

Deep Learning for Instance Retrieval: A Survey

Wei Chen, Yu Liu, Weiping Wang, Erwin M. Bakker, Theodoros Georgiou,
Paul Fieguth, Li Liu, and Michael S. Lew

Abstract—In recent years a vast amount of visual content has been generated and shared from many fields, such as social media platforms, medical imaging, and robotics. This abundance of content creation and sharing has introduced new challenges, particularly that of searching databases for similar content — Content Based Image Retrieval (CBIR) — a long-established research area in which improved efficiency and accuracy are needed for real-time retrieval. Artificial intelligence has made progress in CBIR and has significantly facilitated the process of instance search. In this survey we review recent instance retrieval works that are developed based on deep learning algorithms and techniques, with the survey organized by deep feature extraction, feature embedding and aggregation methods, and network fine-tuning strategies. Our survey considers a wide variety of recent methods, whereby we identify milestone work, reveal connections among various methods and present the commonly used benchmarks, evaluation results, common challenges, and propose promising future directions.

Index Terms—Content based image retrieval, Instance retrieval, Deep learning, Convolutional neural networks, Literature survey

INTRODUCTION

CONTENT Based Image Retrieval (CBIR) is the problem of searching for relevant images in an image gallery by analyzing the visual content (colors, textures, shapes, objects *etc.*), given a query image [1],[2]. CBIR has been a longstanding research topic in the fields of computer vision and multimedia [1],[2]. With the exponential growth of image data, the development of appropriate information systems that efficiently manage such large image collections is of utmost importance, with image searching being one of the most indispensable techniques. Thus there is a nearly endless potential for applications of CBIR, such as person/vehicle reidentification [3],[4], landmark retrieval [5], remote sensing [6], online product searching [7].

Generally, CBIR methods can be grouped into two different tasks [8],[9]: Category level Image Retrieval (CIR) and Instance level Image Retrieval (IIR). The goal of CIR is to find an arbitrary image representative of the same category as the query (*e.g.*, dogs, cars) [10],[11]. By contrast, in the IIR task, a query image of a particular instance (*e.g.*, the Eiffel Tower, my neighbor's dog) is given and the goal is to find images containing the same instance that may be captured under different conditions like different imaging distances, viewing angles, backgrounds, illuminations, and weather conditions (reidentifying exemplars of the same instance) [12],[13]. The focus of this survey is the IIR task¹.

In many real world applications, IIR is usually to find a desired image requiring a search among thousands, millions, or even billions of images. Hence, searching efficiently is as critical

as searching accurately, to which continued efforts have been devoted [12],[14],[15]. To enable accurate and efficient retrieval over a large-scale image collection, *developing compact yet discriminative feature representations* is at the core of IIR.

During the past two decades, startling progress has been witnessed in image representation which mainly consists of two important periods, *i.e.*, feature engineering and deep learning. In the feature engineering era, the field was dominated by various milestone handcrafted image representations like SIFT [23] and Bag of visual Words (BoW) [24]. The deep learning era was reignited in 2012 when a Deep Convolutional Neural Network (DCNN) referred as “AlexNet” [25] won the first place in the ImageNet classification contest with a breakthrough reduction in classification error rate. Since then, the dominating role of SIFT like local descriptors has been replaced by data driven Deep Neural Networks (DNNs) which can learn powerful feature representations with multiple levels of abstraction directly from data. During the past decade, DNNs have set the state of the art in various classical computer vision tasks, including image classification [25],[26], object detection [27], semantic segmentation [28], and image retrieval [17].

Given this period of rapid evolution, the goal of this paper is to provide a comprehensive survey of the recent achievements in IIR. In comparison with existing excellent surveys on traditional image retrieval [13],[18],[19],[20], as contrasted in Table 1, our focus in this paper is reviewing deep learning based methods for IIR, particularly on issues of retrieval accuracy and efficiency.

1.1 Summary of Progress since 2012

After the highly successful image classification implementation of AlexNet [25], significant exploration of DCNNs for instance retrieval tasks was undertaken and representative methods are shown in Figure 1. Building on these methods, more recent progress for IIR can be achieved from the perspectives of off-the-shelf models and fine-tuned models, which form the basis for this survey.

Off-the-shelf models, based on DCNNs with fixed parameters [29],[30],[31], extract features at image scales or patch scales, which correspond to single-pass and multiple-pass schemes, respectively. These methods focus on effective feature use, for which researchers have proposed embedding and aggregation

W. Chen is with Academy of Advanced Technology Research of Hunan, Changsha, China

Y. Liu is with DUTRU International School of Information Science and Engineering, Dalian University of Technology, China.

W. Wang is with Academy of Advanced Technology Research of Hunan, Changsha. E. Bakker, T. Georgiou, and M. Lew are with Leiden Institute of Advanced Computer Science, Leiden University, the Netherlands.

P. Fieguth is with the Department of Systems Design Engineering, University of Waterloo, Canada.

L. Liu is with Academy of Advanced Technology Research of Hunan, Changsha, China, and also with Center for Machine Vision and Signal Analysis, University of Oulu, Finland.

Corresponding author: Li Liu, li.liu@oulu.fi, dreamliu2010@gmail.com

This work was supported by China Scholarship Council (No. 201703170183), the Academy of Finland under grant 331883, Infotech Project FRAGES, and the National Natural Science Foundation of China under Grant 61872379, 62022091, 61825305, and 62102061.

1. If not further specified, “image retrieval”, “IIR”, and “instance retrieval” are considered equivalent and will be used interchangeably.

TABLE 1: A summary and comparison of the primary surveys in the field of image retrieval.

Title	Year	Published in	Main Content
Content-Based Image Retrieval at the End of the Early Years [1]	2000	TPAMI	This paper discusses the steps for image retrieval systems, including image processing, feature extraction, user interaction, and similarity evaluation.
Image Search from Thousands to Billions in 20 Years [16]	2013	TOMM	This paper gives a good presentation of image search achievements from 1970 to 2013, but the methods are not deep learning-based.
Deep Learning for Content-Based Image Retrieval: A Comprehensive Study [17]	2014	ACM MM	This paper introduces supervised metric learning methods for fine-tuning AlexNet. Details of instance-based image retrieval are limited.
Semantic Content-based Image Retrieval: A Comprehensive Study [18]	2015	JVCI	This paper presents a comprehensive study about CBIR using traditional methods; deep learning is introduced as a section with limited details.
Socializing the Semantic Gap: A Comparative Survey on Image Tag Assignment, Refinement, and Retrieval [19]	2016	CSUR	A taxonomy is introduced to structure the growing literature of image retrieval. Deep learning methods for feature learning is introduced as future work.
Recent Advance in Content based Image Retrieval: A Literature Survey [20]	2017	arXiv	This survey presents image retrieval from 2003 to 2016. Neural networks are introduced in a section and mainly discussed as a future direction.
Information Fusion in Content-based Image Retrieval: A Comprehensive Overview [21]	2017	Information Fusion	This paper presents information fusion strategies in CBIR. Deep convolutional networks for feature learning are introduced briefly but not covered thoroughly.
A Survey on Learning to Hash [22]	2018	TPAMI	This paper focuses on hash learning algorithms and introduces the similarity-preserving methods and discusses their relationships.
SIFT Meets CNN: A Decade Survey of Instance Retrieval [13]	2018	TPAMI	This paper presents a comprehensive review of instance retrieval based on SIFT and CNN methods.
Deep Learning for Instance Retrieval: A Survey	2021	Ours	Our survey focuses on deep learning methods. We expand the review with indepth details on IIR, including methods of feature extraction, feature embedding and aggregation, and network fine-tuning.

methods, such as R-MAC [32], CroW [15], and SPoC [12] to promote the discriminativity of the extracted features. Fine-tuned models, based on DCNNs in which parameters are updated [29], are more popular since deep networks themselves have been investigated extensively. To learn better retrieval features, researchers have proposed to improve the network architectures and/or update the pre-stored parameters [31].

This survey will explore recent progress in detail in the context of the three following themes:

(1) *Deep Feature Extraction* (Section 3.1)

One key step in IIR is to make the descriptors as semantic-aware [26],[49] as possible. For this, some recent work focus on the input data of DCNNs, thereby instance features can be extracted from the whole image, *e.g.*, CroW [15], VLAD-CNN [50] or from image patches, *e.g.*, MOP-CNN [30], FAemb [38]. For instance, evaluated on the Holidays dataset [51], patch-level input scheme can improve mAP by 8.29% compared to the results (70.53%) achieved using image-level input [30]. Others focus on exploring different feature extractors, *e.g.*, one layer of a given DCNN, to get the output activations. Initially, fully-connected layers are usually chosen to extract features [52],[53], and then convolutional layers are popularly used [12],[32]. Afterwards, some work leverage the complementarity of different extractors to explore layer-level fusion, such as MoF [36], and model-level fusion, such as DeepIndex [52] to promote the retrieval performance.

(2) *Feature Embedding and Aggregation* (Section 3.2)

Recent work revisit the classical embedding and aggregation methods and apply on deep features. Most work have a preference towards mapping individual vector from convolutional layer [24],[54],[55] and then aggregating into a global feature. The mapping process can be realized by using a pre-trained codebook (*i.e.*, built separately), such as VLAD-CNN [50], DeepIndex [52] or learned as parameters during training (built simultaneously), such as NetVLAD [43], GeM-DSM [56]. Some work aggregate local features into a global one by direct pooling [21] or sophisticated pooling-based methods [12] without the aggregation operation, such as R-MAC [32].

(3) *Network Fine-tuning for Learning Representations* (Section 4)

DCNNs pretrained on source datasets for image classification are influenced by domain shifts when performing retrieval tasks on new datasets. Thus, it is necessary to fine-tune deep networks to the specific domain [39], by using supervised or unsupervised fine-tuning methods. As depicted in Figure 1, recent supervised fine-tuning methods focus on designing objective functions (*e.g.*, Listwise loss [57]) and sample sampling strategies, such as NetVLAD [43], Triplet Network [42]. Unsupervised methods focus on mining the relevance among training samples by using clustering, such as SfM-GeM [46] or manifold learning, such as AILIR [58]. Recently, convolution-free models that only rely on transformer layers have shown competitive performance and been used as a powerful alternative to DCNNs, such as IRT [59], reranking Transformers [60].

1.2 Key Challenges

The goal of each of the preceding three themes is to address the competing objectives of *accuracy* and *efficiency*, with both objectives continuing to present challenges:

A) Accuracy related challenges depend on the input data, the feature extractor, and the way in which the extracted features are processed:

- **Invariance challenge:** The instance in an image may be rotated, translated, or scaled differently, so the final features are affected by these transformations and retrieval accuracy may be degraded [30]. It is necessary to incorporate invariance into the IIR pipeline [61],[62].
- **Distraction challenge:** IIR systems may need to focus on only a certain object, or even only a small portion. DCNNs may be affected by image clutter or background, so multiple-pass schemes have been examined where region proposals are studied before feature extraction.
- **Discriminativity challenge:** Deep features for IIR are needed to be as discriminative as possible to distinguish instances with small differences, leading to many explorations in feature processing. These include feature embedding and aggregation methods, to promote feature discriminativity; and attention mechanisms, to highlight the most relevant regions within the extracted features or to enable deep networks to focus on regions of interest.

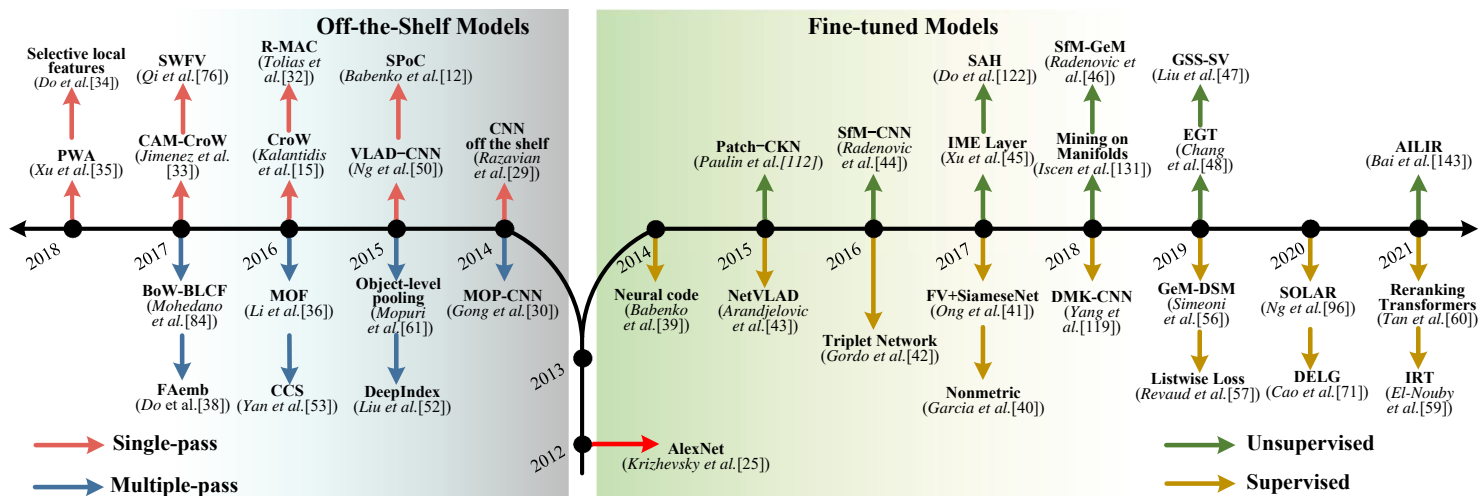


Fig. 1: Representative methods in IIR. Off-the-shelf models have model parameters which are not further updated or tuned when extracting retrieval features. For single-pass schemes, the key step is the feature embedding and aggregation to promote the discriminativity of the extracted image-level activations [15],[33],[34],[35], whereas for multiple-pass schemes the goal is to extract instance features at region scales and eliminate image clutter as much as possible [30],[36],[37],[38]. In contrast, for fine-tuned models, the model parameters are tuned towards the retrieval task and address the issue of domain shifts. For supervised fine-tuning, the key step lies in the design of objective functions and sample sampling strategies [39],[40],[41],[42],[43], while the success of unsupervised fine-tuning is to mine the relevance among training samples [44],[45],[46],[47],[48]. See Sections 3 and 4 for details.

- Fine-tune challenge: DCNNs can be fine-tuned as powerful extractors to capture fine semantic distinctions among instances. These explorations have offered improved accuracy, however the area remains a major challenge.

B) Efficiency related challenges are important, especially for large-scale datasets [63]. Retrieval systems should respond quickly when given a query image. Deep features are high-dimensional and contain semantic-aware information to support higher accuracy, yet is often at the expense of efficiency.

On the one hand, the efficiency is related to the format of features, *i.e.*, real valued or binary. Hash codes have advantages in storage and searching [22],[39], however for hashing methods one needs to carefully consider the loss function design [64],[65], to obtain optimal codes for high retrieval accuracy.

On the other hand, efficiency is also related to the mechanism of feature matching. For example, instead of time-consuming cross-matching between local features, one can choose to use global features to perform an initial ranking and then a post-step re-ranking via the features of top-ranked images.

2 GENERAL FRAMEWORK OF IIR

Figure 2 offers an overview of the general framework for deep-learning-based IIR, involving three main stages.

1) Deep feature extraction: (Section 3.1)

Feature extraction is the first step of IIR and can be realized in a single-pass or multiple-pass way. Single-pass methods take as input the whole image, whereas multiple-pass methods depend on region extraction, as depicted in Figure 4.

The activations from fully-connected layers of a given DCNN can be used as retrieval features whether based on a whole image or on patches. The tensors from convolutional layers can be used when further processed by sophisticated pooling, as shown in Figure 2. Different layers of the same deep network can be combined as a more powerful extractor [36],[66]. Furthermore, it is possible to fuse the activations from layers of different models [67],[68]. Feature extraction is the step to produce vanilla

network activations (*i.e.*, 3D tensors or a single vector), these activations, in most cases, are needed to be further processed.

2) Embedding and aggregation: (Section 3.2)

Feature embedding and aggregation are two essential steps to produce global or local features. Feature embedding maps individual local features into higher-dimensional space whereas feature aggregation summarizes the multiple mapped vectors or all individual features into a global vector. Global features may come from pooling convolutional feature maps directly [69],[70] or using some sophisticated weighting methods [12],[15] (*i.e.*, both without feature embedding). Feature embedding method using a pre-generated codebook can be performed to encode individual convolutional vectors and then aggregated [24],[54],[55]. For local features, the well-embedded representations for all regions of interest are stored individually and used for cross-matching in the reranking stage without aggregation.

3) Feature matching:

Feature matching is a process to measure the feature similarity between images and then return a ranked list. Global matches can be computed efficiently via such as Euclidean distance. For local features [5],[71], the image similarity is usually evaluated by summarizing the similarities across local features, using classical RANSAC [72] or more recent variations [73],[74]. Storing local features separately and then estimating their similarity individually lead to additional memory and search costs [71],[74], therefore in most cases local features are used to re-rank the initial ranking image matched by global features [32],[64],[71],[75].

The three preceding stages for IIR rely on DCNNs as backbone architectures. In almost all cases, pre-stored parameters in these backbones can be fine-tuned (Section 4) to be better suited for instance retrieval and to contribute to better performance.

The detailed categorization of the material of the following sections is shown in Figure 3.

