

HIER: Metric Learning Beyond Class Labels via Hierarchical Regularization

Sungyeon Kim¹ Boseung Jeong¹ Suha Kwak^{1,2}
 Dept. of CSE, POSTECH¹ Graduate School of AI, POSTECH²
 {sungyeon.kim, boseung01, suha.kwak}@postech.ac.kr
<http://cvlab.postech.ac.kr/research/HIER>

Abstract

Supervision for metric learning has long been given in the form of equivalence between human-labeled classes. Although this type of supervision has been a basis of metric learning for decades, we argue that it hinders further advances in the field. In this regard, we propose a new regularization method, dubbed HIER, to discover the latent semantic hierarchy of training data, and to deploy the hierarchy to provide richer and more fine-grained supervision than inter-class separability induced by common metric learning losses. HIER achieves this goal with no annotation for the semantic hierarchy but by learning hierarchical proxies in hyperbolic spaces. The hierarchical proxies are learnable parameters, and each of them is trained to serve as an ancestor of a group of data or other proxies to approximate the semantic hierarchy among them. HIER deals with the proxies along with data in hyperbolic space since the geometric properties of the space are well-suited to represent their hierarchical structure. The efficacy of HIER is evaluated on four standard benchmarks, where it consistently improved the performance of conventional methods when integrated with them, and consequently achieved the best records, surpassing even the existing hyperbolic metric learning technique, in almost all settings.

1. Introduction

Learning a discriminative and generalizable embedding space has been a vital step within many machine learning tasks including content-based image retrieval [18, 26, 36, 37], face verification [20, 33], person re-identification [6, 50], few-shot learning [29, 35, 39], and representation learning [18, 45, 53]. Deep metric learning has aroused lots of attention as an effective tool for this purpose. Its goal is to learn an embedding space where semantically similar samples are close together and dissimilar ones are far apart. Hence, the semantic affinity between samples serves as the main supervision for deep metric learning, and has long been given in the form of equivalence between their human-labeled classes.

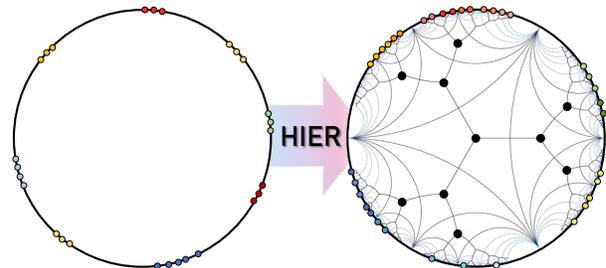


Figure 1. Motivation of HIER. HIER aims to discover an inherent but latent semantic hierarchy of data (colored dots on the boundary) by learning hierarchical proxies (larger black dots) in hyperbolic space. The semantic hierarchy is deployed to provide rich and fine-grained supervision that cannot be derived only by human-labeled classes.

Although this type of supervision has been a basis of metric learning for decades, we argue that it now hinders further advances of the field. The equivalence of human-labeled classes deals with only a tiny subset of possible relations between samples due to the following two reasons. First, the class equivalence is examined at only a single fixed level of semantic hierarchy, although different classes can be semantically similar if they share the same super-class. Second, the equivalence is a binary relation that ignores the degree of semantic affinity between two classes. It is difficult to overcome these two limitations since the semantic hierarchy of data is latent and only human-labeled classes are available from existing datasets in the standard metric learning setting. However, once they are resolved, one can open up the possibility of providing rich supervision beyond human-labeled classes.

In this regard, we propose a new regularization method, called HIERarchical Regularization (HIER), to discover and deploy the latent semantic hierarchy of training data for metric learning. Since the semantic hierarchy can capture not only the predefined set of classes but also their subclasses, super-classes, and affinities between them, our regularizer is expected to provide richer and more fine-grained supervision beyond inter-class separability induced by common metric learning losses. However, it is challenging to establish such a semantic hierarchy of data given their

human-labeled classes only due to the absence of annotation for the semantic hierarchy.

HIER tackles this challenge by learning *hierarchical proxies* in *hyperbolic space*. The hierarchical proxies are learnable parameters, and each of them is trained to serve as an ancestor of a group of data or other proxies to approximate the semantic hierarchy. Also, HIER deals with the proxies in hyperbolic space since the space is well-suited to represent hierarchical structures of the proxies and data. It has been reported in literature that Euclidean space with zero curvature is not optimal for representing data exhibiting such a semantic hierarchy [16]. In contrast, hyperbolic space with constant negative curvature can represent the semantic hierarchy of data effectively using relatively small dimensions since its volume increases exponentially as its Poincaré radius increases [31, 32].

To be specific, HIER is designed as a soft approximation of hierarchical clustering [7, 25, 43], where similar data should be near by one another in a tree-structured hierarchy while dissimilar ones should be placed in separate branches of the hierarchy (See Figure 1). During training, HIER constructs a triplet of samples or proxies, in which two of them are similar and the other is dissimilar based on their hyperbolic distances. Then, the proxy closest to the entire triplet is considered as the *lowest common ancestor* (LCA) of the triplet in a semantic hierarchy; likewise, we identify another proxy as the LCA of only the similar pair in the same manner. Given the triplet and two LCAs (proxies), HIER encourages that each of the LCAs and its associated members of the triplet are close together and that the dissimilar one of the triplet is far apart from the LCA of the similar pair. This allows the hierarchical proxies to approximate a semantic hierarchy of data in the embedding space, without any off-the-shelf module for pseudo hierarchical labeling [51].

The major contribution of our work from the perspective of conventional metric learning is three-fold:

- We study a new and effective way of providing semantic supervision beyond human-labeled classes, which has not been actively studied in metric learning [13, 24, 30, 51, 52].
- HIER consistently improved performance over the state-of-the-art metric learning losses when integrated with them, and consequently achieved the best records in almost all settings on four public benchmarks.
- By taking advantage of hyperbolic space, HIER substantially improved performance particularly on low-dimensional embedding spaces.

Also, when regarding our work as a hyperbolic metric learning method, a remarkable contribution of HIER is that it allows taking full advantage of both hyperbolic embedding

space and the great legacy of conventional metric learning since it can be seamlessly incorporated with any metric learning losses based on spherical embedding spaces. The only prior studies on hyperbolic metric learning [11, 51] are not compatible with conventional losses for metric learning since they learn a distance metric directly on hyperbolic space, in which, unlike the spherical embedding spaces, norms of embedding vectors vary significantly.

2. Related Work

2.1. Deep Metric Learning

Deep metric learning aims to learn an embedding space where data points of the same class are close to each other, and those of different classes are far apart. To this end, a number of studies have proposed various loss functions, which can be categorized into two ways, pair-based and proxy-based losses. Pair-based losses handle the relations between pairs of data points. Contrastive loss aims to minimize the distance between positive pairs and maximize the distance between negative pairs. Triplet loss [33, 44] employs a triplet of anchor, positive and negative samples, and optimizes the embedding space by enforcing the constraint that the distance between the positive pair is smaller than that between the negative pair with a pre-defined margin. Higher-order variants of these losses have also been proposed to capture more complex relations between embedding vectors [36, 37, 46, 47]. On the other hand, proxy-based losses consider the relations between data points and proxies, where the proxy is an additional learnable embedding vector that represents the class of training data. Proxy-NCA [26] associates a data point with proxies and pulls the positive pair of a data point and proxy closer while pushing away the negative pair. Proxy-Anchor [17] associates a proxy with all data points in a batch, allowing it to consider rich relations between data points while reflecting their relative hardness through gradients.

