was pre-trained on the ImageNet 2012 dataset. We also evaluate the effect of the entropy loss on hashing for the DAH. The **DAH-e** is the DAH algorithm where $\eta$ is set to zero, which implies that the target hash values are not driven to align with the source categories. We follow the standard protocol for unsupervised domain adaptation, where all the labeled source data and all the unlabeled target data is used for training.

**Results and Discussion**: The results are reported for the target classification in each of the transfer tasks in Tables 2 and 3, where accuracies denote the percentage of correctly

**Table 2:** Recognition accuracies (%) for domain adaptation experiments on the *Office* dataset. {`Amazon` (A), `Dslr` (D), `Webcam` (W)}. A→W implies A is source and W is target.

| Expt. | A→D | A→W | D→A | D→W | W→A | W→D | Avg. |
|---|---|---|---|---|---|---|---|
| GFK | 48.59 | 52.08 | 41.83 | 89.18 | 49.04 | 93.17 | 62.32 |
| TCA | 51.00 | 49.43 | 48.12 | 93.08 | 48.83 | 96.79 | 64.54 |
| CORAL | 54.42 | 51.70 | 48.26 | 95.97 | 47.27 | 98.59 | 66.04 |
| JDA | 59.24 | 58.62 | 51.35 | 96.86 | 52.34 | 97.79 | 69.37 |
| DAN | 67.04 | 67.80 | 50.36 | 95.85 | 52.33 | 99.40 | 72.13 |
| DANN | 72.89 | 72.70 | 56.25 | 96.48 | 53.20 | 99.40 | 75.15 |
| DAH-e | 66.27 | 66.16 | 55.97 | 94.59 | 53.91 | 96.99 | 72.31 |
| DAH | 66.47 | 68.30 | 55.54 | 96.10 | 53.02 | 98.80 | 73.04 |

classified target data samples. We present results with hash length $d = 64$ bits. The DAH algorithm consistently outperforms the baselines across all the domains for the *Office-Home* dataset. However, DANN marginally surpasses DAH for the *Office* dataset, prompting us to reason that domain adversarial training is more effective than DAH when the categories are fewer in number. Since domain alignment is category agnostic, it is possible that the aligned domains are not classification friendly in the presence of large number of categories. When the number of categories is large, as in *Office-Home*, DAH does best at extracting transferable features to achieve higher accuracies. We also note that DAH delivers better performance than DAH-e; thus, minimizing the entropy on the target data through 7 aids in improved alignment of the source and target samples, which boosts the accuracy.

**Feature Analysis**: We also study the feature representations of the penultimate layer (*fc*7) outputs using t-SNE embeddings as in [12]. Figure 3a depicts the $\mathcal{A}$-distance between domain pairs using Deep (VGG-F), DAN and DAH features. Ben-David et al. [2] defined $\mathcal{A}$-distance as the distance between two domains that can be viewed as the discrepancy between two domains. Although it is difficult to estimate its exact value, an approximate distance measure is given by $2(1 - 2\epsilon)$, where $\epsilon$ is the generalization error for a binary classifier trained to distinguish between the two domains. We used a LIBLINEAR SVM [15] clas-

**Table 3:** Recognition accuracies (%) for domain adaptation experiments on the *Office-Home* dataset. {Art (Ar), Clipart (Cl), Product (Pr), Real-World (Rw)}. Ar→Cl implies Ar is source and Cl is target.

| Expt. | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GFK | 21.60 | 31.72 | 38.83 | 21.63 | 34.94 | 34.20 | 24.52 | 25.73 | 42.92 | 32.88 | 28.96 | 50.89 | 32.40 |
| TCA | 19.93 | 32.08 | 35.71 | 19.00 | 31.36 | 31.74 | 21.92 | 23.64 | 42.12 | 30.74 | 27.15 | 48.68 | 30.34 |
| CORAL | 27.10 | 36.16 | 44.32 | 26.08 | 40.03 | 40.33 | 27.77 | 30.54 | 50.61 | 38.48 | 36.36 | 57.11 | 37.91 |
| JDA | 25.34 | 35.98 | 42.94 | 24.52 | 40.19 | 40.90 | 25.96 | 32.72 | 49.25 | 35.10 | 35.35 | 55.35 | 36.97 |
| DAN | 30.66 | 42.17 | 54.13 | 32.83 | 47.59 | 49.78 | 29.07 | 34.05 | 56.70 | 43.58 | 38.25 | 62.73 | 43.46 |
| DANN | 33.33 | 42.96 | 54.42 | 32.26 | 49.13 | 49.76 | 30.49 | 38.14 | 56.76 | 44.71 | 42.66 | 64.65 | 44.94 |
| DAH-e | 29.23 | 35.71 | 48.29 | 33.79 | 48.23 | 47.49 | 29.87 | 38.76 | 55.63 | 41.16 | 44.99 | 59.07 | 42.69 |
| DAH | 31.64 | 40.75 | 51.73 | 34.69 | 51.93 | 52.79 | 29.91 | 39.63 | 60.71 | 44.99 | 45.13 | 62.54 | 45.54 |

sifier with 5-fold cross-validation to estimate $\epsilon$. Figure 3a indicates that the DAH features have the least discrepancy between the source and target compared to DAN and Deep features. This is also confirmed with the t-SNE embeddings in Figures 3b-3d. The Deep features show very little overlap between the domains and the categories depict minimal clustering. Domain overlap and clustering improves as we move to DAN and DAH features, with DAH providing the best visualizations. This corroborates the efficacy of the DAH algorithm to exploit the feature learning capabilities of deep neural networks to learn representative hash codes to address domain adaptation.

