1	Graph-based Discriminative Features Learning for Fine-grained Image
2	Retrieval
3	Han SUN <sup>1,2*</sup> , Wenxi Lang <sup>1,2</sup> , Can Xu <sup>3</sup> , Ningzhong Liu <sup>1,2</sup> , Huiyu Zhou <sup>4</sup>
4	1 College of Computer Science and Technology, Nanjing University of Aeronautics
5	and Astronautics, Nanjing 211106, China
6	2 MIIT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing
7	211106, China
8	3 College of Computer Science and Engineering, Nanjing University of Science and
9	Technology, Nanjing 210094, China
10	4 School of Informatics, University of Leicester, Leicester LE1 7RH, U.K.
11	Abstract: Fine-grained image retrieval has gradually become a hot topic in computer
12	vision, which aims to retrieve images with the same subcategories from general visual
13	categories. Though fine-grained image retrieval has made a breakthrough with the
14	development of convolutional neural networks, its performance is still limited by the
15	low discriminative feature embedding. To solve this problem, most prior works focus
16	on mining more discriminative features with various strategies. In this paper, we
17	propose a novel graph-based discriminative features learning network for fine-grained
18	image retrieval (GDF-Net). We first design a global fine-grained feature aggregation
19	module, which reconstructs the discriminative features through capturing context
20	correlation based on a K-Nearest Neighbor graph. To reduce storage overhead and
21	speed up retrieval, we further design a semantic hash encoding module, which generates
22	a semantically compact hash code under the guidance of Cauchy quantization loss and

23	bit balance loss. Validated by extensive experiments and ablation studies, our method
24	consistently outperforms state-of-the-art generic retrieval methods as well as fine-
25	grained retrieval methods on three datasets, e.g., CUB Birds, Stanford Dogs and
26	Stanford Cars.
27	
28	Keywords: Fine-grained image retrieval; graph convolutional neural network; deep
29	hashing
30	

31 **1 Introduction** 

As an emerging research topic, Fine-grained image retrieval (FGIR) aims to 32 retrieve images belonging to a certain meta-category (e.g., dogs) and return images with 33 34 the same subcategory (e.g., walker hound or basset) as a query image. Different from 35 the classical image retrieval, for the FGIR task, images that belong to the different subcategories are similar to each other as the discriminative information among them 36 37 is quite imperceptible (e.g., the texture of a dog's head usually determines whether it's walker hound or basset), while images within the same subcategory are also difficult to 38 be distinguished due to the various pose, illumination, background, and shooting angle. 39 To this end, distinguishing and retrieving fine-grained images is still a challenging 40 research focus at present. 41

As a cornerstone, learning discriminative feature representation within and outside
fine-grained categories plays a key role in FGIR. Earlier works[1, 2] mainly rely on the
hand-craft features, while more recent studies[3-7] prefer to extract features from a pre-

trained Convolutional Neural Network (CNN). Although CNN's powerful semantic 45 expression capabilities have significantly improved the feature quality, constructing a 46 47 high-quality feature representation in FGIR is still an open problem, whose main challenge lies in efficiently mining and embedding discriminative local information. 48 Most recent works[8-12] usually design specific localization strategies to perceive 49 object's discriminative parts first. For instance, Selective Convolutional Descriptor 50 Aggregation (SCDA) [13] first localizes the salient foreground region and irrelevant 51 background noise with a pre-trained CNN, and then obtained a more representative 52 53 feature embedding by introducing the flood-fill algorithm to filter the noise. ExchNet[14] and WSDAN[15] obtains parts of fine-grained objects based on the 54 attention mechanism, which further designs channel-wise constraint and data 55 56 enhancement to enhance the discriminativeness of parts, respectively. It seems that all the previous works try to locate and encode objects' discriminative parts and patches 57 independently. However, what has been neglected is the global structure relationship 58 59 among these local discriminative regions.