Recently, loss functions that leverage the underlying hierarchical relations of data beyond class equivalence relations have been proposed for metric learning, such as hierarchical triplet loss [13] and hierarchical proxy-based loss [52]. These methods predefine the hierarchy of data, and adjust the distances between data points to fit the discrete hierarchy. In contrast, HIER aims to learn *continuous hierarchical representations* in a data-driven manner without the need for prior knowledge of the hierarchy or pre-existing algorithms [51]. This is a significant advantage over existing hierarchical metric learning methods [13, 51, 52], which have a discrete nature and require the number of hierarchy levels and clusters per level to be specified. Our method allows for a more robust and adaptable method for hierarchical clustering, automatically adjusting to the latent semantic hierarchy of the data.

2.2. Hyperbolic Embedding

Learning hyperbolic embedding has drawn the attention of many research fields because it can encode data such as text or images into hyperbolic space with semantically rich representation due to its high capacity and tree-like property. Since learning embedding in the hyperbolic space for natural language processing [27] has been successful, there are some attempts to learn hyperbolic embeddings for computer vision tasks such as few-shot learning [12, 16] and deep metric learning [11, 51]; they have demonstrated that hyperbolic embeddings are able to improve the model performance. Instead of building complex network architectures whose all layers operate in hyperbolic space, they proposed hybrid architectures where only the last layer maps inputs in Euclidean space to the hyperbolic embedding vectors, and the remainders still operate in Euclidean space. Yan *et al.* [51] proposed the unsupervised hyperbolic metric learning framework by conducting hierarchical clustering with pre-defined hierarchy thresholds. Ermolov *et al.* [11] utilized the pairwise cross-entropy loss with hyperbolic distances through joint use of vision transformers [3, 9, 41].

Unfortunately, these approaches do not fully take advantage of hyperbolic space for the following reasons. First, a pre-defined set of hierarchy thresholds restricts the hierarchical property of hyperbolic space. Second, the learned embedding space is less generalizable since the latent semantic hierarchy of data in hyperbolic space, which can provide more fine-grained supervision and allow a more generalizable embedding space to be learned, is not considered. Contrary to the existing hyperbolic metric learning methods, our loss can preserve the heritage of metric learning as well as fully take advantage of hyperbolic space thanks to HIER.

2.3. Hierarchical Clustering

Hierarchical clustering is a recursive grouping of a dataset into clusters with an increasingly finer granularity. It has usually been used for data analysis [38, 55], visualization [34], and mining fine-grained relations of data [1]. Recently, cost function based methods [7, 25, 43] have been developed.

Dasgupta [7] first proposed a cost function based method, where the cost function is optimized by pairwise similarities between data points and the leaves of the lowest common ancestor (LCA). Wang and Wang [43] have improved the cost function of [7] by introducing a triplet manner that determines the weights of the costs according to the triplet relations. However, These cost functions cannot be optimized by the stochastic gradient descent method due to its inherent discrete nature. As an effort to overcome this limitation, Monath *et al.* [25] proposed a gradient-based hyperbolic hierarchical clustering (gHHC) that can be applied to continuous tree representations in hyperbolic space, and

improved computational cost by approximating a distribution over LCA for two or three data points.

The aforementioned works have enabled hierarchical clustering for both discrete and continuous trees, but they have the limitation of requiring a ground-truth weighted graph. Therefore, they are not suitable for hyperbolic metric learning benchmarks, where the ground-truth weighted graph is absent. HIER can overcome this problem by leveraging virtual ancestors of data points, named as *hierarchical proxies*, while still allowing for stochastic gradient descent optimization.

3. Proposed Method

This section first presents preliminaries to our work, the Poincaré ball model and hierarchical clustering. Then, the objective of HIER is elaborated on.

3.1. Preliminary

3.1.1 The Poincaré ball Model

Hyperbolic space means a Riemannian manifold with negative curvature, which has multiple conformal models [2]. In this work, we use the Poincaré ball, which is a popular and well-studied space in hyperbolic geometry. The n -dimensional Poincaré ball model $(\mathbb{D}_c^n, g^{\mathbb{D}})$ is defined by the manifold $\mathbb{D}_c^n = \{\mathbf{x} \in \mathbb{R}^n : c\|\mathbf{x}\| < 1\}$ and Riemannian metric $g^{\mathbb{D}} = \lambda_c^2 g^E$, where c is a curvature hyperparameter, $\lambda_c = \frac{2}{1-c\|\mathbf{x}\|^2}$ is the conformal factor, and $g^E = I_n$ is Euclidean metric tensor.

An output of a conventional embedding network, lying in Euclidean space, can be transformed into a point on a Poincaré ball by a mapping function called *exponential map*. The common form of the exponential map is given by

$$\exp_0(\mathbf{v}) = \tanh \sqrt{c}\|\mathbf{v}\| \frac{\mathbf{v}}{\sqrt{c}\|\mathbf{v}\|}. \quad (1)$$

One of the important properties of hyperbolic space is that it is not a vector space, so vectors can be calculated algebraically only by introducing gyrovector spaces [42], a generalized notion of vector spaces. The addition operation in gyrovector spaces, called Möbius addition, between a pair of vectors $\mathbf{u} \in \mathbb{D}_c^n$ and $\mathbf{v} \in \mathbb{D}_c^n$ is defined as

$$\mathbf{u} \oplus_c \mathbf{v} := \frac{(1 + 2c\langle \mathbf{u}, \mathbf{v} \rangle + c\|\mathbf{v}\|^2)\mathbf{u} + (1 - c\|\mathbf{u}\|^2)\mathbf{v}}{1 + 2c\langle \mathbf{u}, \mathbf{v} \rangle + c^2\|\mathbf{u}\|^2\|\mathbf{v}\|^2}. \quad (2)$$

In the Poincaré ball, the hyperbolic distance between the two vectors \mathbf{u} and \mathbf{v} is then formulated as

$$d_H(\mathbf{u}, \mathbf{v}) = \frac{2}{\sqrt{c}} \operatorname{arctanh}(\sqrt{c}\|\mathbf{u} \oplus_c \mathbf{v}\|). \quad (3)$$

As the curvature c of Eq. (3) approaches to 0, the hyperbolic distance becomes equal to the Euclidean distance. As the norm of vectors increases, the hyperbolic distance grows

much faster than the Euclidean distance that increases linearly. Thanks to its inherent characteristic, when the radius of the Poincaré ball increases, its volume increases exponentially and thus enables to represent the semantic hierarchy of data even with low dimensions [31, 32]. However, at the same time, this property leads to serious optimization issues when the hyperbolic distance is incorporated with existing metric learning losses. Specifically, the range of hyperbolic distances is not bounded, and accordingly their scales vary significantly by norms of embedding vectors. Thus, the use of hyperbolic distance makes the criteria of conventional metric learning losses futile. One tentative solution is to normalize the embedding vectors and utilize an existing metric learning loss. However, such a naïve method cannot encode hierarchical structures, which is the main reason for adopting hyperbolic space.

3.1.2 Conventional Hierarchical Clustering

Hierarchical clustering is an algorithm that builds a hierarchy of clusters in a tree-like structure by grouping data with a progressively finer granularity. Intuitively, the optimal clustering constructs a tree where similar data points are close to each other and dissimilar data points are located in separate branches. The Dasgupta cost [7] is an objective function that measures the quality of hierarchical clustering based on this property, and is formulated as follows:

$$\mathcal{C} := \sum_{i,j} w_{i,j} |\text{leaves}(i \vee j)|, \quad (4)$$

where $w_{i,j}$ is the ground-truth weight indicating similarity between two nodes i and j , $i \vee j$ denotes their lowest common ancestor (LCA), and $\text{leaves}(k)$ is a set of leaves of the sub-tree whose root is node k . To minimize the Dasgupta cost, the higher the ground-truth weight between two nodes, the smaller the number of descendants of their LCA should be. This means that similar nodes should be located close to each other at leaves of the tree structure.