### 5.4. Unsupervised Domain Adaptive Hashing

In this section, we study the performance of our algorithm to generate compact and efficient hash codes from the data for classifying unseen test instances, when no labels are available. This problem has been addressed in the literature, with promising empirical results [7, 11, 21]. However, in a real-world setting, labels may be available from a different, but related (source) domain; a strategy to utilize the labeled data from the source domain to learn representative hash codes for the target domain is therefore of immense practical importance. Our work is the first to identify and address this problem. We consider the following scenarios to address this real-world challenge: (i) No labels are available for a given dataset and the hash codes need to be learned in a completely unsupervised manner. We evaluate against baseline unsupervised hashing methods (**ITQ**) [22] and (**KMeans**) [25] and also state-of-the-art methods for unsupervised hashing (**BA**) [7] and (**BDNN**) [11]. (ii) Labeled data is available from a different, but related source domain. A hashing model is trained on the labeled source data and is used to learn hash codes for the target data. We refer to this method as **NoDA**, as no domain adaptation is performed. We used the deep pairwise-supervised hashing (DPSH) algorithm [30] to train a deep network with the source data and applied the network to generate hash codes for the target data. (iii) Labeled data is available from a different, but related source domain and we use our DAH formulation to learn hash codes for the target domain by reducing domain disparity. (iv) Labeled data is available

**Table 4:** Mean average precision @64 bits. For the NoDA and DAH results, Art is the source domain for Clipart, Product and Real-World and Clipart is the source domain for Art. Similarly, Amazon and Webcam are source target pairs.

| Expt. | NoDA | ITQ | KMeans | BA | BDNN | DAH | SuH |
|---|---|---|---|---|---|---|---|
| Amazon | 0.324 | 0.465 | 0.403 | 0.367 | 0.491 | 0.582 | 0.830 |
| Webcam | 0.511 | 0.652 | 0.558 | 0.480 | 0.656 | 0.717 | 0.939 |
| Art | 0.155 | 0.191 | 0.170 | 0.156 | 0.193 | 0.302 | 0.492 |
| Clipart | 0.160 | 0.195 | 0.178 | 0.179 | 0.206 | 0.333 | 0.622 |
| Product | 0.239 | 0.393 | 0.341 | 0.349 | 0.407 | 0.414 | 0.774 |
| Real-World | 0.281 | 0.323 | 0.279 | 0.273 | 0.336 | 0.533 | 0.586 |
| Avg. | 0.278 | 0.370 | 0.322 | 0.301 | 0.382 | 0.480 | 0.707 |

in the target domain. This method falls under supervised hashing (**SuH**) (as it uses labeled data in the target domain to learn hash codes in the same domain) and denotes the upper bound on the performance. It is included to compare the performance of unsupervised hashing algorithms relative to the supervised algorithm. We used the DPSH algorithm [30] to train a deep network on the target data and used it to generate hash codes on a validation subset.

**Results and Discussion**: We applied the precision-recall curves and the mean average precision (mAP) measures to evaluate the efficacy of the hashing methods, similar to previous research [7, 11, 21]. The results are depicted in Figures 4 and 5 (precision-recall curves) and Table 4 (mAP values), where we present hashing with code length $d = 64$ bits. Hashing performance with $d = 16$ bits also follows a similar trend and is presented in the supplementary material. For the sake of brevity, we drop the results with Dslr as it is very similar to Webcam, with little domain difference. We note that the NoDA has the poorest performance due to domain mismatch. This demonstrates that domain disparity needs to be considered before deploying a hashing network to extract hash codes. The unsupervised hashing methods ITQ, KMeans, BA and BDNN perform slightly better compared to NoDA. The proposed DAH algorithm encompasses hash code learning and domain adaptation in a single integrated framework. It is thus able to leverage the labeled data in the source domain in a meaningful manner to learn efficient hash codes for the target domain. This accounts for its improved performance, as is evident in Figures 4 and 5 and Table 4. The supervised hashing technique (SuH) uses labels from the target and therefore depicts the
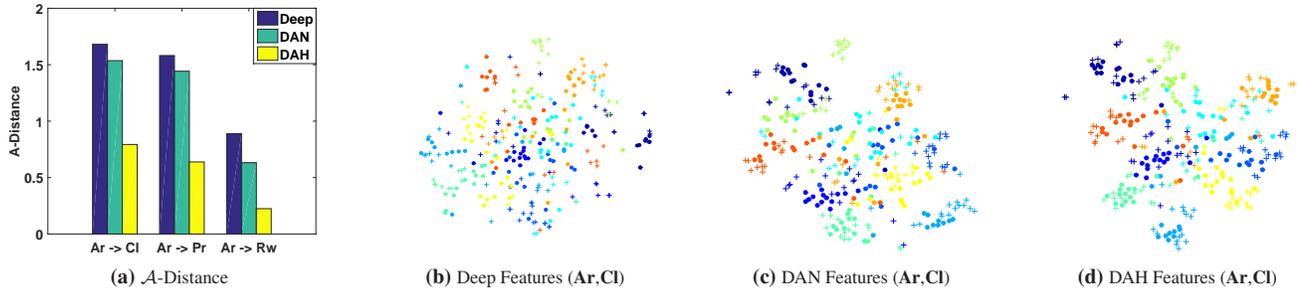
**(a)** $\mathcal{A}$-Distance  **(b)** Deep Features (**Ar,Cl**)  **(c)** DAN Features (**Ar,Cl**)  **(d)** DAH Features (**Ar,Cl**)

**Figure 3:** Feature analysis of $fc7$ layer. (a) $\mathcal{A}$-distances for Deep, DAN and DAH, (b), (c) and (d) t-SNE embeddings for 10 categories from Art (●) and Clipart(+) domains. Best viewed in color.



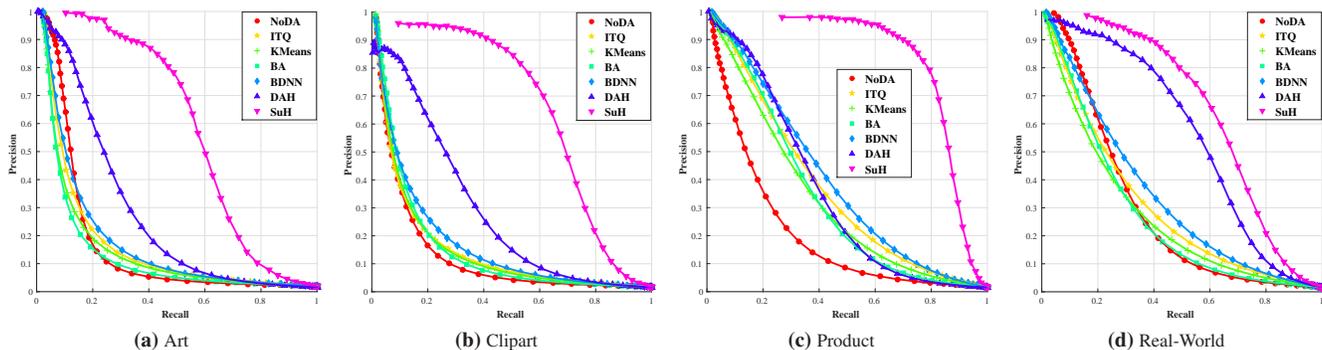**(a)** Art  **(b)** Clipart  **(c)** Product  **(d)** Real-World

