

# Towards Interpretable Deep Metric Learning with Structural Matching

Wenliang Zhao<sup>1,2,3\*</sup>, Yongming Rao<sup>1,2,3\*</sup>, Ziyi Wang<sup>1,2,3</sup>, Jiwen Lu<sup>1,2,3†</sup>, Jie Zhou<sup>1,2,3</sup>

<sup>1</sup>Department of Automation, Tsinghua University, China

<sup>2</sup>State Key Lab of Intelligent Technologies and Systems, China

<sup>3</sup>Beijing National Research Center for Information Science and Technology, China

zhaowl20@mails.tsinghua.edu.cn; raoyongming95@gmail.com;

wziyi20@mails.tsinghua.edu.cn; {lujiwen, jzhou}@tsinghua.edu.cn

## Abstract

How do the neural networks distinguish two images? It is of critical importance to understand the matching mechanism of deep models for developing reliable intelligent systems for many risky visual applications such as surveillance and access control. However, most existing deep metric learning methods match the images by comparing feature vectors, which ignores the spatial structure of images and thus lacks interpretability. In this paper, we present a deep interpretable metric learning (DIML) method for more transparent embedding learning. Unlike conventional metric learning methods based on feature vector comparison, we propose a structural matching strategy that explicitly aligns the spatial embeddings by computing an optimal matching flow between feature maps of the two images. Our method enables deep models to learn metrics in a more human-friendly way, where the similarity of two images can be decomposed to several part-wise similarities and their contributions to the overall similarity. Our method is model-agnostic, which can be applied to off-the-shelf backbone networks and metric learning methods. We evaluate our method on three major benchmarks of deep metric learning including CUB200-2011, Cars196, and Stanford Online Products, and achieve substantial improvements over popular metric learning methods with better interpretability. Code is available at <https://github.com/wl-zhao/DIML>.

## 1. Introduction

Visual similarity plays an important role in a range of vision tasks including image retrieval [33], person identification [4] and image clustering [30]. Recent advances in learning visual similarity are mostly driven by Deep Metric Learning (DML), which leverages deep neural networks to

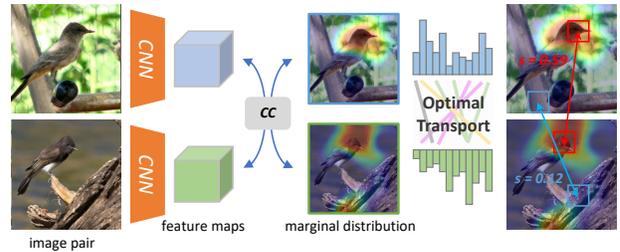


Figure 1: The main idea of the proposed deep interpretable metric learning (DIML) method. Unlike most existing deep metric learning methods match the images by comparing feature vectors, we propose a structural matching strategy that explicitly aligns the spatial embeddings by computing an optimal matching flow between feature maps of the two images to improve the interpretability of visual similarity.

learn an embedding space where the embedding similarity in this space can meaningfully reflect the semantic similarity between samples. A variety of deep metric learning methods have been proposed and have shown strong superiority in learning accurate and generalizable visual similarities on various tasks [7, 42, 16]. Despite the great progress in learning discriminative embeddings, deep metric learning methods with better interpretability have drawn limited attention from the community. Understanding the underlying matching mechanism of deep metric learning models is of critical importance for developing reliable intelligent systems for many risky visual applications such as surveillance [34] and access control [21].

To improve the transparency of deep visual models, many efforts have been made recently by either explaining the existing models [52, 31, 1, 2] or modifying models to achieve better interpretability [48, 49]. For example, visual attribution methods leverage correlation or gradient to find the important regions that have high contributions to the final prediction. [48] and [49] propose to add part constraints

\*Equal contribution.

†Corresponding author.

and tree structures to construct interpretable CNN models respectively. However, these methods are only designed for explaining the reasoning process of how the output of a deep model is produced and did not consider the interaction between samples. Although they achieve promising results on image classification [52, 2], visual question answering [31], and image generation [1], they cannot explain how visual similarity is composed. Therefore, how to improve the interpretability of deep metric learning methods is still an open problem that has barely been visited in previous works.

In this paper, we present a deep interpretable metric learning (DIML) framework as a first step towards more transparent embedding learning. Different from most existing deep metric learning methods that match the images by directly comparing feature vectors, we propose to leverage the spatial structure of images during matching to improve interpretability, as illustrated in Figure 1. More specifically, we measure the similarity of two images by computing an optimal matching flow between the feature maps using the optimal transport theory such that the similarity can be decomposed into several part-wise similarities with different contributions to the overall similarity. Our framework consists of three key components: 1) **Structural Similarity (SS)**. Unlike most existing deep metric learning methods that match the images by comparing feature vectors, we propose a new similarity/distance metric by measuring the similarity of corresponding parts in the feature maps based on the optimal matching flow; 2) **Spatial Cross-Correlation (CC)**. To handle the view variance in the image retrieval problem, we propose to use spatial cross-correlation as the initial marginal distribution to compute the optimal transport plan; 3) **Multi-scale Matching (MM)**. We also devise a multi-scale matching strategy to better incorporate existing metric learning methods and enable us to adaptively adjust the extra computational cost in large-scale search problems. Since our method is model-agnostic and our contribution is orthogonal to previous deep metric learning methods on architectures [14], objective functions [33, 16] and sampling strategies [44, 51], our method can be applied to off-the-shelf backbone networks and metric learning methods even without training. Extensive experimental study on three major benchmarks of deep metric learning including CUB200-2011 [40], Cars196 [17] and Stanford Online Products (SOP) [24] shows that our method enables us to achieve more interpretable metric learning while substantially improving various metric learning methods with or without re-training the models.

## 2. Related Work

**Deep Metric Learning.** Deep metric learning (DML) has drawn increasing attention recently and become one of the primary framework for a range of vision tasks including

image retrieval [33, 16], image clustering [30], person re-identification [4, 27, 3] and face recognition [37, 7, 26]. Previous works on deep metric learning commonly focus on learning more accurate and robust embeddings to better reflect the semantic relations among samples. To achieve this goal, a variety of deep metric learning approaches are proposed to improve the architectures [45, 14], objective functions [10, 30, 5, 24, 33, 16] and sampling strategies [44, 9, 20, 51, 28]. Different from these works, there is a line of deep metric learning research on developing more effective distance or similarity metrics. Except for the commonly used  $\ell^p$  distance and cosine similarity, signal-to-noise ratio (SNR) [46] and hyperbolic geodesic distance [15] have also proven to be effective to reflect the semantic relationships among samples. However, these deep metric learning methods only consider the distance or similarity between feature vectors, which ignores the spatial structure of images and thus lacks interpretability. In this work, we propose to measure the similarity of two images by explicitly leveraging the spatial structures of images such that more accurate and interpretable similarity of two samples can be obtained.

**Explainable & Interpretable Vision Models.** Recent years have witnessed remarkable progress in various computer vision tasks driven by the success of deep learning [18, 12, 19]. Despite the impressive discriminative power, the interpretability is often viewed as an Achilles' heel of deep models. Improving the explainability and interpretability of deep models has attracted increasing attention in recent years. Existing works can be roughly divided into two groups: 1) explaining existing models through visualization and diagnosis of deep representations; 2) modifying deep models to learn disentangled and interpretable representations. For example, Zhou *et al.* [52] proposes a method named Class Activation Mapping (CAM), which identifies discriminative regions in feature maps of CNNs by analyzing the effects on the final classification results. Grad-CAM [31] improves the method by combining both the input features and the gradients of a model's layer. Apart from these methods focusing on explaining and analyzing trained models, interpretable vision models are developed by revising the architectures or training procedure of conventional deep models. Zhang *et al.* [48] design interpretable CNNs by enforcing each filter in a high-level convolutional layer represents a specific object part. [49] combine the CNNs and decision tree to inherit the advantages of the two types of models to construct power yet interpretable image classification models. However, these methods only explain the reasoning process of how the output of a deep model is produced and did not consider the interaction between samples. Therefore, they cannot analyze and explain how the similarity of the two samples is composed. Recently, Williford *et al.* [43] present a study on explainable face recognition, which uses image editing techniques to gener-

