

Structure Similarity Preservation Learning for Asymmetric Image Retrieval

Hui Wu, Min Wang, Wengang Zhou, *Senior Member, IEEE*, Houqiang Li, *Fellow, IEEE*

Abstract—Asymmetric image retrieval is a task that seeks to balance retrieval accuracy and efficiency by leveraging lightweight and large models for the query and gallery sides, respectively. The key to asymmetric image retrieval is realizing feature compatibility between different models. Despite the great progress, most existing approaches either rely on classifiers inherited from gallery models or simply impose constraints at the instance level, ignoring the structure of embedding space. In this work, we propose a simple yet effective structure similarity preserving method to achieve feature compatibility between query and gallery models. Specifically, we first train a product quantizer offline with the image features embedded by the gallery model. The centroid vectors in the quantizer serve as anchor points in the embedding space of the gallery model to characterize its structure. During the training of the query model, anchor points are shared by the query and gallery models. The relationships between image features and centroid vectors are considered as structure similarities and constrained to be consistent. Moreover, our approach makes no assumption about the existence of any labeled training data and thus can be extended to an unlimited amount of data. Comprehensive experiments on large-scale landmark retrieval demonstrate the effectiveness of our approach. Our code is released at: <https://github.com/MCC-WH/SSP>.

Index Terms—Multimedia search, Asymmetric image retrieval

I. INTRODUCTION

In recent years, deep learning-based visual search methods [1]–[12] have achieved great success. In a typical visual search system, a deployed deep representation model is used to embed both query and gallery images into a discriminative embedding space. Usually, the embedding features of the large-scale gallery set are extracted and indexed in advance. During the retrieval stage, query features are extracted online and the retrieval is performed by ranking the distance, *e.g.*, Euclidean distance or cosine similarity, between the gallery features and the input query features.

Conventional methods for image retrieval, as described in previous literature [13], [14], utilize a symmetric image retrieval approach, in which the same deep representation

This work is supported by National Natural Science Foundation of China under Contract 62102128 and 62021001. It was also supported by GPU cluster built by MCC Lab of Information Science and Technology Institution, USTC, and the Supercomputing Center of the USTC.

Hui Wu, Wengang Zhou, and Houqiang Li are with the CAS Key Laboratory of Technology in Geospatial Information Processing and Application System, Department of Electronic Engineering and Information Science, University of Science and Technology of China, Hefei, 230027, China (e-mail: wh241300@mail.ustc.edu.cn; zhgw@ustc.edu.cn; lihq@ustc.edu.cn). Min Wang is with the Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei 230030, China (e-mail: wangmin@iai.ustc.edu.cn).

Corresponding authors: Min Wang and Wengang Zhou.

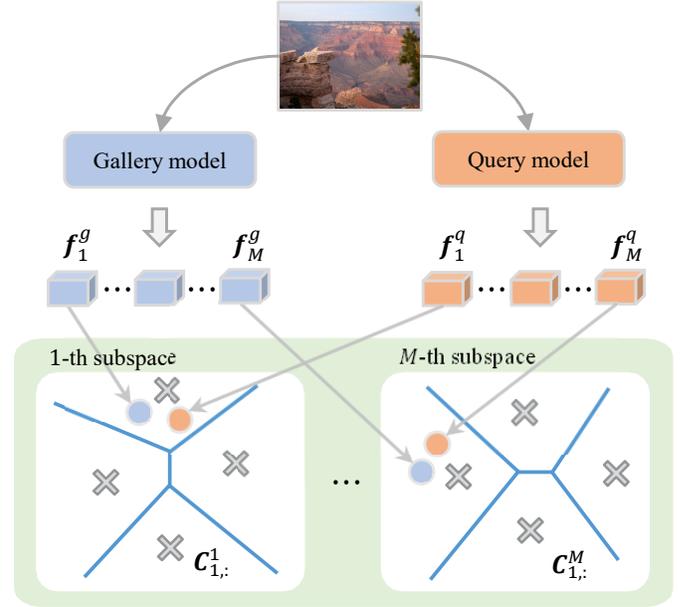


Fig. 1: Illustration of our structure similarity preserving method for *asymmetric image retrieval*. f_i^q (orange) and f_i^g (blue) denote the i -th sub-vector of features from the lightweight query model and the large gallery model, respectively. $C_{i,:}^M$ denotes the i -th centroid vector, *i.e.*, anchor point, in the M -th subspace, which is generated by product quantization. During the learning of the query model, a training image is first embedded by query and gallery models, respectively. Then, we constrain the consistency of the structure similarities between these features and the anchor points to achieve feature compatibility between different models.

model is used to embed both query and gallery images [15], [16]. To obtain high retrieval accuracy, a large powerful model is usually deployed, which is computationally expensive. However, in some real-world applications, gallery images undergo feature extraction offline on resource-rich servers, while query images are processed on resource-constrained end devices, *e.g.*, mobile phones. Due to computational resource constraints, it is difficult to deploy the same large model on these end devices. Lightweight models are better choices due to their low time latency and resource footprint. Thanks to compatible feature learning [17], [18], it is feasible to embed query and gallery images with lightweight and large models separately. This allows enjoying the excellent feature extraction capability of the large model on the server side

while maintaining low resource consumption on the query side. Such an asymmetric setup is denoted as *asymmetric image retrieval* in HVS [16] and AML [15].

As for *asymmetric image retrieval*, it is crucial to ensure that features encoded by the lightweight query model and the large powerful gallery model are compatible with each other. To this end, a straightforward solution is to constrain the features encoded by the two models to be identical, which has shown effectiveness in AML [15]. Another approach, adopted by BCT [17], inherits the classifier of the gallery model to guide the query model. Recently, CSD [19] has considered both feature imitation and neighbor relationship preservation but requires performing multiple retrievals during the training of the query model. Besides, some other works further design advanced restrictions [20]–[24] or advanced network structures [16]. However, all these methods impose constraints only at the instance level or require multiple online retrievals to acquire the nearest neighbor structure, failing to fully preserve the structure of embedding space during the query model learning.

To address above issues, we propose a novel structure similarity preserving approach for ensuring feature compatibility between query and gallery models, which is shown in Figure 1. Our proposed approach involves the extraction of features from an independent dataset using the gallery model, followed by the training of a product quantizer (PQ) in an offline manner. The centroid vectors of the quantizer are then used as anchor points in the embedding space of the gallery model. During query model training, each training image is embedded into features by both query and gallery models, and these features are converted into structure similarities by calculating similarity against anchor points. Our method then constrains the consistency of two structure similarities to optimize the query model. By sharing anchor points between query and gallery models, our method allows the features encoded by two models to align with each other while preserving the embedding space structure.

Compared to previous methods, our approach has two unique advantages. First, we transform features into structure similarities instead of directly performing feature regression. This enables the query model to ignore unimportant feature “details” that are challenging to regress, thereby avoiding over-fitting. Additionally, the centroid vectors of the product quantizer encode the structure information of the gallery embedding space. By sharing these centroids, the embedding spaces of the query and gallery models are closely aligned so that their features are mutually interpretable. Second, our approach leverages the gallery model to derive the structure similarities of the training images, which are further adopted as pseudo-labels to optimize the query model. This eliminates the need for manual annotation of the training dataset, making our approach adaptable and scalable to a variety of real-world scenarios.

To evaluate our method, experiments are conducted on the Revisited Oxford and Paris datasets, with extra 1M distractor images further added for large-scale experiments. Comprehensive experiments with ablation study demonstrate the effectiveness of our method, which achieves the best results

compared to state-of-the-art methods.

II. RELATED WORK

A. Image Retrieval

Given a large corpus, image retrieval aims to efficiently identify images that contain the same object or content as the query image based on feature similarities. Most early retrieval systems rely on bag-of-words representations [25]–[28] with large vocabularies and inverted indexes. In addition, methods for aggregating local features [29], [30] have been explored, including Fisher vectors [31], VLAD [32], and ASMK [28] that produce global descriptors capable of scaling to large databases. To further improve retrieval accuracy, various re-ranking techniques such as spatial verification [25], [27], query expansion [33], and diffusion [34] are further adopted as post-processing steps. Recently, deep learning-based approaches have emerged as promising solutions by framing image retrieval as a metric learning task. Several loss functions [35], [36], pooling methods [4], [14], [37], [38], and training datasets [39], [40] are proposed to enhance the deep representation model.

Although a lot of efforts have been made, optimal retrieval systems typically deploy a large powerful model to process both queries and galleries, which is unaffordable on some resource-constrained end devices. In this work, we focus on *asymmetric image retrieval*, where the query (user) side deploys a lightweight model, while the gallery side deploys a large powerful one.

B. Feature Compatible Learning

It is essential for *asymmetric image retrieval* to encode new (query) features to be interoperable with old (gallery) features. BCT [17] first introduces the problem of *backward-compatible learning* and proposes to inherit the classifier of the gallery model for query model learning. AML [15] performs asymmetric metric learning with the embeddings of anchor and positive/negative samples extracted by query and gallery models, respectively. In HVS [16], both model parameters and architectures are considered simultaneously, and a compatibility-aware neural architecture search method is proposed to search for the optimal query model architecture. LCE [20] proposes a new classifier boundary loss to further improve feature compatibility. However, all these methods impose constraints at the instance level without considering the second-order structural information. The most relevant method to us is CSD [19]. It employs a strategy that incorporates both first-order feature similarity and second-order nearest-neighbor similarity between the gallery and query models. The primary objective is to utilize the gallery model to retrieve the top k nearest neighbor samples from the training dataset. Subsequently, it enforces constraints on the similarity between the feature of the query model and these nearest neighbors to maintain consistency with that of the gallery model. However, it requires multiple time-consuming retrievals during each training iteration to obtain the neighbors of each sample, and using only data samples may not fully capture the structure of the embedding space.