**Figure 4:** Precision-Recall curves @64 bits for the ***Office-Home*** dataset. Comparison of hashing without domain adaptation (**NoDA**), shallow unsupervised hashing (**ITQ**, **KMeans**), state-of-the-art deep unsupervised hashing (**BA**, **BDNN**), unsupervised domain adaptive hashing (**DAH**) and supervised hashing (**SuH**). Best viewed in color.



**(a)** Amazon  **(b)** Webcam

**Figure 5:** Precision-Recall curves @64 bits for the ***Office*** dataset. Comparison of hashing without domain adaptation (**NoDA**), shallow unsupervised hashing (**ITQ**, **KMeans**), state-of-the-art deep unsupervised hashing (**BA**, **BDNN**), unsupervised domain adaptive hashing (**DAH**) and supervised hashing (**SuH**). Best viewed in color.
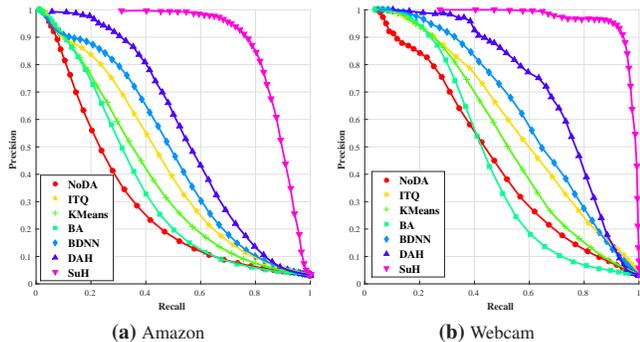
best performance. The proposed DAH framework consistently delivers the best performance relative to SuH when compared with the other hashing procedures. This demonstrates the merit of our framework in learning representative hash codes by utilizing labeled data from a different domain. Such a framework will be immensely useful in a

real-world setting.

## 6. Conclusions

In this paper, we have proposed a novel domain adaptive hashing (DAH) framework which exploits the feature learning capabilities of deep neural networks to learn efficient hash codes for unsupervised domain adaptation. The DAH framework solves two important practical problems: category assignment with weak supervision or insufficient labels (through domain adaptation) and the estimation of hash codes in an unsupervised setting (hash codes for target data). Thus, two practical challenges are addressed through a single integrated framework. This research is the first of its kind to integrate hash code learning with unsupervised domain adaptation. We also introduced a new dataset, ***Office-Home***, which can be used to further research in domain adaptation.

# References

[1] Y. Aytar and A. Zisserman. Tabula rasa: Model transfer for object category detection. In *IEEE ICCV*, 2011. 2

[2] 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(1-2):151–175, 2010. 6

[3] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013. 2

[4] L. Bruzzone and M. Marconcini. Domain adaptation problems: A dasvm classification technique and a circular validation strategy. *IEEE, PAMI*, 32(5):770–787, 2010. 2

[5] Y. Cao, M. Long, J. Wang, Q. Yang, and P. S. Yu. Deep visual-semantic hashing for cross-modal retrieval. In *ACM-SIGKDD*, 2016. 2

[6] Z. Cao, M. Long, and Q. Yang. Transitive hashing network for heterogeneous multimedia retrieval. In *AAAI*, 2016. 2

[7] M. A. Carreira-Perpinán and R. Raziperchikolaei. Hashing with binary autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 557–566, 2015. 2, 7

[8] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*, 2014. 2, 4

[9] R. Chattopadhyay, Q. Sun, W. Fan, I. Davidson, S. Panchanathan, and J. Ye. Multisource domain adaptation and its application to early detection of fatigue. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(4):18, 2012. 2

[10] W.-S. Chu, F. De la Torre, and J. F. Cohn. Selective transfer machine for personalized facial action unit detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3515–3522, 2013. 2

[11] T.-T. Do, A.-D. Doan, and N.-M. Cheung. Learning to hash with binary deep neural network. In *European Conference on Computer Vision*, pages 219–234. Springer, 2016. 2, 7

[12] 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*, pages 647–655, 2014. 2, 6

[13] L. Duan, I. W. Tsang, and D. Xu. Domain transfer multiple kernel learning. *IEEE PAMI*, 34(3):465–479, 2012. 2

[14] V. Erin Liong, J. Lu, G. Wang, P. Moulin, and J. Zhou. Deep hashing for compact binary codes learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2475–2483, 2015. 2

[15] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. Liblinear: A library for large linear classification. *Journal of machine learning research*, 9(Aug):1871–1874, 2008. 5, 6

[16] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *CVPR*, pages 2960–2967, 2013. 1, 2

[17] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016. 1, 2, 3, 6

[18] X. Glorot, A. Bordes, and Y. Bengio. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 513–520, 2011. 1, 2

[19] B. Gong, K. Grauman, and F. Sha. Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation. In *ICML (1)*, pages 222–230, 2013. 2

[20] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *IEEE CVPR*, 2012. 1, 2, 5, 6

[21] Y. Gong and S. Lazebnik. Iterative quantization: A procrustean approach to learning binary codes. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 817–824. IEEE, 2011. 7

[22] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, 2013. 2, 4, 7

[23] R. Gopalan, R. Li, and R. Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *2011 international conference on computer vision*, pages 999–1006. IEEE, 2011. 2

[24] A. Gretton, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, and B. K. Sriperumbudur. Optimal kernel choice for large-scale two-sample tests. In *Advances in neural information processing systems*, pages 1205–1213, 2012. 3, 4, 5, 11