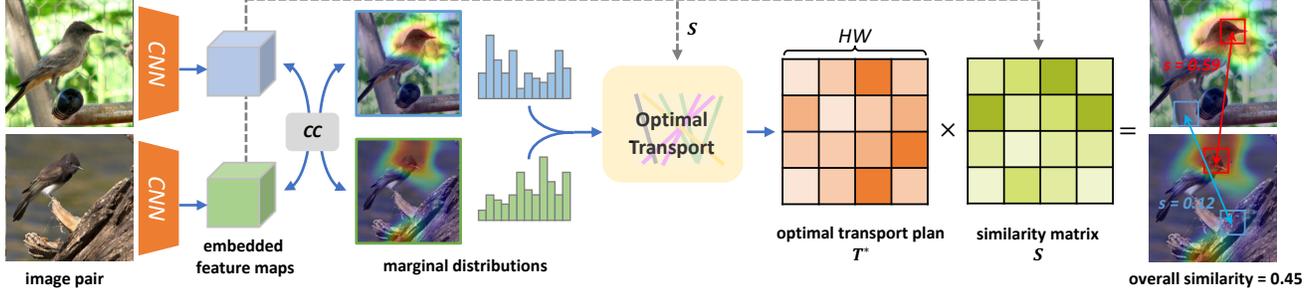


Figure 2: The overall pipeline of our deep interpretable metric learning (DIML) framework. The feature maps extracted from the backbone CNN model are further fed into the cross-correlation module (CC) to compute the marginal distributions that represent the weights of each location. The optimal transport plan then is obtained using the marginal distributions and the similarity matrix. Our framework decomposes the visual similarity to part-wise similarities and their contributions, which enable us to interpret and analyze how a deep model distinguishes two images.

ate a new dataset to evaluate what regions contribute to face matching. Their benchmark requires prior knowledge on face structures and thus is hard to generalize to other image matching problems. Different from these works, we propose to study a new and more generic problem of interpretable deep metric learning and provide a basic solution.

### 3. Approach

#### 3.1. Preliminaries: Deep Metric Learning

Deep metric learning aims to find a distance metric parameterized by deep neural networks to map the input image feature pairs to a distance in  $\mathbb{R}$  that reflects the semantic similarity of the two images defined by labels. Formally, given a set of images  $\mathcal{X} = \{x^k\}_{k=1}^N$  and the corresponding labels  $\mathcal{Y} = \{y^k\}_{k=1}^N$ , deep metric learning introduces the deep neural networks  $f: \mathcal{X} \rightarrow \Phi \subset \mathbb{R}^C$  to map an image to a feature  $\phi^k = f(x^k)$ , where the semantic patterns of the input image are extracted. The mainstreams of deep metric learning aim to learn Mahalanobis distance metrics  $d(\cdot, \cdot)$ , which can be formulated as:

$$d(x^k, x^l) = \|Mf(x^k) - Mf(x^l)\|_2 = \|g(\phi^k) - g(\phi^l)\|_2,$$

where  $g(\phi^k) = M\phi^k := \psi^k \in \Psi$  is a parametrized linear projection from the feature space  $\Phi$  to an embedding space  $\Psi \subset \mathbb{R}^D$ . Following the configuration in the backbone networks like ResNet [12] and Inception [35],  $f$  can be decomposed into  $f = \text{GAP} \circ f_1$ , where  $f_1$  extracts a feature map  $\omega^k = f_1(x^k) \in \mathbb{R}^{H \times W \times C}$  and GAP is the global average pooling. The GAP operation abstracts the feature maps into vectors so as to enable fast similarity calculation.

However, the abstraction on deep features also loses the spatial structures of the images during the embedding process, which makes most deep metric learning methods lack interpretability—deep models can tell us whether the two images are similar but cannot show us the reason. Since it is of

importance to understand the matching mechanism in many risky visual applications, developing a more interpretable deep metric learning method becomes a critical research topic but it has barely been visited in previous works.

#### 3.2. Structural Matching via Optimal Transport

To exploit the spatial structures in images for more interpretable deep metric learning, we devise a new structural matching scheme to compute feature similarity based on optimal transport theory [39].

Our core algorithm is adopted from the optimal transport theory, which aims to seek the *minimal cost transport plan* between two distributions. Given a source distribution  $\mu^s$  and a target distribution  $\mu^t$  that are defined on probability space  $\mathcal{U}$  and  $\mathcal{V}$  respectively, the minimal cost transport plan can be obtained by minimizing the Wasserstein distance between the two distributions:

$$\pi^* = \arg \inf_{\pi \in \Pi(\mu^s, \mu^t)} \int_{\mathcal{U} \times \mathcal{V}} c(u, v) d\pi(u, v), \quad (1)$$

where  $\pi^*$  is the optimal transport plan,  $\Pi(\mu^s, \mu^t)$  is the joint probability distribution with marginals  $\mu^s$  and  $\mu^t$ , and  $c: \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}^+$  is the cost function of transportation.

Different from the above generic formulation, here we only need to consider the discrete distribution matching for image feature maps. Consider two feature maps  $\omega^s, \omega^t \in \mathbb{R}^{H \times W \times C}$  obtained by a backbone (e.g. ResNet50 [12]). We first use the projection layer  $g$  to map each element in the feature maps  $\omega_i^k$  into an embedding space of dimension  $D$  individually<sup>1</sup>:

$$z_i^s = g(\omega_i^s) \in \mathbb{R}^D, \quad z_j^t = g(\omega_j^t) \in \mathbb{R}^D. \quad (2)$$

<sup>1</sup>For the sake of simplicity, we use a single subscript  $i \in [1, HW]$  to index the spatial location. For pre-trained metric learning models, we can directly apply the original projection layer on the elements in the feature maps. Thus, our method does not need any modifications on the parameters.

The cost of transporting one unit of mass from  $i$  to  $j$  is:

$$C_{i,j} = c(i, j) := d(z_i^s, z_j^t), \quad (3)$$

where we use the distance metric  $d(\cdot, \cdot)$  for two vectors (e.g., Euclid distance or cosine distance) as the transport cost function  $c$ . In this discrete case, the transport plan  $\pi$  matching the two distributions also becomes discrete. Given the two corresponding discrete distributions  $\mu^s$  and  $\mu^t$ , the original optimal transport problem is equivalent to:

$$\begin{aligned} T^* &= \arg \min_{T \geq 0} \text{tr}(CT^\top), \\ \text{subject to } T\mathbf{1} &= \mu^s, \quad T^\top \mathbf{1} = \mu^t. \end{aligned} \quad (4)$$

$T^*$  is the optimal matching flow between these two distributions, which can be also viewed as the structural matching plan of the two images.  $T_{i,j}^*$  is the amount of mass that needs to move from  $i$  to  $j$  in order to reach an overall minimum cost, which represents the contribution of location pair  $(i, j)$  to the overall matching.

To efficiently solve the optimization problem in (4), we adopt the Sinkhorn divergence algorithm [6] by introducing an entropic regularizer to enable fast training and inference. More details about the algorithm can be found in Supplementary Material. Note that the iterative algorithm is fully differentiable, which can be easily implemented by using the automatic differentiation library like PyTorch [25] and directly apply the matching process to any deep metric learning pipelines.

**Discussions.** Some closely related works of the proposed structural matching scheme include EMD metric learning [50] and Wasserstein embedding learning [8]. However, different from our method, they usually focus on learning better embeddings for set inputs, which can be naturally solved by the Wasserstein distance metric learning framework. Here our main contribution is not the matching algorithm itself but the introduction of structural matching for learning more interpretable visual similarity.

### 3.3. Deep Interpretable Metric Learning

In Section 3.2, we have already shown how to calculate the distance between two distributions using the optimal transport. In this section, we describe how to perform structural matching in metric learning. Specifically, our method consists of three components: 1) we use optimal transport to calculate the *structural similarity* (SS) of two images; 2) we propose to calculate the spatial *cross-correlation* (CC) to initialize the marginal distributions in Equation (1); 3) we propose *multi-scale matching* (MM) to improve the metric and reduce the computation cost.