3 RETRIEVAL WITH OFF-THE-SHELF DCNN MODELS

Because of their size, DCNNs need to be trained, initially for classification tasks, on exceptionally large-scale datasets and be

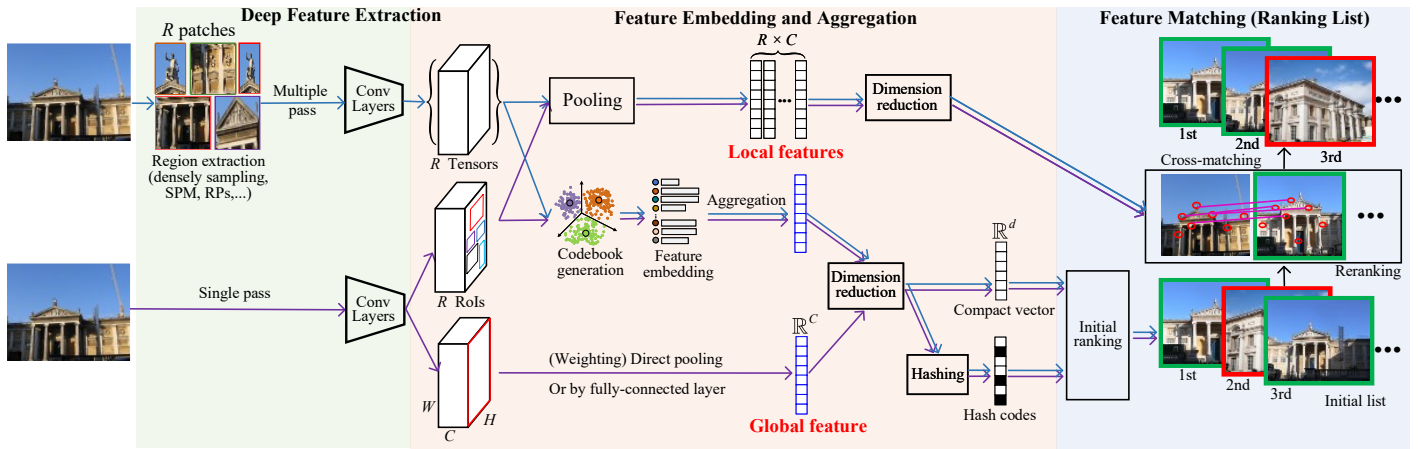


Fig. 2: General framework of IIR, which includes feature extraction on image or image patches, followed by feature embedding and aggregation methods to improve feature discriminativity. Feature matching can be performed by using global features (initial filtering) or use local features to rerank the top-ranked images matched by global features.

able to recognize images from different classes. One possible scheme, then, is that DCNNs effectively trained for classification directly serve as off-the-shelf feature detectors for the image retrieval task, the topic of this survey. That is, one can propose to undertake image retrieval on the basis of DCNNs, trained for classification and with their pre-trained parameters frozen.

do not necessarily extract features well suited to image retrieval. In particular, a classification decision can be successful as long as features remain within classification boundaries, however features from such models may show insufficient capacity for retrieval where feature matching itself is more important than classification. This section will survey the strategies which have been developed to improve the quality of feature representations, particularly based on feature extraction / fusion (Section 3.1) and feature embedding / aggregation (Section 3.2).

Deep Learning for Instance-level Image Retrieval (overall survey)

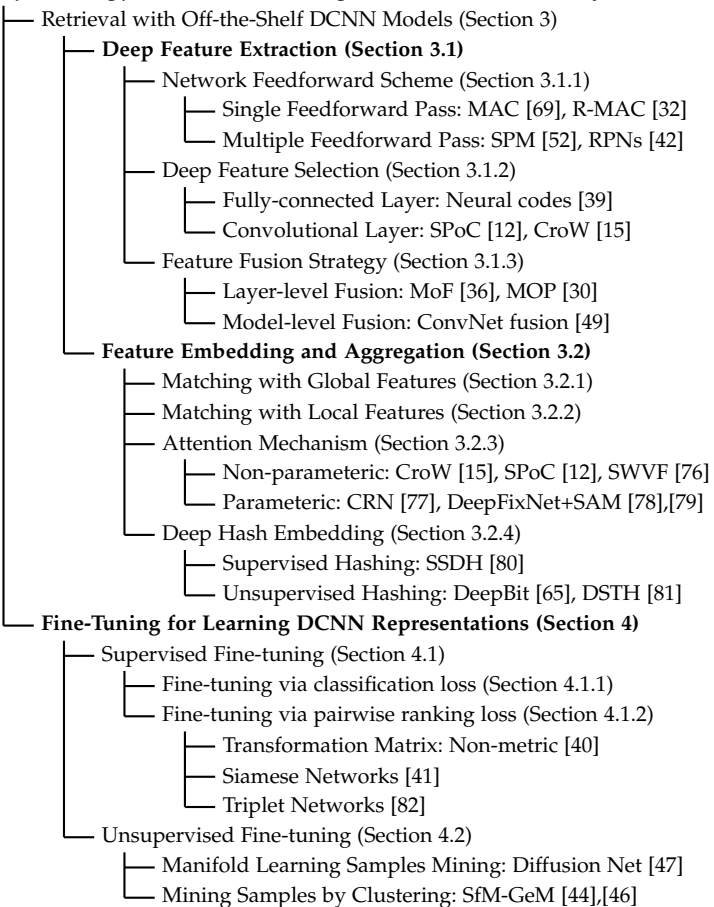


Fig. 3: This survey is organized around three key themes in instance-level image retrieval, shown in bold.

There are limitations to this approach, most fundamentally that there is a model-transfer or domain-shift challenge between tasks [13],[31],[83], meaning that models trained for classification

3.1 Deep Feature Extraction

Feature extraction is about the mechanism by which retrieval features can be extracted from off-the-shelf DCNNs. For an input image x and a network $f(\cdot; \theta)$, we denote its features from a convolutional layer as $\mathbf{A} := f_{conv}(x) \in \mathbb{R}^{H \times W \times C}$ with height H , width W , and channels C while that from a fully-connected layer as $\mathbf{B} := f_{fc}(x) \in \mathbb{R}^{D \times 1}$ with the dimensional D .

3.1.1 Network Feedforward Scheme

Network feedforward schemes focus on how images are fed into a DCNN, which includes single-pass and multiple-pass.

a. Single Feedforward Pass Methods.

Single feedforward pass methods take the whole image and feed it into an off-the-shelf model to extract features. The approach is relatively efficient since the input image is fed only once. For these methods, both the fully-connected layer and last convolutional layer can be used as feature extractors [84].

Early network-based IIR work focused on leveraging DCNNs as a fixed extractor to obtain global features, especially based on the fully-connected layers [29],[39], requiring close to zero engineering effort. However, extracting features in this way affects retrieval accuracy since the extracted features may include background information or activations for irrelevant objects.

The key to single-pass schemes is to embed and aggregate features to improve their discriminativity, such that features of two related images (*i.e.*, including the same object) are more similar than these of two unrelated images [12]. For this purpose, it is possible to first map the features \mathbf{B} into a high-dimensional space and then to aggregate them into a final global feature [30]. Another direction is to treat regions in convolutional features \mathbf{A} as different sub-vectors, such that a combination of sub-vectors of all feature maps are used to represent the input image [15],[32].

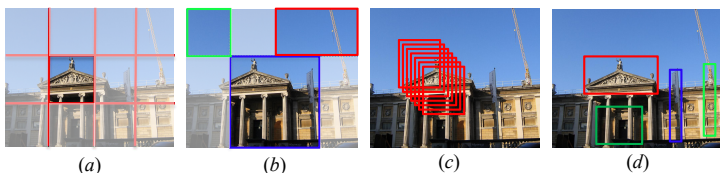


Fig. 4: Image patch generation schemes: (a) Sliding windows [37],[38]; (b) Spatial pyramid modeling [85]; (c) Dense sampling [30],[36]; (d) Region proposals from region proposal networks [27],[42].

b. Multiple Feedforward Pass Methods.

Compared to single-pass schemes, multiple-pass methods are more time-consuming [13] because several patches are generated and then fed into the network, multiple-pass schemes are more helpful for addressing the “invariance challenges” and “distraction challenges” in Section 1.2. Local patches at multiple scales become more robust for image translation, scaling and rotation [30],[61]. Also, these patches are helpful to filter several irrelevant background information.

The representations are usually produced from two stages: patch detection and patch description. Multi-scale image patches are obtained using sliding windows [37],[38] or spatial pyramid model (SPM) [52],[85],[86], as shown in Figure 4. For example, Zheng *et al.* [86] partition an image by using SPM and extract features at increasing scales, thus enabling the integration of global, regional, local contextual information.

Patch detection methods lack retrieval efficiency since irrelevant patches are also detected [32]. For example, Cao *et al.* [87] propose to merge image patches into larger regions with different hyper-parameters, where the hyper-parameter selection is viewed as an optimization problem to maximize the similarity between query and candidate features.

Instead of generating multi-scale image patches randomly or densely, region proposal methods introduce a degree of purpose. Region proposals can be generated using object detectors, such as selective search [61], edge boxes [88],[89], and BING [90]. For example, Yu *et al.* [89] propose fuzzy object matching (FOM) for instance search in which the fuzzy objects are generated from 300 object proposals and then clustered to filter out overlapping proposals. Region proposals can also be learned using such as region proposal networks (RPNs) [27],[42] and convolutional kernel networks (CKNs) [91], and then to apply these networks into end-to-end fine-tuning for learning similarity [92]. This usually requires the datasets provide well-localized bounding boxes as supervision, *e.g.*, the datasets INSTRE [93], Oxford-5k [94], Paris-6k [95], GLD-v2 variant [74]. Also, in the off-the-shelf scenarios, the way that using the bounding boxes to crop the query images and use as input the DCNNs has been shown to provide better retrieval performance since only the information relevant to the instance is extracted [92],[32].

3.1.2 Deep Feature Selection

Feature selection decides the receptive field of the extracted features, *i.e.*, global-level from fully-connected layers and regional-level from convolutional layers.

a. Extracted from Fully-connected Layers

It is straightforward to select a fully-connected layer as a global feature extractor [29],[30],[39]. With PCA dimensionality reduction and normalization [29] image similarity can be measured. Extracting features \mathbf{B} from fully-connected layer leads to two obvious limitations for IIR: including irrelevant information, and a lack of local geometric invariance [30].

With regards to the first limitation, image-level global descriptors may include irrelevant patterns or background clutter, especially when a target instance is only a small portion of an image. It may then be more reasonable to extract region-level features at finer scales, *i.e.*, using multiple passes [30],[61],[64]. For the second limitation, an alternative is to extract multi-scale features on a convolutional layer [69],[62]. Further, it makes the global features incompatible with techniques such as spatial verification and re-ranking. Several methods then choose to leverage intermediate convolutional layers [12],[30],[50],[69].

b. Extracted from Convolutional Layers

The neurons in a convolutional layer are connected only to a local region of the input image, and this smaller receptive field ensures that the produced features \mathbf{A} , usually from the last layer, preserve more local structural information [96],[97] and are more robust to image transformations [12] thereby address the “invariance challenge”. For instance, Razavian *et al.* [69] extract multi-scale features on the last convolutional layer and Morère *et al.* [62] incorporate a series of nested pooling layers into CNN. Both of them provide higher feature invariance. Thus, many image retrieval methods use convolutional layers as feature extractors [33],[50],[69],[98].

Sum/average and max pooling are two simple aggregation methods to produce global features [69]. For a pooled layer, the last convolutional layer usually yields superior accuracy over other shallower or later fully-connected layers [97]. There is no other operation on the feature maps before pooling, so we illustrate these methods as “direct pooling” in Figure 2.

Instead of direct pooling, many sophisticated aggregation methods have been explored, such as channel-wise or spatial-wise feature weighting on the convolutional feature maps [62],[99],[100]. These aggregation methods aim to highlight feature importance [15] or reduce the undesirable influence of bursty descriptors of some regions [34],[101]. For clarity, we illustrate the representative strategies in Figure 5. Note that these feature aggregation methods are usually performed before channel-wise sum/max pooling and does not embed features into a higher dimensional space.

One rationale behind using convolutional features is that each such vector can act as a “dense SIFT” feature [12] since each vector corresponds to a region in the input image. Inspired by this perception, many works leverage embedding methods (*e.g.*, BoW) used for SIFT features [23] on the regional feature vectors and then aggregate them (*e.g.*, by sum pooling) into a global descriptor. Feature embedding methods address the discriminativity challenge via mapping individual features into a high-dimensional space and make them distinguishable [34]. Feature embedding is followed by PCA to reduce feature dimensionality and whitening to down-weight co-occurrence between features.

3.1.3 Feature Fusion Strategies

Fusion studies the complementarity of different features which includes layer-level and model-level fusion explorations.

a. Layer-level Fusion

With layer-level fusion it is possible to fuse multiple fully-connected layers in a deep network [52],[66]. For instance, Liu *et al.* [52] introduce DeepIndex to incorporate multiple global features from different fully connected layers. The activation from the first fully-connected layer is taken as column indexing, and that from the second layer serves as row indexing. Similarly, it is also possible to fuse the activations from multiple

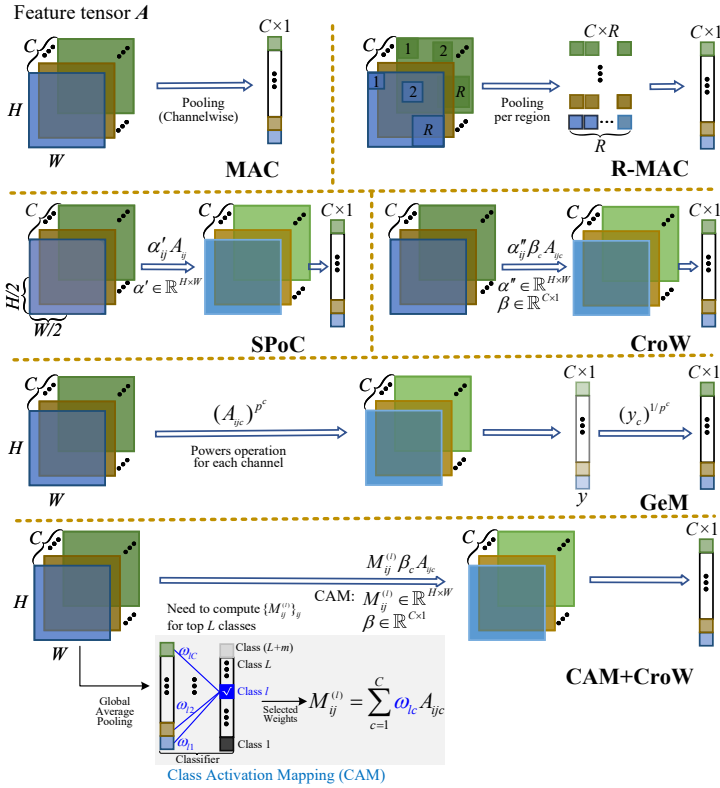


Fig. 5: Representative methods in single-pass methods, focusing on convolutional feature tensor \mathbf{A} . We denote the entry in \mathbf{A} corresponding to channel c , at spatial location (i, j) as A_{ijc} : MAC [69], R-MAC [32], SPoC with the per-channel Gaussian weighting $\alpha'_{ij} A_{ij}$ where $\alpha'_{ij} = \exp\left\{-\frac{(i-H/2)^2 + (j-W/2)^2}{2\sigma^2}\right\}$ [12], CroW with α'' computed by summing all C feature maps at location (i, j) and β computed by summing the $H \times W$ -array at each feature map c [15], GeM with channel-wise powers operation [46], and CAM+CroW by performing $M_{ij}^{(l)} = \sum_{c=1}^C \omega_{lc} A_{ijc}$ where ω_{lc} are weights activated by l -th class [33].

convolutional layers. For instance, Li *et al.* [99] apply the R-MAC encoding scheme on five convolutional layers of VGG-16 and then concatenate them into a multi-scale feature vector.

Features from fully-connected layers retain global high-level semantics, whereas features from convolutional layers can present local low- and intermediate-level cues. Global and local features therefore complement each other when measuring semantic similarity and can, to some extent, guarantee retrieval performance [102],[103]. Such features can be concatenated directly [71],[102], with convolutional features normally filtered by sliding windows or region proposal nets. Direct concatenation can also be replaced by other advanced methods, such as Multi-layer Orderless Fusion (MOF) of Li *et al.* [36], which is inspired by Multi-layer Orderless Pooling (MOP) [30]. However local features cannot play a decisive role in distinguishing subtle feature differences if global and local features are treated identically. Yu *et al.* [102] use a mapping function to assert local features in refining the return ranking lists, via an exponential mapping function for tapping the complementary strengths of convolutional and fully-connected layers. Similarly, Liu *et al.* [4] design two sub-networks on top of convolutional layers to obtain global and local features and then learn to fuse these features, thereby adaptively adjusting the fusion weights. Instead of directly fusing the layer activations, Zhang *et al.* [104] fuse the

index matrices which are generated based on the two feature types extracted from the same CNN, a feature fusion which has low computational complexity.

It is worth considering which layer combinations are better for fusion given their differences and complementarity. Yu *et al.* [102] compare the performance of different combinations between fully-connected and convolutional layers on the Oxford 5k, Holiday, and UKBench datasets. The results show that the combinations including the first fully-connected layer always perform better. Li *et al.* [36] demonstrate that fusing convolutional and fully-connected layers outperforms the fusion of only convolutional layers. Fusing two convolutional layers with one fully-connected layer achieves the best performance on the Holiday and UKBench datasets.

b. Model-level Fusion

It is possible to combine features from different models; such a fusion focuses more on model complementarity, with methods categorized into *intra-model* and *inter-model*.

Intra-model fusion suggests multiple deep models having similar or highly compatible structures, while inter-model fusion involves models with differing structures. For example, Simonyan *et al.* [49] introduce a ConvNet intra-model fusion strategy to improve the feature learning capacity of VGG where VGG-16 and VGG-19 are fused. To attend to different parts of an image object, Wang *et al.* [105] realize the multi-feature fusion by selecting all convolutional layers of VGG-16 to extract image representations, which is demonstrated to be more robust than using only single-layer features.