[25] K. He, F. Wen, and J. Sun. K-means hashing: An affinity-preserving quantization method for learning binary compact codes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2938–2945, 2013. 2, 7

[26] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001. 4

[27] J. Hoffman, E. Rodner, J. Donahue, K. Saenko, and T. Darrell. Efficient learning of domain-invariant image representations. In *ICLR*, 2013. 2

[28] K. Jarrett, K. Kavukcuoglu, Y. Lecun, et al. What is the best multi-stage architecture for object recognition? In *2009 IEEE 12th International Conference on Computer Vision*, pages 2146–2153. IEEE, 2009. 5

[29] Q.-Y. Jiang and W.-J. Li. Deep cross-modal hashing. *arXiv preprint arXiv:1602.02255*, 2016. 2

[30] W.-J. Li, S. Wang, and W.-C. Kang. Feature learning based deep supervised hashing with pairwise labels. In *IJCAI, 2016*, 2016. 4, 7

[31] M. Long, Y. Cao, J. Wang, and M. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, pages 97–105, 2015. 1, 2, 3, 4, 5, 6

[32] M. Long, J. Wang, G. Ding, J. Sun, and P. Yu. Transfer joint matching for unsupervised domain adaptation. In *CVPR*, pages 1410–1417, 2014. 2

[33] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2200–2207, 2013. 1, 2, 5, 6

[34] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. In *NIPS*, 2016. 1, 2, 3, 5

[35] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *CVPR*, pages 94–101. IEEE, 2010. 5

[36] 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 Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. 5

[37] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1717–1724, 2014. 2

[38] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang. Domain adaptation via transfer component analysis. *Neural Networks, IEEE Trans. on*, 22(2):199–210, 2011. 1, 2, 6

[39] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE TKDE*, 22(10):1345–1359, 2010. 2

[40] M. Pantic, M. Valstar, R. Rademaker, and L. Maat. Web-based database for facial expression analysis. In *ICME*. IEEE, 2005. 5

[41] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 32(3):53–69, 2015. 2

[42] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *ECCV*, 2010. 1, 2, 5

[43] S. Shekhar, V. M. Patel, H. V. Nguyen, and R. Chellappa. Generalized domain-adaptive dictionaries. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 361–368, 2013. 1

[44] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. In *ICCV, TASK-CV*, 2015. 1, 2, 6

[45] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE, 2011. 5

[46] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4068–4076, 2015. 1, 2, 3

[47] A. Vedaldi and K. Lenc. Matconvnet – convolutional neural networks for matlab. In *Proceeding of the ACM Int. Conf. on Multimedia*, 2015. 5

[48] J. Wang, H. T. Shen, J. Song, and J. Ji. Hashing for similarity search: A survey. *arXiv preprint arXiv:1408.2927*, 2014. 1

[49] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014. 2, 3

[50] H. Zhu, M. Long, J. Wang, and Y. Cao. Deep hashing network for efficient similarity retrieval. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016. 4

# Supplementary Material

## 7. Loss Function Derivative

In this section we outline the derivative of Equation 8 for the backpropagation algorithm;

$$\min_{\mathcal{U}} \mathcal{J} = \mathcal{L}(\mathcal{U}_s) + \gamma \mathcal{M}(\mathcal{U}_s, \mathcal{U}_t) + \eta \mathcal{H}(\mathcal{U}_s, \mathcal{U}_t), \tag{8}$$

where, $\mathcal{U} \coloneqq \{\mathcal{U}_s \cup \mathcal{U}_t\}$ and $(\gamma, \eta)$ control the importance of domain adaptation (1) and target entropy loss (7) respectively. In the following subsections, we outline the derivative of the individual terms w.r.t. the input $\mathcal{U}$.

### 7.1. Derivative for MK-MMD

$$\mathcal{M}(\mathcal{U}_s, \mathcal{U}_t) = \sum_{l \in \mathcal{F}} d_k^2(\mathcal{U}_s^l, \mathcal{U}_t^l), \tag{1}$$

$$d_k^2(\mathcal{U}_s^l, \mathcal{U}_t^l) = \left\| \mathbb{E}[\phi(\boldsymbol{u}^{s,l})] - \mathbb{E}[\phi(\boldsymbol{u}^{t,l})] \right\|_{\mathscr{H}_k}^2. \tag{2}$$

We implement the linear MK-MMD loss according to [24]. For this derivation, we consider the loss at just one layer. The derivative for the MK-MMD loss at every other layer can be derived in a similar manner. The output of $i^{th}$ source data point at layer $l$ is represented as $\boldsymbol{u}_i$ and the output of the $i^{th}$ target data point is represented as $\boldsymbol{v}_i$. For ease of representation, we drop the superscripts for the source $(s)$, the target $(t)$ and the layer $(l)$. Unlike the conventional MMD loss which is $\mathcal{O}(n^2)$, the MK-MMD loss outlined in [24] is $\mathcal{O}(n)$ and can be estimated online (does not require all the data). The loss is calculated over every batch of data points during the back-propagation. Let $n$ be the number of source data points $\mathcal{U} \coloneqq \{\boldsymbol{u}_i\}_{i=1}^n$ and the number of target data points $\mathcal{V} \coloneqq \{\boldsymbol{v}_i\}_{i=1}^n$ in the batch. We assume equal number of source and target data points in a batch and that $n$ is even. The MK-MMD is defined over a set of 4 data points $\boldsymbol{w}_i = [\boldsymbol{u}_{2i-1}, \boldsymbol{u}_{2i}, \boldsymbol{v}_{2i-1}, \boldsymbol{v}_{2i}], \forall i \in \{1, 2, \ldots, n/2\}$. The MK-MMD is given by,

$$\mathcal{M}(\mathcal{U}, \mathcal{V}) = \sum_{m=1}^{\kappa} \beta_m \frac{1}{n/2} \sum_{i=1}^{n/2} h_m(\boldsymbol{w}_i), \tag{9}$$

where, $\kappa$ is the number of kernels and $\beta_m = 1/\kappa$ is the weight for each kernel and,

$$h_m(\boldsymbol{w}_i) = k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{u}_{2i}) + k_m(\boldsymbol{v}_{2i-1}, \boldsymbol{v}_{2i}) - k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{v}_{2i}) - k_m(\boldsymbol{u}_{2i}, \boldsymbol{v}_{2i-1}), \tag{10}$$

where, $k_m(\boldsymbol{x}, \boldsymbol{y}) = \exp\left(-\frac{||\boldsymbol{x}-\boldsymbol{y}||_2^2}{\sigma_m}\right)$. Re-writing the MK-MMD in terms of the kernels, we have,

$$\mathcal{M}(\mathcal{U}, \mathcal{V}) = \frac{2}{n\kappa} \sum_{m=1}^{\kappa} \sum_{i=1}^{n/2} \left[ k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{u}_{2i}) + k_m(\boldsymbol{v}_{2i-1}, \boldsymbol{v}_{2i}) - k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{v}_{2i}) - k_m(\boldsymbol{u}_{2i}, \boldsymbol{v}_{2i-1}) \right], \tag{11}$$