**Structural Similarity (SS).** Given the marginal distributions  $\mu^s$  and  $\mu^t$  (which we will discuss in detail later) and the cost matrix  $C$ , we can then obtain the optimal transport

$T^*$  by solving (4). Once we have  $T^*$ , we can define the structural similarity of two feature maps  $z^s, z^t \in \mathbb{R}^{HW \times D}$  as follows:

$$s_{\text{struct}}(z^s, z^t) = \sum_{1 \leq i, j \leq HW} s(z_i^s, z_j^t) T_{i,j}^*, \quad (5)$$

where  $s(\cdot, \cdot)$  is a function to calculate the similarity between two vectors. Our structural similarity enables us to investigate the composition of the overall similarity, thus we can easily decompose the similarity and understand how the similarity between different locations in the two images contribute to the overall similarity. Similarly, given any distance function  $d(\cdot, \cdot)$ , we can also derive our structural distance:

$$d_{\text{struct}}(z^s, z^t) = \sum_{1 \leq i, j \leq HW} d(z_i^s, z_j^t) T_{i,j}^*, \quad (6)$$

**Cross-Correlation (CC).** Another important part is the definition of the marginal distributions  $\mu^s$  and  $\mu^t$ . One trivial solution is to initialize them with uniform distributions, i.e.,

$$\mu_i^s = \mu_i^t = \frac{1}{HW}, \forall 1 \leq i \leq HW, \quad (7)$$

which indicates similarity of each location has the identical weight to the overall similarity. In the structural matching algorithm, the marginal distributions should characterize the importance of each spatial location. Simply using uniform distributions implies that we want to match all the features with equal importance, which is not desired in some cases. For example, some image contains background information that may be less useful for matching thus we want to impose lower weights on the background. Another common circumstance is when we want to match two images with different views (e.g., the first image contains the whole object and the second one only contains a part of it), and similarly we only need to focus on the certain part of the first image and treat the rest as background. To find the areas that are most related to the similarity, we propose to calculate the cross-correlation between the two images as the marginal distributions for the matching algorithm. Specifically, we first perform global average pooling to  $z^s, z^t$  and obtain the global feature  $\bar{z}^s, \bar{z}^t$ . We then slide the global feature of one image on the feature map of the other image and calculate point-wise correlation at each spatial location. Formally, the cross-correlation is calculated as:

$$\alpha_i^s = \frac{\langle \bar{z}^s, z_i^t \rangle}{\|\bar{z}^s\| \|z_i^t\|}, \alpha_i^t = \frac{\langle \bar{z}^t, z_i^s \rangle}{\|\bar{z}^t\| \|z_i^s\|}, \quad (8)$$

where  $\langle \cdot, \cdot \rangle$  is the dot product and  $\alpha_i^k \in [-1, 1]$ . After obtaining the cross-correlation, we can use  $\alpha_i^k$  to reflect the importance of  $z_i^k$  in the matching problem. To further reduce the effects of low correlation regions, we discard

the negative value of  $\alpha_i^k$  and normalize it to obtain the final marginal distributions:

$$\gamma_i^k = \max(0, \alpha_i^k), \mu_i^k = \frac{1}{\sum_{i'} \gamma_{i'}^k} \gamma_i^k \quad (9)$$

$$\forall 1 \leq i \leq HW, \quad k \in \{s, t\}.$$

Once we have the marginal distributions  $\mu^{(k)}$ , we can then apply the structural matching algorithm in Section 3.2 to calculate the similarity between two images. We will show in Section 4.3 that cross-correlation is an indispensable component to improve the power of DIML.

**Multi-scale Matching (MM).** Although DIML can capture the structural similarity of two images and can provide results easily understood by humans, it requires more computation ( $\mathcal{O}(H^2W^2)$ ) to solve the optimal transport problem. In the application of image retrieval, there are usually a great number of images in the gallery. Given an image as an anchor, calculate the structural similarity between the anchor and all the images in the gallery is inefficient. To reduce the computational cost, we propose a multi-scale matching method for image retrieval. Let  $z^a \in \mathbb{R}^{H \times W \times D}$  be the feature map of the anchor image and  $z^k \in \mathbb{R}^{H \times W \times D}$ ,  $k = 1, \dots, N$  be the feature maps of all the images in the gallery. In the first scale ( $1 \times 1$ ), we compute the global feature using global average pooling to get  $\bar{z}^a, \bar{z}^k \in \mathbb{R}^D$ , and calculate cosine similarity between  $\bar{z}^a$  and each  $\bar{z}^k$  as conventional DML methods. We can then define a truncation number  $K$  and select the images with top- $K$  similarity score and denote the indices of them as  $\mathcal{I}_K$  to further enhance the similarity with our method. In the second scale ( $H \times W$ ), we calculate the structural similarity between  $z^a$  and each  $z^k$ ,  $k \in \mathcal{I}_K$ . Since  $K$  is fixed, the extra computational cost of DIML can be controlled. By multi-scale matching, we can filter out the obvious dissimilar samples ( $1 \times 1$  scale, cosine similarity) and focus on the hard ones ( $H \times W$  scale, structural similarity). Combining the similarity at two scales can also capture both semantic and spatial information, which is also helpful to improve retrieval precision. We will show later in Section 4.3 that a small  $K$  can yield a significant performance boost.

### 3.4. Implementation

One of the major advantages of DIML is that we can apply DIML to any pre-train model to improve performance with *no need of training*. Besides, we can also incorporate DIML into the training objective. In this section, we will describe how to use DIML in these two scenarios.

**Testing.** Given a pre-trained model, we first calculate the feature maps  $\omega^s, \omega^t \in \mathbb{R}^{H \times W \times C}$  (before the global pooling layer) of the image pair  $x^s, x^t$ . We can then use the algorithm describe in Section 3.3 to compute the structural similarity. However,  $HW$  may be sometimes large in practice (e.g., for

ResNet50 [12],  $H = W = 7$ ). Therefore, we can use ROI Align [11] to pool the feature maps to  $\mathbb{R}^{H' \times W' \times C}$ , where  $H' < H$  and  $W' < W$ . With smaller feature maps, we can then calculate the structural similarity with a relatively lower computational cost. In our implementation, we use  $H' = W' = G$  and we found  $G = 4$  can achieve good trade-off between cost and performance.

**Training.** We can also combining DIML and existing metric learning methods to facilitate training. We now use Margin loss [44] as an example to show how to incorporate DIML into the training objective. The Margin loss [44] is defined as

$$\mathcal{L}_{\text{margin}}(k, l) = \left( \sigma + (-1)^{I(y_k \neq y_l)} (D_{k,l} - \beta) \right)_+, \quad (10)$$

where  $\sigma$  and  $\beta$  are learnable parameters, and  $D_{kl}$  is used to measure the distance between image  $k$  and  $l$ :

$$D_{k,l} = \frac{1}{2} (d_{\text{struct}}(z^k, z^l) + d(\bar{z}^k, \bar{z}^l)) \quad (11)$$

For the implementation details of other metric learning methods, please refer to the Supplementary Material.

## 4. Experiments

To evaluate the performance of our proposed DIML, we conduct experiments on three widely used datasets in the image retrieval research field: CUB200-2011 [40], Cars196 [17], and Standard Online Products (SOP) [24].

### 4.1. Experiment Setups

**Datasets.** We evaluate our method under a zero-shot image retrieval setting, where the training set and test set contain image classes with no intersection. We follow the training/test set splits in previous works [23, 29]:

- CUB200-2011 [40] contains 11,788 images of birds from 200 species. The first 100 classes (5,864 images) are used for training, while other 100 classes (5,924 images) are kept for testing.
- Cars196 [17] contains 16,185 images of cars from 196 classes. We use the first 98 classes (8,054 images) for training and leave the rest 98 classes (8,131 images) for testing.
- SOP [24] contains 120,053 images from 22,634 classes. We use the first 11,318 classes (59,551 images) for training and other 11,316 (60,502 images) for testing.

**A Fair Evaluation Protocol.** Although there are many previous metric learning methods, [23] pointed out that the improvements over time are not significant, due to the unfair comparisons of different methods. Therefore, we try our best

to provide fair results by implementing all the methods under the same evaluation protocol. For all the baseline methods, we use ResNet-50 [12] pre-trained on ImageNet [18] as the backbone. We freeze the BatchNorm layers during training and modify the output channel of the last linear layer to a fixed embedding dimension  $D$ . We use embedding size  $D = 128$  and other implementation settings following [29] for most experiments unless otherwise noted.

**Evaluation Metrics.** Most previous works use Recall@K, Normalized Mutual Information, and F1 score as accuracy metrics. However, as is pointed by [23], NMI and F1 scores sometimes give us wrong pictures of the embedding space. To this end, we adopt the evaluation metrics used in [23]: Precision@1, R-Precision, and MAP@R. For the formal definition of the metrics, see Supplementary Material.

**Implementation.** It is also worth noting that our proposed DIML does not require any training. Therefore, we aim to prove that our method can improve the performance given *any* trained model as the baseline. Therefore, we perform experiments on a wide range of loss functions (Margin [44], Arcface [7], *etc.*) and sampling methods (Distance [44], N-Pair [33], *etc.*) to prove the effectiveness and the generalization ability of our method. For most of the baseline methods, we follow the implementation from [29].