Unfortunately, this cost function cannot be optimized by stochastic gradient descent because of its discrete nature. Its extensions [4, 25] enable gradient-based optimization for clustering, but they still have limitations in that they demand the ground-truth weighted graph, which is not available in metric learning datasets. For these reasons, the Dasgupta cost and its extensions are not directly applicable to deep metric learning.

3.2. Our Hierarchical Metric Learning Objective

We propose HIER, a novel regularization method that discovers and deploys the latent semantic hierarchy of data in hyperbolic space. Since the semantic hierarchy is expected to capture not only predefined classes but also their sub-classes, super-classes, and affinities between these notions, our regularizer provides more fine-grained supervi-

sion beyond inter-class separability induced by common metric learning losses. Also, it is designed to be incorporated with any conventional metric learning losses based on the spherical embedding.

3.2.1 HIER

To produce rich supplementary supervision as a regularizer, HIER aims to discover the latent semantic hierarchy of data with no additional annotation for the hierarchy. Specifically, it is built upon an extension of the Dasgupta cost [25], but is reshaped for unsupervised clustering with no ground-truth weighted graph. The key idea of HIER is to employ *hierarchical proxies* as learnable ancestors of data points in the hierarchy.

Given a triplet $\{x_i, x_j, x_k\}$, in which x_i and x_j are related to each other and x_k is irrelevant, HIER encourages that the pair of related samples have the same LCA and the rest has a different LCA. Note that we do not consider class labels of data when sampling such a triplet since the goal of HIER is to discover the latent semantic hierarchy beyond the predefined classes. Instead, we employ the reciprocal nearest neighbor to ensure the relevance relations of a triplet. To be specific, the set of feasible triplets, denoted by \mathcal{T} , is estimated as follows:

$$\mathcal{T} = \{(x_i, x_j, x_k) \mid (x_j \in R_K(x_i)) \wedge (x_k \notin R_K(x_i))\} \\ \text{where } R_K(x) = \{x' \mid (x' \in N_K(x)) \wedge (x \in N_K(x'))\}, \quad (5)$$

where $x \in \mathbb{D}_c^n$ denotes an embedding vector in hyperbolic space, $N_K(x)$ is the K -nearest neighbors of x , and $R_K(x)$ stands for the K -reciprocal nearest neighbors of x . The probability that a hierarchical proxy $\rho \in P$ is the LCA of embedding vectors x_i and x_j is given by

$$\pi_{ij}(\rho) = \exp\left(-\max\{d_H(x_i, \rho), d_H(x_j, \rho)\}\right), \quad (6)$$

where d_H is the hyperbolic distance defined in Eq. (3). Then, the most likely LCA of x_i and x_j is sampled from the probability distribution π_{ij} by the Gumbel-max trick as follows:

$$\rho_{ij} = \arg \max_{\rho} (\pi_{ij}(\rho) + g_{ij}), \quad (7)$$

where g_{ij} is an *i.i.d.* sample drawn from Gumbel(0, 1); the Gumbel noise g_{ij} allows to avoid falling into local optima. In the same way of Eq. (6) and (7), ρ_{ijk} that is likely to be the LCA of the entire triplet $\{x_i, x_j, x_k\}$ is sampled from the hierarchical proxy set P except for ρ_{ij} .

Using the hyperbolic distances between a sampled triplet and their LCAs, our objective for hierarchical regularization is implemented as a stack of three triplet losses:

$$\mathcal{L}_{\text{HIER}}(t) = [d_H(x_i, \rho_{ij}) - d_H(x_i, \rho_{ijk}) + \delta]_+ \\ + [d_H(x_j, \rho_{ij}) - d_H(x_j, \rho_{ijk}) + \delta]_+ \quad (8) \\ + [d_H(x_k, \rho_{ijk}) - d_H(x_k, \rho_{ij}) + \delta]_+,$$

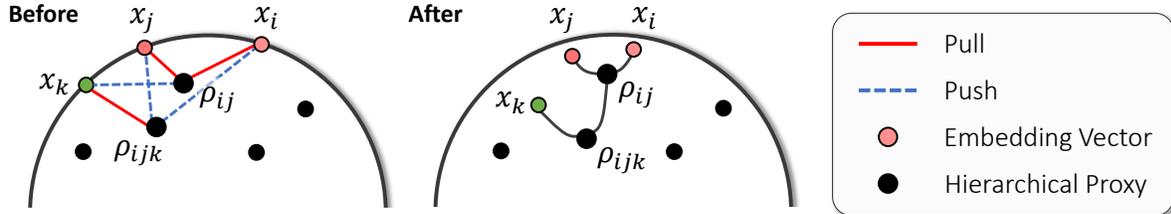


Figure 2. A conceptual illustration of the learning objective in Eq. (8). Each hierarchical proxy is colored in black and different colors indicate different classes. The associations defined by the losses are expressed by edges, where the red solid line means the pulling and the blue dotted line is the pushing. Relevant pairs are pulled into their LCA, and the remaining sample is pulled into LCA of the entire triplet.

where δ is a margin hyperparameter. Figure 2 illustrates the behavior of HIER in terms of the relations between a triplet and their LCAs. HIER forces the relevant samples x_i and x_j to be close to ρ_{ij} , but encourages them to increase the distance from ρ_{ijk} by the margin. Meanwhile, the opposite signal is applied to the irrelevant sample x_k . These operations encourage ρ_{ijk} to be relatively close to the center of the Poincaré ball and ρ_{ij} to be away from the center. Consequently, they become representatives of higher and lower levels of the semantic hierarchy, respectively. In addition, relevant samples are clustered, while unrelated samples are branched out from them. As a result, x_i and x_j belong to a child of ρ_{ij} , and x_k becomes a child of ρ_{ijk} , forming a tree-like hierarchical structure.

3.2.2 Total Objective

Although HIER provides rich and fine-grained supervision induced by the approximate semantic hierarchy of training data, a metric learning loss still plays crucial roles as its supervisory signals are reliable thanks to the use of ground-truth labels while HIER relies on self-supervised hierarchical clustering. The final loss function for metric learning is thus a linear combination of the two terms as follows:

$$\mathcal{L} = \mathcal{L}_{\text{ML}} + \lambda \sum_{t \in \{\mathcal{T}_x, \mathcal{T}_\rho\}} \mathcal{L}_{\text{HIER}}(t), \quad (9)$$

where λ is a weight hyperparameter, \mathcal{L}_{ML} is a metric learning loss, and \mathcal{T}_x and \mathcal{T}_ρ are sets of triplets of samples and hierarchical proxies that satisfy the condition in Eq. (5), respectively. By feeding both samples and proxies as input to the loss, we imply $\mathcal{L}_{\text{HIER}}$ to encourage hierarchical relations between ancestors (*i.e.*, proxies) as well as individual samples.

For the metric learning objective \mathcal{L}_{ML} , either those based on cosine (Euclidean) distances or those of hyperbolic distances can be incorporated. When using a metric learning loss based on spherical embedding, the loss is calculated through the Euclidean metric with ℓ_2 normalized embedding vectors. In this case, the metric learning loss only controls the angles between the normalized embedding vectors and HIER adjusts their hyperbolic distances based on their positions and norms.