We now outline the derivative of 11 w.r.t. source output $\boldsymbol{u}_q$ and target output $\boldsymbol{v}_q$. The derivative is,

$$\frac{\partial \mathcal{M}}{\partial \boldsymbol{u}_q} = \frac{2}{n\kappa} \sum_{m=1}^{\kappa} \sum_{i=1}^{n/2} \left[ \frac{2}{\sigma_m} k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{u}_{2i}).(\boldsymbol{u}_{2i-1} - \boldsymbol{u}_{2i}).(\mathcal{I}\{q = 2i\} - \mathcal{I}\{q = 2i - 1\}) \right.$$
$$\left. + \frac{2}{\sigma_m} k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{v}_{2i}).(\boldsymbol{u}_{2i-1} - \boldsymbol{v}_{2i}).\mathcal{I}\{q = 2i - 1\} + \frac{2}{\sigma_m} k_m(\boldsymbol{u}_{2i}, \boldsymbol{v}_{2i-1}).(\boldsymbol{u}_{2i} - \boldsymbol{v}_{2i-1}).\mathcal{I}\{q = 2i\} \right],$$

(12)

where, $\mathcal{I}\{.\}$ is the indicator function which is 1 if the condition is true, else it is false. The derivative w.r.t. the target data output $\boldsymbol{v}_q$ is,

$$\frac{\partial \mathcal{M}}{\partial \boldsymbol{v}_q} = \frac{2}{n\kappa} \sum_{m=1}^{\kappa} \sum_{i=1}^{n/2} \left[ \frac{2}{\sigma_m} k_m(\boldsymbol{v}_{2i-1}, \boldsymbol{v}_{2i}).(\boldsymbol{v}_{2i-1} - \boldsymbol{v}_{2i}).(\mathcal{I}\{q = 2i\} - \mathcal{I}\{q = 2i - 1\}) \right.$$
$$\left. - \frac{2}{\sigma_m} k_m(\boldsymbol{u}_{2i-1}, \boldsymbol{v}_{2i}).(\boldsymbol{u}_{2i-1} - \boldsymbol{v}_{2i}).\mathcal{I}\{q = 2i\} - \frac{2}{\sigma_m} k_m(\boldsymbol{u}_{2i}, \boldsymbol{v}_{2i-1}).(\boldsymbol{u}_{2i} - \boldsymbol{v}_{2i-1}).\mathcal{I}\{q = 2i - 1\} \right],$$

(13)

## 7.2. Derivative for Supervised Hash Loss

The supervised hash loss is given by,

$$\min_{\mathcal{U}_s} \mathcal{L}(\mathcal{U}_s) = - \sum_{s_{ij} \in \mathcal{S}} \left( s_{ij} \boldsymbol{u}_i^{\top} \boldsymbol{u}_j - \log\left(1 + \exp(\boldsymbol{u}_i^{\top} \boldsymbol{u}_j)\right) \right)$$
$$+ \sum_{i=1}^{n_s} \left|\left| \boldsymbol{u}_i - \text{sgn}(\boldsymbol{u}_i) \right|\right|_2^2.$$

(5)

The partial derivative of 5 w.r.t. source data output $\boldsymbol{u}_p$ is given by,

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{u}_q} = \sum_{s_{ij} \in \mathcal{S}} \left[ I\{i = q\}\left(\sigma(\boldsymbol{u}_i^{\top} \boldsymbol{u}_j) - s_{ij}\right)\boldsymbol{u}_j + I\{j = q\}\left(\sigma(\boldsymbol{u}_i^{\top} \boldsymbol{u}_j) - s_{ij}\right)\boldsymbol{u}_i \right] + 2(\boldsymbol{u}_q - \text{sgn}(\boldsymbol{u}_q))$$

(14)

where, $\sigma(x) = \frac{1}{1+exp(-x)}$. We assume sgn(.) to be a constant and avoid the differentiability issues with sgn(.) at 0. Since the $\mathcal{S}$ is symmetric, we can reduce the derivative to,

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{u}_q} = \sum_{j=1}^{n_s} \left[ 2\left(\sigma(\boldsymbol{u}_q^{\top} \boldsymbol{u}_j) - s_{qj}\right)\boldsymbol{u}_j \right] + 2\left(\boldsymbol{u}_q - \text{sgn}(\boldsymbol{u}_q)\right).$$

(15)

## 7.3. Derivative for Unsupervised Entropy Loss

We outline the derivative of $\frac{d\mathcal{H}}{d\mathcal{U}}$ in the following section, where $\mathcal{H}$ is defined as,

$$\mathcal{H}(\mathcal{U}_s, \mathcal{U}_t) = -\frac{1}{n_t} \sum_{i=1}^{n_t} \sum_{j=1}^{C} p_{ij} \log(p_{ij})$$

(7)

and $p_{ij}$ is the probability of target data output $\boldsymbol{u}_i^t$ belonging to category $j$, given by

$$p_{ij} = \frac{\sum_{k=1}^{K} \exp(\boldsymbol{u}_i^{t\top} \boldsymbol{u}_k^{s_j})}{\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(\boldsymbol{u}_i^{t\top} \boldsymbol{u}_{k'}^{s_l})}$$

(6)