Inter-model fusion is a way to bridge different features given the fact that different deep networks have different receptive fields [52],[68],[79],[97]. For instance, a two-stream attention network [79] is introduced to implement image retrieval where the main network for semantic prediction is VGG-16 while an auxiliary network is used for predicting attention maps. Similarly, considering the importance and necessity of inter-model fusion to bridge the gap between mid-level and high-level features, Liu *et al.* [52] and Zheng *et al.* [97] combine VGG-19 and AlexNet to learn combined features, while Ozaki *et al.* [68] concatenate descriptors from six different models.

Inter-model and intra-model fusion are relevant to model selection. There are some strategies to determine how to combine the features from two models. It is straightforward to fuse all features from the candidate models and then learn a metric based on the concatenated features [52],[79], which is a kind of “*early fusion*” strategy. Alternatively, it is also possible to learn optimal metrics separately for the features from each model, and then to combine these metrics for final retrieval ranking [36],[106], which is a kind of “*late fusion*” strategy.

Discussion. Layer-level fusion and model-level fusion are conditioned on the fact that the associated layers or networks have different feature description capacities. For these fusion strategies, the key question is *what features are the best to be combined?* Some explorations have been made on the basis of off-the-shelf models, such as Xuan *et al.* [107], who illustrates the effect of combining different numbers of features and different sizes within the ensemble. Chen *et al.* [108] analyze the performance of embedded features from off-the-shelf image classification and object detection models with respect to image retrieval.

3.2 Feature Embedding and Aggregation

The primary aim of feature embedding and aggregation is to further promote feature discriminativity, targeting for the

“discriminativity challenge”, and obtain final global and/or local features for retrieving specific instances.

3.2.1 Matching with Global Features

Global features can be extracted from fully-connected layers, followed by dimensionality reduction and normalization [29],[39]. They are easy to implement and there is no further aggregation process. Gong *et al.* [30] extract fully-connected activations for local image patches at three scale levels and embed patch-level activations individually using VLAD. Thus, the final concatenated features significantly tackle the invariance challenge caused by image rotations.

Convolutional features can also be aggregated into compact a global feature. Simple aggregation methods are sum/average or max pooling [69],[70]. Sum/average pooling is less discriminative, because it takes into account all activated convolutional outputs, thereby weakening the effect of highly activated features [34]. As a result, max pooling is particularly well suited for sparse features having a low probability of being active, however max pooling may be inferior to sum/average pooling when image features are whitened [12].

Figure 5 illustrates sophisticated feature aggregation methods using channel-wise or spatial-wise weighting [12],[62]. For example, Babenko *et al.* [12] propose sum-pooling convolutional features (SPoC) to obtain compact descriptors weighted by α' with a Gaussian center prior. Similarly, it is possible to treat regions in feature maps as different sub-vectors [32],[69],[97], thus combinations of R sub-vectors are used to represent the input image, such as R-MAC [32]. Since convolutional features may include repetitive patterns and each vector may correspond to identical regions, the resulting descriptors may be bursty, which makes the final aggregated global feature less distinguishable. As a solution, Pang *et al.* [101] leverage heat diffusion to weigh convolutional features at the aggregation stage, and reduce the undesirable influence of burstiness.

Convolutional features have an interpretation as descriptors of local regions, thus many works leverage embedding methods, including BoW, VLAD, and FV, to encode regional feature vectors and then aggregate them into a global descriptor. Note that BoW and VLAD can be extended by using other metrics, such as a Hamming distance [109]. Here we briefly describe the principle of Euclidean embeddings.

BoW [24] is a widely used feature embedding which leads to a sparse vector of occurrence. Let $\mathbf{a} = \{a_1, a_2, \dots, a_R\}$ be a set of R local features, each of dimensionality d . BoW requires a pre-defined codebook $\mathbf{c} = \{c_1, c_2, \dots, c_K\}$ with K centroids, usually learned offline, to cluster these local descriptors, and maps each descriptor a_t to the nearest centroid c_k . For each centroid, one can count and normalize the number of occurrences as

$$g(c_k) = \frac{1}{R} \sum_{r=1}^R \phi(a_r, c_k) \quad (1)$$

$$\phi(a_r, c_k) = \begin{cases} 1 & \text{if } c_k \text{ is the closest codeword for } a_r \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Thus BoW considers the number of descriptors belonging to each c_k (*i.e.* 0-order feature statistics), so the BoW representation is the concatenation of all mapped vectors:

$$G_{BoW}(\mathbf{a}) = [g(c_1), \dots, g(c_k), \dots, g(c_K)]^\top \quad (3)$$

BoW is simple to implement the encoding of local descriptors, such as convolutional feature maps [36],[84] or fully-connected activations [86],[104], or to encode regional descriptors [110],[111]. Mukherjee *et al.* [111] extract image patches

based on information entropy and feed into a pre-trained VGG-16, then use BoW to embed and aggregate the patch-level descriptors from a fully-connected layer. Embedded BoW vectors are typically high-dimensional and sparse, so not well suited to large-scale datasets in terms of the mentioned efficiency challenge.

VLAD [54] stores the sum of residuals for each visual word. Similar to BoW, it generates K visual word centroids, then each feature a_r is assigned to its nearest visual centroid c_k :

$$g(c_k) = \frac{1}{R} \sum_{r=1}^R \phi(a_r, c_k)(a_r - c_k) \quad (4)$$

The VLAD representation is stacked by the residuals for all centroids, with dimension $(d \times K)$, *i.e.*,

$$G_{VLAD}(\mathbf{a}) = [\dots, g(c_k)^\top, \dots]^\top \quad (5)$$

VLAD captures first-order feature statistics, *i.e.*, $(a_r - c_k)$. Similar to BoW, the performance of VLAD is affected by the number of clusters: more centroids produce larger vectors that are harder to index. For instance-level image retrieval, Gong *et al.* [30] concatenate the activations of a fully-connected layer with VLAD applied to image-level and patch-level inputs [112]. Ng *et al.* [50] replace BoW [24] with VLAD [54], and are the first to encode local features into VLAD representations. This idea inspired another milestone work [43] where, for the first time, VLAD is plugged into the last convolutional layer, which allows end-to-end training via back-propagation.

FV [55] extends BoW by encoding the first and second order statistics. FV clusters the set of local descriptors by a Gaussian Mixture Model (GMM) with K components to generate a dictionary $\mathbf{c} = \{\mu_k; \Sigma_k; w_k\}_{k=1}^K$ made up of mean / covariance / weight triples [113], where the covariance may be simplified by keeping only its diagonal elements. For each local feature a_r , a GMM is given by

$$\gamma_k(a_r) = w_k \times p_k(a_r) / \left(\sum_{k=1}^K w_k p_k(a_r) \right) \quad \text{s.t.} \quad \sum_{k=1}^K w_k = 1 \quad (6)$$

where $p_k(a_r) = \mathcal{N}(a_r, \mu_k, \sigma_k^2)$. All local features are assigned into each component k in the dictionary, which is computed as

$$\begin{aligned} g_{w_k} &= \frac{1}{R\sqrt{w_k}} \sum_{r=1}^R (\gamma_k(a_r) - w_k) \\ g_{u_k} &= \frac{\gamma_k(a_r)}{R\sqrt{w_k}} \sum_{r=1}^R \left(\frac{a_r - \mu_k}{\sigma_k} \right), \\ g_{\sigma_k^2} &= \frac{\gamma_k(a_r)}{R\sqrt{2w_k}} \sum_{r=1}^R \left[\left(\frac{a_r - \mu_k}{\sigma_k} \right)^2 - 1 \right] \end{aligned} \quad (7)$$

The FV representation is produced by concatenating vectors from the K components:

$$G_{FV}(\mathbf{a}) = [g_{w_1}, \dots, g_{w_K}, g_{u_1}, \dots, g_{u_K}, g_{\sigma_1^2}, \dots, g_{\sigma_K^2}]^\top \quad (8)$$

The FV representation defines a kernel from a generative process and captures more statistics than BoW and VLAD. FV vectors do not increase the computational cost significantly but require more memory. Applying FV without memory controls may lead to suboptimal performance [114].

Discussion. Traditionally, pooling-based aggregation methods (*e.g.*, in Figure 5) are directly plugged into deep networks and then the whole model is used end-to-end. The three embedding methods (BoW, VLAD, FV) are initially trained with large

pre-defined vocabularies [52],[115]. One needs to pay attention on their properties before choosing an embedding: BoW and VLAD are computed in the rigid Euclidean space where performance is closely related to the number of centroids, whereas FV can capture higher-order statistics and improves the effectiveness of feature embedding at the expense of a higher memory cost. Further, although vocabularies are usually built separately and pre-trained before encoding deep features, it is necessary to integrate the training of networks and the learning of vocabulary parameters into a unified framework so as to guarantee training and testing efficiency. For example, VLAD is integrated into deep networks where each spatial column feature is used to construct clusters via k-means [50]. This idea led to NetVLAD [43], where deep networks are fine-tuned with the VLAD vectors. The FV method is also combined with deep networks for retrieval tasks [34],[41].

3.2.2 Matching with Local Features

Although matching with global features has high efficiency for both feature extraction and similarity computation, global features are not compatible with spatial verification and correspondence estimation, which are important procedures for instance-level retrieval tasks, motivating work on matching with local features. In terms of the matching process, global features are matched only once while local feature matching is evaluated by summarizing the similarity across all individual local features (*i.e.*, many-to-many matching).

One important aspect of local features is to detect the keypoints for an instance within an image, and then to describe the detected keypoints as a set of local descriptors. Inspired by [116], the common strategies of this whole procedure for IIR can be categorized as *detect-then-describe* and *describe-then-detect*.

In terms of *detect-then-describe*, we regard the descriptors around keypoints as local features, similar to [29],[42]. Coarse regions can be detected, for example, by using the methods depicted in Figure 4, and regions of interest in an image can be detected by using region proposal networks (RPNs) [27],[74]. The extracted coarse regions around the keypoints are fed into a DCNN, followed by feature description. Traditional detectors can also be used to detect fine regions around a keypoint. For instance, Zheng *et al.* [86] employ the popular Hessian-Affine detector [117] to get an affine-invariant local region. Paulin *et al.* [112] and Mishchuk *et al.* [110] detect regions using the Hessian-Affine detector and feed into patch-convolutional kernel networks (Patch-CKNs) [91]. Note that it becomes more convenient for the case where bounding boxes annotations have been provided by datasets (see Section 5.1), and then the image regions can be cropped directly for further reranking [92].

Rather than performing keypoint detection early on, it is possible to postpone the detection stage on the convolutional feature maps, *i.e.*, *describe-then-detect*. One can select regions on the convolutional feature maps to obtain a set of local features [32],[37],[84]; the local maxima of the feature maps are then detected as keypoints [56]. A similar strategy is also used in network fine-tuning [5],[60],[71],[74],[118], where the keypoints on the convolutional feature maps can be selected based on attention scores predicted by an attention network [71],[5], or based on single-head and multi-head attention modules in transformers [60],[59]. This approach to keypoint selection is better for achieving computational efficiency.

After keypoint detection and description, a large number of local features are used in the matching stage to perform instance-level retrieval, and the image similarity is evaluated

by matching across all local features. Local matching techniques include spatial verification and selective match kernels (SMK) [73]. Spatial verification assumes object instances are rigid so that local matches between images can be estimated as an affine transformation using RANdom SAMple Consensus (RANSAC) [72]. One limitation of RANSAC is its high computational complexity of estimating the transformation model when all local descriptors are considered; instead, it is possible to apply RANSAC to a small number of top-ranked local descriptors, such as those selected by approximate nearest neighbor [5]. SMK weighs the contributions of individual matches with a non-linear selective function, but is still memory intensive. Its extension, the Aggregated Selective Match Kernel (ASMK), focuses more on aggregating similarities between local features without explicitly modeling the geometric alignment, which can produce a more compact representation [73],[118]. Recently, Teichmann *et al.* [74] introduced Regional Aggregated Selective Match Kernel (R-ASMK) to combine information from detected regions, boosting image retrieval accuracy compared to the ASMK.

Discussion. Using local descriptors to perform instance retrieval tasks has two limitations. First, the local descriptors for an image are stored individually and independently, which is memory-intensive, and not well-suited for large-scale scenarios. Second, estimating the similarity between the query and database images depends on cross-matching all local descriptor pairs, which incurs additional searching cost and then a low retrieval efficiency. Therefore, most instance retrieval systems using local features follow a two-stage paradigm: initial filtering and re-ranking [64],[71],[75],[92],[119], as in Figure 2. The initial filtering stage is to employ a global descriptor to select a set of candidate matching images, thereby reducing the solution space; the re-ranking stage is to use local descriptors to re-rank the top-ranked images from the global descriptor.

3.2.3 Attention Mechanism

Attention mechanism can be regarded as a kind of feature aggregation, whose aim is to highlight the most relevant feature parts. It can effectively address the “*distraction challenge*” and also promote feature discriminativity [98], realized by computing an attention map. Approaches to obtaining attention maps can be categorized into non-parametric and parametric groups, as shown in Figure 6, where the main difference is whether the importance weights in the attention map are learnable.

Non-parametric weighting is a straightforward method to highlight feature importance, and the corresponding attention maps can be obtained by channel-wise or spatial-wise pooling, as in Figure 6 (a,b). For spatial-wise pooling, Kalantidis *et al.* [15] propose an effective CroW method to weight and pool feature maps, which concentrate on weighting activations at different spatial locations, without considering the relations between these activations. In contrast, Ng *et al.* [96] explore the correlations among activations at different spatial locations on the convolutional feature maps.

Channel-wise weighting methods are also popular non-parametric attention mechanisms [35],[100]. Xu *et al.* [35] rank the weighted feature maps to build “probabilistic proposals” to select regional features. Jimenez *et al.* [33] combine CroW and R-MAC to propose Classes Activation Maps (CAM) to weigh the feature map per class. Xiang *et al.* [100] employ a Gram matrix to analyze the correlations between different channels and then obtain channel sensitivity information to tune the importance of each feature map. Channel-wise and spatial-wise weighting methods are usually integrated into a deep model to highlight feature importance [15],[105].

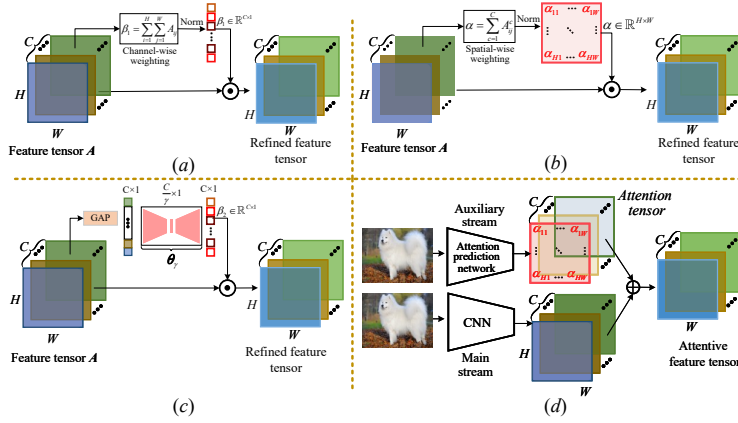


Fig. 6: Illustration of attention mechanisms. (a)-(b) Non-parametric schemes: The attention is based on convolutional feature A . Channel-wise attention in (a) produces a C -dimensional importance vector β_1 [15],[35]; Spatial-wise attention in (b) computes a 2-dimensional attention map α [15],[33],[96]. (c)-(d) Parametric schemes: The attention weights are learned by a trainable network. In (c), β_2 are provided by a sub-network with parameters θ , [71],[77],[98],[120]. In (d), the attention maps, as a tensor, are predicted by some auxiliary saliency extraction models from the input image directly [78],[79],[121].

Parametric attention maps, shown in Figure 6 (c,d), can be learned via deep networks, where the input can be either image patches or feature maps [77],[100],[103], approaches which are commonly used in supervised metric learning [98]. Kim *et al.* [77] make the first attempt to propose a shallow network (CRN) to take as input the feature maps of convolutional layers and outputs a weighted mask indicating the importance of spatial regions in the feature maps. The resulting mask modulates feature aggregation to create a global representation of the input image. Noh *et al.* [5] design a 2-layer CNN with a softplus output layer to compute scores which indicate the importance of different image regions. Inspired by R-MAC, Kim *et al.* [120] employ a pre-trained ResNet101 to train a context-aware attention network using multi-scale feature maps.

Apart from using feature maps as inputs, a whole image can be used to learn feature importance, for which specific networks are needed [78],[79],[121]. Mohedano [78] explores different saliency models, including DeepFixNet and Saliency Attentive Model. Yang *et al.* [79] and Wei *et al.* [121] introduce a two-stream network for image retrieval in which the auxiliary stream, DeepFixNet, is used specifically for predicting attention maps, which are then fused with the feature maps produced by the main network. For image retrieval, attention mechanisms can be combined with supervised metric learning [96].

3.2.4 Hashing Embedding

Real-valued features extracted by deep networks are typically high-dimensional, and therefore are not well-suited to retrieval efficiency. As a result, there is significant motivation to transform deep features into more compact codes. Since their computational and storage efficiency are beneficial for the “*efficiency challenge*”, hashing algorithms have been widely used for global [62],[80] and local descriptors [64],[85],[86].