4. Experiments

In this section, we evaluate our method on four benchmark datasets for deep metric learning and compare it with the state of the art.

4.1. Experimental Setup

Datasets. On four benchmark datasets, namely CUB-200-2011 (CUB) [49], Cars-196 (Cars) [19], Stanford Online Product (SOP) [37], and In-shop Clothes Retrieval (In-Shop) [21], models are evaluated and compared. We split the datasets into train and test sets, directly following the standard protocol presented in [17].

Evaluation metric. We measure the performance on the datasets by Recall@ k , the fraction of queries that have at least one relevant sample in their k -nearest neighbors on a learned hyperbolic embedding space.

Embedding network architectures. For fair comparisons with previous work, we utilize ResNet50 [14] and vision transformer architecture with three types of pretraining scheme (ViT-S [9], DeiT-S [41] and DINO [3]). All encoder is pretrained for ImageNet classification [8]. In ViT variants, the linear projection layer for patch embedding is frozen during training. We adjust the size of the last FC layer according to the dimensionality of embedding vectors. We note that L_2 normalization is not used for hyperbolic embedding. Instead, the exponential mapping layer (See Eq. (1)) is appended to the last embedding layer.

Implementation details. Our embedding models are optimized by AdamW [22] with the learning rate value 10^{-5} for ViT-S and DeiT-S, and 5×10^{-6} for DINO models. For experiments using ResNet50, we follow the training setting of [17]. Training images are randomly resized and cropped to 224×224 and randomly flipped horizontally while test images are resized to 256×256 and then center cropped. For hyperbolic embedding, we use curvature parameter $c = 0.1$ and clipping radius $r = 2.3$ following previous work [11]. In proxy anchor loss, we use a high learning rate for proxies by scaling 1×10^4 times. For all of our experiments, we maintain a consistent set of hyperparameters, including $K = 20$ for the number of nearest neighbors, $|P| = 512$ for the number of hierarchical proxies, $\lambda = 1$ for the loss

Methods	Arch.	CUB			Cars			SOP			In-Shop		
		R@1	R@2	R@4	R@1	R@2	R@4	R@1	R@10	R@100	R@1	R@10	R@20
<i>Backbone architecture: CNN</i>													
NSoftmax [54]	R ¹²⁸	56.5	69.6	79.9	81.6	88.7	93.4	75.2	88.7	95.2	86.6	96.8	97.8
MIC [30]	R ¹²⁸	66.1	76.8	85.6	82.6	89.1	93.2	77.2	89.4	94.6	88.2	97.0	-
XBM [48]	R ¹²⁸	-	-	-	-	-	-	80.6	91.6	96.2	91.3	97.8	98.4
XBM [48]	B ⁵¹²	65.8	75.9	84.0	82.0	88.7	93.1	79.5	90.8	96.1	89.9	97.6	98.4
HTL [13]	B ⁵¹²	57.1	68.8	78.7	81.4	88.0	92.7	74.8	88.3	94.8	80.9	94.3	95.8
MS [46]	B ⁵¹²	65.7	77.0	86.3	84.1	90.4	94.0	78.2	90.5	96.0	89.7	97.9	98.5
SoftTriple [28]	B ⁵¹²	65.4	76.4	84.5	84.5	90.7	94.5	78.6	86.6	91.8	-	-	-
PA [17]	B ⁵¹²	68.4	79.2	86.8	86.1	91.7	95.0	79.1	90.8	96.2	91.5	98.1	98.8
NSoftmax [54]	R ⁵¹²	61.3	73.9	83.5	84.2	90.4	94.4	78.2	90.6	96.2	86.6	97.5	98.4
†ProxyNCA++ [40]	R ⁵¹²	69.0	79.8	87.3	86.5	92.5	95.7	80.7	92.0	96.7	90.4	98.1	98.8
Hyp [11]	R ⁵¹²	65.5	76.2	84.9	81.9	88.8	93.1	79.9	91.5	96.5	90.1	98.0	98.7
HIER (ours)	R ⁵¹²	70.1	79.4	86.9	88.2	93.0	95.6	80.2	91.5	96.6	92.4	98.2	98.8
<i>Backbone architecture: ViT</i>													
IRT _R [10]	De ¹²⁸	72.6	81.9	88.7	-	-	-	83.4	93.0	97.0	91.1	98.1	98.6
Hyp [11]	De ¹²⁸	74.7	84.5	90.1	82.1	89.1	93.4	83.0	93.4	97.5	90.9	97.9	98.6
HIER (ours)	De ¹²⁸	75.2	84.2	90.0	85.1	91.1	95.1	82.5	92.7	97.0	91.0	98.0	98.6
Hyp [11]	DN ¹²⁸	78.3	86.0	91.2	86.0	91.9	95.2	84.6	94.1	97.7	92.6	98.4	99.0
HIER (ours)	DN ¹²⁸	78.5	86.7	91.5	88.4	93.3	95.9	84.9	94.2	97.5	92.6	98.4	98.9
Hyp [11]	V ¹²⁸	84.0	90.2	94.2	82.7	89.7	93.9	85.5	94.9	98.1	92.7	98.4	98.9
HIER (ours)	V ¹²⁸	84.2	90.1	93.7	86.4	91.9	95.1	85.6	94.6	97.8	92.7	98.4	98.9
IRT _R [10]	De ³⁸⁴	76.6	85.0	91.1	-	-	-	84.2	93.7	97.3	91.9	98.1	98.9
DeiT-S [41]	De ³⁸⁴	70.6	81.3	88.7	52.8	65.1	76.2	58.3	73.9	85.9	37.9	64.7	72.1
Hyp [11]	De ³⁸⁴	77.8	86.6	91.9	86.4	92.2	95.5	83.3	93.5	97.4	90.5	97.8	98.5
HIER (ours)	De ³⁸⁴	78.7	86.8	92.0	88.9	93.9	96.6	83.0	93.1	97.2	90.6	98.1	98.6
DINO [3]	DN ³⁸⁴	70.8	81.1	88.8	42.9	53.9	64.2	63.4	78.1	88.3	46.1	71.1	77.5
Hyp [11]	DN ³⁸⁴	80.9	87.6	92.4	89.2	94.1	96.7	85.1	94.4	97.8	92.4	98.4	98.9
HIER (ours)	DN ³⁸⁴	81.1	88.2	93.3	91.3	95.2	97.1	85.7	94.6	97.8	92.5	98.6	99.0
ViT-S [9]	V ³⁸⁴	83.1	90.4	94.4	47.8	60.2	72.2	62.1	77.7	89.0	43.2	70.2	76.7
Hyp [11]	V ³⁸⁴	85.6	91.4	94.8	86.5	92.1	95.3	85.9	94.9	98.1	92.5	98.3	98.8
HIER (ours)	V ³⁸⁴	85.7	91.3	94.4	88.3	93.2	96.1	86.1	95.0	98.0	92.8	98.4	99.0

Table 1. Performance of metric learning methods on the four datasets. Their network architectures are denoted by abbreviations, R-ResNet50 [14], B-Inception with BatchNorm [15], De-DeiT [41], DN-DINO [3] and V-ViT [9]. Note that ViT [9] is pretrained on ImageNet-21k [8]. Superscripts denote their embedding dimensions and † indicates models using larger input images.

weight of HIER, and $\delta = 0.1$ for the margin.

4.2. Quantitative Results

We compare the performance of the proposed method with the existing state of the arts on four standard datasets. In these experiments, we adopt proxy anchor loss [17] as a metric learning objective \mathcal{L}_{ML} . For a fair comparison with other methods, we divide the previous work in terms of backbone architecture and embedding dimension, which is summarized in Table 1.