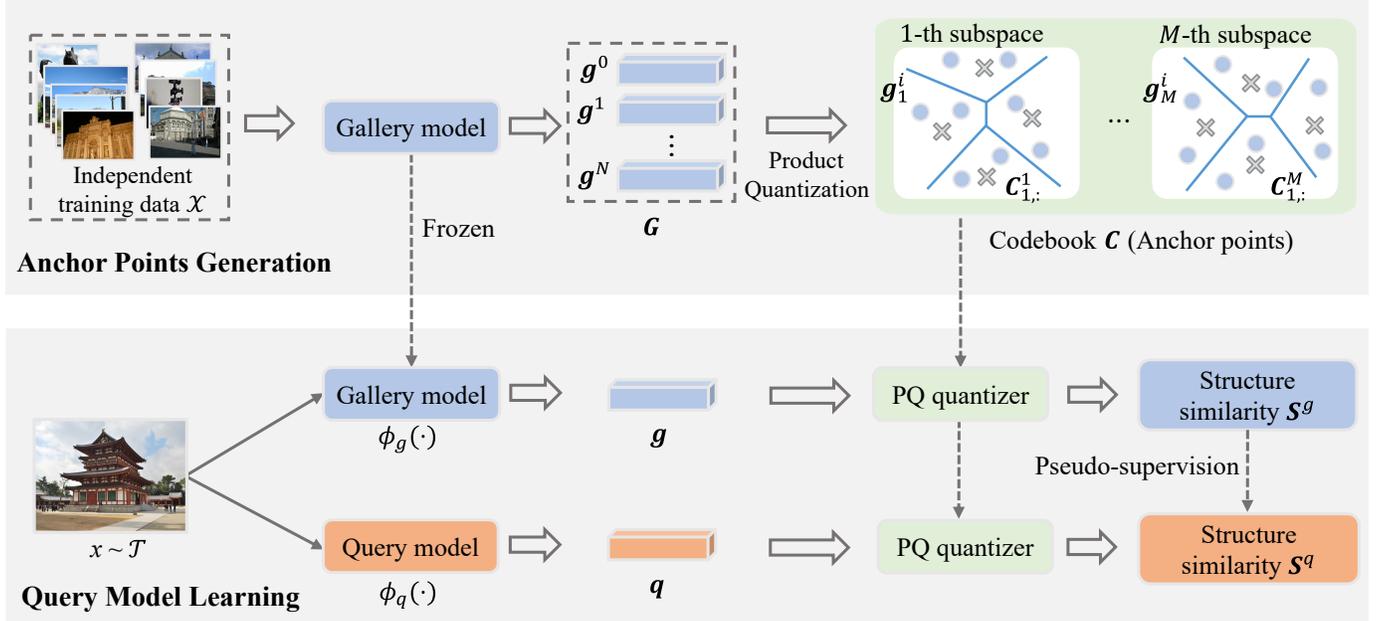


Fig. 2: An overview of our framework. A well trained gallery model $\phi_g(\cdot)$ is first applied to extract the image features $[g^1, g^2, \dots, g^N]$ of the training data \mathcal{X} . Then, these features are utilized to train a product quantizer [41] (Section III-B1), whose codebook \mathcal{C} serves as the anchor points in the embedding space of $\phi_g(\cdot)$. During query model learning, an image x of training dataset \mathcal{T} is mapped into two feature vectors g and q by the query and gallery models, respectively. Then, the similarities between the query/gallery feature (q/g) and the anchor points are regarded as structure similarities (S^q/S^g). Finally, we constrain the consistency between structure similarities S^q and S^g to optimize the query model (Section III-B2). This ensures the embedding spaces of query and gallery models are well aligned, which is essential for *asymmetric retrieval*.

Differently, our method generates a large number of anchor points in the embedding space with product quantization [41]. These anchor points are densely distributed in the embedding space of the gallery model, which carve its structure more delicately. Then, the gallery feature is converted into structure similarity by computing similarities against anchor points, which further serve as the pseudo-label to optimize the query model. Notably, the anchor points are shared by both models, so that feature compatibility is achieved.

C. Lightweight Network

With the evolution of the model architecture, deep convolutional neural networks (CNNs) have made tremendous advances in various computer vision tasks. In real-world applications, computational complexity is another important consideration in addition to accuracy. Typically, it is expected to achieve the best accuracy with a limited computational budget, which is determined by the target computing platforms. The immediate need to deploy high-performance deep neural networks on a range of resource-constrained end devices has motivated a series of studies on efficient model design, including SqueezeNets [42], MobileNets [43], [44], ShuffleNets [45], [46] and EfficientNets [47]. All these methods aim at designing lightweight architecture to achieve better speed-accuracy trade-offs.

In this work, we focus on *asymmetric image retrieval*, where query features are extracted on some resource-constrained end

devices. Our approach employs the various lightweight models mentioned above as query models.

D. Knowledge Transfer

Knowledge transfer aims at learning a student model by transferring knowledge from a pretrained teacher model. It is first introduced by Hinton *et al.* [48], where the student model learns from real labels and soft predicted class logits by the teacher. FitNet [49] distills knowledge through intermediate features and Euclidean distance is used to measure the distance between them. After that, PKT [50] models the knowledge of the teacher model as a probability distribution and uses KL divergence to measure the distance. In RKD [51] and DARK [52], geometric relationships between multiple examples, such as angles and distances, are used as knowledge to guide students learning. CRD [53] combines contrastive learning and knowledge distillation and uses contrastive loss to transfer knowledge between different modalities. Other approaches use multi-stage information to transfer knowledge. AT [54] uses multi-layer attention maps to transfer knowledge. FSP [55] generates FSP matrices from layer features and uses them to guide the learning process of a small model.

However, these methods only transfer knowledge between models but do not consider feature compatibility between them. Thus, they fail to meet the needs of *asymmetric image retrieval*. In our approach, centroid vectors of a product quantizer, which is trained using the image features extracted by the gallery model, serve as the anchor points for both

models. Thus, the embedding spaces of the query and gallery model are constrained to be aligned during knowledge transfer.

III. OUR APPROACH

In this section, we first give a formulation of asymmetric image retrieval. After that, we elaborate our structure similarity preserving framework.

A. Problem Formulation

Let $\phi_g(\cdot)$ and $\phi_q(\cdot)$ denote the gallery and query models, respectively. For a visual retrieval system, the gallery model $\phi_g(\cdot)$ is first trained and then used to map the gallery images \mathcal{R} into feature vectors. During testing, the query model $\phi_q(\cdot)$ processes queries \mathcal{Q} , and the retrieval is reduced to the nearest neighbor search in the embedding space. Some evaluation metric, e.g., mean Average Precision (mAP), is used to evaluate the performance of a retrieval system, which is abbreviated as $\mathcal{M}(\phi_q(\cdot), \phi_g(\cdot))$ for simplicity. In a conventional symmetric retrieval system, the query model is usually the same as the gallery model, i.e., $\phi_q(\cdot) = \phi_g(\cdot)$. It typically deploys a large powerful model to achieve high retrieval accuracy, which cannot be satisfied in resource-constrained scenarios.

As for an asymmetric retrieval system, the model $\phi_g(\cdot)$ is well-trained and fixed. To facilitate resource-constrained application scenarios, it needs to learn a compatible lightweight query model $\phi_q(\cdot)$ which is significantly smaller than $\phi_g(\cdot)$ in terms of parameter size and computational complexity. The core of *asymmetric image retrieval* is that the feature embeddings of query and gallery models are mutually interpretable. In other words, we expect that *asymmetric image retrieval* achieves a retrieval accuracy similar to that of *symmetric retrieval*, i.e., $\mathcal{M}(\phi_q(\cdot), \phi_g(\cdot)) \approx \mathcal{M}(\phi_g(\cdot), \phi_g(\cdot))$ so that the balance between performance and efficiency is achieved.

B. Structure Similarity Preservation Learning

In this work, we propose a structure similarity preserving framework, which is shown in Figure 2. A product quantizer is first trained with the features extracted by the gallery model. The centroids of the quantizer serve as the anchor points to characterize the space structure. During the training of the query model, the gallery model is frozen. Each training sample is mapped into two embeddings by the query and gallery models, respectively. Then, two embeddings are converted to structure similarities by calculating the similarities against the centroids. Finally, our method restricts the consistency between two structure similarities to optimize the query model. Since the anchor points generated by the gallery model are shared with the query model, their embedding space is well-aligned after training.

1) *Anchor Points Generation*: To achieve a comprehensive characterization of the embedding space, our approach requires selecting representative anchor points in the embedding space of the gallery model. These anchor points are fixed reference points in the embedding space which are used to convert query and gallery features into structure similarities. A straightforward approach is to use flat k -means clustering to generate

a series of anchor points. However, our method requires a large number of anchor points to delicately characterize the space structure. If k -means clustering is adopted, the required training samples and computational complexity are several times the number of centroids. When the number of centroids is large, the cost of clustering is unaffordable. To this end, our approach employs product quantization (PQ) [41] to efficiently expand the number of anchor points at a lower cost.