Hash functions can be plugged as a layer into deep networks, so that hash codes and deep networks can be simultaneously trained and optimized, either supervised [80] or unsupervised [65]. During hash function training, the hash codes of originally similar images are embedded as closely as possible, and the hash

codes of dissimilar images are as separated as possible. d -dim hash codes from a hash function $h(\cdot)$ for an image x can be formulated as $b_x = h(x) = h(f(x; \theta)) \in \{+1, -1\}^d$. Because hash codes are non-differentiable their optimization is difficult, so $h(\cdot)$ can be relaxed to be differentiable by using *tanh* or *sigmoid* functions [22].

When binarizing real-valued features, it is crucial to preserve image similarity and to improve hash code quality [22]. These two aspects are at the heart of hashing algorithms to maximize retrieval accuracy.

a. Hash Functions to Preserve Image Similarity

Preserving similarity seeks to minimize the inconsistencies between real-valued features and corresponding hash codes, for which a variety of strategies have been adopted.

Loss functions can significantly influence similarity preservation, which includes both supervised and unsupervised methods. With class labels available, many loss functions are designed to learn hash codes in a Hamming space. As a straightforward method, one can optimize either the difference between matrices computed from the binary codes and their supervision labels [62],[81] or the difference between the hash codes and real-valued deep features [64],[65]. Song *et al.* [64] propose to learn hash codes for regional features in which each local feature is converted to a set of binary codes by multiplying a hash function and the raw RoI features, then the differences between RoI features and hash codes are characterized by an L_2 loss. Do *et al.* [122] regularize hash codes with a reconstruction loss, which ensure that codes can be reconstructed to their inputs so that similar/dissimilar inputs are mapped to similar/dissimilar hash codes. Lin *et al.* [65] learn hash codes and address the “*invariance challenge*” by introducing an objective function which characterize the difference between the binary codes which are computed from the original image and the geometric transformed one.

b. Improving Hash Function Quality

A good hash function seeks to have binary codes uniformly distributed; that is, maximally filling and using the hash code space, normally on the basis of bit uncorrelation and bit balance [22],[65]. Bit uncorrelation implies that different bits are as independent as possible, so that a given set of bits can aggregate more information within a given code length [65]. Bit balance means that each bit should have a 50% chance of being +1 or -1, thereby maximizing code variance and information [22]. Morère *et al.* [62] use the uniform distribution $U(0,1)$ to build a regularization term to make hash codes distribute evenly where the codes are learned by a Restricted Boltzmann Machine layer. Likewise, Lin *et al.* [65] optimize the mean of learned hash codes to be close to 0.5 to prevent any bit bias towards zero or one.

4 RETRIEVAL VIA LEARNING DCNN REPRESENTATIONS

The off-the-shelf DCNNs pre-trained on source datasets for classification are quite robust to inter-class variability. However, in most cases, deep features extracted based on off-the-shelf models may not be sufficient for accurate retrieval, even with the strategies discussed in Section 3. In order for models to be more effective for retrieval, a common practice is network fine-tuning, *i.e.*, updating the pre-stored parameters [31]. Fine-tuning methods have been studied extensively to learn better features, whose primary aim is to address the “*fine-tune challenge*”. A standard dataset with clear and well-defined ground-truth labels is indispensable for the supervised fine-tuning and subsequently pair-wise supervisory information is incorporated into ranking

loss to update networks by regularizing on retrieval representations, otherwise it is necessary to develop unsupervised fine-tuned methods. After network fine-tuning, features can be organized as global or local to perform retrieval.

For the most feature strategies we presented in Section 3, including feature extraction, feature embedding and feature aggregation. Note that fine-tuning does not contradict or render irrelevant these feature processing methods; indeed, these strategies are complementary and can be equivalently incorporated as part of network fine-tuning. To this end, this section will survey the strategies which have been developed, based on the patch-level, image-level, or class-level supervision, to fine-tune deep networks for better instance retrieval.

4.1 Supervised Fine-tuning

The way to realize supervised fine-tuning can be determined by the given class labels or pairwise supervisory information.

4.1.1 Fine-tuning via Classification Loss

When class labels of a new dataset are available (*e.g.*, INSTRE [93], GLDV2 [5],[63]), it is preferable to begin with a previously-trained DCNN, trained on a separate dataset, with the backbone DCNN typically chosen from one of AlexNet, VGG, GoogLeNet, or ResNet. The DCNN can then be fine-tuned, as shown in Figure 7 (a), by optimizing its parameters on the basis of a cross entropy loss

$$\mathcal{L}_{CE}(\hat{p}_i, y_i) = -\sum_i^c (y_i \times \log(\hat{p}_i)) \quad (9)$$

Here y_i and \hat{p}_i are the ground-truth labels and the predicted logits, respectively, and c is the total number of categories. The milestone work in such fine-tuning is [39], in which AlexNet is re-trained on the Landmarks dataset. According to the class labels, the image-level features are required to compute the logits. Thus, the descriptors extracted from local regions on convolutional feature maps [5],[71] or image patch inputs [74] are further needed to be aggregated.

A classification-based fine-tuning method enables to enforce higher similarity for intra-class samples and diversity for inter-class samples. Cao *et al.* [71] employ the ArcFace loss [126], which uses the margin-adjusted cosine similarity in the form of softmax loss, to induce smaller intra-class variance and show excellent results for instance retrieval. Recently, Boudiaf *et al.* [127] claim that cross entropy loss can minimize intra-class distances while maximizing inter-class distances. Cross entropy loss is, in essence, maximizing a common mutual information between the retrieval features and the ground-truth labels. Therefore, it can be regarded as an upper bound on a new pairwise loss, which has a structure similar to various pairwise ranking losses, of which representatives are introduced below.

4.1.2 Fine-tuning via Pairwise Ranking Loss

With affinity information (*e.g.*, samples from the same group) indicating similar and dissimilar pairs, fine-tuning methods based on pairwise ranking loss learn an optimal metric which minimizes or maximizes the distance of pairs to maintain their similarity. Network fine-tuning via ranking loss involves two types of information [17]:

- 1) A pair-wise constraint, corresponding to a Siamese network as in Figure 7 (c), in which input images are paired with either a positive or negative sample;

- 2) A triplet constraint, associated with triplet networks as in Figure 7 (e), in which anchor images are paired with both similar and dissimilar samples [17].

These pairwise ranking loss based methods are categorized into globally supervised approaches (Figure 7 (c,d)) and locally supervised approaches (Figure 7 (g,h)), where the former ones learn a metric on global features by satisfying all constraints, whereas the latter ones focus on local areas by only satisfying the given local constraints (*e.g.* region proposals).

To be specific, consider a triplet set $X = \{(x_a, x_p, x_n)\}$ in a mini-batch, where (x_a, x_p) indicates a similar pair and (x_a, x_n) a dissimilar pair. Features $f(x; \theta)$ of one image are extracted by a network $f(\cdot)$ with parameters θ , for which we can represent the affinity information for each similar or dissimilar pair as

$$\mathcal{D}_{ij} = \mathcal{D}(x_i, x_j) = \|f(x_i; \theta) - f(x_j; \theta)\|_2^2 \quad (10)$$

a. Refining with Transformation Matrix.

Learning the similarity among input samples can be implemented by optimizing the weights of a linear transformation matrix [40]. It transforms the concatenated feature pairs into a common latent space using a transformation matrix $\mathbf{W} \in \mathbb{R}^{2d \times 1}$, where d is the final feature dimension. The similarity score of these pairs are predicted via a sub-network $\mathcal{S}_W(x_i, x_j) = f_W(f(x_i; \theta) \cup f(x_j; \theta); \mathbf{W})$ [40]. In other words, the sub-network f_W predicts how similar the feature pairs are. Given the affinity information of feature pairs $\mathcal{S}_{ij} = \mathcal{S}(x_i, x_j) \in \{0, 1\}$, the binary labels 0 and 1 indicate the similar (positive) or dissimilar (negative) pairs, respectively. The training of function f_W can be achieved by using a regression loss:

$$\mathcal{L}_W(x_i, x_j) = |\mathcal{S}_W(x_i, x_j) - \mathcal{S}_{ij}(\text{sim}(x_i, x_j) + m) - (1 - \mathcal{S}_{ij})(\text{sim}(x_i, x_j) - m)| \quad (11)$$

where $\text{sim}(x_i, x_j)$ can be the cosine function for guiding the training of \mathbf{W} and m is a margin. By optimizing the regression loss and updating \mathbf{W} , deep networks maximize the similarity of similar pairs and minimize that of dissimilar pairs. It is worth noting that the pre-stored parameters in the deep models are frozen when optimizing \mathbf{W} . The pipeline of this approach is depicted in Figure 7 (b).

b. Fine-tuning with Siamese Networks.

Siamese networks represent important options in implementing metric learning for fine-tuning, as in Figure 7 (c) and Figure 8 (a). It is a structure composed of two branches that share the same weights across layers. Siamese networks are trained on paired data, consisting of an image pair (x_i, x_j) such that $\mathcal{S}(x_i, x_j) \in \{0, 1\}$. A Siamese loss is formulated as

$$\mathcal{L}_{Siam}(x_i, x_j) = \frac{1}{2} \mathcal{S}(x_i, x_j) \mathcal{D}(x_i, x_j) + \frac{1}{2} (1 - \mathcal{S}(x_i, x_j)) \max(0, m - \mathcal{D}(x_i, x_j)) \quad (12)$$

Siamese loss has recently been reaffirmed as a very effective metric in category-level image retrieval, outperforming many more sophisticated losses if implemented carefully [123]. Enabled by the standard Siamese network, this objective function is used to learn the similarity between semantically relevant samples under different scenarios [46],[124]. For example, Radenović *et al.* [46] employ a Siamese network on matching and non-matching global feature pairs which are aggregated by GeM-based pooling. The deep network fine-tuned by the Siamese loss generalizes better and converges at higher retrieval performance. Ong *et al.*

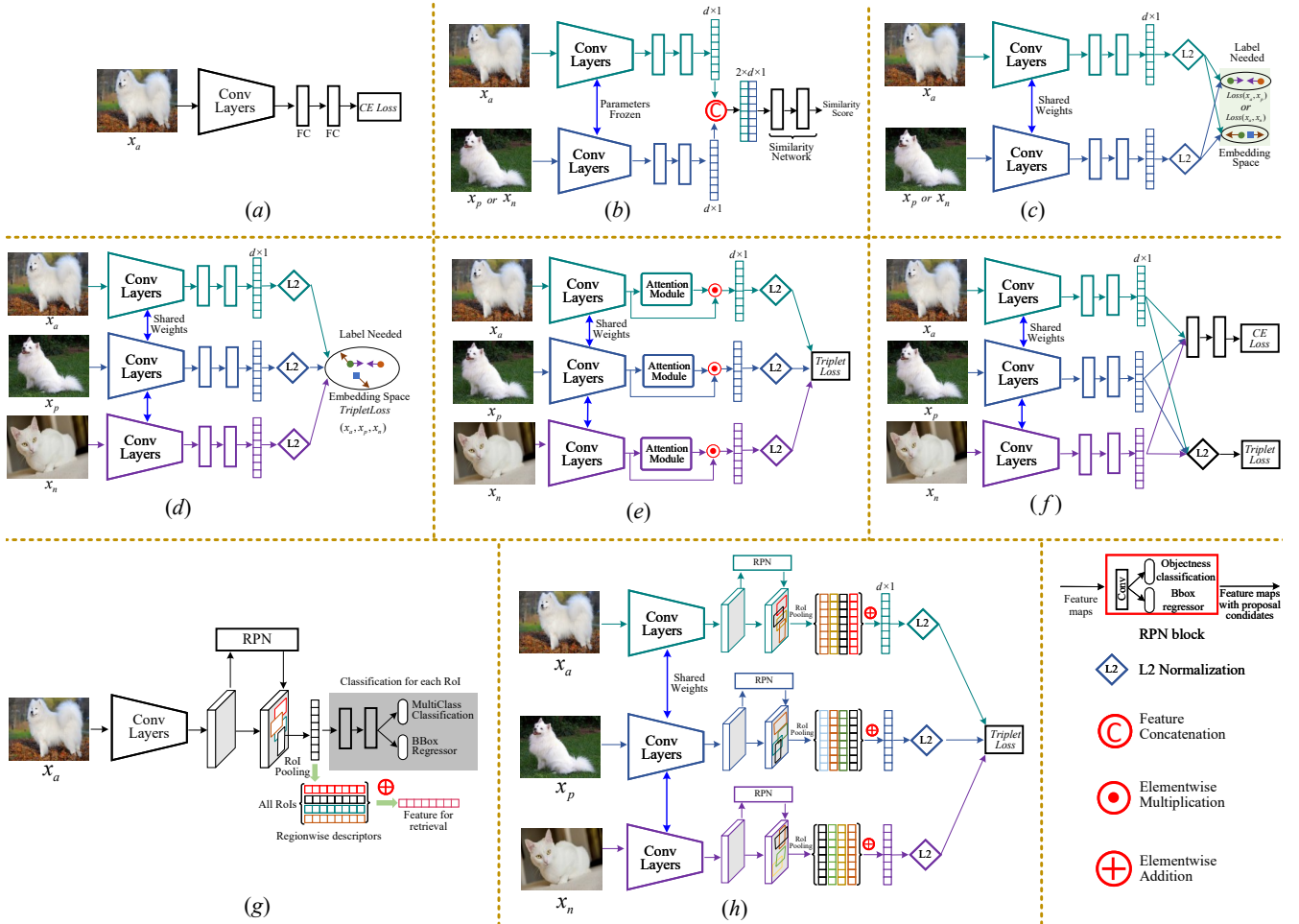


Fig. 7: Schemes of supervised fine-tuning. Anchor, positive, and negative images are indicated by x_a , x_p , x_n , respectively. (a) classification loss [39]; (b) similarity learning by using a transformation matrix [40]; (c) Siamese loss [46],[41],[123],[124]; (d) triplet loss [42]; (e) an attention block into DCNNs to highlight regions [64]; (f) combining classification loss and pairwise ranking loss [100],[125]; (g) region proposal networks (RPNs) to locate the RoI and highlight specific regions or instances [92]; (h) inserting the RPNs of (g) into DCNNs, such that the RPNs extract regions or instances at the convolutional layer [42],[82].

[41] leverage the Siamese network to learn image features which are then fed into the Fisher Vector model for further encoding. Siamese networks can also be applied to hashing learning in which the Euclidean distance $\mathcal{D}(\cdot)$ in Eq. 12 is computed for binary codes [128].

An implicit drawback of the Siamese loss is that it may penalize similar image pairs even if the margin between these pairs is small or zero [128], if the constraint is too strong and unbalanced. At the same time, it is hard to map the features of similar pairs to the same point when images contain complex contents or scenes. To tackle this limitation, Cao *et al.* [129] adopt a double-margin Siamese loss [128] to relax the penalty for similar pairs by setting a margin m_1 instead of zero, in which case the original single-margin Siamese loss is re-formulated as

$$\mathcal{L}_{\mathcal{D}_{Siam}}(x_i, x_j) = \frac{1}{2} \mathcal{S}(x_i, x_j) \max(0, \mathcal{D}(x_i, x_j) - m_1) + \frac{1}{2} (1 - \mathcal{S}(x_i, x_j)) \max(0, m_2 - \mathcal{D}(x_i, x_j)) \quad (13)$$

where $m_1 > 0$ and $m_2 > 0$ are the margins affecting the similar and dissimilar pairs, respectively, as in Figure 8 (b), meaning that the double margin Siamese loss only applies a contrastive force when the distance of a similar pair is larger than m_1 . The mAP metric of retrieval is improved when using the double margin Siamese loss [128].

More recently, transformers have been trained under the regularization of cross entropy [60] and Siamese loss [59] for instance-level retrieval and achieved competitive performance, positioning it as an alternative to convolutional architectures. As observed by [59], the transformer-based architecture is less impacted than convolutional networks by feature collapse since each input feature is projected to different sub-spaces before the multi-headed attention. Moreover, the transformer backbone operates as a learned aggregation operator, thereby avoiding the design of sophisticated feature aggregation methods.

c. Fine-tuning with Triplet Networks.

Triplet networks optimize similar and dissimilar pairs simultaneously. As shown in Figure 7 (d) and Figure 8 (c), the plain triplet networks adopt a ranking loss for training:

$$\mathcal{L}_{Triplet}(x_a, x_p, x_n) = \max(0, m + \mathcal{D}(x_a, x_p) - \mathcal{D}(x_a, x_n)) \quad (14)$$

which indicates that the distance of an anchor-negative pair $\mathcal{D}(x_a, x_n)$ should be larger than that of an anchor-positive pair $\mathcal{D}(x_a, x_p)$ by a certain margin m .

Given the datasets that provide bounding box annotations, such as INSTRE, Oxford-5k, Paris-6k, and their variants, the bounding box annotations are used as patch-level supervision to train a region detector which enables the final DCNNs to locate specific regions or objects. As an example, region proposal

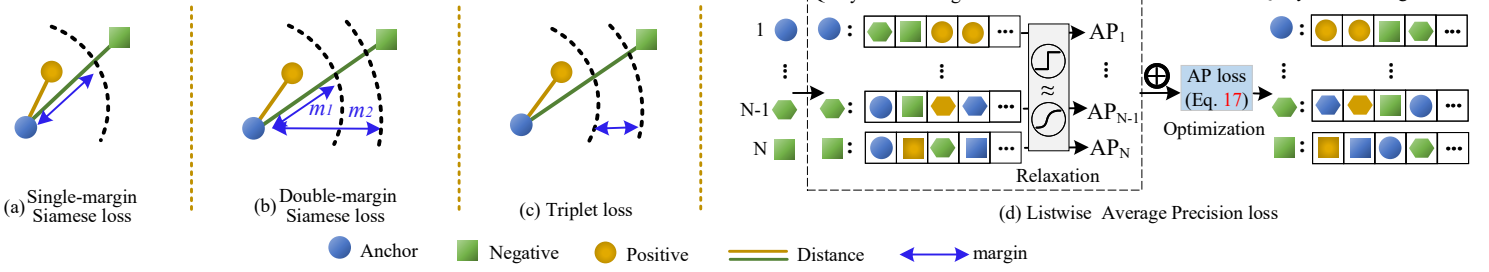


Fig. 8: Illustrations of different losses for network fine-tuning. The same shape with different colors denotes images that include the same instance. (a)-(c) have been introduced in the text [46],[128],[42]. (d) Listwise AP loss considers a mini-batch of N features simultaneously and directly optimizes the Average-Precision computed from these features [57],[130].

networks (RPNs) [27] is fine-tuned and subsequently plugged into DCNNs and trained end-to-end [92], as shown in Figure 7 (g). RPNs yield the regressed bounding box coordinates of objects and are trained by the multi-class classification loss. Once fine-tuned, RPNs can produce regional features for each detected region by RoI pooling and perform better instance search.