## 4.2. Main Results

We first evaluate our method by applying DIML to a wide range of metric learning methods. We measure the performance using the evaluation metrics aforementioned: Precision@1 (P@1), R-Precision (RP) and MAP@R (M@R), and the results<sup>2</sup> are shown in Table 1. For all the experiments, we set the truncation number  $K = 100$  and feature map size  $G = 4$ . We observe that our method can improve the performance for *all* the models on *all* the three benchmarks, without any extra training. Especially, we find on Cars196 dataset, the performances of all the methods are enhanced profoundly after equipped with our DIML.

## 4.3. Ablation Study and Analysis

In this section, we will evaluate our DIML in various settings and provide detailed analyses through experiments and visualization.

**Effects of different components.** The DIML consists of three components: structural similarity (SS), cross-correlation (CC), and multi-scale matching (MM). We will analyze the effect of each one, as is shown in Table 2. We start with two baseline methods Margin [44] and Multi-Similarity [42], and add the three components gradually. First, we adopt structural similarity instead of standard cosine similarity (where we use uniform distribution in Equation (7) for  $\mu^s$  and  $\mu^t$ ). We find that SS can improve the

<sup>2</sup>For more results, please refer to the Supplementary Material.

Table 1: **Applying DIML to various deep metric learning methods.** Experimental results show that our method can improve the performance of all the methods consistently.

| Method           | CUB200-2011  |              |              | Cars196      |              |              | SOP          |              |              |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                  | P@1          | RP           | M@R          | P@1          | RP           | M@R          | P@1          | RP           | M@R          |
| Contrastive [10] | 61.77        | 34.25        | 23.24        | 67.45        | 30.01        | 18.61        | 73.27        | 40.92        | 37.5         |
| + DIML           | <b>64.43</b> | <b>35.16</b> | <b>24.29</b> | <b>73.35</b> | <b>30.76</b> | <b>20.13</b> | <b>74.47</b> | <b>41.58</b> | <b>38.29</b> |
| Triplet-R [30]   | 58.34        | 32.00        | 20.93        | 62.73        | 26.95        | 15.24        | 66.60        | 33.62        | 30.26        |
| + DIML           | <b>60.60</b> | <b>32.63</b> | <b>21.74</b> | <b>67.92</b> | <b>27.65</b> | <b>16.72</b> | <b>68.73</b> | <b>35.04</b> | <b>31.79</b> |
| Triplet-S [30]   | 59.28        | 32.77        | 21.79        | 67.00        | 30.0         | 18.23        | 73.67        | 40.45        | 37.14        |
| + DIML           | <b>62.85</b> | <b>33.87</b> | <b>23.04</b> | <b>72.06</b> | <b>30.89</b> | <b>20.04</b> | <b>75.14</b> | <b>41.68</b> | <b>38.42</b> |
| Triplet-H [29]   | 61.39        | 33.21        | 22.15        | 71.76        | 32.53        | 20.83        | 73.28        | 39.98        | 36.56        |
| + DIML           | <b>62.02</b> | <b>33.50</b> | <b>22.49</b> | <b>74.75</b> | <b>32.94</b> | <b>21.76</b> | <b>73.62</b> | <b>40.14</b> | <b>36.79</b> |
| Triplet-D [44]   | 61.99        | 33.92        | 22.90        | 73.07        | 32.18        | 20.81        | 77.34        | 44.25        | 40.80        |
| + DIML           | <b>63.40</b> | <b>34.49</b> | <b>23.59</b> | <b>77.31</b> | <b>33.05</b> | <b>22.61</b> | <b>77.81</b> | <b>44.82</b> | <b>41.39</b> |
| NPair [33]       | 60.30        | 33.53        | 22.27        | 69.52        | 32.24        | 20.25        | 76.47        | 43.48        | 39.94        |
| + DIML           | <b>62.17</b> | <b>34.02</b> | <b>22.85</b> | <b>74.65</b> | <b>32.91</b> | <b>21.67</b> | <b>76.86</b> | <b>43.87</b> | <b>40.38</b> |
| Angular [41]     | 61.36        | 34.17        | 23.00        | 70.93        | 32.97        | 21.31        | 73.79        | 41.42        | 37.90        |
| + DIML           | <b>63.77</b> | <b>35.09</b> | <b>24.06</b> | <b>74.72</b> | <b>33.80</b> | <b>22.83</b> | <b>74.91</b> | <b>42.17</b> | <b>38.73</b> |
| GenLifted [13]   | 58.27        | 32.86        | 21.83        | 66.88        | 30.96        | 19.00        | 74.84        | 42.28        | 38.66        |
| + DIML           | <b>61.07</b> | <b>33.82</b> | <b>22.98</b> | <b>72.95</b> | <b>31.93</b> | <b>20.88</b> | <b>75.92</b> | <b>43.08</b> | <b>39.55</b> |
| ProxyNCA [22]    | 62.76        | 35.05        | 24.03        | 71.05        | 31.62        | 20.55        | 74.70        | 41.32        | 37.96        |
| + DIML           | <b>64.75</b> | <b>36.02</b> | <b>25.10</b> | <b>74.86</b> | <b>32.43</b> | <b>22.00</b> | <b>76.17</b> | <b>42.65</b> | <b>39.36</b> |
| Histogram [38]   | 59.96        | 33.07        | 22.15        | 69.49        | 31.62        | 19.76        | 71.15        | 38.06        | 34.70        |
| + DIML           | <b>62.69</b> | <b>33.80</b> | <b>23.00</b> | <b>74.50</b> | <b>32.36</b> | <b>21.26</b> | <b>72.06</b> | <b>38.57</b> | <b>35.30</b> |
| Quadruplet [5]   | 61.53        | 34.05        | 22.93        | 69.64        | 31.40        | 19.67        | 77.02        | 44.27        | 40.88        |
| + DIML           | <b>62.80</b> | <b>34.65</b> | <b>23.62</b> | <b>75.66</b> | <b>32.35</b> | <b>21.69</b> | <b>78.08</b> | <b>45.16</b> | <b>41.79</b> |
| SNR [46]         | 62.00        | 34.72        | 23.59        | 72.95        | 32.72        | 21.28        | 77.82        | 44.98        | 41.51        |
| + DIML           | <b>64.55</b> | <b>35.25</b> | <b>24.27</b> | <b>77.57</b> | <b>33.54</b> | <b>23.02</b> | <b>78.50</b> | <b>45.65</b> | <b>42.24</b> |
| Softmax [47]     | 61.06        | 32.7         | 21.55        | 72.61        | 31.17        | 19.88        | 77.02        | 43.47        | 40.25        |
| + DIML           | <b>63.30</b> | <b>33.71</b> | <b>22.64</b> | <b>76.39</b> | <b>32.06</b> | <b>21.49</b> | <b>78.17</b> | <b>44.62</b> | <b>41.43</b> |
| Margin [44]      | 62.47        | 34.12        | 23.14        | 72.18        | 32.00        | 20.82        | 78.39        | 45.64        | 42.34        |
| + DIML           | <b>65.16</b> | <b>35.37</b> | <b>24.51</b> | <b>76.62</b> | <b>32.85</b> | <b>22.48</b> | <b>79.26</b> | <b>46.44</b> | <b>43.19</b> |
| Arcface [7]      | 61.39        | 33.70        | 22.4         | 73.37        | 31.90        | 20.52        | 77.55        | 44.44        | 41.07        |
| + DIML           | <b>64.72</b> | <b>34.88</b> | <b>23.72</b> | <b>77.24</b> | <b>32.88</b> | <b>22.34</b> | <b>78.52</b> | <b>45.45</b> | <b>42.10</b> |
| MS [42]          | 62.56        | 32.74        | 21.99        | 74.81        | 32.72        | 21.60        | 77.90        | 44.97        | 41.54        |
| + DIML           | <b>64.89</b> | <b>33.99</b> | <b>23.34</b> | <b>78.44</b> | <b>33.57</b> | <b>23.31</b> | <b>78.53</b> | <b>45.59</b> | <b>42.22</b> |
| ProxyAnchor [16] | 65.24        | 35.81        | 24.87        | 82.36        | 36.00        | 25.85        | 79.10        | 46.31        | 42.91        |
| + DIML           | <b>66.46</b> | <b>36.49</b> | <b>25.58</b> | <b>86.13</b> | <b>37.90</b> | <b>28.11</b> | <b>79.22</b> | <b>46.43</b> | <b>43.04</b> |

performance on all the datasets *except* for SOP (as is highlighted by underline). It is mainly because that the SS algorithm aims to match every part of a source image to a target image. However, the views vary a lot in SOP dataset, which hinders the application of SS. Second, we show that multi-scale matching can make use of the semantic information and improve the performance on all three datasets. Finally, we replace the uniform distribution with the one calculated by cross-correlation. We find the marginal distributions obtained in this way are helpful to explore the important area of the images and can further improve the performance.