Suppose there exists some training data $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ for anchor points generation. Since the gallery model $\phi_g(\cdot)$ is frozen, we first employ it to extract the features $\mathbf{G} = [\mathbf{g}^1; \mathbf{g}^2; \dots; \mathbf{g}^N] \in \mathbb{R}^{N \times d}$ of images in \mathcal{X} offline:

$$\mathbf{g}^i = \phi_g(x_i) \in \mathbb{R}^d, \quad i = 1, 2, \dots, N. \quad (1)$$

Then, each feature vector $\mathbf{g}^i \in \mathbf{G}$ is split into M distinct sub-vectors $u_j(\mathbf{g}^i) \in \mathbb{R}^{d^*}$, $j = 1, 2, \dots, M$:

$$\underbrace{\mathbf{g}_1^i, \dots, \mathbf{g}_{d^*}^i}_{u_1(\mathbf{g}^i)}, \dots, \underbrace{\mathbf{g}_{d-d^*+1}^i, \dots, \mathbf{g}_d^i}_{u_M(\mathbf{g}^i)}, \quad (2)$$

where \mathbf{g}_j^i denotes the j -th feature dimension of \mathbf{g}^i , $d^* = d/M$ and d is a multiple of M . After that, we perform k -means clustering on each sub-vector set $[u_j(\mathbf{g}^1); u_j(\mathbf{g}^2); \dots; u_j(\mathbf{g}^N)] \in \mathbb{R}^{N \times d^*}$, $j = 1, 2, \dots, M$, individually to obtain the corresponding sub-codebook $\mathbf{C}^j \in \mathbb{R}^{K \times d^*}$, where K is the number of centroids. Then, the anchor points in the gallery space are defined as the Cartesian product of sub-codebooks:

$$\mathbf{C} = \mathbf{C}^1 \times \mathbf{C}^2 \times \dots \times \mathbf{C}^M \in \mathbb{R}^{K^M \times d}, \quad (3)$$

in which any centroid vector is formed by concatenating M different sub-centroid vectors.

Compared with k -means clustering, PQ has two distinctive advantages. First, it is easy to generate a large number of anchor points \mathbf{C} . The total number of anchor points is K^M . Second, instead of storing the huge anchor points directly, it only needs to store $M \times K$ sub-centroids. During training, we also adopt the splitting mechanism to calculate the similarity by segments, instead of directly computing the similarities between feature vectors and all anchor points, which greatly reduces our training overhead. The complete learning procedure is summarized in Algorithm 1.

2) *Query Model Learning: Structure Similarity Calculation*. During the query model learning, the feature vectors of query and gallery models are first converted into structure similarities by calculating the similarities against anchor points. Given an image x in training dataset \mathcal{T} . Let \mathbf{g} and \mathbf{q} be its feature vectors extracted by the gallery and query models, respectively:

$$\mathbf{g} = \phi_g(x) \in \mathbb{R}^d, \quad \mathbf{q} = \phi_q(x) \in \mathbb{R}^d. \quad (4)$$

We first split them into M sub-vectors:

$$\begin{aligned} \mathbf{g} &\rightarrow u_1(\mathbf{g}), u_2(\mathbf{g}), \dots, u_M(\mathbf{g}), \\ \mathbf{q} &\rightarrow u_1(\mathbf{q}), u_2(\mathbf{q}), \dots, u_M(\mathbf{q}). \end{aligned} \quad (5)$$

Then, we calculate the structure similarities $\mathbf{S}_{i,:}^g$ and $\mathbf{S}_{i,:}^q$ by computing for each sub-vectors $u_i(\mathbf{g})$ and $u_i(\mathbf{q})$ the simi-

Algorithm 1: Anchor points generation

Input: Training data $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$; gallery model $\phi_g(\cdot)$; number of subvectors M ; number of centroids per subvector K

Output: Anchor points $\mathbf{C} = \mathbf{C}^1 \times \mathbf{C}^2 \times \dots \times \mathbf{C}^M \in \mathbb{R}^{K^M \times d}$

- 1 **for** each image x_i in training data \mathcal{X} **do**
- 2 Extract image feature with gallery model according to Equation (1);
- 3 Split image feature into M sub-vectors $u_j(\mathbf{g}^i) \in \mathbb{R}^{d^*}$, $j = 1, 2, \dots, M$ according to Equation (2);
- 4 **end**
- 5 **for** each sub-vectors set $[u_j(\mathbf{g}^1); u_j(\mathbf{g}^2); \dots; u_j(\mathbf{g}^N)] \in \mathbb{R}^{N \times d^*}$ **do**
- 6 Perform k -means clustering with K centroids;
- 7 Obtain corresponding sub-codebook $\mathbf{C}^j \in \mathbb{R}^{K \times d^*}$;
- 8 **end**

larities against the corresponding K centroid vectors in the pretrained quantizer:

$$\begin{aligned} \mathbf{S}_{i,:}^g &= [s(u_i(\mathbf{g}), \mathbf{C}_{1,:}^i), \dots, s(u_i(\mathbf{g}), \mathbf{C}_{K,:}^i)] \in \mathbb{R}^K, \\ \mathbf{S}_{i,:}^q &= [s(u_i(\mathbf{q}), \mathbf{C}_{1,:}^i), \dots, s(u_i(\mathbf{q}), \mathbf{C}_{K,:}^i)] \in \mathbb{R}^K, \end{aligned} \quad (6)$$

where $\mathbf{C}_{i,:}^j$ denotes the i -th centroid vector in the j -th subspace and $s(\cdot, \cdot)$ is the similarity metric. In this work, cosine similarity is considered, and $s(\cdot, \cdot)$ is formulated as:

$$s(u_i(\mathbf{g}), \mathbf{C}_{i,:}^j) = \frac{\mathbf{C}_{i,:}^j u_i(\mathbf{g})^T}{\|\mathbf{C}_{i,:}^j\|_2 \|u_i(\mathbf{g})\|_2}. \quad (7)$$

Finally, we impose consistency constraints \mathcal{L}_c on the structure similarities \mathbf{S}^g and \mathbf{S}^q to optimize $\phi_q(\cdot)$ so that the feature embedding \mathbf{q} shares the same structure similarity as \mathbf{g} in the embedding space of the gallery model. Notably, the anchor points are shared between the query and gallery models, thus their embedding spaces are well aligned.

Structure Similarity Preserving Constraint. For *asymmetric image retrieval*, a desirable query model $\phi_q(\cdot)$ not only maintains feature compatibility but also preserves the structure similarity of \mathbf{g} in the embedding space of gallery model. To this end, our method constrains the consistency between two structure similarity \mathbf{S}_i^g and \mathbf{S}_i^q for the corresponding sub-vector pair $u_i(\mathbf{g})$ and $u_i(\mathbf{q})$.

Specifically, Kullback–Leibler (KL) divergence is adopted to measure the distance between $\mathbf{S}_{i,:}^g$ and $\mathbf{S}_{i,:}^q$. First, $\mathbf{S}_{i,:}^g$ is converted into the form of probability distribution:

$$\mathbf{p}_{i,:}^g := \left[\frac{\exp(\mathbf{S}_{i,1}^g/\tau_g)}{\sum_{l=1}^K \exp(\mathbf{S}_{i,l}^g/\tau_g)}, \dots, \frac{\exp(\mathbf{S}_{i,K}^g/\tau_g)}{\sum_{l=1}^K \exp(\mathbf{S}_{i,l}^g/\tau_g)} \right], \quad (8)$$

where τ_g is a temperature value used for controlling the sharpness of the assignments. Similarly, the probability distribution

corresponding to the i -th subvector of the query feature \mathbf{q} is formulated as:

$$\mathbf{p}_{i,:}^q := \left[\frac{\exp(\mathbf{S}_{i,1}^q/\tau_q)}{\sum_{l=1}^K \exp(\mathbf{S}_{i,l}^q/\tau_q)}, \dots, \frac{\exp(\mathbf{S}_{i,K}^q/\tau_q)}{\sum_{l=1}^K \exp(\mathbf{S}_{i,l}^q/\tau_q)} \right]. \quad (9)$$

Then, the structure similarity preserving constraint is defined as the KL divergence between two probabilities over the same sub-centroid vectors:

$$\mathcal{L}_{\text{KL}}^i = \text{KL}(\mathbf{p}_{i,:}^g \parallel \mathbf{p}_{i,:}^q) = \sum_{l=1}^K \mathbf{p}_{i,l}^g \log \frac{\mathbf{p}_{i,l}^g}{\mathbf{p}_{i,l}^q}, \quad (10)$$

which consists of the cross-entropy of $\mathbf{p}_{i,:}^g$ and $\mathbf{p}_{i,:}^q$, and the entropy of $\mathbf{p}_{i,:}^g$. The latter is independent of the feature of the query model and thus does not affect the training. The final objective function is defined as the summation of all consistency losses corresponding to the M distinct sub-vectors:

$$\mathcal{L}_{\text{final}} = \sum_{i=1}^M \mathcal{L}_{\text{KL}}^i, \quad (11)$$

which is used for optimizing query model end to end.