Further, local supervised metric learning has been explored based on the fact that RPNs [27] enable deep models to learn regional features for particular instance objects [125],[64],[42],[82]. RPNs used in the triplet formulation are shown in Figure 7 (h). Firstly, regression loss (RPNs loss) is used to minimize the regressed bounding box relative to ground-truth. Then, the regional features for all detected RoIs are aggregated into a global one and L2-normalized for the triplet loss. Note that, in some cases, jointly training an RPN loss and triplet loss leads to unstable results, a problem addressed in [42] by first training a CNN to produce R-MAC using a rigid grid, after which the parameters in convolutional layers are fixed and RPNs are trained to replace the rigid grid.

Attention mechanisms can also be combined with metric learning for fine-tuning [64], as in Figure 7 (e), where the attention module is typically end-to-end trainable and takes as input the convolutional feature maps. Song *et al.* [64] introduce a convolutional attention layer to explore spatial-semantic information, highlighting regions in images to significantly improve the discrimination for image retrieval.

Recent studies [100],[125] have jointly optimized the triplet loss and classification loss to further improve network capacity, as shown in Figure 7 (f). The overall joint function is

$$\mathcal{L}_{Joint} = \lambda_1 \cdot \mathcal{L}_{Triplet}(x_{i,a}, x_{i,p}, x_{i,n}) + \lambda_2 \cdot \mathcal{L}_{CE}(\hat{y}_i, y_i) \quad (15)$$

where the cross entropy loss (CE loss) \mathcal{L}_{CE} is defined in Eq. (9) and the triplet loss $\mathcal{L}_{Triplet}$ in Eq. (14). λ_1 and λ_2 are hyper-parameters tuning the tradeoff between the two loss functions.

4.1.3 Discussion

In some cases, pairwise ranking loss cannot effectively learn the variations between samples and still suffers from a weaker generalization capability if the training set is not ordered correctly. Therefore, pairwise ranking loss requires careful sample mining and weighting strategies to obtain the most informative pairs, especially when considering mini-batches. The hard negative mining strategy is commonly used [43],[46],[118], however further sophisticated mining strategies have recently been developed. Mishchuk *et al.* [110] calculate a pair-wise distance matrix on all mini-batch samples to select two closest negative and one anchor-positive pair to form a triplet. Instead of traversing all possible two-tuple or three-tuple combinations, it is possible to consider all positive samples in one cluster and negative samples

together. Liu *et al.* [4] introduce a group-group loss to decrease the intra-group distance and increase the inter-group distance. Considering all samples may be beneficial for stabilizing optimization and promoting generalization due to a larger data diversity, however the extra computational cost remains an issue to be addressed.

Substantial research has been devoted to pair-wise ranking loss, while cross entropy loss, mainly used for classification, has been largely overlooked. Recently, Boudiaf *et al.* [127] claim that cross entropy loss can match and even surpass the pair-wise ranking loss when carefully tuned on fine-grained category-level retrieval tasks. In fact, the greatest improvements have come from enhanced training schemes (*e.g.*, data augmentation, learning rate polices, batch normalization freeze) rather than intrinsic properties of pairwise ranking loss. Further, although several sophisticated ranking losses have been explored and validated for category-level retrieval, Musgrave *et al.* [123] revisited these losses and found that most of them perform on par to vanilla Siamese loss and triplet loss, so there is merit to consider these losses also for instance-level image retrieval tasks.

Both cross entropy loss and pair-wise ranking loss regularize on the embedded features and the corresponding labels so as to maximize their mutual information [127]. Their effectiveness is not guaranteed to give retrieval results that also optimize mAP [57]. To tackle this limitation one can directly optimize the average precision (AP) metric using the listwise AP loss,

$$\mathcal{L}_{mAP} = 1 - \frac{1}{N} \sum_{i=1}^N \text{AP}(x_i^\top X_N, Y_i) \quad (16)$$

which optimizes the global ranking of thousands of images simultaneously, instead of only a few images at a time. Here Y_i is the binary label to evaluate the relevance between batch images. $X_N = \{x_1, x_2, \dots, x_j, \dots, x_N\}$ denotes the features of all images, where each x_i is used as a potential query to rank the remaining batch images. Each similarity score $x_i^\top x_j$ can be measured by a cosine function.

It is demonstrated that training with AP-based loss improves retrieval performance [57],[130]. However average precision, as a metric, is normally non-differentiable. To directly optimize the AP loss during back-propagation, the key is that the indicator function for AP computing needs to be relaxed using methods such as triangular kernel-based soft assignment [57] or sigmoid function [130], as shown in Figure 8 (d).

4.2 Unsupervised Fine-tuning

Supervised network fine-tuning becomes infeasible when there is insufficient supervisory information, normally because of cost or unavailability. Therefore unsupervised fine-tuning methods

for image retrieval are quite necessary, but less studied [131]. For unsupervised fine-tuning, two directions are to mine relevance among features via manifold learning, and via clustering techniques, each discussed below.

4.2.1 Mining Samples with Manifold Learning

Manifold learning focuses on capturing intrinsic correlations on a manifold structure to mine or deduce relevance, as illustrated in Figure 9. Initial similarities between the extracted global features [132] or local features [14],[133] are used to construct an affinity matrix, which is then re-evaluated and updated using manifold learning [134]. According to the manifold similarity in the updated affinity matrix, positive and hard negative samples are selected for metric learning using pairwise ranking loss based functions such as pair loss [42],[133] or triplet loss [131],[135]. Note that this is different from the aforementioned methods for pairwise ranking loss based fine-tuning methods, where the hard positive and negative samples are explicitly selected from an ordered dataset according to the given affinity information.

It is important to capture the geometry of the manifold of deep features, generally involving two steps [134], known as diffusion. First, the affinity matrix (Figure 9) is interpreted as a weighted kNN graph, where each vector is represented by a node, and edges are defined by the pairwise affinities of two connected nodes. Then, the pairwise affinities are re-evaluated in the context of all other elements by diffusing the similarity values through the graph [48],[131],[133],[135], with recent strategies proposed such as regularized diffusion (RDP) [58] and regional diffusion [133]. For more details on diffusion methods refer to survey [134].

Most algorithms follow the two steps of [134]; the differences among methods lie primarily in three aspects:

- 1) **Similarity initialization**, which affects the subsequent kNN graph construction in an affinity matrix. Usually, an inner product [48] or Euclidean distance [45] is directly computed for the affinities. A Gaussian kernel function can be used [134],[135], or consider regional similarity from image patches [133].
- 2) **Transition matrix definition**, a row-stochastic matrix [134], determines the probabilities of transitioning from one node to another in the graph. These probabilities are proportional to the affinities between nodes, which can be measured by Geodesic distance (*e.g.* the summation of weights of relevant edges).
- 3) **Iteration scheme**, to re-evaluate and update the values in the affinity matrix by the manifold similarity until some convergence is achieved. Most algorithms are iteration-based [131],[134], as illustrated in Figure 9.

Diffusion process algorithms are indispensable for unsupervised fine-tuning. Better image similarity is guaranteed when it is improved based on initialization (*e.g.* regional similarity [133] or higher order information [45]). Diffusion is normally iterative and is computationally demanding [135], a limitation which cannot meet the efficiency requirements of image retrieval. To reduce the computational complexity, Bai *et al.* [58] propose a regularized diffusion process, facilitated by an efficient iteration-based solver. Zhao *et al.* [135] regard the diffusion process as a non-linear kernel mapping function, which is then modelled by a deep neural network. Other studies replace the diffusion process on a kNN graph with a diffusion network [47], which is derived from graph convolution networks [136], an end-to-end trainable framework which allows efficient computation during training and testing.

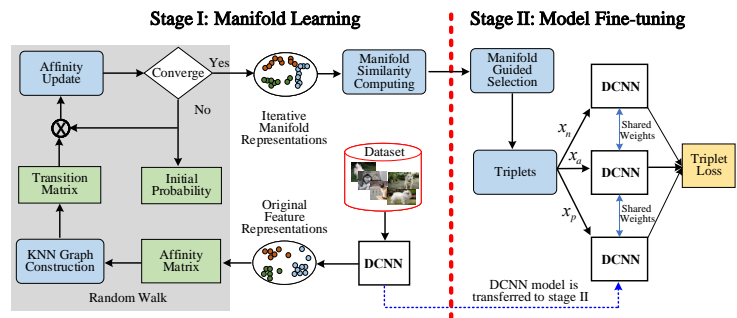


Fig. 9: Paradigm of manifold learning for unsupervised metric learning, based on triplet loss [131],[135].

Once the manifold space is learned, samples are mined by computing geodesic distances based on the Floyd-Warshall algorithm or by comparing the set difference [131]. The selected samples are fed into deep networks to perform fine-tuning.

4.2.2 Mining Samples by Clustering

Clustering is used to explore proximity information that has been studied in instance-level retrieval [44],[137],[138],[139],[140]. The rationale behind these methods is that samples in a cluster are likely to satisfy a degree of similarity.

One class of methods for clustering deep features is via k-means. Given k cluster centroids, during each training epoch a deep network alternates between two steps: first, a soft assignment between the feature representations and the cluster centroids; second, the cluster centroids are refined and, at the same time, the deep network is updated by learning from current high confidence assignments using a certain regularization. These two steps are repeated until a convergence criterion is met, at which point the cluster assignments are used as pseudo-labels [138],[139]. Alternatively, the pseudo-labels can be calculated from the samples in a cluster, *e.g.*, the mean values. For example, Tzelepi *et al.* [137] compute k nearest feature representations with respect to a query feature and then compute their mean vectors, which is used as a target for the query feature. In this case, fine-tuning is performed by minimizing the squared distance between each query feature and the mean of its k nearest features. Liu *et al.* [81] propose a self-taught hashing algorithm using a kNN graph construction to generate pseudo labels that are used to analyze and guide network training. Shen *et al.* [141] and Radenović *et al.* [44],[46] use Structure-from-Motion (SfM) for each image cluster to explore sample reconstructions to select images for triplet loss. Clustering methods depend on the Euclidean distance, making it difficult to reveal the intrinsic relationship between objects.

There are further techniques for instance retrieval, such as by using AutoEncoder [122],[142], generative adversarial networks (GANs) [143], convolutional kernel networks [112],[144], and graph convolutional networks [47]. For these methods, they focus on devising novel unsupervised frameworks to realize unsupervised learning, instead of iterative similarity diffusion or cluster refinement on feature space. For example, instead of performing iterative traversal on a set of nearest neighbors defined by kNN graph, Liu *et al.* [47] employ graph convolutional networks [136] to directly encode the neighbor information into image descriptors and then train the deep models to learn a new feature space. This method is demonstrated to significantly improve retrieval accuracy while maintaining efficiency. GANs are also explored, for the first time, for instance-level retrieval

in an unsupervised fashion [143]. The generator retrieves images that contain similar instances as a given image, while the discriminator judges whether the retrieved images have the specified instance which appeared in the query image. During training, the discriminator and the generator play a min-max game via an adversarial reward which is computed based on the cosine distance between the query image and the images retrieved by the generator.

5 STATE OF THE ART PERFORMANCE

5.1 Datasets

To demonstrate the effectiveness of methods, we choose the following commonly-used datasets for performance comparison:

UKBench (UKB) [145] consists of 10,200 images of objects. This dataset has 2,550 groups of images, each group having four images of the same object from different viewpoints or illumination conditions, which can be regarded as a kind of class-level supervision information. All images can be used as a query.

Holidays [51] consists of 1,491 images collected from personal holiday albums. Most images are scene-related. The dataset comprises 500 groups of similar images with a single query image for each group. The dataset also provides position information of the interest regions for each image.

Oxford-5k [94] consists of 5,062 images for 11 Oxford buildings. Each building is associated with five hand-drawn bounding box queries. According to the relevance level, each image of the same building is assigned a label *Good* (i.e., *positive*), *OK* (i.e., *positive*), *Junk*, or *Bad* (i.e., *negative*). *Junk* images can be discarded or regarded as *negative* examples [146],[54]. To build a tuple for each given query, one can select a positive example whose label corresponds to *Good* or *OK* in the same category, and select one negative example from each of the remaining building categories. Furthermore, an additional disjoint set of 100,000 distractor images is added to obtain Oxford-105k.

Paris-6k [95] includes 6,412 images and is categorized into 12 groups by architecture. The supervision information can be used like that of Oxford-5k. Likewise, an additional disjoint set of 100,000 distractor images is added to obtain Paris-106k.

INSTRE [93] consists of 28,543 images from 250 different object classes, including three disjoint subsets²: INSTRE-S1, INSTRE-S2, INSTRE-M. INSTRE dataset has bounding box annotations, providing single-labelled and double-labelled class information for single- and multiple-object retrieval, respectively. One can use the class information to build a tuple, with two positive examples from the same class and one negative from one of the remaining classes. The performance evaluation on INSTRE in our experiments follows the protocol in [133].

Google Landmarks Dataset (GLD) [5],[63] consists of GLD-v1 and GLD-v2. GLD-v2 is mainly recommended to use and it has the advantage of stability where all images have permissive licenses [71]. GLD-v2 is divided into three subsets: (i) 118k query images with ground-truth annotations, (ii) 4.1M training images of 203k landmarks with labels, and (iii) 762k index images of 101k landmarks. Due to its large scale, GLD-v2 provides class-level ground-truth which can be used to build training tuples. Due to its image diversity, it may produce clutter images for each landmark so it is necessary to introduce pre-processing methods to select the more relevant images [147]. Finally, the training set is cleaned by removing these clutters, consisting of a subset “GLD-v2-clean” containing 1.6M images of 81k landmarks. Since Google landmarks dataset stills lack bounding box for objects of

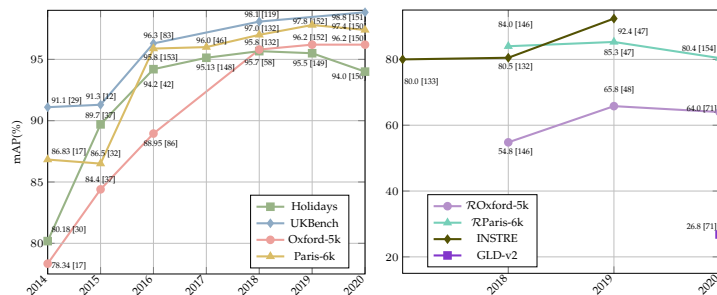


Fig. 10: Performance improved from 2014 to 2020.

interest, Teichmann *et al.* [74] provide a new dataset of landmark bounding boxes, based on GLD. This patch-level supervision information can help locate the most relevant regions.

Note that, additional queries and distractor images have been added into Oxford-5k and Paris-6k, producing the Revisited Oxford (*ROxford*) and Revisited Paris (*RParis*) datasets where each image of the same building is assigned a label *Easy*, *Hard*, *Unclear*, or *Negative* [146]. Different label combinations are used as *positive* according to the difficulty level of different setups. During testing, if there are no positive images for a query, then that query is excluded from the evaluation. For details, we refer the reader to [146]. We undertake partial comparisons under the *hard* evaluation protocol on these revisited datasets.

5.2 Evaluation Metrics

Average precision (AP) refers to the coverage area under the precision-recall (PR) curve. A larger AP implies a higher PR curve and better retrieval accuracy. AP can be calculated as $AP = \frac{1}{R} \sum_{k=1}^N P(k) \cdot rel(k)$, where R denotes the number of relevant results for the query image from the total number N of images. $P(k)$ is the precision of the top k retrieved images, and $rel(k)$ is an indicator function equal to 1 if the item within rank k is a relevant image and 0 otherwise. Mean average precision (mAP) is adopted for the evaluation over all query images, i.e., $mAP = \frac{1}{Q} \sum_{q=1}^Q AP(q)$, where Q is the number of query images.

The **N-S score** is a metric used for UKBench [145]; the N-S score is the average for the top-4 precision over the dataset.

5.3 Performance Comparison and Analysis

Overview. Figure 10 summarizes the performance over 6 datasets from 2014 to 2020. Early on, the powerful feature extraction of DCNNs led to rapid improvements. Subsequent key ideas have been to extract instance features at the region level to reduce image clutter [30], and to improve feature discriminativity by using methods including feature fusion [86],[148],[151], feature aggregation [32],[69], and feature embedding [86]. Fine-tuning is an important strategy to improve performance by tuning deep networks specific for learning instance features [58],[150]. For instance, the accuracy increases steadily from 78.34% [17] to 96.2% [152] on the Oxford-5k dataset when manifold learning is used to fine-tune deep networks. The mAP on *RParis-6k* and *ROxford-5k* is smaller than *Paris-6k* and *Oxford-5k*, leaving room for improvement.