**Effects of truncation number.** To show how the truncation number  $K$  affect our DIML, we test our method on Margin [44] and Multi-Similarity [42] with  $K$  increasing from 0 to 500 (Figure 3). Note that  $K = 0$  means no structural similarity is used, which is identical to the baseline. We find that even a small  $K$  will bring considerable improvement on the performance (especially for the P@1 metric). Generally, the retrieval accuracy grows with  $K$  increasing and saturates

Table 2: **Effects of the three components in our DIML:** Structural Similarity (SS), Multi-scale Matching (MM) and Cross Correlation (CC). We show that our method can enhance the performance of the baseline methods by combining the three components together.

| Baseline    | Components |    |    | CUB200-2011  |              | Cars196      |              | SOP          |              |              |              |              |
|-------------|------------|----|----|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|             | SS         | MM | CC | P@1          | M@R          | P@1          | M@R          | P@1          | M@R          |              |              |              |
| Margin [44] | ✓          |    |    | 62.47        | 23.14        | 72.18        | 20.82        | 78.39        | 42.34        |              |              |              |
|             | ✓          | ✓  |    | 63.64        | 22.52        | 74.86        | 21.24        | 77.30        | 41.02        |              |              |              |
|             | ✓          | ✓  | ✓  | <b>64.96</b> | <b>23.87</b> | <b>76.02</b> | <b>22.02</b> | <b>78.53</b> | <b>42.45</b> |              |              |              |
| MS [42]     | ✓          |    |    | 62.56        | 21.99        | 74.81        | 21.60        | 77.90        | 41.54        |              |              |              |
|             | ✓          | ✓  |    | 63.52        | 21.66        | 75.63        | 21.10        | 75.81        | 39.38        |              |              |              |
|             | ✓          | ✓  | ✓  | <b>64.40</b> | <b>22.83</b> | <b>77.39</b> | <b>22.77</b> | <b>77.87</b> | <b>41.55</b> |              |              |              |
|             |            |    |    |              |              |              | <b>64.89</b> | <b>23.34</b> | <b>78.44</b> | <b>23.31</b> | <b>78.53</b> | <b>42.22</b> |

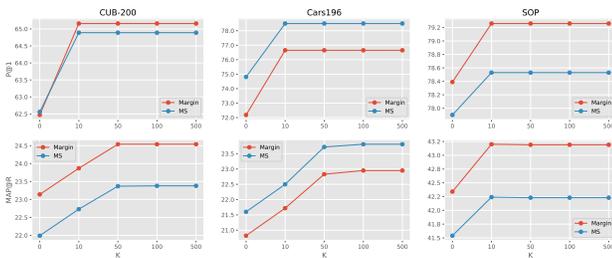


Figure 3: **Comparisons of different truncation number.** We test for different truncation number  $K$  ranging from 0 to 500. Experimental results show that a small  $K$  can already bring considerable performance improvement.

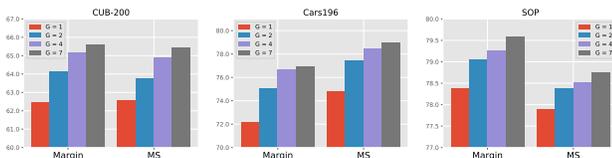


Figure 4: **Effects of the size of feature map.** Generally, the performance of our DIML is better with higher  $G$ . DIML with  $G = 4$  yields good results within relatively low computational costs.

before  $K$  reaches 100. This phenomenon indicates that with a fixed and relatively small  $K$ , we can already enjoy a significant performance boost with constant extra computational cost and no extra training cost.

**Effects of feature map size.** We then perform an ablation study on the feature map size  $G$ . In our experiments, we use ResNet50 [12] as our backbone, and the size of the feature map before the pooling layer is  $7 \times 7$ . Hence, we need to pool the feature map to a smaller one ( $G \times G$ ) to reduce computational complexity. Specifically, we let  $G \leq 7$  and evaluate for the cases where  $G = 1, 2, 4, 7$ . The results are shown in Table 4 and we observe that the performance of

Table 3: **Effects of training.** Our method can substantially improve the baseline model with or without training.

| Baseline    | Setting |       | CUB200-2011  |              | Cars196      |              | SOP          |              |
|-------------|---------|-------|--------------|--------------|--------------|--------------|--------------|--------------|
|             | test    | train | P@1          | M@R          | P@1          | M@R          | P@1          | M@R          |
| Margin [44] | ✓       |       | 62.47        | 23.14        | 72.18        | 20.82        | 78.39        | 42.34        |
|             | ✓       | ✓     | <b>65.16</b> | <b>24.54</b> | <b>76.65</b> | <b>22.95</b> | <b>79.26</b> | <b>43.19</b> |
| MS [42]     | ✓       |       | 62.56        | 21.99        | 74.81        | 21.60        | 77.90        | 41.54        |
|             | ✓       | ✓     | <b>64.89</b> | <b>23.38</b> | <b>78.50</b> | <b>23.81</b> | <b>78.53</b> | <b>42.23</b> |
|             |         |       | <b>65.36</b> | <b>24.90</b> | <b>75.61</b> | <b>22.34</b> | <b>78.81</b> | <b>42.89</b> |
|             |         |       | <b>65.72</b> | <b>24.37</b> | <b>78.90</b> | <b>23.80</b> | <b>79.00</b> | <b>42.96</b> |

Table 4: **Effects of embedding size.** Our DIML is robust to the changing of the embedding size  $D$  and can improve the performance of the baseline methods consistently.

| $D$ | Method             | CUB200-2011  |              |              | Cars196      |              |              |
|-----|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|
|     |                    | P@1          | RP           | M@R          | P@1          | RP           | M@R          |
| 64  | Margin             | 59.39        | 32.59        | 21.53        | 69.31        | 30.98        | 19.67        |
|     | Margin [44] + DIML | <b>62.98</b> | <b>33.88</b> | <b>22.95</b> | <b>74.44</b> | <b>31.96</b> | <b>21.60</b> |
|     | MS [42]            | 58.52        | 31.35        | 20.23        | 71.67        | 30.90        | 19.57        |
|     | MS + DIML          | <b>61.73</b> | <b>32.61</b> | <b>21.62</b> | <b>76.94</b> | <b>31.95</b> | <b>21.72</b> |
|     | ProxyAnchor [16]   | 62.56        | 34.61        | 23.45        | 78.08        | 34.35        | 23.95        |
|     | ProxyAnchor + DIML | <b>65.01</b> | <b>35.53</b> | <b>24.40</b> | <b>83.11</b> | <b>36.49</b> | <b>26.55</b> |
| 512 | Margin [44]        | 64.92        | 35.94        | 24.92        | 73.68        | 32.03        | 21.09        |
|     | Margin + DIML      | <b>66.91</b> | <b>36.82</b> | <b>25.89</b> | <b>76.67</b> | <b>32.62</b> | <b>22.20</b> |
|     | MS [42]            | 65.92        | 35.14        | 24.17        | 76.85        | 33.93        | 22.78        |
|     | MS + DIML          | <b>68.15</b> | <b>36.04</b> | <b>25.14</b> | <b>79.74</b> | <b>34.68</b> | <b>24.01</b> |
|     | ProxyAnchor[16]    | 67.30        | 37.40        | 26.38        | 84.75        | 37.56        | 27.66        |
|     | ProxyAnchor + DIML | <b>67.93</b> | <b>37.92</b> | <b>26.88</b> | <b>87.01</b> | <b>39.03</b> | <b>29.39</b> |

our method is better with larger  $G$  in general. We can also find  $G = 4$  is a good trade-off between performance and computational complexity.

**Effects of training.** Besides the default setting where we use DIML to test on any pre-trained model, we can also incorporate the structural similarity into the training objectives (see Section 3.4 for details). In Table 3, we compare the performance in three scenarios: (1) without DIML testing or training (same as baseline) (2) with DIML testing only (2) with DIML testing and training. The results are listed in Table 3. We find that for most cases, using DIML to test a pre-trained model can already improve the baseline by a significant margin. Besides, it is also useful sometimes to apply DIML in the training stage.