Algorithm 2: Query model learning

Input: Training set \mathcal{T} ; well-trained gallery model $\phi_g(\cdot)$; random initialized query model $\phi_q(\cdot)$; anchor points \mathbf{C}

Output: Query model $\phi_q(\cdot)$ compatible with $\phi_g(\cdot)$

- 1 **for** each image x in training set \mathcal{T} **do**
- 2 Extract image feature with gallery and query models according to Equation (4);
- 3 Split \mathbf{g} and \mathbf{q} into M sub-vectors according to Equation (5);
- 4 Calculate the structure similarities $\mathbf{S}_{i,:}^g$ and $\mathbf{S}_{i,:}^q$ according to Equation (6);
- 5 Impose consistency constraints \mathcal{L}_c on the structure similarities \mathbf{S}^g and \mathbf{S}^q to optimize $\phi_q(\cdot)$ according to Equation (12);
- 6 **end**

Soft-assignment vs. Hard-assignment. The centroids of the product quantizer serve as the anchor points in the embedding space of the gallery model. By quantizing the feature, we convert the feature regression into an assignment prediction task. When setting temperature $\tau_g = 0$, the probability $\mathbf{p}_{i,:}^g$ in Equation (8) will be a one-hot vector with the only 1 at index $j = \arg \max_j (\mathbf{S}_{i,j}^g)$. Thus, Equation (8) is simplified as

$$\mathcal{L}_{\text{KL}}^i = \sum_{l=1}^K \mathbf{p}_{i,l}^g \log \frac{\mathbf{p}_{i,l}^g}{\mathbf{p}_{i,l}^q} = \log \frac{1}{\mathbf{p}_{i,j}^q}. \quad (12)$$

Optimizing this loss encourages the query model to regress the anchor point, to which feature \mathbf{g} is quantized. It avoids the query model regressing the feature “details” of the gallery model. However, the relationships between the feature vector and anchor points carry discriminative knowledge, and simply ignoring them may lead to inferior performance. Thus, we set $\tau_g > 0$ to use the soft assignments as the prediction target. The

QUERY MODEL $\phi_q(\cdot)$	FLOPS (G)		PARAM(M)	
	ABS	%	ABS	%
ResNet101 [56]	42.85	100.0	42.50	100.0
ShuffleNetV2 (0.5 \times) [46]	0.84	1.96	2.44	5.74
ShuffleNetV2 [46]	1.44	3.36	3.35	7.88
MobileNetV2 [44]	2.50	5.83	4.85	11.41
EfficientNetB0 [47]	2.86	6.67	6.63	15.60
EfficientNetB1 [47]	3.92	9.15	9.13	21.49
EfficientNetB2 [47]	4.50	10.51	10.58	24.90
EfficientNetB3 [47]	6.24	14.57	13.84	32.56

TABLE I: The computational complexity and the number of parameters of all lightweight query models adopted in this work are listed. ABS: absolute number. %: relative to ResNet101 [56]. (0.5 \times) denotes a model with 0.5 \times width. FLOPS are calculated when the input image size is 362 \times 362.

overall learning process of the query model is summarized in Algorithm 2.

IV. EXPERIMENTS

A. Experimental Setup

Dataset. SfM-120k [14] and the clean version of Google landmark v2 (GLDv2) [39] are adopted as training set \mathcal{T} . SfM-120k includes images selected from 3D reconstructions of landmarks and city scenes. Following the common setting [15], we use 91,642 images from 551 3D models for training and the remaining images of 162 3D models for validation. The clean version of GLDv2 [39] consists of 1,580,470 images from 81,313 categories. We randomly select 80% images as the training set and let the rest as the validation set. $\mathcal{R}1M$ [40] is used as extra images \mathcal{X} for anchor points generation. It is collected from Yahoo Flickr Creative Commons 100 Million (YFCC100m) dataset [57] and contains 1M distractor images. **Query and Gallery Models.** ResNet101 trained by DELG [58] and GeM [14] are deployed as gallery models, which are denoted as R101-DELG and R101-GeM in this work, respectively. As for lightweight query models, ShuffleNets [46], MobileNets [44] and EfficientNets [47] are chosen. To adapt the model for image retrieval tasks, only the feature extractor of the model is kept and the other layers are both removed. Then, GeM pooling [14] is applied on the last convolutional feature map, followed by another whitening layer, which is implemented by a fully-connected layer. The whitening layer is initialized in the embedding space of the gallery models and kept frozen during the training of the query model. In Table I, we list the number of parameters and the computational complexity (in FLOPS) of the lightweight models adopted in this work.

Evaluation Datasets and Metrics. The revisited Oxford5k [25] and Paris6k [59] datasets are used for evaluation, which are denoted as $\mathcal{R}Oxf$ and $\mathcal{R}Par$ [40]. All datasets describe specific landmarks of buildings under a variety of different observation conditions, each with 70 query images, and 4,993 and 6,322 gallery images, respectively. We follow the common setting [15] to report mAP under the Medium and Hard settings for two datasets. Large-scale experiment results

are further reported with the $\mathcal{R}1M$ (1M distractor images) [40] dataset added to the database.

Implementation Details. Under the *asymmetric image retrieval* setting, gallery models typically use very deep models (e.g., ResNet101 [56]). It is expensive to extract the embeddings of training images during training in terms of computation and memory, especially with such large gallery models. In addition, the gallery model is not optimized during the query model learning. Therefore, our method first extracts all the embedding of training images with the large gallery models offline and caches them in memory.

When SfM-120k is adopted as the training set \mathcal{T} , we follow the settings of AML [15]. The query models are trained on an NVIDIA RTX 3090 GPU for 10 epochs with a batch size of 64. When GLDv2 is adopted, the image size is set to 512 \times 512. Following the setting of DELG [58], random cropping, random color jittering, and random horizontal flipping are used as data augmentation. We train the query models on four NVIDIA RTX 3090 GPUs with a batch size of 256 for 5 epochs. All models are optimized using Adam with an initial learning rate of 10^{-3} and a weight decay of 10^{-6} . A linear decay scheduler is employed to gradually decay the learning rate to 0 when the desired number of steps is reached. τ_g and τ_q were set to 0.1 and 1.0, respectively. For the anchor points generation, the number of centroids K in each subspace is set to 256. The number of subspaces M is set to 64 and 32 when R101-GeM and R101-DELG are adopted as the gallery model, respectively.

During the testing phase, images are resized to a maximum size of 1024 \times 1024 pixels while maintaining their original aspect ratio. Image features are extracted at three scales, namely, $1/\sqrt{2}$, 1, $\sqrt{2}$. We apply L_2 normalization to each scale independently, followed by averaging the features of the three scales and applying another L_2 normalization. Under the *asymmetric image retrieval* setting, we leverage the lightweight query model $\phi_q(\cdot)$ to extract the features of queries and perform retrieval in the gallery, whose features are extracted by a large gallery model $\phi_g(\cdot)$.

B. Comparison with State-of-the-art Methods

mAP Comparison. In Table II, we provide a comprehensive comparison of our proposed approach with state-of-the-art methods on various benchmark datasets. To evaluate the effectiveness of our method under different scenarios, we conduct experiments using different query models, gallery models, and training datasets. We compare the performance of two lightweight query models, MobileNetV2 [44] and EfficientNetB3 [47], two large gallery models with varying performance, R101-GeM [14] and R101-DELG [58], and two training datasets with different sizes, SfM-120k [14] and GLDv2 [39]. The first six rows in Table II illustrate the performance of our method when using both large models (R101-GeM and R101-DELG) and small models (MobileNetV2 and EfficientNetB3) under a symmetrical setting.

We first evaluate our approach using R101-GeM as the gallery model and SfM-120k as the training set. Our method outperforms the most effective solution to *asymmetric image*

METHOD	QUERY MODEL	GALLERY MODEL	TRAINING SET \mathcal{T}	MEDIUM				HARD			
				\mathcal{ROxf}	$\mathcal{ROxf}+\mathcal{R1M}$	\mathcal{RPar}	$\mathcal{RPar}+\mathcal{R1M}$	\mathcal{ROxf}	$\mathcal{ROxf}+\mathcal{R1M}$	\mathcal{RPar}	$\mathcal{RPar}+\mathcal{R1M}$
GeM [†] [14]	R101-GeM	R101-GeM		65.43	45.23	76.75	52.34	40.13	19.92	55.24	24.77
GeM [†] [14]	EfficientNetB3	EfficientNetB3	SfM-120k	54.22	37.10	71.21	44.67	27.53	17.49	48.00	18.45
GeM [†] [14]	MobileNetV2	MobileNetV2		58.81	40.02	67.87	42.25	33.41	17.71	40.97	16.59
DELG [†] [58]	R101-DELG	R101-DELG		78.55	66.02	88.58	73.65	60.89	41.75	76.05	51.46
DELG [†] [58]	EfficientNetB3	EfficientNetB3	GLDv2	66.64	49.67	81.78	61.10	43.82	24.89	63.90	32.34
DELG [†] [58]	MobileNetV2	MobileNetV2		62.42	42.21	77.91	55.09	36.56	18.64	57.96	28.81
Contr* [15]				47.10	18.00	61.50	28.80	21.80	6.30	37.70	8.80
Reg [15]	MobileNetV2	R101-GeM	SfM-120k	49.20	26.50	65.00	34.60	23.30	7.80	40.70	12.70
CSD [19]				63.59	40.29	76.05	43.08	38.51	17.93	52.67	17.43
Ours				63.98	41.07	76.54	45.40	37.91	19.22	53.59	19.02
Contr* [15]				45.20	24.70	63.70	32.80	19.60	12.20	40.90	12.50
Reg [15]	EfficientNetB3	R101-GeM	SfM-120k	52.90	29.70	65.20	39.00	27.80	10.40	42.40	16.00
CSD [19]				64.49	43.39	76.11	45.58	39.06	19.12	53.64	19.78
Ours				65.14	43.95	76.87	48.22	39.38	20.01	54.50	20.64
Contr* [15]				66.42	45.76	83.13	53.10	45.99	23.34	66.79	30.24
Reg [15]	MobileNetV2	R101-DELG	GLDv2	72.75	56.03	85.81	65.23	53.07	32.21	69.96	39.29
HVS [17]				74.39	58.24	86.86	67.44	54.68	34.77	72.42	43.39
LCE [20]				75.45	58.03	87.24	67.30	54.95	33.88	73.03	43.01
CSD [19]				75.94	59.45	87.27	68.52	57.51	36.41	73.45	44.31
Ours				77.88	60.26	88.34	70.23	60.05	37.29	75.08	46.16
Contr* [15]				69.45	49.70	83.81	59.36	46.19	26.49	68.15	35.24
Reg [15]	EfficientNetB3	R101-DELG	GLDv2	74.60	59.88	86.09	67.69	53.41	33.31	72.21	42.63
HVS [17]				76.41	62.72	87.07	71.54	56.13	36.86	74.53	49.09
LCE [20]				75.89	61.90	86.63	70.98	55.21	36.53	73.62	48.94
CSD [19]				77.64	64.29	87.95	72.90	59.32	39.84	75.11	49.13
Ours				79.46	63.22	89.14	73.07	62.17	39.05	76.88	49.54