We report results using off-the-shelf models (Table 2) and fine-tuning networks (Table 3). In Table 2, single-pass and multiple-pass are analyzed, while supervised and unsupervised fine-tuning are compared in Table 3. Since there are many aspects that vary across the different methods, making them not directly comparable, we mainly draw some general claims or trends based on the collected results.

TABLE 2: Performance evaluation of off-the-shelf DCNN models. “●” indicates that the models or layers are combined to learn features; “PCA_w” indicates PCA with whitening on the extracted features to improve robustness; “MP” means Max Pooling; “SP” means Sum Pooling. The CNN-M network with “*” has an architecture similar to that of AlexNet. “-” means that the results were not reported.

Type	Method	Backbone DCNN	Output Layer	Embed. Aggre.	Feat. Dim	Holidays	UKB	Oxford5k (+100k)	Paris6k (+100k)	Brief Conclusions and Highlights
Single-pass	Neural codes [39]	AlexNet	FC6	PCA	128	74.7	3.42 (N-S)	43.3 (38.6)	–	Compressed neural codes of different layers are explored. AlexNet is also fine-tuned for retrieval.
	SPoC [12]	VGG16	Conv5	SPoC + PCA _w	256	80.2	3.65 (N-S)	58.9 (57.8)	–	Exploring Gaussian weighting scheme <i>i.e.</i> , the centering prior, to improve the discrimination of features.
	CroW [15]	VGG16	Conv5	CroW + PCA _w	256	85.1	–	68.4 (63.7)	76.5 (69.1)	The spatialwise and channelwise weighting mechanisms are utilized to highlight crucial convolutional features.
	R-MAC [32]	VGG16	Conv5	R-MAC + PCA _w	512	–	–	66.9 (61.6)	– (75.7)	Sliding windows with different scales on convolutional feature maps to encode multiple image regions.
	Multi-layer CNN [102]	VGG16	FC6 ● Conv4~5	SP	4096	91.4	3.68 (N-S)	61.5 (–)	–	Layer-level feature fusion and the complementary properties of different layers are explored.
	BLCF [84]	VGG16	Conv5	BoW + PCA _w	25k	–	–	73.9 (59.3)	82.0 (64.8)	Both global features and local features are explored, demonstrating that local features have higher accuracy.
Multiple-pass	OLDFP [61]	AlexNet	FC6	MP + PCA _w	512	88.5	3.81 (N-S)	60.7 (–)	66.2 (–)	Exploring the impact of proposal number. Patches are extracted by RPNs (see Figure 4 (d)) and the features are encoded in an orderless way.
	MOP-CNN [30]	AlexNet	FC7	VLAD + PCA _w	2048	80.2	–	–	–	Image patches are extracted densely, as shown in Figure 4 (c). Multi-scale patch features are further embedded into VLAD descriptors.
	CNNaug-ss [29]	Overfeat [155]	FC	PCA _w	15k	84.3	91.1 (mAP)	68.0 (–)	79.5 (–)	Image patches are extracted densely, as shown in Figure 4 (c). Image regions at different locations with different sizes are included.
	MOF [36]	CNN-M* [156]	FC7 ● Conv	SP or MP + BoW	20k	76.8	3.00 (N-S)	–	–	Exploring layer-level fusion scheme. Image patches are extracted using spatial pyramid modeling, as shown in Figure 4 (b).
	Multi-scale CNN [69]	VGG16	Conv5	SP or MP + PCA _w	32k	89.6	95.1 (mAP)	84.3 (–)	87.9 (–)	Image patches are extracted in a dense manner, as shown in Figure 4 (c). Geometric invariance is considered when aggregating patch features.
	LDD [115]	VGG19	Conv5	BoW + PCA _w	500k	84.6	–	83.3 (–)	87.2 (–)	Image patches are obtained using a uniform square mesh, as shown in Figure 4 (a). Patch features are encoded into BoW descriptors.
	DeepIndex [52]	AlexNet ● VGG19	FC6-7 ● FC17-18	BoW + PCA	512	81.7	3.32 (N-S)	–	75.4 (–)	Exploring layer-level and model-level fusion methods. Image patches are extracted using spatial pyramid modeling, as shown in Figure 4 (b).
	CCS [53]	GoogLeNet	Conv	VLAD + PCA _w	128	84.1	3.81 (N-S)	64.8 (–)	76.8 (–)	Object proposals are extracted by RPNs, as shown in Figure 4 (d). Object level and point level feature concatenation schemes are explored.

TABLE 3: Performance evaluation of methods in which DCNN models are fine-tuned, in a supervised or an unsupervised manner. “CE Loss” means the models are fine-tuned using the classification-based loss function in the form of Eq. 9. “Siamese Loss” is in the form of Eq. 12. “Regression Loss” is in the form of Eq. 11. “Triplet Loss” is in the form of Eq. 14.

Type	Method	Backbone DCNN	Training set	Output Layer	Embed. Aggre.	Loss Function	Feat. Dim	Holidays	UKB	Oxford5k (+100k)	Paris6k (+100k)	Brief Conclusions and Highlights
Supervised Fine-tuning	Neural codes [39]	AlexNet	Landmarks [39]	FC6	PCA	CE Loss	128	78.9	3.29 (N-S)	55.7 (52.3)	–	The first work which fine-tunes deep networks for IIR. Compressed neural codes and different layers are explored.
	Nonmetric [40]	VGG16	Landmarks	Conv5	PCA _w	Regression Loss	512	–	–	88.2 (82.1)	88.2 (82.9)	Similarity learning of similar and dissimilar pairs is performed by a neural network, optimized using regression loss.
	SIAM-FV [41]	VGG16	Landmarks	Conv5	FV + PCA _w	Siamese Loss	512	–	–	81.5 (76.6)	82.4 (–)	Fisher Vector is integrated on top of VGG and is trained with VGG simultaneously.
	Faster R-CNN [92]	VGG16	Oxford5k Paris6k	Conv5	MP / SP	Regression Loss	512	–	–	75.1 (–)	80.7 (–)	RPN is fine-tuned, based on bounding box coordinates and class scores for specific region query which is region-targeted.
	SIFT-CNN [124]	VGG16	Holidays UKB	Conv5	SP	Siamese Loss	512	88.4	3.91 (N-S)	–	–	SIFT features are used as supervisory information for mining positive and negative samples.
	Quartet-Net [129]	VGG16	GeoPair [129]	FC6	PCA	Siamese Loss	128	71.2	87.5 (mAP)	48.5 (–)	48.8 (–)	Quartet-net learning is explored to improve feature discrimination where double-margin contrastive loss is used.
	NetVLAD [43]	VGG16	Tokyo Time Machine	VLAD Layer	PCA _w	Triplet Loss	256	79.9	–	62.5 (–)	72.0 (–)	VLAD is integrated at the last convolutional layer of VGG16 network as a plugged layer.
	Deep Retrieval [82]	ResNet-101	Landmarks	Conv5 Block	MP + PCA _w	Triplet Loss	2048	90.3	–	86.1 (82.8)	94.5 (90.6)	Dataset is cleaned automatically. Features are encoded by R-MAC. RPN is used to extract the most relevant regions.
DELFI [5]	ResNet-101	GLDv1	Conv4 Block	Attention + PCA _w	CE Loss	2048	–	–	83.8 (82.6)	85.0 (81.7)	Exploring the FCN to extract region-level features and construct feature pyramids of different sizes.	
Unsupervised Fine-tuning	MoM [131]	VGG16	Flickr 7M	Conv5	MP + PCA _w	Siamese Loss	64	87.5	–	78.2 (72.6)	85.1 (78.0)	Exploring manifold learning for mining dis/similar samples. Features are tested globally and regionally.
	GeM [46]	VGG16	Flickr 7M	Conv5	GeM Pooling	Siamese Loss	512	83.1	–	82.0 (76.9)	79.7 (72.6)	Fine-tuning CNNs on an unordered dataset. Samples are selected from an automated 3D reconstruction system.
	SfM-CNN [44]	VGG16	Flickr 7M	Conv5	PCA _w	Siamese Loss	512	82.5	–	77.0 (69.2)	83.8 (76.4)	Employing Structure-from-Motion to select positive and negative samples from unordered images.
	MDP-CNN [135]	ResNet-101	Landmarks	Conv5 Block	SP	Triplet Loss	2048	–	–	85.4 (85.1)	96.3 (94.7)	Exploring global feature structure by modeling the manifold learning to select positive and negative pairs.
	IME-CNN [45]	ResNet-101	Oxford105k Paris106k	IME Layer	MP	Regression Loss	2048	–	–	92.0 (87.2)	96.6 (93.3)	Graph-based manifold learning is explored to mine the matching and non-matching pairs in unordered datasets.

Evaluation for single feedforward pass. In general, we observe that fully-connected layers used as feature extractors may give a lower accuracy (*e.g.*, 74.7% on Holidays in [39]), compared to using convolutional layers in Table 2. For the case where the same VGG net is used, the way to embed or aggregate features is critical. The methods shown in Figure 5 improve the discrimination of convolutional feature maps and perform differently in Table 2, 66.9% of R-MAC [95] and 58.9% of SPoC [12] on Oxford-5k, differences which we see as critical factors for further analysis. If embedded by a BoW model, the results are competitive on Oxford-5k and Paris-6k (73.9% and 82.0%, respectively), while its codebook size is 25k, which may affect retrieval efficiency. Moreover, layer-level feature fusion improves retrieval accuracy. Yu *et al.* [102] combine three layers (mAP of 91.4% on Holidays), outperforming the performance of non-fusion method [12] (mAP of 80.2%).

Evaluation for multiple feedforward pass. Results for the

methods of Figure 4 are reported in Table 2. Among them, extracting image patches densely using VGG [49] has the highest performance on the 4 datasets [29], and rigid grid with BoW encoding [115] is competitive (mAP of 87.2% on Paris-6k). These two methods consider more patches, even background information, when used for feature extraction. Instead of generating patches densely, region proposals and spatial pyramid modeling introduce a degree of purpose and efficiency in processing image objects. Spatial information is better maintained using multiple-pass schemes than with single-pass. For example, a shallower network (AlexNet) and region proposal networks in [61] have a UKBench N-Score of 3.81, higher than using deeper networks [12],[39],[102]. Besides feeding image patches into the same network, model-level fusion also exploits complementary spatial information to improve accuracy. For instance, as reported in [52], which combines AlexNet and VGG, the results on Holidays (81.7% of mAP) and UKBench (3.32 of N-Score) are better than

these in [36] (76.75% and 3.00, respectively).

Evaluation for supervised fine-tuning. Compared to off-the-shelf models, fine-tuning deep networks usually improves accuracy, see Table 3. For instance, the result on Oxford-5k [32] by using a pre-trained VGG is improved from 66.9% to 81.5% in [41] when a single-margin Siamese loss is used. Similar trends can be also observed on the Paris-6k dataset. For classification-based fine-tuning, its performance may be improved by using powerful DCNNs and feature enhancement methods such as the attention mechanism in [5], with an mAP increased from 55.7% in [39] to 83.8% in [5] on Oxford-5k. As for pairwise ranking loss fine-tuning, in some cases the loss used for fine-tuning is essential for performance improvement. For example, RPN is re-trained using regression loss on Oxford-5k and Paris-6k (75.1% and 80.7%, respectively) [92]. Its results are lower than the results from [40] (88.2% and 88.2%, respectively) where a transformation matrix is used to learn visual similarity. However, when RPN is trained by using triplet loss such as [82], the effectiveness of retrieval is improved significantly where the results are 86.1% (on Oxford-5k) and 94.5% (on Paris-6k). Feature embedding methods are important for retrieval accuracy; Ong *et al.* [41] embedded *Conv5* feature maps by Fisher Vector and achieved an mAP of 81.5% on Oxford-5k, while embedding feature maps by using VLAD achieves an mAP of 62.5% on this dataset [43],[44].

Evaluation for unsupervised fine-tuning. Compared to supervised fine-tuning, unsupervised fine-tuning methods are relatively less explored. The difficulty for unsupervised fine-tuning is to mine sample relevance without ground-truth labels. In general, unsupervised fine-tuning methods should be expected to have lower performance than supervised. For instance, supervised fine-tuning using Siamese loss [124] achieves an mAP 88.4% on Holidays, while unsupervised fine-tuning using the same loss function in [44],[46],[131] achieves 82.5%, 83.1%, and 87.5%, respectively. However, unsupervised fine-tuning methods can achieve a similar accuracy, even outperform the supervised fine-tuning, if a suitable feature embedding method is used. For instance, Zhao *et al.* [135] explore global feature structure modeling the manifold learning, producing an mAP of 85.4% (on Oxford-5k) and 96.3% (on Paris-6k), which is similar to supervised results [82] of 86.1% (on Oxford-5k) and 94.5% (on Paris-6k). As another example, the precision of ResNet-101 fine-tuned by cross entropy loss achieves 83.8% on Oxford-5k [5], while the precision is further improved to 92.0% when an IME layer is used to embed features and fine-tuned in an unsupervised way [45]. Note that fine-tuning strategies are related to the type of the target retrieval datasets. As demonstrated in Table 4 and [71], fine-tuning on different datasets may produce a different final retrieval performance.

Network depth. We compare the efficacy of DCNNs by depth, following the fine-tuning protocols³ in [46]. For fair comparisons, all convolutional features from these backbone DCNNs are aggregated by MAC [69], and fine-tuned by using the same loss function with the same learning rate, thus the adopted methods are the same except for the DCNN depth. We use the default feature dimension (*i.e.* AlexNet (256), VGG (512), GoogLeNet (1024), ResNet-50/101 (2048)). The results are reported in Figure 11 (a). We observe that the deeper networks consistently lead to better accuracy due to extracting more discriminative features.

Feature aggregation methods. The methods of embedding convolutional feature maps were illustrated in Figure 5. We use the off-the-shelf VGG (without updating parameters) on the Oxford and Paris datasets. The results are reported in Figure 11 (b).

TABLE 4: Evaluations of training sets and retrieval reranking. Numerical results are cited from [71].

Conditions	Global	Local reranking	Training set		ROxf	RPar	GLD-v2 testing	
			GLD-v1	GLD-v2-clean				
ResNet-50	Case 1	✓	✗	✓	✗	45.1	63.4	20.4
	Case 2	✓	✗	✗	✓	51.0	71.5	24.1
	Case 3	✓	✓	✓	✗	54.2	64.9	22.3
	Case 4	✓	✓	✗	✓	57.9	71.0	24.3
ResNet-101	Case 5	✓	✗	✓	✗	51.2	64.7	21.7
	Case 6	✓	✗	✗	✓	55.6	72.4	26.0
	Case 7	✓	✓	✓	✗	59.3	65.5	24.3
	Case 8	✓	✓	✗	✓	64.0	72.8	26.8

We observe that the different ways to aggregate the same off-the-shelf DCNN leads to differences in retrieval performance. These reported results provide a reference for feature aggregation when one uses convolutional layers for performing retrieval tasks.

Global feature dimension. We add fully-connected layers on the top of pooled convolutional features of ResNet-50 to obtain global descriptors with their dimensions varying from 32 to 8192. The results of 5 datasets are shown in Figure 11 (c). It is expected that higher-dimension features usually capture more semantics and are helpful for retrieval. The performance tends to be stable when the dimension is very large.

Number of image regions. We compare the retrieval performance when different number of regions are fed and other components are kept the same, as depicted in Figure 11 (d). Convolutional features of each region are pooled as 2048-dim regional features by MAC and then aggregated into a global one. Note that the final memory requirement is identical for the case that a holistic image is used as input (*i.e.*, regarded as the case where only one region is used). Regional inputs on an image are extracted with a 40% overlap of neighboring regions and the number varying from 1 to 41. For Oxford-5k, the best result is given by the case where 9 image regions are used. For the rest datasets, 3 image regions give the best results. Finally, more regions extracted from one image decline the retrieval mAP. A reason is that features of background or irrelevant regions have also been aggregated, and negatively affect the performance.

Fine-tuning datasets and retrieval reranking. We compare performance on \mathcal{R} Oxford-5k, \mathcal{R} Paris-6k, and GLD-v2, aiming at comparing the role of different fine-tuning training sets and the effectiveness of retrieval reranking. Table 4 lists 8 experimental scenarios using two network backbones, as in [71].

Since GLD-v2 provides class-level ground-truth, its including images show large context diversity and may pose challenges to the network fine-tuning. Thus, the pre-processing steps, as proposed in [147],[63], are necessary to select the more coherent images, referring to the GLD-v2-clean subset. As a result, when using the global features only, this cleaned version of the training set improves the performance, as observed in Cases 1/5 and Cases 2/6 for ResNet-50/ResNet-101, respectively. As an important postprocessing strategy, reranking further boosts the accuracy after the initial filtering step by using global features.