The experimental results demonstrate the effectiveness of our proposed method. Compared with existing CNN-based methods, our method using ResNet50 as the backbone achieves superior performance on all datasets except for SOP. Specifically, it outperforms ProxyNCA++ by a significant margin of 1.7% on the Cars dataset and 2% on the In-Shop dataset. Furthermore, we reimplement Hyp using ResNet50 and observe that its performance lags behind ours. For ViT-based models, HIER consistently out-

performs the state of the art in terms of R@1 for almost all settings. Specifically, our method with 128 embedding dimensions achieves a large margin improvement of 3.7%, 2.4%, and 3.0% in R@1 for V¹²⁸, DN¹²⁸, and De¹²⁸ models on the Cars dataset, respectively, and surpasses the performance of 512-dimensional CNN-based methods. These improvements suggest that our method allows the embedding space to be more discriminative and generalized with smaller dimensionality, thanks to utilizing hyperbolic space and considering latent semantic hierarchy.

4.3. Qualitative Results

To demonstrate that HIER actually represents a latent semantic hierarchy of data, we visualize the learned embedding vectors which are projected to a 2-dimensional Poincaré ball. For visualization, we use UMAP [23] with hyperboloid distance metric as a dimensional reduction technique. Figure 3 shows that a tree-like hierarchical struc-

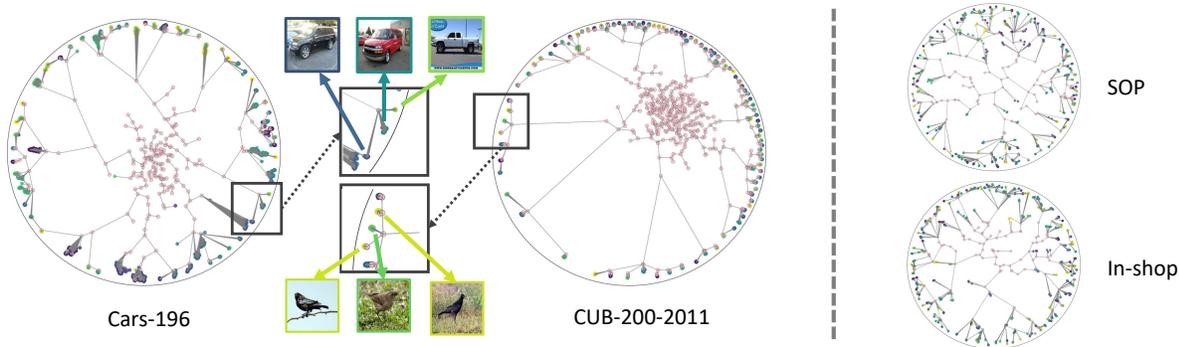


Figure 3. UMAP visualization of our embedding space learned on the train split of Cars, CUB, SOP, and In-shop. Pink ones indicate hierarchical proxies and other colors represent distinct classes. The gray lines are the ancestor-descendant relations between data points.

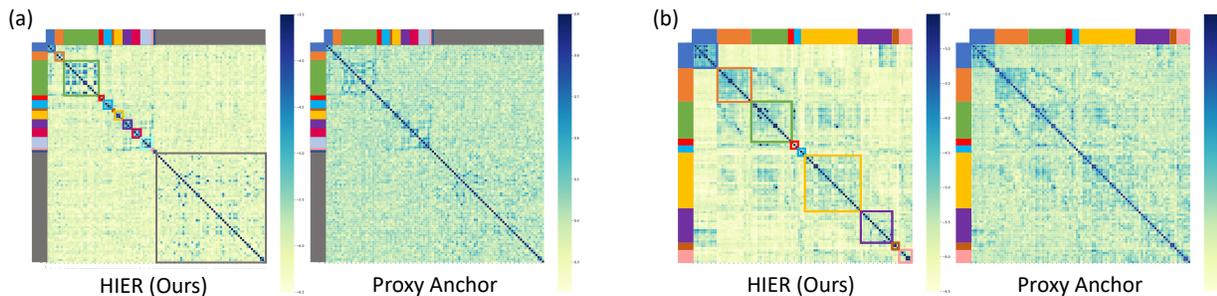


Figure 4. Class-to-class affinity matrices of proxy anchor and ours on the CUB (a) and Cars (b) datasets, which show the inter-class similarity. The different colors (13 colors of CUB and 9 colors of Cars) of the sidebar indicate distinct actual super-classes at the order level referring to the hierarchy labels of [5].

ture between data and hierarchical proxies is constructed in the embedding space learned by HIER. The middle figures are enlarged figures of part of the embedding space, which shows that samples in a sub-tree share common semantics, although not of the same class. This visualization proves that HIER discovers and deploys latent and meaningful hierarchies between data in the embedding space.

4.4. Analysis on Semantic Hierarchy of HIER

To demonstrate that HIER can capture a meaningful semantic hierarchy, we perform two qualitative experiments that can show whether the embedding space reflects the super-class and sub-class relations between data.

In Figure 4, we present the affinity matrix between classes on the embedding space trained with our method, and compare it with that of the proxy anchor loss [17]. The similarity measures of affinity matrices for proxy anchor loss and ours are set to cosine similarity and negative hyperbolic distance, respectively. Since classes with the same super-class share common attributes, they have a high degree of similarity in both methods. However, in the case of proxy anchor loss, classes belonging to different super-classes often have high similarity, whereas HIER can clearly distinguish between classes sharing the same super-class and classes that are not.

In addition, to investigate the sub-class characteristics,

we present the nearest neighbor samples of hierarchical proxies that are close to the boundary of the Poincaré ball (*i.e.*, representing a relatively low level of semantic hierarchy). As shown in Figure 5, each hierarchical proxy on the CUB dataset is associated with pose variation of birds and backgrounds of images (*e.g.*, a bird floating on water, a bird flying in the sky, and a bird standing on a branch). On the Cars dataset, hierarchical proxies represent the viewpoints of the cars. (*e.g.*, front, half side, and side views of the car). These results suggest that HIER also captures sub-class relations, thereby enabling the model to learn intra-class variation and features shared across different classes.

4.5. Ablation Studies

Impact of HIER and hyperbolic space. We investigate the effect of HIER loss and hyperbolic space by comparing the performance of three models trained using two different metric learning losses as \mathcal{L}_{ML} in Eq. (9): without HIER loss and with HIER loss in spherical space (*i.e.*, $HIER_{sph}$), and with HIER loss in hyperbolic space (*i.e.*, HIER (ours)). The results are summarized in Table 2. Overall, using the HIER loss leads to performance improvements across four datasets regardless of the kinds of metric learning loss. However, in spherical space, the improvement is not significant, possibly due to the limitations of spherical space in effectively representing the continuous latent semantic hi-

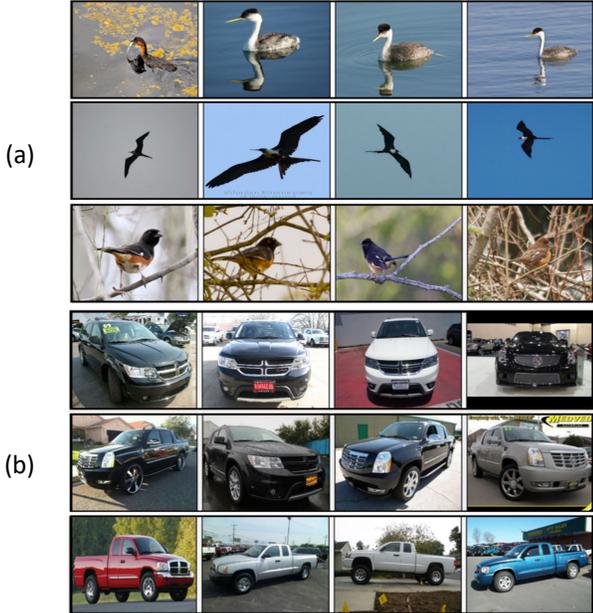


Figure 5. Top-4 neighbors of hierarchical proxies that are close to the boundary of the Poincaré ball on the CUB (a) and Cars (b) datasets. The samples in each row are the nearest neighbors of a hierarchical proxy.

erarchy of data discovered by the HIER loss. In contrast, applying the HIER loss in hyperbolic space results in significant performance improvements across all datasets. These results demonstrate that the HIER loss can effectively discover and deploy the latent semantic hierarchy of data in hyperbolic space.