TABLE II: mAP comparison against existing methods on the full benchmark. Black bold: best results under the same setting. \dagger : our re-implementation. The first six rows illustrate the performance of our method when using both large models (R101-GeM and R101-DELG) and small models (MobileNetV2 and EfficientNetB3) under a symmetrical setting.

retrieval in AML [15], *i.e.*, direct feature regression (Reg), by a large margin. Furthermore, our approach achieves consistently superior or comparable performance compared to CSD [19], which takes neighbor similarity into consideration. For instance, when MobileNetV2 is deployed as the query model, our method still outperforms CSD in most settings, with an mAP improvement of 1.59% on the $\mathcal{RPar} + \mathcal{R1M}$ dataset. It is worth mentioning that in CSD, retrieved real data points are used to calculate neighbor similarity, whereas our method generates a significant number of anchor points in the embedding space of the gallery model, enabling a more detailed characterization of the spatial structure than that obtained using real data points.

Next, we evaluate our method with R101-DELG as the gallery model and GLDv2 as the training set. Our approach achieves better performance than the best previous method in most cases, regardless of whether the query model is MobileNetV2 or EfficientNetB3. *e.g.*, when MobileNetV2 is deployed as the query model, our method outperforms CSD by 2.94% and 2.54% on the \mathcal{ROxf} dataset with Medium and Hard protocols, respectively. Similarly, on the \mathcal{RPar} dataset with Medium and Hard protocols, our method outperforms CSD by 1.07% and 1.63%, respectively. All these results convincingly

demonstrate the superiority of our approach.

While we acknowledge that our approach may perform less favorably than symmetric retrieval when large models are deployed on both the query and gallery sides, it is important to highlight the efficiency of our approach when a smaller model is deployed on the query side. With only 5.8% of the computational FLOPS required by a model like ResNet101 [56], our approach achieves a remarkable 90% performance. Additionally, it is clear from Table II that symmetric retrieval performance suffers when small models are deployed on both sides of the query and the database, and our asymmetric retrieval approach better balances performance and computational complexity in this case.

Discussion about Training Overhead. Both our method and CSD introduce additional time overhead during training. As for CSD [19], it performs retrieval in additional databases to obtain nearest neighbors during each iteration of the training process. In the CSD paper, the authors take the training set as an additional database. When dealing with GLDv2 [39] with 1, 264, 376 images, the time overhead of this extra retrieval step becomes non-negligible. For instance, a single retrieval operation in this scenario incurs a latency of 0.105 seconds, and training 5 epochs with multi-threading acceleration takes

GALLERY MODEL	QUANTIZER TYPE	TRAINING SET \mathcal{T}	MEDIUM		HARD	
			\mathcal{ROxf}	\mathcal{RPar}	\mathcal{ROxf}	\mathcal{RPar}
R101-GeM [14]	k -means _{1,024}	SfM-120k	31.1	51.9	11.6	23.8
	Spectral _{1,024}		30.1	54.5	10.3	25.9
	k -means _{4,096}		40.9	58.3	18.3	30.7
	Spectral _{4,096}		43.7	57.9	20.1	29.4
	k -means _{16,384}		46.7	56.9	21.6	28.5
	Spectral _{16,384}		47.2	59.1	21.8	32.7
	k -means _{65,536}		52.8	60.8	26.9	33.2
PQ _{64 256}	63.9	76.5	37.9	53.5		
R101-DELG [58]	k -means _{1,024}	GLDv2	38.3	56.5	18.9	36.4
	Spectral _{1,024}		35.8	54.2	17.1	37.6
	k -means _{4,096}		59.3	73.9	37.2	53.6
	Spectral _{4,096}		57.1	73.0	36.1	52.9
	k -means _{16,384}		63.4	79.0	39.7	59.7
	Spectral _{16,384}		62.7	79.6	38.4	60.5
	k -means _{65,536}		67.1	81.1	46.0	62.4
PQ _{32 256}	77.8	88.3	60.0	75.0		

TABLE III: mAP (asymmetric) comparison of **different quantizers**. MobileNetV2 [44] is used as query model. “ k -means _{i} ” means a flatten k -means quantizer with i centroids. “Spectral _{i} ” means that we cluster the data with i centroids using spectral clustering. PQ_{32||256} denotes that we split the feature vector into 32 subvectors, with each subvector quantized to 256 centroids. R101-GeM and R101-DELG denote the ResNet101 trained by GeM [14] and DELG [58], respectively.

GALLERY MODEL	SUBSPACE NUMBER M	TRAINING SET \mathcal{T}	MEDIUM		HARD	
			\mathcal{ROxf}	\mathcal{RPar}	\mathcal{ROxf}	\mathcal{RPar}
R101-GeM [14]	2	SfM-120k	51.1	66.6	28.5	37.9
	4		53.5	68.5	29.6	41.1
	8		55.1	69.3	30.5	42.6
	16		56.4	70.9	32.7	45.1
	32		62.0	74.6	35.8	50.9
	64		63.9	76.5	37.9	53.5
R101-DELG [58]	2	GLDv2	68.4	80.8	46.4	64.4
	4		70.2	83.4	49.8	68.2
	8		73.9	87.0	54.8	72.7
	16		77.6	87.2	59.2	73.8
	32		77.8	88.3	60.0	75.0
	64		78.2	88.0	59.8	74.7

TABLE IV: mAP (asymmetric) **comparison of different number M of subspaces**. τ_q and τ_g are set as 1.0 and 0.1, respectively. MobileNetV2 [44] is used as query model. R101-GeM and R101-DELG denote the ResNet101 trained by GeM [14] and DELG [58], respectively.

approximately 24 hours.

In contrast, our approach introduces online computation of structural similarity, which indeed adds some time overhead. Specifically, our method divides features into 64 subvectors, each quantized to 256 clustering centers. The time overhead for computing structural similarity is approximately 16 ms. Note that even with this extra time overhead, based on the structural similarity, our approach costs much less training time than CSD with linear retrieval step. In detail, the total time overhead for training 5 epochs amounts to about 11 hours. Therefore, our method achieves better retrieval performance and costs less training time than CSD.

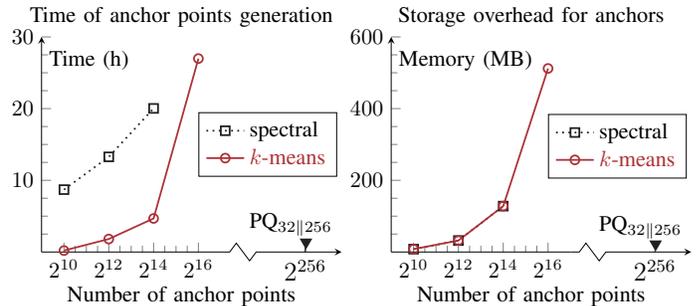


Fig. 3: Comparison of the **time overhead** and **storage complexity** of different anchor generation methods. “ k -means” means a flatten k -means quantizer and “spectral” denotes for spectral clustering. PQ_{32||256} denotes that we split the feature vector into 32 subvectors, with each subvector quantized to 256 centroids.

C. Ablation Study

Comparison with different clustering options. Our method generates anchor points in the embedding space of the gallery model, whose number is related to the granularity of the space division. In this experiment, we compare product quantization with the flat k -means and spectral clustering [60], [61].

As shown in Table III, adopting flatten k -means or spectral clustering leads to severe performance degradation, which is mainly due to the coarse granularity of the space partition. The performance gradually increases as the number of anchor images increases, which shows the need for a large number of anchor points. However, as shown in Figure 3, when the number of required centroids is large, k -means and spectral clustering lead to heavy computations and unaffordable time overheads, making it difficult to further scale up the number of anchor points. This limits the granularity of the partitioning for the space, which makes the relationship between features and anchor points fail to reflect the structure of the space well. Besides, when adopting spectral and k -means clustering, we need to save a large number of centroid vectors. Thus, product quantization is more suitable.

Number of Subspaces. Table IV shows the mAP of our method with different numbers M of subspaces. As M increases, the performance increases at all settings. When M is small, the number of anchor points is small. The division granularity of the gallery embedding space is too coarse to delicately characterize its structure. In contrast, when M is large, the number of equivalent anchor points is large, e.g., when M is 32, it reaches 256^{32} . By constraining the consistency of structure similarities between feature embeddings of the same training sample, the embedding spaces of the query and gallery models are well aligned.

Similarity Type. As shown in Table V, we explore two types of similarities, including negative Euclidean distance and Cosine similarity. When negative Euclidean distance is adopted as the similarity strategy, $s(\cdot, \cdot)$ is formulated as:

$$s(u_i(\mathbf{g}), \mathbf{C}_{i,:}^j) = -\|\mathbf{C}_{i,:}^j - u_i(\mathbf{g})\|_2. \quad (13)$$

“Cosine similarity” leads to better performance. The negative Euclidean distance ranges from 0 to $-\infty$, and the probabilities

GALLERY MODEL	SIMILARITY TYPE	TRAINING SET \mathcal{T}	MEDIUM		HARD	
			\mathcal{R}_{Oxf}	\mathcal{R}_{Par}	\mathcal{R}_{Oxf}	\mathcal{R}_{Par}
R101-GeM [14]	Equation (13)	SfM-120k	56.4	69.7	32.1	44.0
	Equation (7)		63.9	76.5	37.9	53.5
R101-DELG [58]	Equation (13)	GLDv2	65.4	79.6	42.9	59.9
	Equation (7)		77.8	88.3	60.0	75.0

TABLE V: mAP (asymmetric) comparison of **different similarity types**. MobileNetV2 [44] is used as query model. R101-GeM and R101-DELG denote the ResNet101 trained by GeM and DELG, respectively.