6 CONCLUSIONS AND OUTLOOKS

As a comprehensive yet timely survey on instance retrieval using deep learning, this paper has discussed the main challenges, presented a taxonomy of recent developments according to their roles in implementing instance retrieval, highlighted the recent representative methods and analyzed their merits and demerits, discussed the datasets, evaluation protocols, and SOTA performance. Nowadays the exponentially increasing amount of image

3. <https://github.com/filipradenovic/cnnimageretrieval-pytorch>

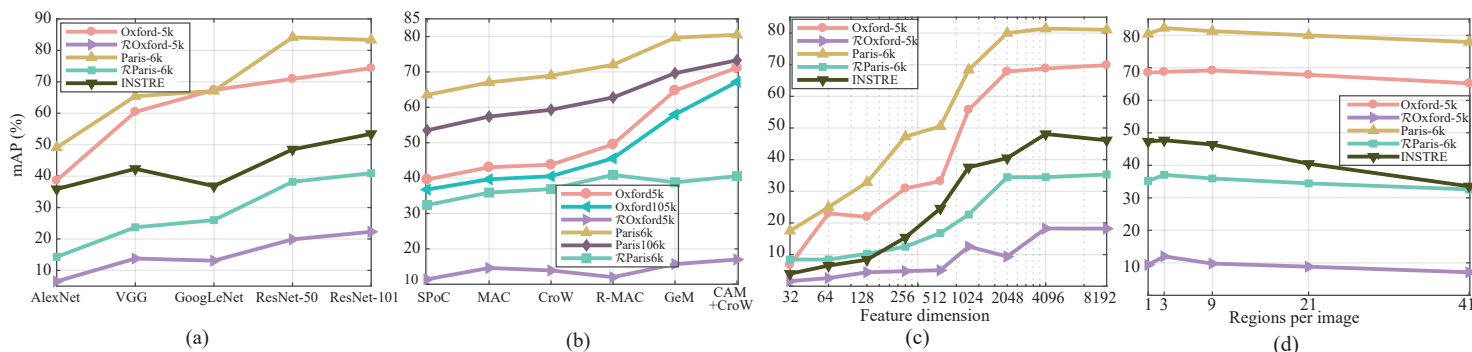


Fig. 11: (a) The effectiveness of different DCNNs; (b) Comparison of the feature aggregation methods in Figure 5; (c) The impact of global feature dimension by using ResNet-50; (d) Performance comparison when aggregating different numbers of image regions.

and video data due to surveillance, e-commerce, medical images, handheld devices, robotics, *etc.*, offers an endless potential for applications of instance retrieval. Although significant progress has been made, as discussed in Section 1.2, the main challenges in instance retrieval have not been fully addressed. Below we identify a number of promising directions for future research.

(1) Accurate and robust feature representations. One of the main challenges in instance retrieval is the large intra-class variations due to changes in viewpoint, scale, illumination, weather condition and background clutter *etc.*, as we discussed in Section 1.2. However, DCNN representations have very little invariance, even though trained with lots of augmented data [146]. Fortunately, before deep learning, with instance retrieval there are lots of important ideas in handling such intra-class variations like local interest point detectors and local invariant descriptors. Therefore, it is worth enabling DCNN to learn more accurate and robust representations via leveraging such traditional wisdom to design better DCNNs. In addition, unlike most objects in existing benchmarks which are rigid, planar and textured, textureless objects, 3D objects, transparent objects, reflective surfaces, *etc.* are still very hard to find.

In addition, pursuing accuracy alone is not sufficient, as instance retrieval systems should be able to resist potential adversarial attacks. Recently, deep networks have been proven to be fooled rather easily by adversarial examples [157], *i.e.*, images added with intentionally designed yet nearly imperceptible perturbations, which raises serious safety and robustness concerns. However, adversarial robustness in instance retrieval [157],[158] has received very little attention, and should merit further effort.

(2) Compact and efficient deep representations. In instance retrieval, searching efficiently is as critical as searching accurately, especially for the pervasive mobile or wearable devices with very limited computing resources. However, existing methods adopt large scale, energy hungry DCNNs that are very difficult to be deployed in mobile devices. Hence, there has been pressing needs to develop compact, efficient, yet reusable deep representations tailored to the resource limited devices, like using binary neural networks [64],[85],[86].

(3) Learning with fewer labels. Deep learning require a large amount of high-quality labeled dataset to achieve high accuracy. The presence of labels errors or the limited amount of labeled data can greatly degrade DCNN’s accuracy. However, collecting massive amounts of accurately labeled data is costly. In practical scenarios, datasets like GLDv2 [5],[63] have long-tailed distributions, and noisy labels. Thus, to address such limitations, few shot learning [159], self-supervised learning [160], imbalanced learning [63], noisy label aware learning [161] *etc.* should be paid more attention in instance retrieval in the future.

(4) Continual learning for instance retrieval. In specific, the current IIR methods make restrictive assumptions, such as the training data being enough and stationary, retraining from scratch being possible when new data becomes available, which is problematic in realistic conditions. Our living world is continuously varying, and in general data distributions are often non-stationary, new data may be added, and previously unseen classes may be encountered. Thus, continual learning plays a vital role in continuously updating the IIR systems. The key issues are how to retain and utilize the previously learned knowledge, how to update the retrieval system as new images becomes available, and how to learn and improve over time.

(5) Privacy-aware instance retrieval. Most IIR systems concentrate on improving the accuracy or efficiency performance, and the higher performance might come at the cost of users’ privacy. Therefore, in some cases, such as personalized search systems, the privacy protection problem is also an important issue to be considered. Deep models should be privacy-aware and protect users’ personalized searching experience to avoid their worries about using such IIR systems.

(6) Video instance retrieval. Searching a specific instance in an image cannot always meet the requirements in some scenarios such as the video surveillance system in the field of searching criminals. Currently, with the rapid growth of video data, retrieving a certain object, place, or action in videos has become more and more important and highly necessary in the future. For video instance retrieval, 3D-CNNs models need to be built to learn video’s spatio-temporal representations to compute the semantic similarity of instances.

REFERENCES

- [1] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, “Content-based image retrieval at the end of the early years,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 12, pp. 1349–1380, 2000.
- [2] M. S. Lew, N. Sebe, C. Djeraba, and R. Jain, “Content-based multimedia information retrieval: State of the art and challenges,” *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 2, no. 1, pp. 1–19, 2006.
- [3] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, “Scalable person re-identification: A benchmark,” in *ICCV*, 2015, pp. 1116–1124.
- [4] X. Liu, S. Zhang, X. Wang, R. Hong, and Q. Tian, “Group-group loss-based global-regional feature learning for vehicle re-identification,” *IEEE Trans. Image Process.*, vol. 29, pp. 2638–2652, 2019.
- [5] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, “Largescale image retrieval with attentive deep local features,” in *ICCV*, 2017.
- [6] U. Chaudhuri, B. Banerjee, and A. Bhattacharya, “Siamese graph convolutional network for content based remote sensing image retrieval,” *Comput. Vis. Image Underst.*, vol. 184, pp. 22–30, 2019.
- [7] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, “Deepfashion: Powering robust clothes recognition and retrieval with rich annotations,” in *CVPR*, 2016, pp. 1096–1104.
- [8] A. Gordo and D. Larlus, “Beyond instance-level image retrieval: Leveraging captions to learn a global visual representation for semantic retrieval,” in *CVPR*, 2017, pp. 6589–6598.

- [9] B. Barz and J. Denzler, "Content-based image retrieval and the semantic gap in the deep learning era," in *ICPR*, 2021, pp. 245–260.
- [10] X. Wang, X. Han, W. Huang, D. Dong, and M. R. Scott, "Multi-similarity loss with general pair weighting for deep metric learning," in *CVPR*, 2019, pp. 5022–5030.
- [11] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *CVPR*, 2014, pp. 1386–1393.
- [12] A. Babenko and V. Lempitsky, "Aggregating local deep features for image retrieval," in *ICCV*, 2015, pp. 1269–1277.
- [13] L. Zheng, Y. Yang, and Q. Tian, "SIFT meets CNN: A decade survey of instance retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 5, pp. 1224–1244, 2018.
- [14] W. Zhang, C.-W. Ngo, and X. Cao, "Hyperlink-aware object retrieval," *IEEE Trans. Image Process.*, vol. 25, no. 9, pp. 4186–4198, 2016.
- [15] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *ECCV*, 2016.
- [16] L. Zhang and Y. Rui, "Image search from thousands to billions in 20 years," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 9, no. 1s, p. 36, 2013.
- [17] J. Wan, D. Wang, S. C. H. Hoi, P. Wu, J. Zhu, Y. Zhang, and J. Li, "Deep learning for content-based image retrieval: A comprehensive study," in *ACM MM*, 2014, pp. 157–166.
- [18] A. Alzu'bi, A. Amira, and N. Ramzan, "Semantic content-based image retrieval: A comprehensive study," *J. Vis. Commun. Image Represent.*, vol. 32, pp. 20–54, 2015.
- [19] X. Li, T. Uricchio, L. Ballan, M. Bertini, C. G. Snoek, and A. D. Bimbo, "Socializing the semantic gap: A comparative survey on image tag assignment, refinement, and retrieval," *ACM Comput. Surv. (CSUR)*, vol. 49, no. 1, pp. 1–39, 2016.
- [20] W. Zhou, H. Li, and Q. Tian, "Recent advance in content-based image retrieval: A literature survey," *arXiv preprint arXiv:1706.06064*, 2017.
- [21] L. Piras and G. Giacinto, "Information fusion in content based image retrieval: A comprehensive overview," *Inf. Fusion*, vol. 37, pp. 50–60, 2017.
- [22] J. Wang, T. Zhang, N. Sebe, H. T. Shen *et al.*, "A survey on learning to hash," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 769–790, 2018.
- [23] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.
- [24] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *CVPR*, 2003, pp. 1470–1477.
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NeurIPS*, 2012.
- [26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [27] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015, pp. 91–99.
- [28] S. Minaee, Y. Y. Boykov, F. Porikli, A. J. Plaza, N. Kehtarnavaz, and D. Terzopoulos, "Image segmentation using deep learning: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.
- [29] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features off-the-shelf: an astounding baseline for recognition," in *CVPR workshops*, 2014, pp. 806–813.
- [30] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *ECCV*, 2014.
- [31] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *NeurIPS*, 2014, pp. 3320–3328.
- [32] G. Toliás, R. Sivic, and H. Jégou, "Particular object retrieval with integral max-pooling of CNN activations," in *ICLR*, 2015, pp. 1–12.
- [33] A. Jiménez, J. M. Alvarez, and X. Giró Nieto, "Class-weighted convolutional features for visual instance search," in *BMVC*, 2017, pp. 1–12.
- [34] T.-T. Do, T. Hoang, D.-K. L. Tan, H. Le, T. V. Nguyen, and N.-M. Cheung, "From selective deep convolutional features to compact binary representations for image retrieval," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 15, no. 2, pp. 1–22, 2018.
- [35] J. Xu, C. Wang, C. Qi, C. Shi, and B. Xiao, "Unsupervised part-based weighting aggregation of deep convolutional features for image retrieval," in *AAAI*, 2018, pp. 7436–7443.
- [36] Y. Li, X. Kong, L. Zheng, and Q. Tian, "Exploiting hierarchical activations of neural network for image retrieval," in *ACM MM*, 2016.
- [37] A. Sharif Razavian, J. Sullivan, A. Maki, and S. Carlsson, "A baseline for visual instance retrieval with deep convolutional networks," in *ICLR*, 2015.
- [38] T.-T. Do and N.-M. Cheung, "Embedding based on function approximation for large scale image search," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 3, pp. 626–638, 2017.
- [39] A. Babenko, A. Slesarev, A. Chigorin, and V. Lempitsky, "Neural codes for image retrieval," in *ECCV*, 2014, pp. 584–599.
- [40] N. Garcia and G. Vogiatzis, "Learning non-metric visual similarity for image retrieval," *Image Vis. Comput.*, vol. 82, pp. 18–25, 2019.
- [41] E.-J. Ong, S. Husain, and M. Bober, "Siamese network of deep fisher-vector descriptors for image retrieval," *arXiv preprint arXiv:1702.00338*, 2017.
- [42] A. Gordo, J. Almazán, J. Revaud, and D. Larlus, "Deep image retrieval: Learning global representations for image search," in *ECCV*, 2016.
- [43] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *CVPR*, 2015, pp. 5297–5307.
- [44] F. Radenović, G. Toliás, and O. Chum, "CNN image retrieval learns from BoW: Unsupervised fine-tuning with hard examples," in *ECCV*, 2016, pp. 3–20.
- [45] J. Xu, C. Wang, C. Qi, C. Shi, and B. Xiao, "Iterative manifold embedding layer learned by incomplete data for large-scale image retrieval," *IEEE Trans. Multimedia*, vol. 21, no. 6, pp. 1551–1562, 2017.
- [46] F. Radenović, G. Toliás, and O. Chum, "Fine-tuning CNN image retrieval with no human annotation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 7, pp. 1655–1668, 2018.
- [47] C. Liu, G. Yu, M. Volkovs, C. Chang, H. Rai, J. Ma, and S. K. Gorti, "Guided similarity separation for image retrieval," in *NeurIPS*, 2019.
- [48] C. Chang, G. Yu, C. Liu, and M. Volkovs, "Explore-exploit graph traversal for image retrieval," in *CVPR*, 2019, pp. 9423–9431.
- [49] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [50] J. Yue-Hei Ng, F. Yang, and L. S. Davis, "Exploiting local features from deep networks for image retrieval," in *CVPR workshops*, 2015.
- [51] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *ECCV*, 2008.
- [52] Y. Liu, Y. Guo, S. Wu, and M. S. Lew, "Deepindex for accurate and efficient image retrieval," in *ICMR*, 2015, pp. 43–50.
- [53] K. Yan, Y. Wang, D. Liang, T. Huang, and Y. Tian, "CNN vs. SIFT for image retrieval: Alternative or complementary?" in *ACM MM*, 2016.
- [54] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *CVPR*, 2010.
- [55] F. Perronnin and C. Dance, "Fisher kernels on visual vocabularies for image categorization," in *CVPR*, 2007, pp. 1–8.
- [56] O. Siméoni, Y. Avrithis, and O. Chum, "Local features and visual words emerge in activations," in *CVPR*, 2019, pp. 11 651–11 660.
- [57] J. Revaud, J. Almazán, R. S. Rezende, and C. R. d. Souza, "Learning with average precision: Training image retrieval with a listwise loss," in *ICCV*, 2019, pp. 5107–5116.
- [58] B. Song, X. Bai, Q. Tian, and L. J. Latecki, "Regularized diffusion process on bidirectional context for object retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 5, pp. 1213–1226, 2018.
- [59] A. El-Nouby, N. Neverova, I. Laptev, and H. Jégou, "Training vision transformers for image retrieval," *arXiv preprint arXiv:2102.05644*, 2021.
- [60] F. Tan, J. Yuan, and V. Ordonez, "Instance-level image retrieval using reranking transformers," *ICCV*, 2021.
- [61] K. Reddy Mopuri and R. Venkatesh Babu, "Object level deep feature pooling for compact image representation," in *CVPR Workshops*, 2015.
- [62] O. Morere, J. Lin, A. Veillard, L.-Y. Duan, V. Chandrasekhar, and T. Poggio, "Nested invariance pooling and RBM hashing for image instance retrieval," in *ICMR*, 2017, pp. 260–268.
- [63] T. Weyand, A. Araujo, B. Cao, and J. Sim, "Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval," in *CVPR*, 2020, pp. 2575–2584.
- [64] J. Song, T. He, L. Gao, X. Xu, and H. T. Shen, "Deep region hashing for efficient large-scale instance search from images," in *AAAI*, 2018.
- [65] K. Lin, J. Lu, C.-S. Chen, J. Zhou, and M.-T. Sun, "Unsupervised deep learning of compact binary descriptors," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.
- [66] W. Zhou, H. Li, J. Sun, and Q. Tian, "Collaborative index embedding for image retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 5, pp. 1154–1166, 2017.
- [67] H. Liu, Y. Tian, Y. Yang, L. Pang, and T. Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *CVPR*, 2016.
- [68] K. Ozaki and S. Yokoo, "Large-scale landmark retrieval/recognition under a noisy and diverse dataset," in *CVPR Workshop*, 2019.
- [69] A. S. Razavian, J. Sullivan, S. Carlsson, and A. Maki, "Visual instance retrieval with deep convolutional networks," *ITE Trans. Media Technol. Appl.*, vol. 4, no. 3, pp. 251–258, 2016.
- [70] S. Pang, J. Xue, J. Zhu, L. Zhu, and Q. Tian, "Unifying sum and weighted aggregations for efficient yet effective image representation computation," *IEEE Trans. Image Process.*, vol. 28, no. 2, pp. 841–852, 2018.
- [71] B. Cao, A. Araujo, and J. Sim, "Unifying deep local and global features for efficient image search," in *ECCV*, 2020, pp. 726–743.
- [72] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