Existing works[16-18] have shown that the global structure benefits to improve the performance of image classification/recognition/segmentation tasks. More intuitively, as the Fig 1 shown, the local discriminative information (e.g., blue back, white abdomen, white and blue striped wings) can't support the model to distinguish and retrieve the Florida Jay from database accurately. However, if we further consider the size of the white area on the wings or the relative spatial relationship between the blue back and white area on the wings by manual search, the result will be significantly

improved. In theory, we consider that local spatial context and region correlations is 67 rather helpful in distinguishing fine-grained images. On the one hand, by mining the 68 69 context of the discriminative regions, more detailed information of the object can be encoded, which alleviates the local embedding issue caused by receptive field of 70 71 convolutional neural network; on the other hand, by dynamically formulating the correlation among different independent discriminative regions, the final feature can 72 not only retain the discriminativeness but also describe object's global structural 73 74 relationship. All of the above issues motivate us to incorporate the local spatial context 75 of independent discriminative regions and the correlation among these independent discriminative regions into the fine-grained image feature extracting. 76



Fig. 1 The top 5 images returned. (a)based on our reconstructed feature through capturing the global spatial context

correlation; (b) based on the independent local discriminative feature

77

```
Based on the above considerations, in this paper, we propose a graph-based
discriminative features learning network for fine-grained image retrieval (GDF-Net),
which consists of the Global Fine-grained Feature Aggregation Module (GFFAM) and
Semantic Hash Encoding Module (SHECM). The GFFAM first learns attention maps
```

that indicating objects' discriminative parts. For each part, the GFFAM regards it as the 82 node and constructs a graph convolutional network. By adopting a propagation strategy 83 84 inspired by the K-Nearest Neighbor, the GFFAM finally learns the dynamic connection relationship among different nodes to aggregate the local discriminative features into 85 global fine-grained features. In order to improve the retrieval speed, we also design a 86 SHECM, which restricts the network to learn compact and high semantic hash codes 87 based on the global fine-grained features. While ensuring the semantics of hash 88 encoding based on classification loss, the SHECM further introduces Cauchy cross-89 90 entropy loss and bit balance loss, which can improve the compactness of the hash code and realize efficient Hamming space retrieval. In summary, the contributions of this 91 92 paper are as follows:

93 (1) We propose an adaptive local feature fusion method based on a graph neural
94 network. In the dissemination of the graph, the correlation among the discriminative
95 part features can be captured effectively, which helps guide the fusion of different fine96 grained part features to improve the final feature quality.

97 (2) We propose an extra semantic hash encoding module in the fine-grained image
98 retrieval task, which outputs a compact hash code to improve the retrieval performance
99 by combining the classification loss with Cauchy cross-entropy loss and bit balance
100 loss.

(3) We evaluate the proposed model on three fine-grained benchmarks. Extensive
experiments indicate our method achieves the best performance in retrieval accuracy.

## 103 **2 Relate work**

## 104 **2.1 Fine-grained image retrieval**

Fine-grained image retrieval is an emerging research hotspot that has attracted more and more research attention in recent years. Unlike general image retrieval tasks, fine-grained image retrieval aims to distinguish subcategory, which has two characteristics between database images and query images: (1) small inter-class variance. Differences among most classes are subtle and difficult to distinguish. (2) large intra-class variance, which is due to factors such as poses and viewpoints.

In the early days, FGIR relied on the use of manual features[19], and as deep learning method evolved, more and more FGIR methods were proposed. And these deep methods can be roughly divided into supervised and unsupervised methods.

For unsupervised retrieval methods, SCDA[13] proposes selective convolutional descriptor aggregation method, which first localizes objects in fine-grained images and retains useful deep descriptors for fine-grained image retrieval.

117 In supervised methods, CRL-WSL[20] presents a unified framework for efficient learning of discriminative features with centralized ranking loss and segmentation of 118 119 target contours using learning saliency regions. DCL-NC[19] improves CRL-WSL[20] 120 by adding a normalize-scale layer and decorrelated ranking loss. However, these two methods encode high dimensionality and may encounter problems of slow query speed 121 and storage redundancy in practical large-scale image retrieval. ExchNet[14] uses hash 122 123 learning to provide compact binary codes for fine-grained images, and designs the operation of feature swapping to make local features aligned. [21] adds a hash layer 124

before the classification layer as a feature-switching layer to guide the classificationand achieves state-of-the-art results.

## 127 **2.2** Graph convolutional networks in computer vision

Graph convolutional network(GCN)[22], a neural network that applies graph structures to learn the representation of a target node by iteratively propagating neighborhood information, has been widely used in various computer vision tasks[23-131 26].