Impact of embedding dimension. The embedding dimension is a seminal hyperparameter in image retrieval benchmarks due to the trade-off between efficacy and efficiency. Therefore, we examine the effect of the dimension of embedding vectors in our method, compared to Hyp [11]. To this end, we evaluate our method on the In-shop [21] dataset with ViT [9] with d -dimensional embedding vectors, where $d \in \{32, 64, 128, 256, 384\}$. Figure 6 demonstrates that our method consistently outpaces Hyp with diverse embedding dimensions. It is worth noting that the performance of our method increases consistently and stably according to the growth of the embedding dimension, while the performance of Hyp drops when the dimension is higher than 128.

Impact of hyperparameters. We investigate the impact of four hyperparameters on the performance of our method. These hyperparameters include the number of nearest neighbors K in Eq. (5), the number of hierarchical proxies $|P|$, the loss weight of HIER λ in Eq. (9), and the margin δ in Eq. (8). We use the DeiT backbone and measure the performance on the Cars dataset. As shown in Figure 7, the results demonstrate that remarkably robust to changes in these hyperparameters, suggesting that it can perform well regardless of their specific values.

Methods	CUB	Cars	SOP	In-Shop
PA	74.7	84.3	82.3	90.4
PA + HIER _{sph}	75.1	84.8	82.4	90.6
PA + HIER (ours)	75.2	85.1	82.5	91.0
MS	75.4	83.5	80.0	88.1
MS + HIER _{sph}	75.6	83.6	80.0	88.0
MS + HIER (ours)	75.8	84.3	80.0	88.2

Table 2. Accuracy in Recall@1 of ours with two metric learning losses [17, 46], and their variants on the four datasets. The network architecture is DeiT [41] with 128 embedding dimensions. HIER_{sph} denotes HIER over spherical space.

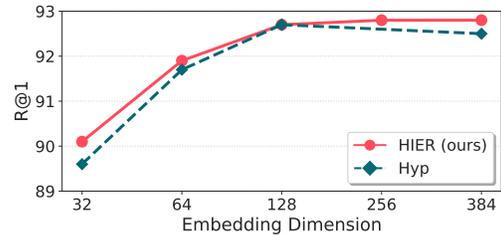


Figure 6. Comparison HIER and Hyp [11] varying embedding dimension on the In-shop dataset using ViT as backbone network.

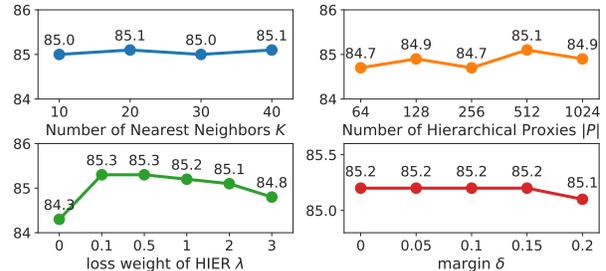


Figure 7. Recall@1 versus four hyperparameters of HIER on the Cars dataset using DeiT [41] with 128 dimensions. Note that $\lambda = 0$ denotes not using HIER loss.

5. Conclusion

In this paper, we have presented HIER, a regularization method for deep metric learning. HIER discovers the latent semantic hierarchy of training data in a self-supervised learning fashion, and deploys the hierarchy to provide richer and more fine-grained supervision that the inter-class separability induced by common metric learning losses based on human-labeled classes. HIER consistently improve the performance of existing techniques when integrated with them, and consequently achieved the best records, surpassing even the prior art of hyperbolic metric learning. Further, since it can be seamlessly incorporated with any existing metric learning losses, HIER enables taking full advantage of both recent hyperbolic embedding and the great legacy of hyper-spherical embedding.

Acknowledgement. This work was supported by the NRF grant and the IITP grant funded by Ministry of Science and ICT, Korea (NRF-2021R1A2C3012728–20%, IITP-2019-0-01906–10%, IITP-2020-0-00842–50%, IITP-2022-0-00926–20%).

References

- [1] Peter F Brown, Vincent J Della Pietra, Peter V Desouza, Jennifer C Lai, and Robert L Mercer. Class-based n-gram models of natural language. *Computational linguistics*, 18(4):467–480, 1992. **3**
- [2] James W Cannon, William J Floyd, Richard Kenyon, Walter R Parry, et al. Hyperbolic geometry. *Flavors of geometry*, 1997. **3**
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. **3, 5, 6, 11**
- [4] Ines Chami, Albert Gu, Vaggos Chatziafratis, and Christopher Ré. From trees to continuous embeddings and back: Hyperbolic hierarchical clustering. *Proc. Neural Information Processing Systems (NeurIPS)*, 2020. **4**
- [5] Jingzhou Chen, Peng Wang, Jian Liu, and Yuntao Qian. Label relation graphs enhanced hierarchical residual network for hierarchical multi-granularity classification. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. **7**
- [6] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: A deep quadruplet network for person re-identification. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. **1**
- [7] Sanjoy Dasgupta. A cost function for similarity-based hierarchical clustering. In *Proc. ACM symposium on Theory of Computing*, 2016. **2, 3, 4**
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: a large-scale hierarchical image database. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. **5, 6**
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. International Conference on Learning Representations (ICLR)*, 2021. **3, 5, 6, 8, 11**
- [10] Alaaeldin El-Nouby, Natalia Neverova, Ivan Laptev, and Hervé Jégou. Training vision transformers for image retrieval. *arXiv preprint arXiv:2102.05644*, 2021. **6**
- [11] Aleksandr Ermolov, Leyla Mirvakhabova, Valentin Khruikov, Nicu Sebe, and Ivan Oseledets. Hyperbolic vision transformers: Combining improvements in metric learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. **2, 3, 5, 6, 8**
- [12] Zhi Gao, Yuwei Wu, Yunde Jia, and Mehrtash Harandi. Curvature generation in curved spaces for few-shot learning. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. **3**
- [13] Weifeng Ge, Weilin Huang, Dengke Dong, and Matthew R. Scott. Deep metric learning with hierarchical triplet loss. In *Proc. European Conference on Computer Vision (ECCV)*, 2018. **2, 6**
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. **5, 6**
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proc. International Conference on Machine Learning (ICML)*, 2015. **6**
- [16] Valentin Khruikov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. Hyperbolic image embeddings. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. **2, 3**
- [17] Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. **2, 5, 6, 7, 8, 11**
- [18] Sungyeon Kim, Minkyoo Seo, Ivan Laptev, Minsu Cho, and Suha Kwak. Deep metric learning beyond binary supervision. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. **1**
- [19] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 554–561, 2013. **5, 11**
- [20] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. SpheroFace: Deep hypersphere embedding for face recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. **1**
- [21] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. DeepFashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. **5, 8, 11**
- [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. International Conference on Learning Representations (ICLR)*, 2019. **5**
- [23] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. Umap: Uniform manifold approximation and projection. *The Journal of Open Source Software*, 2018. **6**
- [24] Timo Milbich, Karsten Roth, Homanga Bharadhwaj, Samarth Sinha, Yoshua Bengio, Björn Ommer, and Joseph Paul Cohen. Diva: Diverse visual feature aggregation for deep metric learning. In *Proc. European Conference on Computer Vision (ECCV)*, 2020. **2**
- [25] Nicholas Monath, Manzil Zaheer, Daniel Silva, Andrew McCallum, and Amr Ahmed. Gradient-based hierarchical clustering using continuous representations of trees in hyperbolic space. In *Proc. ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019. **2, 3, 4**
- [26] Yair Movshovitz-Attias, Alexander Toshev, Thomas K Leung, Sergey Ioffe, and Saurabh Singh. No fuss distance metric learning using proxies. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017. **1, 2**
- [27] Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. *Advances in neural information processing systems*, 30, 2017. **3**
- [28] Qi Qian, Lei Shang, Baigui Sun, Juhua Hu, Hao Li, and Rong Jin. Softtriplet loss: Deep metric learning without triplet sampling. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. **6, 11, 12**
- [29] Limeng Qiao, Yemin Shi, Jia Li, Yaowei Wang, Tiejun