- [73] G. Tolias, Y. Avrithis, and H. Jégou, "Image search with selective match kernels: aggregation across single and multiple images," *Int. J. Comput. Vis.*, vol. 116, no. 3, pp. 247–261, 2016.
- [74] M. Teichmann, A. Araujo, M. Zhu, and J. Sim, "Detect-to-retrieve: Efficient regional aggregation for image search," in *CVPR*, 2019.
- [75] S. Pang, J. Ma, J. Zhu, J. Xue, and Q. Tian, "Improving object retrieval quality by integration of similarity propagation and query expansion," *IEEE Trans. Multimedia*, vol. 21, no. 3, pp. 760–770, 2018.
- [76] C. Qi, C. Shi, J. Xu, C. Wang, and B. Xiao, "Spatial weighted fisher vector for image retrieval," in *ICME*, 2017, pp. 463–468.
- [77] H. J. Kim, E. Dunn, and J.-M. Frahm, "Learned contextual feature reweighting for image geo-localization," in *CVPR*, 2017, pp. 3251–3260.
- [78] E. Mohamedano, K. McGuinness, X. Giró-i Nieto, and N. E. O'Connor, "Saliency weighted convolutional features for instance search," in *CBMI*, 2018, pp. 1–6.
- [79] F. Yang, J. Li, S. Wei, Q. Zheng, T. Liu, and Y. Zhao, "Two-stream attentive CNNs for image retrieval," in *ACM MM*, 2017, pp. 1513–1521.
- [80] H.-F. Yang, K. Lin, and C.-S. Chen, "Supervised learning of semantics-preserving hash via deep convolutional neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 2, pp. 437–451, 2018.
- [81] Y. Liu, J. Song, K. Zhou, L. Yan, L. Liu, F. Zou, and L. Shao, "Deep self-taught hashing for image retrieval," *IEEE Trans. Cybern.*, vol. 49, no. 6, pp. 2229–2241, 2018.
- [82] A. Gordo, J. Almazan, J. Revaud, and D. Larlus, "End-to-end learning of deep visual representations for image retrieval," *Int. J. Comput. Vis.*, vol. 124, no. 2, pp. 237–254, 2017.
- [83] H. Azizpour, A. S. Razavian, J. Sullivan, A. Maki, and S. Carlsson, "Factors of transferability for a generic convnet representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 9, pp. 1790–1802, 2016.
- [84] E. Mohamedano, K. McGuinness, N. E. O'Connor, A. Salvador, F. Marqués, and X. Giro-i Nieto, "Bags of local convolutional features for scalable instance search," in *ICMR*, 2016, pp. 327–331.
- [85] W. Zhao, H. Luo, J. Peng, and J. Fan, "Spatial pyramid deep hashing for large-scale image retrieval," *Neurocomputing*, vol. 243, pp. 166–173, 2017.
- [86] L. Zheng, S. Wang, J. Wang, and Q. Tian, "Accurate image search with multi-scale contextual evidences," *Int. J. Comput. Vis.*, vol. 120, no. 1, pp. 1–13, 2016.
- [87] J. Cao, L. Liu, P. Wang, Z. Huang, C. Shen, and H. T. Shen, "Where to focus: Query adaptive matching for instance retrieval using convolutional feature maps," *arXiv preprint arXiv:1606.06811*, 2016.
- [88] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *ECCV*, 2014, pp. 391–405.
- [89] T. Yu, Y. Wu, S. D. Bhattacharjee, and J. Yuan, "Efficient object instance search using fuzzy objects matching," in *AAAI*, 2017, pp. 4320–4326.
- [90] S. Sun, W. Zhou, Q. Tian, and H. Li, "Scalable object retrieval with compact image representation from generic object regions," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 12, no. 2, pp. 1–21, 2015.
- [91] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid, "Convolutional kernel networks," in *NeurIPS*, 2014, pp. 2627–2635.
- [92] A. Salvador, X. Giró-i Nieto, F. Marqués, and S. Satoh, "Faster R-CNN features for instance search," in *CVPR Workshops*, 2016, pp. 9–16.
- [93] S. Wang and S. Jiang, "INSTRE: a new benchmark for instance-level object retrieval and recognition," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 11, no. 3, pp. 37:1–37:21, 2015.
- [94] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *CVPR*, 2007, pp. 1–8.
- [95] —, "Lost in quantization: Improving particular object retrieval in large scale image databases," in *CVPR*, 2008, pp. 1–8.
- [96] T. Ng, V. Balntas, Y. Tian, and K. Mikolajczyk, "SOLAR: Second-order loss and attention for image retrieval," in *ECCV*, 2020, pp. 253–270.
- [97] L. Zheng, Y. Zhao, S. Wang, J. Wang, and Q. Tian, "Good practice in CNN feature transfer," *CoRR*, vol. abs/1604.00133, 2016.
- [98] Y. Lou, Y. Bai, S. Wang, and L.-Y. Duan, "Multi-scale context attention network for image retrieval," in *ACM MM*, 2018, pp. 1128–1136.
- [99] Y. Li, Y. Xu, J. Wang, Z. Miao, and Y. Zhang, "MS-RMAC: Multiscale regional maximum activation of convolutions for image retrieval," *IEEE Signal Process. Lett.*, vol. 24, no. 5, pp. 609–613, 2017.
- [100] X. Xiang, Z. Wang, Z. Zhao, and F. Su, "Multiple saliency and channel sensitivity network for aggregated convolutional feature," in *AAAI*, vol. 33, no. 01, 2019, pp. 9013–9020.
- [101] 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 Trans. Multimedia*, vol. 21, no. 6, pp. 1513–1523, 2018.
- [102] W. Yu, K. Yang, H. Yao, X. Sun, and P. Xu, "Exploiting the complementary strengths of multi-layer CNN features for image retrieval," *Neurocomputing*, vol. 237, pp. 235–241, 2017.
- [103] M. Yang, D. He, M. Fan, B. Shi, X. Xue, F. Li, E. Ding, and J. Huang, "Dolg: Single-stage image retrieval with deep orthogonal fusion of local and global features," in *ICCV*, 2021, pp. 11 772–11 781.
- [104] Z. Zhang, Y. Xie, W. Zhang, and Q. Tian, "Effective image retrieval via multilinear multi-index fusion," *IEEE Trans. Multimedia*, vol. 21, no. 11, pp. 2878–2890, 2019.
- [105] Q. Wang, J. Lai, Z. Yang, K. Xu, P. Kan, W. Liu, and L. Lei, "Improving cross-dimensional weighting pooling with multi-scale feature fusion for image retrieval," *Neurocomputing*, vol. 363, pp. 17–26, 2019.
- [106] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, and Q. Tian, "Query-adaptive late fusion for image search and person re-identification," in *CVPR*, 2015, pp. 1741–1750.
- [107] H. Xuan, R. Souvenir, and R. Pless, "Deep randomized ensembles for metric learning," in *ECCV*, 2018, pp. 723–734.
- [108] B.-C. Chen, L. S. Davis, and S.-N. Lim, "An analysis of object embeddings for image retrieval," *arXiv preprint arXiv:1905.11903*.
- [109] F. Wang, W.-L. Zhao, C.-W. Ngo, and B. Merialdo, "A hamming embedding kernel with informative bag-of-visual words for video semantic indexing," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 10, no. 3, pp. 1–20, 2014.
- [110] A. Mishchuk, D. Mishkin, F. Radenović, and J. Matas, "Working hard to know your neighbor's margins: Local descriptor learning loss," in *NeurIPS*, 2017, pp. 4827–4838.
- [111] A. Mukherjee, J. Sil, A. Sahu, and A. S. Chowdhury, "A bag of constrained informative deep visual words for image retrieval," *Pattern Recognition Letters*, vol. 129, pp. 158–165, 2020.
- [112] M. Paulin, M. Douze, Z. Harchaoui, J. Mairal, F. Perronin, and C. Schmid, "Local convolutional features with unsupervised training for image retrieval," in *ICCV*, 2015, pp. 91–99.
- [113] H. Jegou, F. Perronnin, M. Douze, J. Sánchez, P. Perez, and C. Schmid, "Aggregating local image descriptors into compact codes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 9, pp. 1704–1716, 2012.
- [114] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," *Int. J. Comput. Vis.*, vol. 105, no. 3, pp. 222–245, 2013.
- [115] J. Cao, Z. Huang, and H. T. Shen, "Local deep descriptors in bag-of-words for image retrieval," in *ACM MM*, 2017, pp. 52–58.
- [116] M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, "D2-net: A trainable cnn for joint description and detection of local features," in *CVPR*, 2019, pp. 8092–8101.
- [117] K. Mikolajczyk and C. Schmid, "Scale & affine invariant interest point detectors," *Int. J. Comput. Vis.*, vol. 60, no. 1, pp. 63–86, 2004.
- [118] G. Tolias, T. Jenicek, and O. Chum, "Learning and aggregating deep local descriptors for instance-level recognition," in *ECCV*, 2020.
- [119] J. Yang, J. Liang, H. Shen, K. Wang, P. L. Rosin, and M.-H. Yang, "Dynamic match kernel with deep convolutional features for image retrieval," *IEEE Trans. Image Process.*, vol. 27, no. 11, pp. 5288–5302, 2018.
- [120] J. Kim and S.-E. Yoon, "Regional attention based deep feature for image retrieval," in *BMVC*, 2018, pp. 209–223.
- [121] S. Wei, L. Liao, J. Li, Q. Zheng, F. Yang, and Y. Zhao, "Saliency inside: Learning attentive CNNs for content-based image retrieval," *IEEE Trans. Image Process.*, vol. 28, no. 9, pp. 4580–4593, 2019.
- [122] T.-T. Do, D.-K. Le Tan, T. T. Pham, and N.-M. Cheung, "Simultaneous feature aggregating and hashing for large-scale image search," in *CVPR*, 2017, pp. 6618–6627.
- [123] K. Musgrave, S. Belongie, and S.-N. Lim, "A metric learning reality check," in *ECCV*, 2020, pp. 681–699.
- [124] Y. Lv, W. Zhou, Q. Tian, S. Sun, and H. Li, "Retrieval oriented deep feature learning with complementary supervision mining," *IEEE Trans. Image Process.*, vol. 27, no. 10, pp. 4945–4957, 2018.
- [125] W. Min, S. Mei, Z. Li, and S. Jiang, "A two-stage triplet network training framework for image retrieval," *IEEE Trans. Multimedia*, vol. 22, no. 12, pp. 3128–3138, 2020.
- [126] J. Deng, J. Guo, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," *CVPR*, pp. 4685–4694, 2019.
- [127] M. Boudiaf, J. Rony, I. M. Ziko, E. Granger, M. Pedersoli, P. Piantanida, and I. B. Ayed, "A unifying mutual information view of metric learning: cross-entropy vs. pairwise losses," in *ECCV*, 2020, pp. 548–564.
- [128] J. Lin, O. Moree, A. Veillard, L.-Y. Duan, H. Goh, and V. Chandrasekhar, "Deephash for image instance retrieval: Getting regularization, depth and fine-tuning right," in *ICMR*, 2017, pp. 133–141.
- [129] J. Cao, Z. Huang, P. Wang, C. Li, X. Sun, and H. T. Shen, "Quartet-net learning for visual instance retrieval," in *ACM MM*, 2016, pp. 456–460.
- [130] A. Brown, W. Xie, V. Kalogeiton, and A. Zisserman, "Smooth-AP: Smoothing the path towards large-scale image retrieval," in *ECCV*, 2020, pp. 677–694.
- [131] A. Iscen, G. Tolias, Y. Avrithis, and O. Chum, "Mining on manifolds: Metric learning without labels," in *CVPR*, 2018, pp. 7642–7651.
- [132] A. Iscen, Y. Avrithis, G. Tolias, T. Furon, and O. Chum, "Fast spectral ranking for similarity search," in *CVPR*, 2018, pp. 7632–7641.
- [133] A. Iscen, G. Tolias, Y. Avrithis, T. Furon, and O. Chum, "Efficient diffusion on region manifolds: Recovering small objects with compact CNN representations," in *CVPR*, 2017, pp. 2077–2086.

- [134] M. Donoser and H. Bischof, "Diffusion processes for retrieval revisited," in *CVPR*, 2013, pp. 1320–1327.
- [135] Y. Zhao, L. Wang, L. Zhou, Y. Shi, and Y. Gao, "Modelling diffusion process by deep neural networks for image retrieval," in *BMVC*, 2018.
- [136] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.
- [137] T. Maria and T. Anastasios, "Deep convolutional image retrieval: A general framework," *Signal Process. Image Commun.*, vol. 63, pp. 30–43, 2018.
- [138] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, 2018, pp. 132–149.
- [139] W. Jiang, Y. Wu, C. Jing, T. Yu, and Y. Jia, "Unsupervised deep quantization for object instance search," *Neurocomputing*, vol. 362, pp. 60–71, 2019.
- [140] B. Ke, J. Shao, Z. Huang, and H. T. Shen, "Feature reconstruction by laplacian eigenmaps for efficient instance search," in *ICMR*, 2018.
- [141] T. Shen, Z. Luo, L. Zhou, R. Zhang, S. Zhu, T. Fang, and L. Quan, "Matchable image retrieval by learning from surface reconstruction," in *ACCV*, 2018, pp. 415–431.
- [142] J. Hu, R. Ji, H. Liu, S. Zhang, C. Deng, and Q. Tian, "Towards visual feature translation," in *CVPR*, 2019, pp. 3004–3013.
- [143] C. Bai, H. Li, J. Zhang, L. Huang, and L. Zhang, "Unsupervised adversarial instance-level image retrieval," *IEEE Trans. Multimedia*, 2021.
- [144] M. Paulin, J. Mairal, M. Douze, Z. Harchaoui, F. Perronnin, and C. Schmid, "Convolutional patch representations for image retrieval: an unsupervised approach," *Int. J. Comput. Vis.*, vol. 121, no. 1, pp. 149–168, 2017.
- [145] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *CVPR*, 2006, pp. 2161–2168.
- [146] F. Radenovic, A. Iscen, G. Toliás, Y. Avrithis, and O. Chum, "Revisiting oxford and paris: Large-scale image retrieval benchmarking," in *CVPR*, 2018.
- [147] S. Yokoo, K. Ozaki, E. Simo-Serra, and S. Iizuka, "Two-stage discriminative re-ranking for large-scale landmark retrieval," in *CVPR Workshops*, 2020, pp. 1012–1013.
- [148] A. Alzu'bi, A. Amira, and N. Ramzan, "Content-based image retrieval with compact deep convolutional features," *Neurocomputing*, vol. 249, pp. 95–105, 2017.
- [149] S. S. Husain and M. Bober, "REMAP: Multi-layer entropy-guided pooling of dense CNN features for image retrieval," *IEEE Trans. Image Process.*, vol. 28, no. 10, pp. 5201–5213, 2019.
- [150] L. T. Alemu and M. Pelillo, "Multi-feature fusion for image retrieval using constrained dominant sets," *Image Vis Comput*, vol. 94, p. 103862, 2020.
- [151] L. P. Valem and D. C. G. Pedronette, "Graph-based selective rank fusion for unsupervised image retrieval," *Pattern Recognit Lett*, 2020.
- [152] F. Yang, R. Hinami, Y. Matsui, S. Ly, and S. Satoh, "Efficient image retrieval via decoupling diffusion into online and offline processing," in *AAAI*, vol. 33, 2019, pp. 9087–9094.
- [153] H.-F. Yang, K. Lin, and C.-S. Chen, "Cross-batch reference learning for deep classification and retrieval," in *ACM MM*, 2016, pp. 1237–1246.
- [154] J. Ouyang, W. Zhou, M. Wang, Q. Tian, and H. Li, "Collaborative image relevance learning for visual re-ranking," *IEEE Trans. Multimedia*, 2020.
- [155] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," in *ICLR*, 2014.
- [156] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," in *BMVC*, 2014.
- [157] J. Li, R. Ji, H. Liu, X. Hong, Y. Gao, and Q. Tian, "Universal perturbation attack against image retrieval," in *ICCV*, 2019, pp. 4899–4908.
- [158] G. Toliás, F. Radenovic, and O. Chum, "Targeted mismatch adversarial attack: Query with a flower to retrieve the tower," in *ICCV*, 2019.
- [159] E. Triantafillou, R. Zemel, and R. Urtasun, "Few-shot learning through an information retrieval lens," *NeurIPS*, vol. 30, 2017.
- [160] Z. Deng, Y. Zhong, S. Guo, and W. Huang, "InsCLR: Improving instance retrieval with self-supervision," in *AAAI*, 2022.
- [161] S. Ibrahim, A. Sors, R. S. de Rezende, and S. Clinchant, "Learning with label noise for image retrieval by selecting interactions," in *WACV*, 2022, pp. 2181–2190.

Wei Chen received the Ph.D. degree at Leiden University, in 2021. Before starting with PhD study in Leiden University, he received his Master degree from the National University of Defense Technology, China, in 2016. His research interest focuses on cross-modal retrieval with deep learning methods, and also in the context of incremental learning. He has published papers in international conferences and journal including *CVPR*, *ACM MM*, *PR*, *Neurocomputing*, and *IEEE TMM* etc.



and received the best paper award at MMM2017.



Weiping Wang received the Ph.D. degree in systems engineering from the National University of Defense Technology, Changsha, China. He is a Professor at Academy of Advanced Technology Research of Hunan. He has more than 200 papers published on journals and conferences including *IEEE Transactions multimedia*, *Transactions on Vehicular Technology*, etc. His current research interests focus on intelligent decision experimentation, multi-modal knowledge extraction and knowledge integrated computation.



Erwin M. Bakker is co-director of the LIACS Media Lab at Leiden University. He has published widely in the fields of image retrieval, audio analysis and retrieval and bioinformatics. He was closely involved with the start of the International Conference on Image and Video Retrieval (CIVR). Moreover, he regularly serves as a program committee member or organizing committee member for scientific multimedia and human-computer interaction conferences and workshops.



Theodoros Georgiou received the Ph.D. degree at Leiden University, in 2021. His research interest focuses on deep learning methods applied on higher than two dimensional data. Before starting with PhD study in Leiden University, he received his Master degree from the Leiden University in 2016. He has published papers in international conferences and journal including *ICPR*, *CBMI*, *WCCI* and *IJMIR*.



processing, hierarchical algorithms, data fusion, machine learning, and the interdisciplinary applications of such methods. He has published textbooks on *Statistical Image Processing* and on *Complex Systems theory*.

Li Liu received the Ph.D. degree in information and communication engineering from the National University of Defense Technology, China, in 2012. She is currently a Professor with the College of System Engineering. She has held visiting appointments at the University of Waterloo, Canada, at the Chinese University of Hong Kong, and at the University of Oulu, Finland. Her current research interests include computer vision, pattern recognition and machine learning. Her papers have currently over 7800 citations in Google Scholar.



Michael S. Lew is the head of the Deep Learning and Computer Vision Research Group and a full Professor at LIACS, Leiden University. He has published over a dozen books and 190 peer-reviewed scientific articles in the areas of image retrieval, computer vision, and deep learning. Notably, he had the most cited paper in the *ACM Transactions on Multimedia* and one of the top 10 most cited articles in the history (out of more than 16,000 articles) of the *ACM SIGMM*. He was also a founding member of the advisory committee for the TRECVID video retrieval evaluation project, chair of the steering committee for the *ACM International Conference on Multimedia Retrieval* and a member of the *ACM SIGMM Executive Committee*.



and received the best paper award at MMM2017.