For ease of representation, we will denote the target output $\boldsymbol{u}_i^t$ as $\boldsymbol{v}_i$ and drop the superscript $t$. Similarly, we will denote the $k^{th}$ source data point in the $j^{th}$ category $\boldsymbol{u}_k^{s_j}$ as $\boldsymbol{u}_k^j$, by dropping the domain superscript. We define the probability $p_{ij}$ with the news terms as,

$$p_{ij} = \frac{\sum_{k=1}^{K} \exp(\boldsymbol{v}_i^{\top} \boldsymbol{u}_k^j)}{\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(\boldsymbol{v}_i^{\top} \boldsymbol{u}_{k'}^l)}$$

(16)

Further, we simplify by replacing $\exp(\boldsymbol{v}_i^\top \boldsymbol{u}_k^j)$ with $\exp(i, jk)$. Equation 16 can now be represented as,

$$p_{ij} = \frac{\sum_{k=1}^{K} \exp(i, jk)}{\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(i, lk')} \tag{17}$$

We drop the outer summations (along with the -ve sign) and will reintroduce it at a later time. The entropy loss can be re-phrased using $\log(\frac{a}{b}) = \log(a) - \log(b)$ as,

$$\mathcal{H}_{ij} = \frac{\sum_{k=1}^{K} \exp(i, jk)}{\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(i, lk')} \log\left(\sum_{k=1}^{K} \exp(i, jk)\right) \tag{18}$$

$$- \frac{\sum_{k=1}^{K} \exp(i, jk)}{\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(i, lk')} \log\left(\sum_{l=1}^{C} \sum_{k'=1}^{K} \exp(i, lk')\right) \tag{19}$$

We need to estimate both, $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i}$ for the target and $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}$ for the source. We refer to $\partial \boldsymbol{u}_q^p$ for a consistent reference to source data. The derivative $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}$ for 18 is,

$$\left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}\right]_{18} = \frac{\boldsymbol{v}_i}{\sum_{l,k'} \exp(i, lk')} \left[\sum_k I\{{}_{k=q}^{j=p;}\} \exp(i, jk).\log\left(\sum_k \exp(i, jk)\right) + \sum_k I\{{}_{k=q}^{j=p;}\} \exp(i, jk)\right.$$

$$\left. - p_{ij} \exp(i, pq) \log\left(\sum_k \exp(i, jk)\right)\right], \tag{20}$$

where, $I\{.\}$ is an indicator function which is 1 only when both the conditions within are true, else it is 0. The derivative $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}$ for 19 is,

$$\left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}\right]_{19} = -\frac{\boldsymbol{v}_i}{\sum_{l,k'} \exp(i, lk')} \left[\sum_k I\{{}_{k=q}^{j=p;}\} \exp(i, jk).\log\left(\sum_{l,k'} \exp(i, lk')\right) + p_{ij} \exp(i, pq)\right.$$

$$\left. - p_{ij} \exp(i, pq) \log\left(\sum_{l,k'} \exp(i, lk')\right)\right] \tag{21}$$

Expressing $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p} = \left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}\right]_{18} + \left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p}\right]_{19}$, and defining $\bar{p}_{ijk} = \frac{\exp(i, jk)}{\sum_{l,k'} \exp(i, lk')}$ the derivative w.r.t. the source is,

$$\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{u}_q^p} = \boldsymbol{v}_i \left[\sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk}.\log\left(\sum_k \exp(i, jk)\right) + \sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk}\right.$$

$$- p_{ij} \bar{p}_{ipq} \log\left(\sum_k \exp(i, jk)\right) - \sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk}.\log\left(\sum_{l,k'} \exp(i, lk')\right)$$

$$\left. - p_{ij} \bar{p}_{ipq} + p_{ij} \bar{p}_{ipq} \log\left(\sum_{l,k'} \exp(i, lk')\right)\right] \tag{22}$$

$$= \boldsymbol{v}_i \left[\sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk} \log(p_{ij}) - p_{ij} \bar{p}_{ipq} \log(p_{ij}) + \sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk} - p_{ij} \bar{p}_{ipq}\right] \tag{23}$$

$$= \boldsymbol{v}_i \left(\log(p_{ij}) + 1\right) \left[\sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk} - p_{ij} \bar{p}_{ipq}\right] \tag{24}$$

The derivative of $\mathcal{H}$ w.r.t the **source** output $\boldsymbol{u}_q^p$ is given by,

$$\frac{\partial \mathcal{H}}{\partial \boldsymbol{u}_q^p} = -\frac{1}{n_t} \sum_{i=1}^{n_t} \sum_{j=1}^{C} \boldsymbol{v}_i \left(\log(p_{ij}) + 1\right) \left[\sum_k I\{{}_{k=q}^{j=p;}\} \bar{p}_{ijk} - p_{ij} \bar{p}_{ipq}\right] \tag{25}$$

We now outline the derivative $\frac{\partial \mathcal{H}}{\partial \boldsymbol{v}_i}$ for 18 as,

$$\left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i}\right]_{18} = \frac{1}{\sum_{l,k'} \exp(i, lk')} \left[\log\left(\sum_k \exp(i, jk)\right) \sum_k \exp(i, jk) \boldsymbol{u}_k^j + \sum_k \exp(i, jk) \boldsymbol{u}_k^j\right.$$

$$\left. - \frac{1}{\sum_{l,k'} \exp(i, lk')} \sum_k \exp(i, jk) \log\left(\sum_k \exp(i, jk)\right) \sum_{l,k'} \exp(i, lk') \boldsymbol{u}_{k'}^l\right], \tag{26}$$

and the derivative $\frac{\partial \mathcal{H}}{\partial \boldsymbol{v}_i}$ for 19 as,

$$\left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i}\right]_{19} = -\frac{1}{\sum_{l,k'} \exp(i,lk')}\Big[\log\big(\textstyle\sum_{l,k'} \exp(i,lk')\big)\sum_k \exp(i,jk)\boldsymbol{u}_k^j + \frac{\sum_k \exp(i,jk)}{\sum_{l,k'} \exp(i,lk')}\sum_{l,k'}\exp(i,lk')\boldsymbol{u}_{k'}^l$$
$$-\frac{1}{\sum_{l,k'}\exp(i,lk')}\sum_k \exp(i,jk)\log\big(\textstyle\sum_{l,k'}\exp(i,lk')\big)\sum_{l,k'}\exp(i,lk')\boldsymbol{u}_{k'}^l\Big], \quad (27)$$