GALLERY MODEL	TEMPERATURE τ_g	TRAINING SET \mathcal{T}	MEDIUM		HARD	
			\mathcal{R}_{Oxf}	\mathcal{R}_{Par}	\mathcal{R}_{Oxf}	\mathcal{R}_{Par}
R101-GeM [14]	0.00	SfM-120k	62.7	73.8	35.8	49.0
	0.01		65.0	74.5	38.9	50.7
	0.1		63.9	76.5	37.9	53.5
	0.2		63.5	74.6	36.8	49.7
	0.5		58.7	71.4	31.8	46.0
R101-DELG [58]	0.00	GLDv2	72.6	86.1	53.2	71.4
	0.01		75.6	86.9	57.4	73.3
	0.1		77.8	88.3	60.0	75.0
	0.2		75.0	87.5	55.1	73.9
	0.5		62.8	79.3	45.6	60.7

TABLE VI: Analysis about the **temperature** τ_g . MobileNetV2 [44] is used as query model and τ_g is set to 1.0. R101-GeM and R101-DELG denote the ResNet101 trained by GeM and DELG, respectively.

$p_{i,:}^g$ and $p_{i,:}^q$ obtained after the softmax function do not reflect well the relationship between the feature vectors and the anchor points. On the contrary, the cosine similarity ranges from -1 to 1, which makes the final probability distribution more discriminative.

Soft vs. Hard Assignment. In Table VI, we demonstrate the effect of temperature τ_g in Equation (8). The results for the hard assignment case are denoted by $\tau_g = 0.00$. Choosing a small τ_g , which makes the probability p_i^g sharper (closer to hard assignment), leads to better performance. However, in the extreme case of hard assignment, the performance decreases. Hard assignment ignores the relationship between features and anchor points, which characterizes the space structure and contains more useful knowledge.

Scalability. The structure similarity generated by the gallery model serves as a pseudo-label to supervise the learning of the query model. While the gallery model can be trained using labeled data if available, the query model is exclusively trained through pseudo-labels generated by the gallery model. This means that it is feasible to leverage the vast amount of unlabeled data during the training phase of the query model. In Table VII, we divide the GLDv2 dataset [39] into 10 random splits and train the query models using different amounts of data. In both settings, *e.g.*, using R101-GeM and R101-DELG as the gallery models, the performance gradually improves as the number of training data increases. Our approach does not use any annotations of the training data but only exploits the knowledge provided by the gallery model. Thus, it is possible for our method to improve the performance of the query model

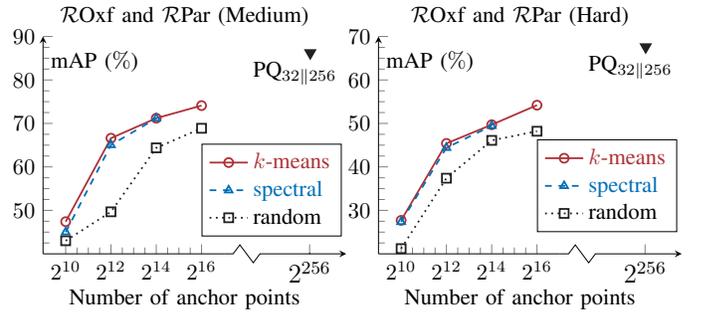


Fig. 4: **Comparison of mAP (asymmetric retrieval) of different methods to generate anchor points.** The mAP is the average of two difficulty settings, medium (left) and high (right). R101-DELG [58] and MobileNetV2 [44] are deployed as the gallery and query models, respectively. “random” denotes for random selecting image features as anchor points. “k-means” means a flatten k -means quantizer and “spectral” refers to spectral clustering of the data. $PQ_{32||256}$ denotes that we split the feature vector into 32 subvectors, with each subvector quantized to 256 centroids.

using a large amount of unlabeled data.

We further validate this in Table VIII, which shows that convincing results are achieved when ImageNet [62] is adopted as the training set. Our approach does not directly regress features but uses the relationships between features and anchor points as knowledge, which somewhat weakens the effect of image data distribution bias. Notably, we achieve remarkable performance when we adopt the distractor set $\mathcal{R}1M$ as the training set, which further illustrates that our approach is able to utilize the available unlabeled data to train query models.

Different methods to generate anchor points. In this section, we explore different ways to generate anchor points such as random selection from gallery features and spectral clustering. As shown in Figure 4, randomly selecting anchor points yields the least favorable results. This observed performance drop may be attributed to the fact that randomly chosen anchor points do not reflect the density of the data distribution, making it difficult for the structural similarity to accurately portray the structural information in the gallery space.

Various Lightweight Models. In this section, we experiment with more lightweight models, whose computational complexity (in FLOPS) is shown in Table I, as query models $\phi_q(\cdot)$. In Figure 6, *symmetric* means that the query and gallery images are both processed using $\phi_q(\cdot)$, while *asymmetric* means that the query and gallery images are processed using $\phi_q(\cdot)$ and $\phi_g(\cdot)$, respectively. The performance becomes better as the model parameters and FLOPS increase under both *asymmetric* and *symmetric* settings. Notably, the performance improvement of the *asymmetric* setting over the *symmetric* setting is more obvious when the number of model parameters is small, *e.g.*, ShuffleNetV2 (0.5 \times), which indicates the advantage of *asymmetric image retrieval* in resource-constrained scenarios. In practical scenarios, it needs to compromise the computational complexity and retrieval accuracy to select an appropriate query model.



Fig. 5: Some successes and failures of our approach on $\mathcal{R}Oxford5k$. In each group, the notation format “query model \rightarrow gallery model” on the top means embedding queries with the query model and retrieving in a gallery set embedded by the gallery model. The image with the orange border on the left is a query, the first row shows the results of symmetric retrieval when the large gallery model deployed on both query and gallery sides. The second row shows the results of our method under asymmetric setting. The third row shows the results of symmetric retrieval when the lightweight query model is deployed on both sides. Images with green borders are true positive images and images with red borders are false positive images.

GALLERY MODEL	TRAINING SET \mathcal{T}	IMAGE NUMBERS	MEDIUM		HARD	
			$\mathcal{R}Oxf$	$\mathcal{R}Par$	$\mathcal{R}Oxf$	$\mathcal{R}Par$
R101-GeM [14]	GLDv2 ($\times 0.1$)	128,078	63.5	75.6	37.2	53.6
	GLDv2 ($\times 0.2$)	256,156	64.1	76.8	38.5	54.3
	GLDv2 ($\times 0.3$)	384,234	64.2	77.6	39.2	56.7
	GLDv2 ($\times 0.4$)	512,312	64.9	77.3	40.1	56.2
R101-DELG [58]	GLDv2 ($\times 0.1$)	128,078	75.0	86.3	55.1	71.9
	GLDv2 ($\times 0.2$)	256,156	76.7	87.0	58.1	72.9
	GLDv2 ($\times 0.3$)	384,234	77.1	87.3	59.3	73.4
	GLDv2 ($\times 0.4$)	512,312	77.0	87.6	59.4	74.1

TABLE VII: mAP (*asymmetric*) comparison of **different dataset size**. ($\times x$) denotes the small dataset formed by randomly selecting x proportion of images from the full GLDv2 dataset [39]. MobileNetV2 [44] is used as query model. R101-GeM and R101-DELG denote the ResNet101 trained by GeM and DELG, respectively.

Visualization of Retrieval Results. In Figure 5, we show some examples of the success and failure of our approach. According to the retrieval results of queries 1 and 2, deploying large gallery models on both the query and gallery sides results in the highest retrieval performance. However, it is important to acknowledge that in resource-constrained scenarios, such as on mobile devices, deploying large models may not be feasible. When lightweight models are used on both sides, there is a noticeable degradation in retrieval performance. Our approach takes an asymmetric approach, where lightweight,

GALLERY MODEL	TRAINING SET \mathcal{T}	IMAGE NUMBERS	MEDIUM		HARD	
			$\mathcal{R}Oxf$	$\mathcal{R}Par$	$\mathcal{R}Oxf$	$\mathcal{R}Par$
R101-GeM [14]	SfM-120k	91,642	63.9	76.5	37.9	53.5
	GLDv2	1,280,787	65.2	77.5	40.1	56.8
	$\mathcal{R}1M$	1,001,001	64.4	76.8	39.6	55.6
	ImageNet	1,281,167	57.8	74.4	31.9	52.6
R101-DELG [58]	SfM-120k	91,642	75.4	84.4	55.1	68.3
	GLDv2	1,280,787	77.8	88.3	60.0	75.0
	$\mathcal{R}1M$	1,001,001	75.5	86.1	56.9	72.5
	ImageNet	1,281,167	56.6	76.4	38.3	59.0

TABLE VIII: mAP (*asymmetric*) comparison of **different training datasets**. MobileNetV2 [44] is used as query model. R101-GeM and R101-DELG denote the ResNet101 trained by GeM and DELG, respectively.

smaller models are deployed on the query side while high-performing large models are used on the gallery side. Additionally, we train the query model to be compatible with the gallery model, striking a balance between computational complexity and retrieval performance.