In the field of object detection, Structural-RNN[23] designs a dual graph neural 132 133 network that combines temporal and spatial relationships among entities to construct graphs that perform graph convolution from different views of the raw data. In the field 134 of image classification, Lee et al.[24] introduces knowledge graphs to zero-sample 135 136 learning tasks and achieves some improvement in multi-label tasks. In the field of semantic segmentation, Graph Long Short Time Memory(Graph LSTM)[25] network 137 proposes a generalization of LSTM from sequential or multidimensional data to general 138 139 graph-structured data. Wang at el.[27] uses graph convolutional neural networks on point clouds and proposes to collect features of edges by edge convolution, which 140 contains neighborhood information, and global geometric features can be learned by 141 stacking. 142

For the task of FGIR, to the best of our knowledge, few works consider introducing the graph convolutional neural networks, while it used to not only express discriminative features but also capture the complex spatial correlations among different local features in our GFFAM.

## 147 **2.3 Hashing method**

Current hashing methods can be classified into two categories, containing dataindependent hashing and data-dependent hashing. Data-independent hashing refers to binary hash codes constructed by random projection or by hand, such as locally sensitive hash functions[28], and LSH methods usually require a long code length to guarantee retrieval performance.

Based on whether supervised information is used or not, data-dependent hashing can be classified as unsupervised and supervised hashing. For example, [29] and [30] are classical unsupervised hashing methods. Although unsupervised methods are less restrictive and more practical, exploiting the available supervised information implies better performance.

158 Thanks to the powerful representation capability of deep neural networks, deep models have been used for supervised hashing, proposing the synthesis of feature 159 learning and hash coding into end-to-end deep hashing methods[31-37]. Convolutional 160 161 neural network hashing (CNNH)[31] and deep pairwise supervised hashing (DPSH)[34] are representative approaches. HashNet[33] proposes to continuously approximate the 162 sign function with the tanh activation function, thus solving the optimization 163 complexity problem and greatly improving the retrieval performance. However, 164 existing hashing methods are based on a linear scan of the hash to maximize the 165 retrieval performance, thus incurring expensive overhead. Therefore, Deep Cauchy 166 Hashing (DCH) [34] proposes a pairwise cross-entropy loss based on Cauchy 167 distribution, which learns almost lossless hash codes by significantly penalizing similar 168

image pairs with Hamming distance greater than a given Hamming radius threshold.
Other recent work includes greedy hash [38], Asymmetric Deep Supervised Hashing
(ADSH) [39], etc.

# 172 **3 Proposed Method**

## 173 **3.1 Overview**

We propose an end-to-end Graph-propagation based discriminative features 174 learning network for Fine-grained image retrieval. As shown in Fig 2, it consists of two 175 key modules: the Global Fine-grained Feature Aggregation Module (GFFAM) and 176 Semantic Hash Encoding Module (SHECM). The GFFAM explores the 177 interdependencies among feature vectors based on a graph convolutional network, 178 which will guide the fusion of independent discriminative elements to enhance the 179 180 global fine-grained features (Section 3.2). The global features are combined in the SHECM to learn a hash embedding for the final retrieval, which we will explain in 181 Section 3.3. 182



Fig.2 Framework of the proposed GDF-Net

### **3.2 Global Fine-grained Feature Aggregation Module**

In order to ensure the discriminative ability of feature while taking both global 185 image-level context and interaction among local features into account, we construct a 186 graph convolution network to aggregate independent local features into global fine-187 grained features iteratively. The proposed GFFAM includes the following two stages. 188 The first stage is that extracts local discriminative features which server as the node of 189 the graph convolution network. In the second stage, the graph updates the 190 representation of each node and propagates them to the top level for aggregation. It 191 192 should be noted that, by calculating the link relationship among different nodes in the graph, the global context information and the interaction among local discriminative 193 features are finally encoded together as the node features propagate along the 194 195 connection direction.

## **3.2.1 Node representation based on local feature**

The discriminative local features are quite important for fine-grained tasks. In 197 198 constructing k-nearest neighbor graph, local features of fine-grained object are used as nodes of the graph. Therefore, we first extract the fine-grained parts with discriminative 199 features. Specifically, in our mission with a weakly supervised image-level label, the 200 location annotation (e.g., bounding boxes or key points) which can indicate instances' 201 important parts is not available during the model training and testing, so we extract the 202 discriminative local features based on attention generation strategy by calculating the 203 204 category response of different parts.