Expressing $\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i} = \left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i}\right]_{18} + \left[\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i}\right]_{19}$, we get,

$$\frac{\partial \mathcal{H}_{ij}}{\partial \boldsymbol{v}_i} = \frac{1}{\sum_{l,k'}\exp(i,lk')}\Big[\log\big(\textstyle\sum_k \exp(i,jk)\big)\sum_k \exp(i,jk)\boldsymbol{u}_k^j - \log\big(\textstyle\sum_{l,k'}\exp(i,lk')\big)\sum_k \exp(i,jk)\boldsymbol{u}_k^j$$
$$+ \sum_k \exp(i,jk)\boldsymbol{u}_k^j - p_{ij}\sum_{l,k'}\exp(i,lk')\boldsymbol{u}_{k'}^l$$
$$- p_{ij}\log\big(\textstyle\sum_k \exp(i,jk)\big)\sum_{l,k'}\exp(i,lk')\boldsymbol{u}_{k'}^l + p_{ij}\log\big(\textstyle\sum_{l,k'}\exp(i,lk')\big)\sum_{l,k'}\exp(i,lk')\boldsymbol{u}_{k'}^l\Big] \quad (28)$$
$$= \Big[\log\big(\textstyle\sum_k \exp(i,jk)\big)\sum_k \bar{p}_{ijk}\boldsymbol{u}_k^j - \log\big(\textstyle\sum_{l,k'}\exp(i,lk')\big)\sum_k \bar{p}_{ijk}\boldsymbol{u}_k^j$$
$$+ \sum_k \bar{p}_{ijk}\boldsymbol{u}_k^j - p_{ij}\sum_{l,k'}\bar{p}_{ijk'}\boldsymbol{u}_{k'}^l$$
$$- p_{ij}\log\big(\textstyle\sum_k \exp(i,jk)\big)\sum_{l,k'}\bar{p}_{ijk'}\boldsymbol{u}_{k'}^l + p_{ij}\log\big(\textstyle\sum_{l,k'}\exp(i,lk')\big)\sum_{l,k'}\bar{p}_{ijk'}\boldsymbol{u}_{k'}^l\Big] \quad (29)$$
$$= \big(\log(p_{ij}) + 1\big)\sum_k \bar{p}_{ijk}\boldsymbol{u}_k^j - \big(\log(p_{ij}) + 1\big)p_{ij}\sum_{l,k'}\bar{p}_{ijk'}\boldsymbol{u}_{k'}^l \quad (30)$$
$$= \big(\log(p_{ij}) + 1\big)\big(\sum_k \bar{p}_{ijk}\boldsymbol{u}_k^j - p_{ij}\sum_{l,k'}\bar{p}_{ijk'}\boldsymbol{u}_{k'}^l\big) \quad (31)$$

The derivative of $\mathcal{H}$ w.r.t. **target** output $\boldsymbol{v}_q$ is given by,

$$\frac{\partial \mathcal{H}}{\partial \boldsymbol{v}_q} = -\frac{1}{n_t}\sum_{j=1}^C \big(\log(p_{qj}) + 1\big)\big(\sum_k \bar{p}_{qjk}\boldsymbol{u}_k^j - p_{qj}\sum_{l,k'}\bar{p}_{qjk'}\boldsymbol{u}_{k'}^l\big) \quad (32)$$

The derivative of $\mathcal{H}$ w.r.t. the source outputs is given by 25 and w.r.t. the target outputs is given by 32.

## 8. Unsupervised Domain Adaptation: Additional Results

In the main paper we had presented results for unsupervised domain adaptation based object recognition with $d = 64$ bits. Here, we outline the classification results with $d = 16$ (DAH-16) and $d = 128$ (DAH-128) bits for the *Office-Home* dataset in Table 5. We also present the (DAH-64), DAN and DANN results for comparison. There is an increase in the average recognition accuracy for $d = 128$ bits compared to $d = 64$ bits because of the increased capacity in representation. As expected, $d = 16$ has a lower recognition accuracy.

**Table 5:** Recognition accuracies (%) for domain adaptation experiments on the *Office-Home* dataset. {Art (Ar), Clipart (Cl), Product (Pr), Real-World (Rw)}. Ar→Cl implies Ar is source and Cl is target.

| Expt. | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DAN | 30.66 | 42.17 | 54.13 | 32.83 | 47.59 | 49.78 | 29.07 | 34.05 | 56.70 | 43.58 | 38.25 | 62.73 | 43.46 |
| DANN | 33.33 | 42.96 | 54.42 | 32.26 | 49.13 | 49.76 | 30.49 | 38.14 | 56.76 | 44.71 | 42.66 | 64.65 | 44.94 |
| DAH-16 | 23.83 | 30.32 | 40.14 | 25.67 | 38.79 | 33.26 | 20.11 | 27.72 | 40.90 | 32.63 | 25.54 | 37.46 | 31.36 |
| DAH-64 | 31.64 | 40.75 | 51.73 | 34.69 | 51.93 | 52.79 | 29.91 | 39.63 | 60.71 | 44.99 | 45.13 | 62.54 | 45.54 |
| DAH-128 | 32.58 | 40.64 | 52.40 | 35.72 | 52.80 | 52.12 | 30.94 | 41.31 | 59.31 | 45.65 | 46.67 | 64.97 | 46.26 |

## 9. Unsupervised Domain Adaptive Hashing: Additional Results

We provide the unsupervised domain adaptive hashing results for $d = 16$ and $d = 128$ bits in Figures 6 and 7 respectively. In Tables 6 and 7, we outline the corresponding mAP values. The notations are along the lines outlined in the main paper. We observe similar trends for both $d = 16$ and $d = 128$ bits compared to $d = 64$ bits. It is interesting to note that with increase in bit size $d$, the mAP does not necessarily increase. Table 7 ($d = 64$) has its mAP values lower than those for $d = 64$ (see main paper) for all the hashing methods. This indicates that merely increasing the hash code length does not always improve mAP scores. Also, the mAP values for Real-World for $d = 128$ bits has DAH performing better than SuH. This indicates that in some cases domain adaptation helps in learning a better generalized model.