**Effects of embedding size.** Our proposed DIML is also robust across different embedding sizes. Apart from the results in Table 1 where  $D = 128$ , we perform experiments with  $D = 64/512$  for Margin [44], Multi-Similarity [42] and Proxy Anchor [16] and the results are summarized in Table 4. We demonstrate that our method can boost the performance of the three methods consistently no matter how the embedding size  $D$  varies.

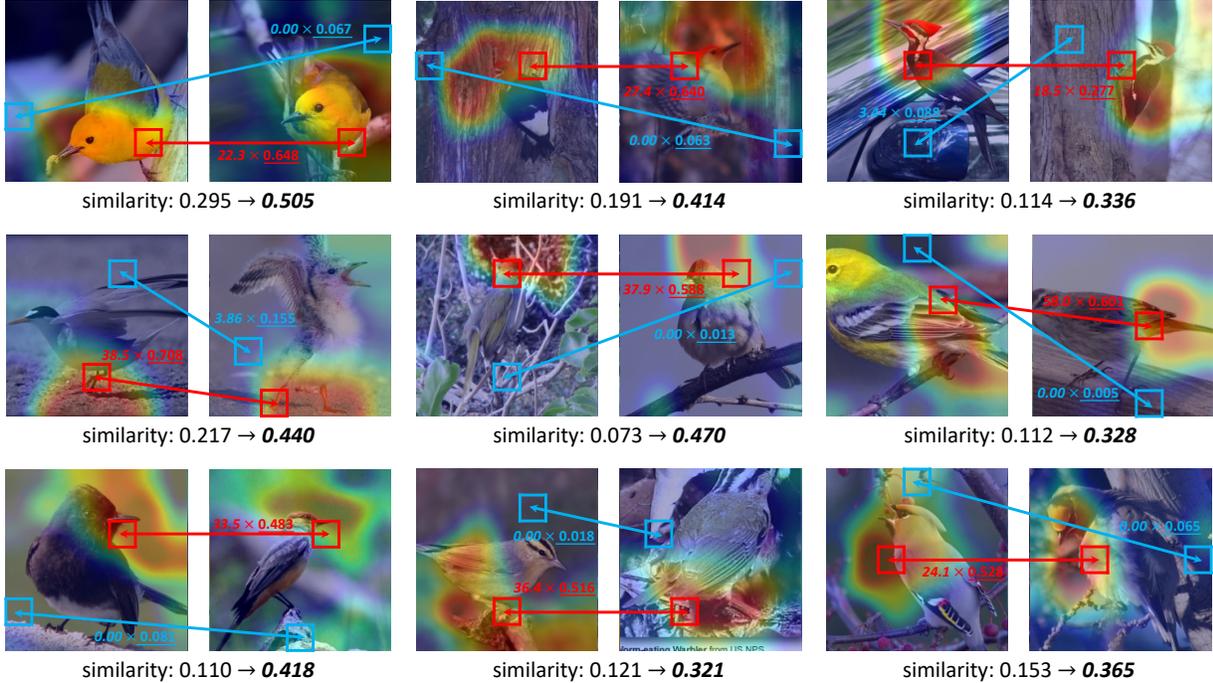


Figure 5: **Visualization.** We use heatmaps to show the marginal distributions obtained by cross-correlation (CC). We also illustrate two most representative part-wise similarity and their contributions to the overall similarity in the form of  $[G^4 T_{i,j}^*] \times [S_{i,j}]$ , where  $G$  is the grid size,  $T_{i,j}^*$  and  $S_{i,j}$  are the matching flow and similarity between location  $i$  and  $j$  respectively. We also show the overall similarity changes after applying our method to the baseline model (cosine similarity  $\rightarrow$  **structural similarity**). All image pairs are positive pairs.

#### 4.4. Visualization

To better understand how our method works, we provide some visualizations for CUB200-2011 [40] in Figure 5, where each pair of images is from the same category. First, we visualize the marginal distributions  $\mu^s$  and  $\mu^t$  (calculated by cross-correlation) through heatmaps and find that they can focus on some discriminative parts in the images (e.g., head, foot, etc.). Second, we show the optimal transport flow  $T_{i,j}^*$  and the similarity  $S_{i,j} = s(z_i^s, z_j^t)$  for some pair of spatial location  $(i, j)$ . Since the sum of the values in  $T^*$  equals 1 and each  $T_{i,j}^*$  is relatively small, we use a re-scaled version  $\hat{T}_{i,j}^* = G^4 T_{i,j}^*$  such that an uniform transport flow yields  $\hat{T}_{i,j}^* = 1, \forall i, j$ . We draw arrows between the location pairs that make a large (or small) contribution to the final structural similarity in red (or blue). The formula along with an arrow takes the form of  $\hat{T}_{i,j}^* \times S_{i,j}$ . We observe that our method can match similar parts and assign a higher  $T_{i,j}^*$  to the pair while enforcing lower  $T_{i,j}^*$  to the parts that are less informative to determine the similarity between the two images. Finally, we demonstrate that by re-weighting the similarity matrix  $S$  with the optimal transport matrix  $T^*$ , our proposed structural similarity (shown in **bold** text) is higher than the standard cosine similarity (shown in *light* text) by a large margin.

#### 5. Conclusion

In this paper, we have presented the deep interpretable metric learning (DIML) method for more transparent embedding learning. We proposed a structural matching strategy that explicitly aligns the spatial embeddings by computing an optimal matching flow between feature maps of the two images. We evaluated our method on three major benchmarks of deep metric learning including CUB200-2011, Cars196 and Stanford Online Products, and achieved substantial improvements over popular metric learning methods with better interpretability. Our method enables deep models to learn metrics in a more human-friendly way, which can be used to inspect and understand the visual similarity of any two samples or applied to any deep metric learning methods with the proposed multi-scale matching strategy to improve image retrieval performance with controllable cost.

#### Acknowledgements

This work was supported in part by the National Natural Science Foundation of China under Grant 61822603, Grant U1813218, and Grant U1713214, in part by a grant from the Beijing Academy of Artificial Intelligence (BAAI), and in part by a grant from the Institute for Guo Qiang, Tsinghua University.

## A. Implementation of DIML

### A.1. The Sinkhorn Algorithm

The Sinkhorn algorithm [6] modifies the original optimal transport problem (Eq.4) into the following one:

$$\begin{aligned} T^* &= \arg \min_{T \geq 0} \text{tr}(CT^\top) + \lambda \text{tr}(T(\log(T) - \mathbf{1}\mathbf{1}^\top)^\top), \\ \text{subject to } T\mathbf{1} &= \mu^s, \quad T^\top \mathbf{1} = \mu^t, \end{aligned} \quad (12)$$

where  $\lambda$  is a non-negative regularization parameter. By adding the entropic regularizer, the Equation (12) becomes a convex problem, which can be solved with Sinkhorn-Knopp algorithm [32]. Starting from an initial matrix  $K = \exp(-C/\lambda)$ , the problem can be solved by iteratively projecting onto the marginal constraints until convergence:

$$\mathbf{a} \leftarrow \mu^s / K\mathbf{b}, \quad \mathbf{b} \leftarrow \mu^t / K^\top \mathbf{a}. \quad (13)$$

After converged, we can obtain the optimal transport plan:

$$T^* = \text{diag}(\mathbf{a})K\text{diag}(\mathbf{b}). \quad (14)$$

### A.2. Testing

In all of our experiments, we use ResNet50 [12] as our backbone. Therefore, the size of the feature map before the pooling layer is  $7 \times 7$ . To reduce computational costs, we first use ROI Align [11] to pool the feature map to  $G \times G$  and  $G = 4$  in most of our experiments unless otherwise noted. According to the multi-scale matching algorithm, for each image as a query, we first sort the images in the gallery using the standard cosine similarity to obtain the indices of top- $K$  candidates  $\mathcal{I}_K$  (we use  $K = 100$  in most of the experiments). We then calculate the proposed structural similarity of all the images in  $\mathcal{I}_K$ . To combine both global and structural information, we use the sum of the cosine similarity and the structural similarity for the top- $K$  images to compute their ranks. The regularization parameter  $\lambda$  in Equation (6) is set to 0.05.