Given the image X and its feature  $F \in R^{H \times W \times N}$ , where H, W and N 205 represent the height, width and number of channels of the feature layers, respectively. 206 We generate Mattention maps for each image and the  $A_k \in \mathbb{R}^{H \times W}$  means the attention 207 map of part k, which may be the wing of a bird or the head of a dog. We use a simple 208 convolution function f() to transform the feature map to attention map, which can be 209 described as a spatial attention mapping indicating the response between original 210 channel-wise features and category labels. In detail, the structure of f() consists of a 211 convolution layer with kernel  $1 \times 1$ , a Batch Normalization layer and a Relu activation 212 213 layer. The calculation method of the attention map is shown in formula (1).

214 
$$A = f(F) = \bigcup_{k=1}^{M} A_k$$
 (1)

With the attention map of different local regions, we then extract the discriminative local features from these parts. Similar to the Bilinear pooling[8, 44], the *M* part features corresponding to these local regions can be calculated by formula (2).

218 
$$f_k = g(A_k \odot F), \ k = 1, 2, \cdots, M$$
 (2)

219 Where  $f_k$  represents the  $k_{th}$  local feature and the  $\odot$  represents the element-wise 220 multiplication of feature map F and the  $k_{th}$  attention map, while g() is the global 221 average pooling operation.

222 **3.2.2 Local feature enhancement** 

For the learned M spatial attention maps, a common phenomenon is that multiple attentions maps may focus on object's similar regions, which will greatly inhibit the diversity of the discriminative features. We further introduce an additional random dropping strategy to force the network to search for other informative local regions. Specifically, for each training image, we first randomly select one attention map  $A_k$ from A. To improve the convergence rate of the model, the min-max normalization is adopted to smooth the value of  $A_k$  to the range of [0, 1], as shown in equation (3).

230 
$$A_k^* = \frac{A_k - \min(A_k)}{\max(A_k) - \min(A_k)}$$
(3)

Here, the  $A_k^*$  represents the  $k_{th}$  attention map after augmentation. We then construct a drop mask  $M_d$  through setting the value of elements larger than the threshold of  $T_d \in [0, 1]$  to 0 and the value of other elements to 1, as:

234 
$$M_d(i,j) = \begin{cases} 0, \ A_k^*(i,j) > T_d \\ 1, \ otherwise \end{cases}$$
(4)

where the  $A_k^{*}(i,j)$  represents the value of element in the  $i_{th}$  row and  $j_{th}$  column 235 belonging to the  $k_{th}$  local feature selected, and the  $M_d(i,j)$  represents the value of 236 drop mask in the corresponding position. Here, the threshold  $T_d$  is set to 0.5. With the 237 drop mask  $M_d$  and the original image, the new masking image  $X_d$  can be obtained 238 by an element-wised multiplication, which will be fed into the network once again and 239 learn M new part features  $f_k^{\ d}$ . With the attention drop, we encourage the attention 240 map to propose other discriminative parts and finally improve the accuracy of 241 localization and the quality of feature. 242



244

Fig.3 Example of k-nearest neighbor graph (green edge indicates the selected neighbor node).

## 245 **3.2.3 Interdependency mining based on Graph-Propagation**

With M local discriminative features  $f_1, f_2, ..., f_M$ , we construct a directed graph 246 G = (V, E) to capture the contextual relationship among these discrete local features, 247 which process can be described as calculating *k*-nearest neighbor of each graph node. 248 Since the graph contains self-loops where each node also points to itself, the feature of 249 each graph node can be updated by aggregating the features of its k-nearest neighbor 250 node. Under the supervision of category label, we finally obtain M reconstructed 251 252 features with the propagation of graph, which can express the contextual structural relationship among different local features while retaining high semantic discrimination. 253 Specially, given vertices  $V = \{1, 2, ..., M\}$  and edges  $E \subseteq V \times V$ , the edge 254