- Huang, and Yonghong Tian. Transductive episodic-wise adaptive metric for few-shot learning. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. 1
- [30] Karsten Roth, Biagio Brattoli, and Bjorn Ommer. Mic: Mining interclass characteristics for improved metric learning. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. 2, 6
- [31] Frederic Sala, Chris De Sa, Albert Gu, and Christopher Ré. Representation tradeoffs for hyperbolic embeddings. In *Proc. International Conference on Machine Learning (ICML)*, 2018. 2, 4
- [32] Rik Sarkar. Low distortion delaunay embedding of trees in hyperbolic plane. In *International Symposium on Graph Drawing*, pages 355–366. Springer, 2011. 2, 4
- [33] Florian Schroff, Dmitry Kalenichenko, and James Philbin. FaceNet: A unified embedding for face recognition and clustering. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 1, 2
- [34] Jinwook Seo and Ben Shneiderman. Interactively exploring hierarchical clustering results [gene identification]. *Computer*, 35(7):80–86, 2002. 3
- [35] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Proc. Neural Information Processing Systems (NeurIPS)*, 2017. 1
- [36] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Proc. Neural Information Processing Systems (NeurIPS)*, 2016. 1, 2
- [37] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 5, 11
- [38] Therese Sørliie, Charles M Perou, Robert Tibshirani, Turid Aas, Stephanie Geisler, Hilde Johnsen, Trevor Hastie, Michael B Eisen, Matt Van De Rijn, Stefanie S Jeffrey, et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proceedings of the National Academy of Sciences*, 2001. 3
- [39] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [40] Eu Wern Teh, Terrance DeVries, and Graham W Taylor. Proxynca++: Revisiting and revitalizing proxy neighborhood component analysis. In *European Conference on Computer Vision (ECCV)*. Springer, 2020. 6
- [41] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *Proc. International Conference on Machine Learning (ICML)*, 2021. 3, 5, 6, 8, 11, 12
- [42] Abraham Albert Ungar. A gyrovector space approach to hyperbolic geometry. *Synthesis Lectures on Mathematics and Statistics*, 2008. 3
- [43] Dingkan Wang and Yusu Wang. An improved cost function for hierarchical cluster trees. *arXiv preprint arXiv:1812.02715*, 2018. 2, 3
- [44] Jiang Wang, Yang Song, T. Leung, C. Rosenberg, Jingbin Wang, J. Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 2
- [45] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2015. 1
- [46] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 6, 8
- [47] Xinshao Wang, Yang Hua, Elyor Kodirov, Guosheng Hu, Romain Garnier, and Neil M Robertson. Ranked list loss for deep metric learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [48] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R Scott. Cross-batch memory for embedding learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6388–6397, 2020. 6
- [49] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 5, 11
- [50] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. Joint detection and identification feature learning for person search. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
- [51] Jiexi Yan, Lei Luo, Cheng Deng, and Heng Huang. Unsupervised hyperbolic metric learning. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3
- [52] Zhibo Yang, Muhammet Bastan, Xinliang Zhu, Douglas Gray, and Dimitris Samaras. Hierarchical proxy-based loss for deep metric learning. 2022. 2
- [53] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 1
- [54] Andrew Zhai and Hao-Yu Wu. Classification is a strong baseline for deep metric learning. *arXiv preprint arXiv:1811.12649*, 2018. 6
- [55] Ying Zhao and George Karypis. Evaluation of hierarchical clustering algorithms for document datasets. In *Proc. Conference on Information and Knowledge Management*, 2002. 3

A. Appendix

This supplementary material provides further analyses, additional experimental results, and implementation details that are left out from the main paper due to the space limit.

A.1. Effect of Deploying Semantic Hierarchy

In this section, we investigate the effect of deploying semantic hierarchies by utilizing proxies. In particular, we quantitatively compare our method with SoftTriple loss [28] which can model intra-class variance of data by utilizing multiple proxies per class. Table 4 shows the performance of two proxy-based methods and the number of proxies they used. In SoftTriple loss, 10 proxies are assigned per class, and our method uses 512 hierarchical proxies in addition to proxies for proxy anchor loss.

The result shows that SoftTriple loss lags significantly in performance despite using a much larger amount of proxies compared to our method. The main difference between these methods is how they handle proxies. SoftTriple assigns multiple proxies per class, so their proxies can only represent sub-classes of data and model intra-class variance. On the other hand, proxies in our method can represent sub-classes or super-classes as well as predefined classes, thus allowing more flexibility in modeling more complex relations of data beyond intra-class variance.

A.2. Embedding Space Visualization

We further visualize the embedding space at the beginning of training, such as Epoch 1, 3, 5, 7, and 10. The results of visualization presented in Figure 8 show that the earlier embedding space does not construct a hierarchical structure between the embedding vectors and hierarchical proxies, while this structure is gradually constructed as the training epoch grows. Specifically, as training progresses, the hierarchical proxies have ancestor-descendant relations between sample or other proxies, where these proxies can be regarded as predefined classes, sub-classes, and super-classes. As a consequence, our method can discover the latent semantic hierarchy of training data, which provides rich and granular supervision beyond human-labeled classes.

A.3. More Qualitative Results

We further verify the effect of HIER in a qualitative perspective. To this end, we present the qualitative results of our method, compared to those of a model optimized by only a metric learning loss, \mathcal{L}_{ML} in Eq. (9) of the main paper. We take proxy anchor loss [17] as the metric learning loss. Figure 9 presents the qualitative results for the four public benchmark datasets, CUB [49], Cars [19], SOP [37] and In-Shop [21]. These results demonstrate that our method is robust against small inter-class variance (CUB and SOP), viewpoint variation and distinct color (Cars), and large intra-class variations and viewpoint changes (In-Shop) by discovering and deploying a latent semantic hierarchy of data, while the proxy anchor still suffers from those problems.

We note that all the results in the figure are obtained from the fully unseen class samples. It implies that it allows the proposed hierarchical regularizer to discover the latent semantic hierarchy of data with no additional annotation for the hierarchy; our method allows the embedding space to be generalized well.