However, it is worth noting that there are instances where retrieval results remain sub-optimal even when large gallery models are deployed on both sides. These cases represent challenges that our approach, or any method, may face. The small model trained by our method is constrained to maintain structural similarity with the gallery model, and there are scenarios where it may struggle with certain query images.

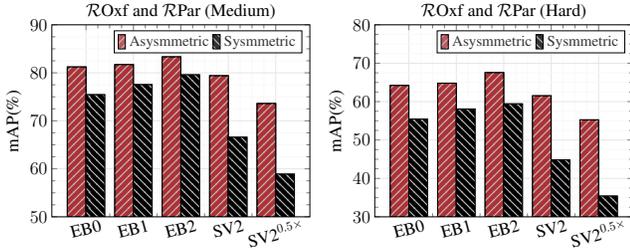


Fig. 6: Analysis of **different model variants**. We use R101-DELG [58] as gallery model $\phi_g(\cdot)$ and compare different architectures as query models $\phi_q(\cdot)$. *Symmetric*: Query and gallery images are both embedded by $\phi_q(\cdot)$; *Asymmetric*: Query and gallery images are embedded by $\phi_q(\cdot)$ and $\phi_g(\cdot)$, respectively. EB0: EfficientNetB0 [47]; EB1: EfficientNetB1 [47]; EB2: EfficientNetB2 [47]; SV2: ShuffleNetV2 [46]; SV2^{0.5x}: ShuffleNetV2 (0.5 \times) [46].

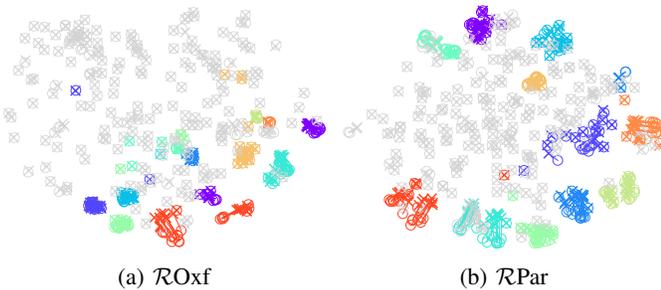


Fig. 7: **T-SNE embeddings of ROxf and RPar datasets**. MobileNetV2 [44] and R101-GeM [14] are used as query and gallery models. Different colors represent different buildings and gray denotes distractor images. We randomly select 10 samples for each building category and 100 in distractors. \circ and \times denote gallery and query models, respectively. A line connects the two representatives of each example.

Qualitative Results. Figure 7 shows the embeddings of some ROxf and RPar images, each processed by a gallery and a query model. For *asymmetric image retrieval*, it is crucial to keep the feature compatibility between query and model models. During training, anchor points are shared by both query and gallery models. We restrict the similarities between

METHOD	RETRIEVAL LATENCY (MS)	MEMORY (MB)	
		ROxf + R1M	RPar + R1M
PQ ₈ ₂₅₆	29.69	7.68	7.69
PQ ₁₆ ₂₅₆	30.57	15.35	15.37
PQ ₃₂ ₂₅₆	32.97	30.70	30.75
PQ ₆₄ ₂₅₆	34.26	61.40	61.50
PQ ₁₂₈ ₂₅₆	41.15	122.80	122.90
PQ ₂₅₆ ₂₅₆	51.87	245.60	245.98
No quantization	152.12	7,782.40	7,783.70

TABLE IX: **Time and memory cost.** We report average search time on a single thread CPU (Intel Xeon CPU E5-2640 v4 @ 2.40GHz) and memory consumption for the gallery sets with 1M distractor images.

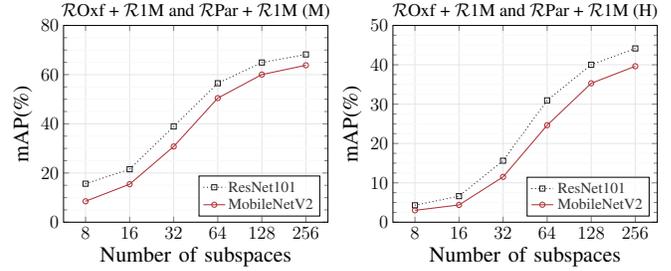


Fig. 8: mAP comparison when PQ [41] is used during **asymmetric and symmetric retrieval**. Horizontal axis denotes that we split the feature into M subvectors. The number of centroids per subvector is set to 256. R101-DELG [58] is deployed as the gallery model. We report the average performance of our method on ROxf + R1M and RPar + R1M datasets with Medium (M) and Hard (H) protocols [40].

two features of the same training sample and anchor points to be consistent, which keeps the structure similarity.

Memory vs. Search Accuracy. In Figure 8, we adopt PQ during the online retrieval, with the corresponding retrieval latency and memory consumption shown in Table IX. PQ is parametrized by the number of sub-vectors M and the number of quantizers per sub-vector K , producing a code of length $M \times \log_2 K$. As M increases, the accuracy of retrieval gradually approximates the direct feature comparison. When $M = 256$, quantization saves 96.8% of memory and 65.9% of retrieval latency with slight performance degradation. In real-world applications, we choose the appropriate M to achieve the performance-memory trade-off.

V. CONCLUSION

In this paper, we propose a structure similarity preserving approach to achieve feature consistency between query and gallery models for asymmetric retrieval. First, we employ product quantization to generate a large number of anchor points in the embedding space of the gallery to characterize its space structure. Then, these anchor points are shared between query and gallery models. The relationships between each training sample and anchor points are considered as structure similarity and constrained to be consistent across different models. This allows the query model to focus less on the feature “details” of the gallery model and more on the overall space structure. Besides, our method does not utilize any annotation from training set, and it is possible for the proposed method to utilize large-scale unlabeled training data, even from different domains. This shows the generalizability of our approach. Extensive experiments show that our method achieves better performance than state-of-the-art asymmetric retrieval methods.

REFERENCES

- [1] O. Simeoni, Y. Avrithis, and O. Chum, “Local features and visual words emerge in activations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1
- [2] A. Babenko, A. Slesarev, A. Chigorin, and V. Lempitsky, “Neural codes for image retrieval,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 584–599. 1