**Table 6:** Mean average precision @16 bits. For the NoDA and DAH results, `Art` is the source domain for `Clipart`, `Product` and `Real-World` and `Clipart` is the source domain for `Art`.

| Expt. | NoDA | ITQ | KMeans | BA | BDNN | DAH | SuH |
|---|---|---|---|---|---|---|---|
| Art | 0.102 | 0.147 | 0.133 | 0.131 | 0.151 | 0.207 | 0.381 |
| Clipart | 0.110 | 0.120 | 0.116 | 0.123 | 0.138 | 0.211 | 0.412 |
| Product | 0.134 | 0.253 | 0.241 | 0.253 | 0.313 | 0.257 | 0.459 |
| Real-World | 0.193 | 0.225 | 0.195 | 0.216 | 0.248 | 0.371 | 0.400 |
| Avg. | 0.135 | 0.186 | 0.171 | 0.181 | 0.212 | 0.262 | 0.413 |

**Table 7:** Mean average precision @128 bits. For the NoDA and DAH results, `Art` is the source domain for `Clipart`, `Product` and `Real-World` and `Clipart` is the source domain for `Art`.

| Expt. | NoDA | ITQ | KMeans | BA | BDNN | DAH | SuH |
|---|---|---|---|---|---|---|---|
| Art | 0.154 | 0.202 | 0.175 | 0.148 | 0.207 | 0.314 | 0.444 |
| Clipart | 0.186 | 0.210 | 0.196 | 0.187 | 0.213 | 0.350 | 0.346 |
| Product | 0.279 | 0.416 | 0.356 | 0.336 | 0.432 | 0.424 | 0.792 |
| Real-World | 0.308 | 0.343 | 0.289 | 0.258 | 0.348 | 0.544 | 0.458 |
| Avg. | 0.232 | 0.293 | 0.254 | 0.232 | 0.300 | 0.408 | 0.510 |



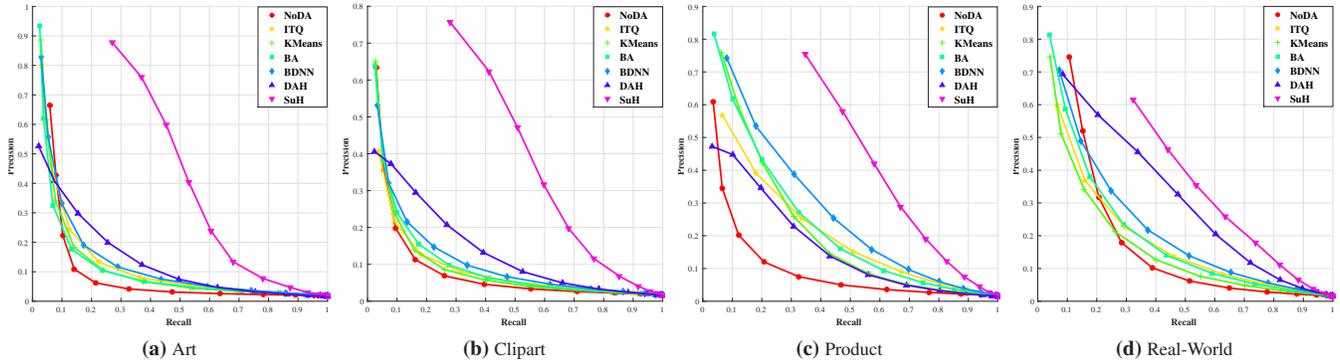**(a)** Art  **(b)** Clipart  **(c)** Product  **(d)** Real-World

**Figure 6:** Precision-Recall curves @16 bits for the ***Office-Home*** dataset. Comparison of hashing without domain adaptation (**NoDA**), shallow unsupervised hashing (**ITQ**, **KMeans**), state-of-the-art deep unsupervised hashing (**BA**, **BDNN**), unsupervised domain adaptive hashing (**DAH**) and supervised hashing (**SuH**). Best viewed in color.



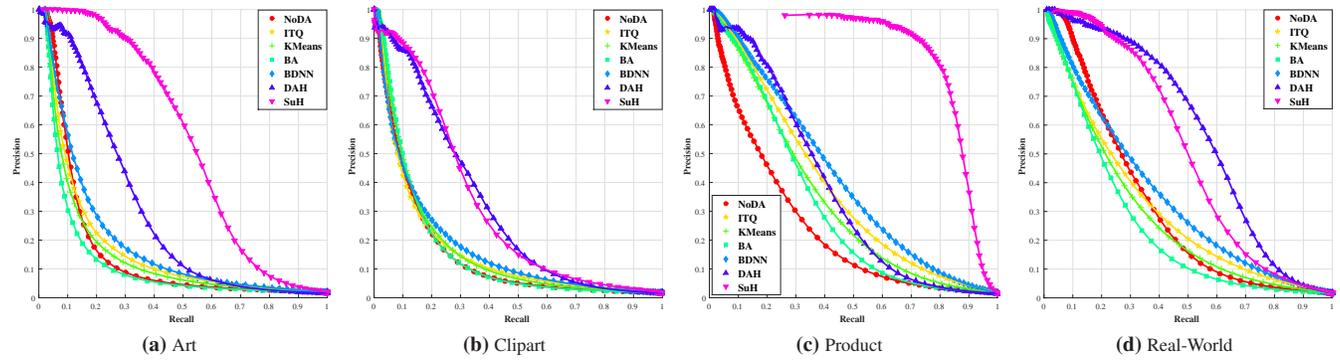**(a)** Art  **(b)** Clipart  **(c)** Product  **(d)** Real-World

**Figure 7:** Precision-Recall curves @128 bits for the ***Office-Home*** dataset. Comparison of hashing without domain adaptation (**NoDA**), shallow unsupervised hashing (**ITQ**, **KMeans**), state-of-the-art deep unsupervised hashing (**BA**, **BDNN**), unsupervised domain adaptive hashing (**DAH**) and supervised hashing (**SuH**). Best viewed in color.