### A.3. Training

Incorporating DIML into the training objectives is quite straightforward. Generally, the loss functions in metric learning can be roughly categorized into distance-based methods (e.g., Contrastive [10], Triplet [9], Margin [44]) and similarity-based methods (e.g., Multi-Similarity [42], Arcface [7], N-Pair [33]). For distance-based methods, we replace the original distance function  $d$  with the average of  $d$  and our structural distance  $d_{\text{struct}}$ ; For similarity-based methods, we replace the original similarity function  $s$  with the average of  $s$  and our structural similarity  $s_{\text{struct}}$ . In this section, we will use several loss functions as examples to demonstrate how to apply DIML during training.

**Margin [44]** The Margin loss [44] is defined as

$$\mathcal{L}_{\text{margin}}(k, l) = \left( \sigma + (-1)^{I(y^k \neq y^l)} (D_{k,l} - \beta) \right)_+, \quad (15)$$

where  $\sigma$  and  $\beta$  are learnable parameters, and  $D_{kl}$  is used to measure the distance between image  $k$  and  $l$ :

$$D_{k,l} = \frac{1}{2} (d_{\text{struct}}(z^k, z^l) + d(\bar{z}^k, \bar{z}^l)), \quad (16)$$

where  $d$  is Euclid distance and  $d_{\text{struct}}$  is derived from  $d$  using Equation (10).

**Multi-Similarity [42]** The original Multi-Similarity is defined as:

$$\begin{aligned} s^*(k, l) &= \begin{cases} s(k, l), & s(k, l) > \min_{p \in \mathcal{P}_k} s(k, p) - \epsilon \\ s(k, l), & s(k, l) < \max_{n \in \mathcal{N}_k} s(k, n) + \epsilon, \\ 0, & \text{otherwise} \end{cases} \\ \mathcal{L}_{\text{MS}} &= \frac{1}{B} \sum_{k \in \mathcal{B}} \left[ \frac{1}{\alpha} \log \left[ 1 + \sum_{p \in \mathcal{P}_k} \exp(-\alpha (s^*(k, p) - \lambda)) \right] \right. \\ &\quad \left. + \frac{1}{\beta} \log \left[ 1 + \sum_{n \in \mathcal{N}_k} \exp(\beta (s^*(k, n) - \lambda)) \right] \right], \end{aligned} \quad (17)$$

where  $s(k, l) = s(\psi^k, \psi^l)$  is the cosine similarity of the embeddings  $\psi^k, \psi^l$  of the two images. To utilize DIML, we can replace  $s$  with

$$s(k, l) \leftarrow \frac{1}{2} (s(\bar{z}^k, \bar{z}^l) + s_{\text{struct}}(z^k, z^l)). \quad (19)$$

Note that in our notation both  $\psi^k$  and  $\bar{z}^k$  represent the same embedding in  $\mathbb{R}^D$ .

**ProxyNCA [22]** It is also worth mentioning there are slight difference when applying DIML to proxy-based methods during training. Taking ProxyNCA [22] as example, the original objective is

$$\mathcal{L}_{\text{proxy}} = -\frac{1}{B} \sum_{k \in \mathcal{B}} \log \left( \frac{\exp(-d(\psi^k, \eta^{y^k}))}{\sum_{c \in \mathcal{C} \setminus \{y^k\}} \exp(-d(\psi^k, \eta^c))} \right), \quad (20)$$

where  $d$  is Euclid distance and  $\eta^c \in \mathbb{R}^D$  is the proxy for the  $c$ -th class. To use DIML, we need to use proxies with the size  $\mathbb{R}^{H \times W \times D}$ , denoted as  $\{\rho^c, c \in \mathcal{C}\}$ . Then, we can replace the  $d(\psi^k, \eta^c)$  with

$$d(\psi^k, \eta^c) \leftarrow \frac{1}{2} (d(\psi^k, \rho^c) + d_{\text{struct}}(z^k, \rho^c)), \quad (21)$$

where we also note that  $\text{GAP}(\rho^c) = \eta^c$ .

Table 5: **Comparisons of different truncation numbers.** We test for different truncation number  $K$  ranging from 0 to 500. Experimental results show that a small  $K$  can already bring considerable performance improvement.

| Baseline             | $K$ | CUB-200      |              |              | Cars196      |              |              | SOP          |              |              |
|----------------------|-----|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                      |     | P@1          | RP           | M@R          | P@1          | RP           | M@R          | P@1          | RP           | M@R          |
| Margin[44]           | 0   | 62.47        | 34.12        | 23.14        | 72.18        | 32.00        | 20.82        | 78.39        | 45.64        | 42.34        |
|                      | 10  | <b>65.16</b> | 34.56        | 23.87        | <b>76.65</b> | 32.52        | 21.72        | <b>79.26</b> | <b>46.44</b> | <b>43.20</b> |
|                      | 50  | 65.16        | 35.43        | <b>24.54</b> | 76.65        | 33.64        | 22.83        | 79.26        | 46.44        | 43.19        |
|                      | 100 | 65.16        | <b>35.48</b> | 24.54        | 76.65        | <b>33.93</b> | <b>22.95</b> | 79.26        | 46.44        | 43.19        |
|                      | 500 | 65.16        | 35.48        | 24.54        | 76.65        | 33.93        | 22.95        | 79.26        | 46.44        | 43.19        |
| Multi-Similarity[42] | 0   | 62.56        | 32.74        | 21.99        | 74.81        | 32.72        | 21.60        | 77.90        | 44.97        | 41.54        |
|                      | 10  | <b>64.89</b> | 33.21        | 22.73        | <b>78.50</b> | 33.26        | 22.50        | <b>78.53</b> | <b>45.60</b> | <b>42.24</b> |
|                      | 50  | 64.89        | 34.04        | 23.37        | 78.50        | 34.46        | 23.72        | 78.53        | 45.60        | 42.23        |
|                      | 100 | 64.89        | <b>34.12</b> | <b>23.38</b> | 78.50        | <b>34.70</b> | <b>23.81</b> | 78.53        | 45.60        | 42.23        |
|                      | 500 | 64.89        | 34.12        | 23.38        | 78.50        | 34.70        | 23.81        | 78.53        | 45.60        | 42.23        |

## B. Experimental Details

### B.1. Evaluation Metrics

We implement the same evaluation metrics as [23], including Precision at 1 (P@1), R-Precision (RP), and Mean Average Precision at R (MAP@R).

**P@1** is also known as Recall@1 in metric learning. Given a sample  $x^q$  and feature encoder  $\phi(\cdot)$ , the set of  $k$  nearest neighbors of  $x^q$  is calculated as the precision of  $k$  nearest neighbors:

$$\mathcal{N}_q^k = \arg \min_{\mathcal{N} \subset \mathcal{X}_{\text{test}}, |\mathcal{N}|=k} \sum_{x^f \in \mathcal{N}} d_e(\phi(x^q), \phi(x^f)) \quad (22)$$

where  $d_e(\cdot, \cdot)$  is the euclidean distance. Then P@ $k$  can be measured as

$$\text{P@}k = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x^q \in \mathcal{X}_{\text{test}}} \frac{1}{k} \sum_{x^i \in \mathcal{N}_q^k} \begin{cases} 1, & y^i = y^q, \\ 0, & \text{otherwise} \end{cases}, \quad (23)$$

where  $y^i$  is the class label of sample  $x^i$ . We only report P@1 in our experiments, i.e.  $k = 1$ .

**R-precision** is defined in [23]. Specifically, for each sample  $x^q$ , let  $R$  be the number of images that are the same class with  $x^q$  and R-precision is simply defined as P@ $R$  (see Equation 23). However, R-precision does not consider the ranking of correct retrievals, so it is not informative enough. To tackle this problem, [23] introduced Mean Average Precision at R.

**MAP@R** is similar to mean average precision, but limit the number of nearest neighbors to R. So it replaces *precision* in MAP calculation with *R-precision*:

$$\text{MAP@}R = \frac{1}{R} \sum_{i=1}^R P(i), \quad (24)$$

where

$$P(i) = \begin{cases} \text{P@}i, & \text{if the } i\text{-th retrieval is correct;} \\ 0, & \text{otherwise.} \end{cases} \quad (25)$$

MAP@R is more informative than P@1 and it can be computed directly from the embedding space without clustering as post-processing.

### B.2. Experimental Setups

For most of the baseline methods, we follow the implementation and the hyper-parameters in [36]. For Proxy Anchor [16], we use their original implementation but set the hyper-parameters as [36] (batch size 112, embedding size 128, etc.). Besides various loss functions, we also experiment with different sampling methods. In Table 1 of the original paper, we use suffixes to represent the sampling methods (-R: Random; -D: Distance [44]; -S Semihard [30]; -H: Softhard [29]).