- [3] H. Wu, M. Wang, W. Zhou, and H. Li, "Learning deep local features with multiple dynamic attentions for large-scale image retrieval," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 11416–11425. **1**
- [4] G. Toliás, R. Sicre, and H. Jégou, "Particular Object Retrieval With Integral Max-Pooling of CNN Activations," in *International Conference on Learning Representations (ICLR)*, 2016, pp. 1–12. **1, 2**
- [5] P. Lu, G. Huang, H. Lin, W. Yang, G. Guo, and Y. Fu, "Domain-aware se network for sketch-based image retrieval with multiplicative euclidean margin softmax," in *Proceedings of the ACM international conference on Multimedia (MM)*, 2021, pp. 3418–3426. **1**
- [6] J. Tian, X. Xu, Z. Wang, F. Shen, and X. Liu, "Relationship-preserving knowledge distillation for zero-shot sketch based image retrieval," in *Proceedings of the ACM international conference on Multimedia (MM)*, 2021, pp. 3418–3426. **1**
- [7] N. Jiang, B. Sheng, P. Li, and T.-Y. Lee, "Photohelper: Portrait photographing guidance via deep feature retrieval and fusion," *IEEE Transactions on Multimedia (TMM)*, pp. 1–1, 2022. **1**
- [8] S. Pang, J. Ma, J. Xue, J. Zhu, and V. Ordóñez, "Deep feature aggregation and image re-ranking with heat diffusion for image retrieval," *IEEE Transactions on Multimedia (TMM)*, vol. 21, no. 6, pp. 1513–1523, 2019. **1**
- [9] I. González-Díaz, M. Birinci, F. Díaz-de María, and E. J. Delp, "Neighborhood matching for image retrieval," *IEEE Transactions on Multimedia (TMM)*, vol. 19, no. 3, pp. 544–558, 2017. **1**
- [10] J. Ouyang, W. Zhou, M. Wang, Q. Tian, and H. Li, "Collaborative image relevance learning for visual re-ranking," *IEEE Transactions on Multimedia (TMM)*, vol. 23, pp. 3646–3656, 2021. **1**
- [11] L. Zheng, S. Wang, Z. Liu, and Q. Tian, "Fast image retrieval: Query pruning and early termination," *IEEE Transactions on Multimedia (TMM)*, vol. 17, no. 5, pp. 648–659, 2015. **1**
- [12] R. Zhou, X. Chang, L. Shi, Y.-D. Shen, Y. Yang, and F. Nie, "Person reidentification via multi-feature fusion with adaptive graph learning," *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, vol. 31, no. 5, pp. 1592–1601, 2020. **1**
- [13] G. Toliás, T. Jeníček, and O. Chum, "Learning and aggregating deep local descriptors for instance-level recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 460–477. **1**
- [14] F. Radenović, G. Toliás, and O. Chum, "Fine-tuning cnn image retrieval with no human annotation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 41, no. 7, pp. 1655–1668, 2019. **1, 2, 6, 7, 8, 9, 10, 11**
- [15] M. Budnik and Y. Avrithis, "Asymmetric metric learning for knowledge transfer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8228–8238. **1, 2, 6, 7**
- [16] R. Duggal, H. Zhou, S. Yang, Y. Xiong, W. Xia, Z. Tu, and S. Soatto, "Compatibility-aware heterogeneous visual search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 10723–10732. **1, 2**
- [17] Y. Shen, Y. Xiong, W. Xia, and S. Soatto, "Towards backward-compatible representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. **1, 2, 7**
- [18] B. Zhang, Y. Ge, Y. Shen, S. Su, F. Wu, C. Yuan, X. Xu, Y. Wang, and Y. Shan, "Towards universal backward-compatible representation learning," in *Proceedings of the International Joint Conference on Artificial Intelligence, IJCAI*, 2022, pp. 1615–1621. **1**
- [19] H. Wu, M. Wang, W. Zhou, H. Li, and Q. Tian, "Contextual similarity distillation for asymmetric image retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 9489–9498. **2, 7**
- [20] Q. Meng, C. Zhang, X. Xu, and F. Zhou, "Learning compatible embeddings," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 9939–9948. **2, 7**
- [21] B. Zhang, Y. Ge, Y. Shen, Y. Li, C. Yuan, X. Xu, Y. Wang, and Y. Shan, "Hot-refresh model upgrades with regression-alleviating compatible training in image retrieval," in *International Conference on Learning Representations (ICLR)*, 2022. **2**
- [22] B. Zhang, S. Su, Y. Ge, X. Xu, Y. Wang, C. Yuan, M. Z. Shou, and Y. Shan, "Darwinian model upgrades: Model evolving with selective compatibility," *arXiv preprint arXiv:2210.06954*, 2022. **2**
- [23] Y. Bai, J. Jiao, Y. Lou, S. Wu, J. Liu, X. Feng, and L.-Y. Duan, "Dual-tuning: Joint prototype transfer and structure regularization for compatible feature learning," *IEEE Transactions on Multimedia (TMM)*, pp. 1–13, 2022. **2**
- [24] S. Wu, L. Chen, Y. Lou, Y. Bai, T. Bai, M. Deng, and L.-Y. Duan, "Neighborhood consensus contrastive learning for backward-compatible representation," *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 2722–2730, 2022. **2**
- [25] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007, pp. 1–8. **2, 6**
- [26] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2003, pp. 1470–1470. **2**
- [27] W. Zhou, Y. Lu, H. Li, Y. Song, and Q. Tian, "Spatial coding for large scale partial-duplicate web image search," in *Proceedings of the ACM international conference on Multimedia (MM)*, 2010, pp. 511–520. **2**
- [28] G. Toliás, Y. Avrithis, and H. Jégou, "To aggregate or not to aggregate: Selective match kernels for image search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 1401–1408. **2**
- [29] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision (IJCV)*, pp. 91–110, 2004. **2**
- [30] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded up robust features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2006, pp. 404–417. **2**
- [31] F. Perronnin, Y. Liu, J. Sánchez, and H. Poirier, "Large-scale image retrieval with compressed fisher vectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010, pp. 3384–3391. **2**
- [32] H. Jégou, M. Douze, J. Sánchez, P. Pérez, and C. Schmid, "Aggregating local image descriptors into compact codes," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 34, no. 9, pp. 1704–1716, 2012. **2**
- [33] O. Chum, J. Philbin, J. Sivic, M. Isard, and A. Zisserman, "Total Recall: Automatic query expansion with a generative feature model for object retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007, pp. 1–8. **2**
- [34] Y. Fang, W. Zhou, Y. Lu, J. Tang, Q. Tian, and H. Li, "Cascaded feature augmentation with diffusion for image retrieval," in *Proceedings of the ACM international conference on Multimedia (MM)*, 2018, pp. 1644–1652. **2**
- [35] J. Revaud, J. Almazán, R. S. Rezende, and C. R. d. Souza, "Learning with average precision: Training image retrieval with a listwise loss," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5107–5116. **2**
- [36] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. **2**
- [37] A. Babenko and V. Lempitsky, "Aggregating deep convolutional features for image retrieval," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1269–1277. **2**
- [38] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 685–701. **2**
- [39] T. Weyand, A. Araujo, B. Cao, and J. Sim, "Google landmarks dataset v2 a large-scale benchmark for instance-level recognition and retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2575–2584. **2, 6, 7, 9, 10**
- [40] F. Radenović, A. Iscen, G. Toliás, Y. Avrithis, and O. Chum, "Revisiting oxford and paris: Large-scale image retrieval benchmarking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5706–5715. **2, 6, 11**
- [41] H. Jégou, M. Douze, and C. Schmid, "Product quantization for nearest neighbor search," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 33, no. 1, pp. 117–128, 2011. **3, 4, 11**
- [42] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "SqueezeNet: AlexNet-level accuracy with 50× fewer parameters and <0.5 mb model size," *arXiv preprint arXiv:1602.07360*, 2016. **3**
- [43] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017. **3**
- [44] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. **3, 6, 8, 9, 10, 11**

- [45] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6848–6856. [3](#)
- [46] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "ShuffleNet v2: Practical guidelines for efficient cnn architecture design," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 116–131. [3](#), [6](#), [11](#)
- [47] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning (ICML)*, 2019, pp. 6105–6114. [3](#), [6](#), [11](#)
- [48] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015. [3](#)
- [49] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *International Conference on Learning Representations (ICLR)*, 2015. [3](#)
- [50] N. Passalis and A. Tefas, "Learning deep representations with probabilistic knowledge transfer," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 268–284. [3](#)
- [51] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [3](#)
- [52] Y. Chen, N. Wang, and Z. Zhang, "DarkRank: Accelerating deep metric learning via cross sample similarities transfer," *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2018. [3](#)
- [53] Y. Tian, D. Krishnan, and P. Isola, "Contrastive representation distillation," in *International Conference on Learning Representations (ICLR)*, 2019. [3](#)
- [54] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *International Conference on Learning Representations (ICLR)*, 2016. [3](#)
- [55] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [3](#)
- [56] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [6](#), [7](#)
- [57] B. Thomee, D. A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L.-J. Li, "Yfcc100m: The new data in multimedia research," *Communications of the ACM (CACM)*, pp. 64–73, 2016. [6](#)
- [58] B. Cao, A. Araujo, and J. Sim, "Unifying deep local and global features for image search," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 726–743. [6](#), [7](#), [8](#), [9](#), [10](#), [11](#)
- [59] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Lost in quantization: Improving particular object retrieval in large scale image databases," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008, pp. 1–8. [6](#)
- [60] Z. Li, F. Nie, X. Chang, Y. Yang, C. Zhang, and N. Sebe, "Dynamic affinity graph construction for spectral clustering using multiple features," *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, vol. 29, no. 12, pp. 6323–6332, 2018. [8](#)
- [61] Z. Li, F. Nie, X. Chang, L. Nie, H. Zhang, and Y. Yang, "Rank-constrained spectral clustering with flexible embedding," *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, vol. 29, no. 12, pp. 6073–6082, 2018. [8](#)
- [62] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255. [9](#)



Hui Wu is currently pursuing the Ph.D. degree in information and communication engineering with the Department of Data Science, from the University of Science and Technology of China.

His research interests include image retrieval, multimedia information retrieval and computer vision.



Min Wang received the B.E., and Ph.D degrees in electronic information engineering from University of Science and Technology of China (USTC), in 2014 and 2019, respectively. She is working in Institute of Artificial Intelligence, Hefei Comprehensive National Science Center.

Her current research interests include binary hashing, multimedia information retrieval and computer vision.



Wengang Zhou (S'20) received the B.E. degree in electronic information engineering from Wuhan University, China, in 2006, and the Ph.D. degree in electronic engineering and information science from the University of Science and Technology of China (USTC), China, in 2011. From September 2011 to September 2013, he worked as a postdoc researcher in Computer Science Department at the University of Texas at San Antonio. He is currently a Professor at the EEIS Department, USTC.

His research interests include multimedia information retrieval, computer vision, and computer game. In those fields, he has published over 100 papers in IEEE/ACM Transactions and CCF Tier-A International Conferences. He is the winner of National Science Funds of China (NSFC) for Excellent Young Scientists. He is the recipient of the Best Paper Award for ICIMCS 2012. He received the award for the Excellent Ph.D Supervisor of Chinese Society of Image and Graphics (CSIG) in 2021, and the award for the Excellent Ph.D Supervisor of Chinese Academy of Sciences (CAS) in 2022. He won the First Class Wu-Wenjun Award for Progress in Artificial Intelligence Technology in 2021. He served as the publication chair of IEEE ICME 2021 and won 2021 ICME Outstanding Service Award. He is currently an Associate Editor and a Lead Guest Editor of IEEE Transactions on Multimedia.



Houqiang Li (S'12, F'21) received the B.S., M.Eng., and Ph.D. degrees in electronic engineering from the University of Science and Technology of China, Hefei, China, in 1992, 1997, and 2000, respectively, where he is currently a Professor with the Department of Electronic Engineering and Information Science.

His research interests include image/video coding, image/video analysis, computer vision, reinforcement learning, etc.. He has authored and co-authored over 200 papers in journals and conferences. He is the winner of National Science Funds (NSFC) for Distinguished Young Scientists, the Distinguished Professor of Changjiang Scholars Program of China, and the Leading Scientist of Ten Thousand Talent Program of China. He is the associate editor (AE) of IEEE TMM, and served as the AE of IEEE TCSVT from 2010 to 2013. He served as the General Co-Chair of ICME 2021 and the TPC Co-Chair of VCIP 2010. He received the second class award of China National Award for Technological Invention in 2019, the second class award of China National Award for Natural Sciences in 2015, and the first class prize of Science and Technology Award of Anhui Province in 2012. He received the award for the Excellent Ph.D Supervisor of Chinese Academy of Sciences (CAS) for four times from 2013 to 2016. He was the recipient of the Best Paper Award for VCIP 2012, the recipient of the Best Paper Award for ICIMCS 2012, and the recipient of the Best Paper Award for ACM MUM in 2011.