Hyperparameters	CUB	Cars	SOP	In-Shop
total epochs	50	50	150	150
warm-up epochs	1	1	5	5
LR of the last layer	$\times 1$	$\times 1$	$\times 10^2$	$\times 10^2$
weight decay	$1e-2$	$1e-2$	$1e-4$	$1e-4$
margin δ	0.1	0.1	0.1	0.1

Table 3. Additional details of hyperparameters for network optimization. LR denotes the learning rate, and its value denotes how many times higher than the original learning rate.

A.4. Additional Implementation Details

Since each dataset that we utilize to evaluate our method has a different characteristic, we set different hyperparameters for network optimization according to each dataset. The summary of the settings of hyperparameter for four datasets are presented in Table 3. For all backbone network variants (*i.e.*, ViT-S [9], DeiT-S [41], and DINO [3]), our model is trained with 50 epochs on CUB [49] and Cars [19], and 150 epochs on SOP [37] and In-Shop [21]. As a warm-up strategy, the last layer which consists of a linear layer followed by the exponential mapping layer is only trained, while the pretrained backbone network is not updated; the warm-up strategy is applied for 1 epoch on CUB and Cars, and 5 epochs on SOP and In-Shop. On the other hand, we use a high learning rate for the last layer (*i.e.*, embedding layer) by scaling 10^2 times for SOP and In-Shop. Furthermore, we set the weight decay factor as $1e-2$ for CUB and Cars, and $1e-4$ for SOP and In-Shop. In all experiments, the margin parameter δ of HIER loss in Eq. (8) is fixed at 0.1.

Methods	CUB			Cars			SOP			In-Shop		
	#Proxies	R@1	R@2	#Proxies	R@1	R@2	#Proxies	R@1	R@10	#Proxies	R@1	R@10
SoftTriple [28]	1,000	72.7	82.7	980	83.2	90.2	113,180	80.9	91.3	39,970	88.5	97.3
PA + HIER (ours)	612	75.2	84.2	610	85.1	91.2	11,830	82.5	92.7	4,509	91.0	98.0

Table 4. Comparison between methods using proxies in terms of performance and the number of proxies on the four benchmark datasets. In these experiments, the backbone network is initialized by weights of DeiT [41].

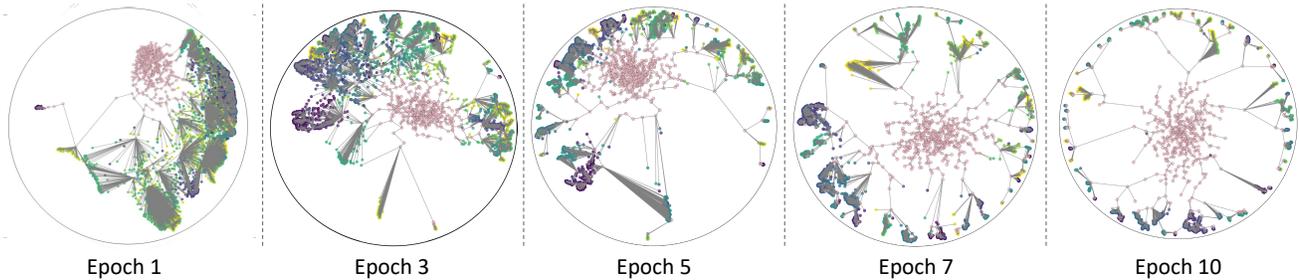


Figure 8. UMAP visualizations of our embedding space learned on the train split of the Cars dataset at different epochs. Pink points indicate hierarchical proxies, and other colors represent distinct classes. The gray line indicates the ancestor-descendant relation between the hierarchical proxy and data points.

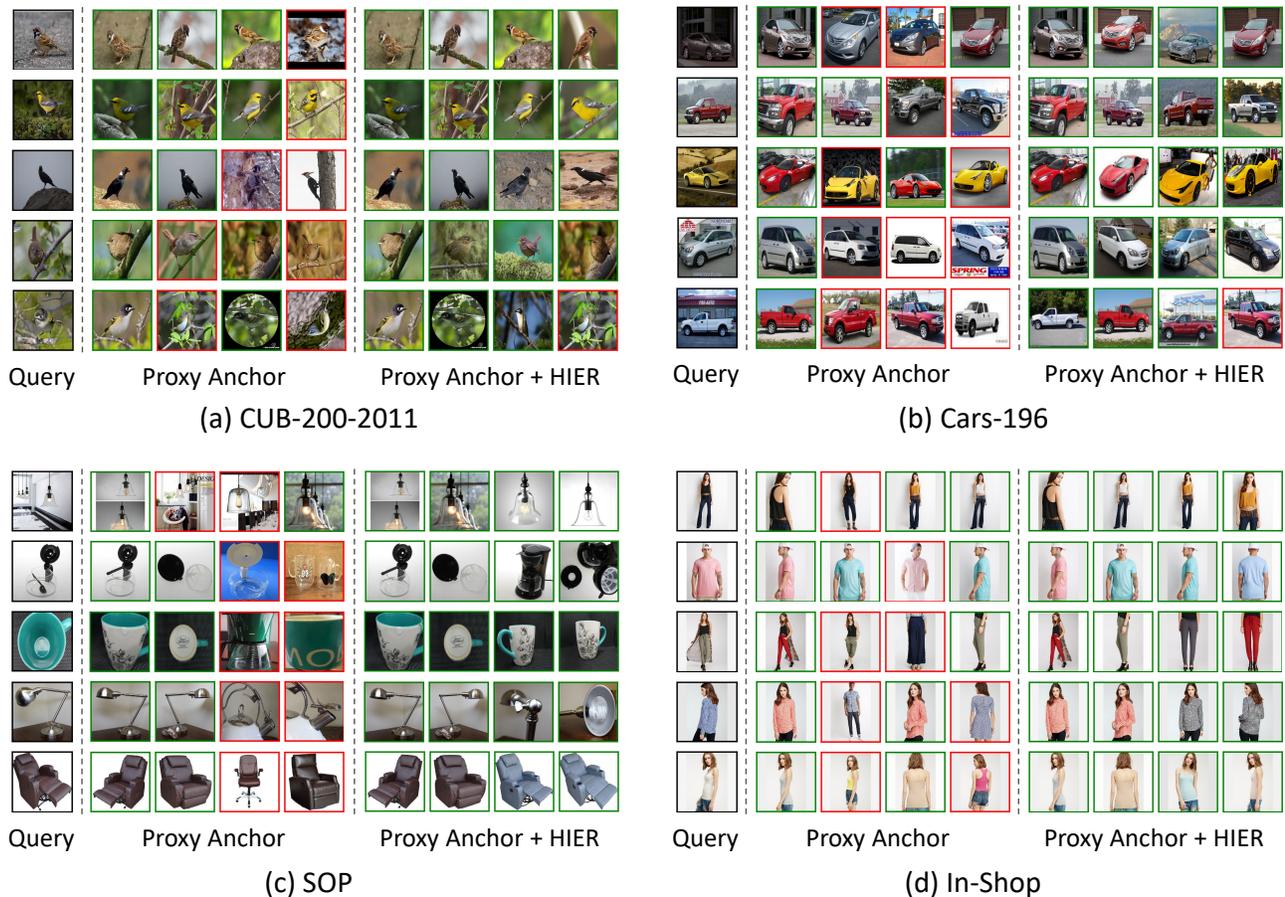


Figure 9. Qualitative results of ours and proxy anchor on the four public benchmark datasets, CUB (a), Cars (b), SOP (c), and In-Shop (d). Queries and the top 4 retrieval results of our method are presented. The true and false matches are colored in green and red, respectively.