## C. Detailed Results

In the original paper, we have demonstrated the effects of truncation number  $K$  and feature map size  $G$  using charts. In this section, we provide the original numerical results that were used to plot those charts in Table 5 and Table 6.

## References

- [1] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B Tenenbaum, William T Freeman, and Antonio Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. *ICLR*, 2019. 1, 2
- [2] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. *CVPR*, 2021. 1, 2

Table 6: **Effects of the size of feature map.** Generally, the performance of our DIMLs is better with higher  $G$ . DIML with  $G = 4$  yields good results within relatively low computational costs.

| Baseline              | $G$ | CUB-200      |              |              | Cars196      |              |              | SOP          |              |              |
|-----------------------|-----|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                       |     | P@1          | RP           | M@R          | P@1          | RP           | M@R          | P@1          | RP           | M@R          |
| Margin[44]            | 1   | 62.47        | 34.12        | 23.14        | 72.18        | 32.00        | 20.82        | 78.39        | 45.64        | 42.34        |
|                       | 2   | 64.15        | 34.79        | 23.83        | 75.04        | 32.59        | 21.85        | 79.06        | 46.29        | 43.03        |
|                       | 4   | 65.16        | 35.48        | 24.54        | 76.65        | <b>33.93</b> | <b>22.95</b> | 79.26        | 46.44        | 43.19        |
|                       | 7   | <b>65.58</b> | <b>35.58</b> | <b>24.79</b> | <b>76.96</b> | 32.93        | 22.66        | <b>79.59</b> | <b>46.83</b> | <b>43.62</b> |
| Multi-Similarity [42] | 1   | 62.56        | 32.74        | 21.99        | 74.81        | 32.72        | 21.60        | 77.90        | 44.97        | 41.54        |
|                       | 2   | 63.77        | 33.33        | 22.60        | 77.45        | 33.25        | 22.60        | 78.39        | 45.56        | 42.15        |
|                       | 4   | 64.89        | 34.12        | 23.38        | 78.50        | <b>34.70</b> | <b>23.81</b> | 78.53        | 45.60        | 42.23        |
|                       | 7   | <b>65.45</b> | <b>34.15</b> | <b>23.55</b> | <b>78.93</b> | 33.64        | 23.50        | <b>78.76</b> | <b>45.90</b> | <b>42.57</b> |

- [3] Guangyi Chen, Yongming Rao, Jiwen Lu, and Jie Zhou. Temporal coherence or temporal motion: Which is more critical for video-based person re-identification? In *ECCV*, pages 660–676, 2020. 2
- [4] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *CVPR*, pages 403–412, 2017. 1, 2
- [5] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *CVPR*, pages 403–412, 2017. 2, 6
- [6] Marco Cuturi. Sinkhorn distances: lightspeed computation of optimal transport. In *NIPS*, volume 2, page 4, 2013. 4, 9
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, pages 4690–4699, 2019. 1, 2, 6, 9
- [8] Charlie Frogner, Farzaneh Mirzazadeh, and Justin Solomon. Learning embeddings into entropic wasserstein spaces. *ICLR*, 2019. 4
- [9] Weifeng Ge. Deep metric learning with hierarchical triplet loss. In *ECCV*, pages 269–285, 2018. 2, 9
- [10] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, volume 2, pages 1735–1742, 2006. 2, 6, 9
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 5, 9
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 2, 3, 5, 6, 7, 9
- [13] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. 6
- [14] Pierre Jacob, David Picard, Aymeric Histace, and Edouard Klein. Metric learning with horde: High-order regularizer for deep embeddings. In *ICCV*, pages 6539–6548, 2019. 2
- [15] Valentin Khruikov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. Hyperbolic image embeddings. In *CVPR*, pages 6418–6428, 2020. 2
- [16] Sungeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric learning. In *CVPR*, pages 3238–3247, 2020. 1, 2, 6, 7, 10
- [17] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCVW*, pages 554–561, 2013. 2, 5
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 2, 6
- [19] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 2
- [20] Xudong Lin, Yueqi Duan, Qiyuan Dong, Jiwen Lu, and Jie Zhou. Deep variational metric learning. In *ECCV*, pages 689–704, 2018. 2
- [21] Iacopo Masi, Yue Wu, Tal Hassner, and Prem Natarajan. Deep face recognition: A survey. In *SIBGRAPI*, pages 471–478. IEEE, 2018. 1
- [22] Yair Movshovitz-Attias, Alexander Toshev, Thomas K Leung, Sergey Ioffe, and Saurabh Singh. No fuss distance metric learning using proxies. In *ICCV*, pages 360–368, 2017. 6, 9
- [23] Kevin Musgrave, Serge Belongie, and Ser-Nam Lim. A metric learning reality check. In *ECCV*, pages 681–699. Springer, 2020. 5, 6, 10
- [24] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, pages 4004–4012, 2016. 2, 5
- [25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 2019. 4
- [26] Yongming Rao, Jiwen Lu, and Jie Zhou. Attention-aware deep reinforcement learning for video face recognition. In *ICCV*, pages 3931–3940, 2017. 2
- [27] Yongming Rao, Jiwen Lu, and Jie Zhou. Learning discriminative aggregation network for video-based face recognition and person re-identification. *IJCV*, 127(6):701–718, 2019. 2
- [28] Karsten Roth, Timo Milbich, and Bjorn Ommer. Pads: Policy-adapted sampling for visual similarity learning. In *CVPR*, pages 6568–6577, 2020. 2
- [29] Karsten Roth, Timo Milbich, Samarth Sinha, Prateek Gupta, Bjorn Ommer, and Joseph Paul Cohen. Revisiting train-

- ing strategies and generalization performance in deep metric learning. In *ICML*, pages 8242–8252. PMLR, 2020. 5, 6, 10
- [30] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 1, 2, 6, 10
- [31] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, pages 618–626, 2017. 1, 2
- [32] Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *The American Mathematical Monthly*, 74(4):402–405, 1967. 9
- [33] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, pages 1857–1865, 2016. 1, 2, 6, 9
- [34] G Sreenu and MA Saleem Durai. Intelligent video surveillance: a review through deep learning techniques for crowd analysis. *Journal of Big Data*, 6(1):1–27, 2019. 1
- [35] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015. 3
- [36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016. 10
- [37] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014. 2
- [38] Evgeniya Ustinova and Victor S. Lempitsky. Learning deep embeddings with histogram loss. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *NeurIPS*, pages 4170–4178, 2016. 6
- [39] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008. 3
- [40] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 2, 5, 8
- [41] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. In *ICCV*, pages 2593–2601, 2017. 6
- [42] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *CVPR*, pages 5022–5030, 2019. 1, 6, 7, 9, 10, 11
- [43] Jonathan R Williford, Brandon B May, and Jeffrey Byrne. Explainable face recognition. In *ECCV*, pages 248–263. Springer, 2020. 2
- [44] Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *ICCV*, pages 2840–2848, 2017. 2, 5, 6, 7, 9, 10, 11
- [45] Hong Xuan, Richard Souvenir, and Robert Pless. Deep randomized ensembles for metric learning. In *ECCV*, pages 723–734, 2018. 2
- [46] Tongtong Yuan, Weihong Deng, Jian Tang, Yinan Tang, and Binghui Chen. Signal-to-noise ratio: A robust distance metric for deep metric learning. In *CVPR*, pages 4815–4824, 2019. 2, 6
- [47] Andrew Zhai and Hao-Yu Wu. Classification is a strong baseline for deep metric learning. In *BMVC*, page 91. BMVA Press, 2019. 6
- [48] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural networks. In *CVPR*, pages 8827–8836, 2018. 1, 2
- [49] Quanshi Zhang, Yu Yang, Haotian Ma, and Ying Nian Wu. Interpreting cnns via decision trees. In *CVPR*, pages 6261–6270, 2019. 1, 2
- [50] Zizhao Zhang, Yubo Zhang, Xibin Zhao, and Yue Gao. Emd metric learning. In *AAAI*, volume 32, 2018. 4
- [51] Wenzhao Zheng, Zhaodong Chen, Jiwen Lu, and Jie Zhou. Hardness-aware deep metric learning. In *CVPR*, pages 72–81, 2019. 2
- [52] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, pages 2921–2929, 2016. 1, 2