254 Specially, given vertices  $V = \{1, 2, ..., M\}$  and edges  $E \subseteq V \times V$ , the edge 255 feature from  $i_{th}$  vertex to  $j_{th}$  vertex can be defined as:

$$e_{ij} = h_{\theta;\sigma}(f_i, f_i - f_j) \tag{5}$$

257  $h_{\theta;\sigma}$  () is an asymmetric edge function implemented through a shared Multilayer 258 Perceptron (MLP), which consists of a convolution layer with 1 × 1 kernel, a Batch 259 Normalization layer and a Relu layer, as:

260

$$h_{\theta;\sigma}(f_i, f_i - f_j) = Relu(\theta \cdot f_i + \sigma(f_i - f_j))$$
(6)

Where  $\theta$  and  $\sigma$  represent parameters of the network. As shown in equation (6), the neighborhood information from different nodes captured by  $f_i - f_j$  can be combined with the global information captured by  $f_i$  gradually. Similar to image convolution, the output of  $i_{th}$  vertex can be obtained by applying a global max pooling (GMP) operation on all edge features associated with the  $i_{th}$  vertex, as:

$$f_i^r = \underset{j:(i,j)\in E}{GMP}e_{ij} \tag{7}$$

## 267 **3.3 Semantic Hash Encoding Module**

In the most of deep hashing methods, the hash layer is designed to encode the features and output a binary hash code, in which '1' indicates the category possesses a certain feature while '0' means it lacks this feature. Once we obtain the high discriminative features which can express the fine-grained categories, what we need to do is just encoding these features into compact hash coding to improve the retrieval efficiency. As shown in Fig.2, different from existing works[45], we design and add a semantic hash coding module before the classification layer.

For the *M* reconstructed context features  $f_i^r$ , the SHECM outputs a *B* bit hash code which can be computed as the formula (8).

277 
$$H_i = tanh((W^H)^T f_i^r + \delta^H), \ i = 1, 2, \cdots, M$$
(8)

278 Where  $H_i \in \mathbb{R}^B$  is the output of the hash layer with the  $f_i^r$ . The  $\delta^H \in B$  and  $W^H \in \mathbb{R}^{M \times B}$  represents the bias and the weights of the hash layer, respectively. The *tanh()* 280 represents the active function which can be described as the formula (9).

281 
$$\tanh(f^r) = \frac{e^{x} - e^{-x}}{e^{x} + e^{-x}}$$
 (9)

The final hash codes can be obtained according to the formula (10). Since the value range of the tanh() is [-1, 1], the  $B_i = 1$  if the  $H_i \ge 0$  and  $B_i = -1$  otherwise. Since the gradient of the sign function at the non-zero point may be zero, the problem of gradient vanishing arises, we only map the real-valued output  $H_i$  to hash code when test the model.

$$B_i = sgn(H_i) \tag{10}$$

288 **3.4 Loss function** 

## 289 **3.4.1 Loss of GFFAM**

When learning the discriminative features in section 3.2.1, it's necessary that each attention map  $A_k$  belonging to same category can point to objects' similar part region. Inspired by the center loss of face recognition[37, 46], we first introduce a loss  $L_{ctr}$  to learn a feature center  $c_k$  for each local discriminative feature  $f_k$ .  $L_{ctr}$  penalizes the variances of features which comes from same parts of different objects with the same category label, which can be formulated by the equation (11).

296 
$$L_{ctr} = \sum_{k=1}^{M+1} ||f_k - c_k||^2$$
(11)

Where  $c_k$  is the part's feature center of  $A_k$  and can be initialized from zero and updated by moving average  $c_k = (1 - \mu)c_k + \mu f_k$ . Here,  $\mu$  controls the update rate of  $c_k$  and  $L_{ctr}$  loss applies only to the original image.

Since the whole learning process of GFFAM is supervised by the category label, we adopt the cross-entropy loss  $L_{ce}$  to constraint the distance among the predicted category and the real image-label  $Y^*$ . For M original local feature  $f_k$ , we simply stack them together and then feed in a SoftMax layer to predict its category probability  $Y_{ori}$ . The loss can be calculated as:

305 
$$L_{ori}(Y_{ori}, Y^*) = L_{ce}(Softmax(\begin{bmatrix} f_1\\f_2\\\vdots\\f_m \end{bmatrix}), Y^*) \qquad (12)$$

Similarly, we also predict the category probability  $Y_{drop}$  and  $Y_{recn}$  for the M dropping feature  $f_k^{\ d}$  and M reconstructed context feature  $f_k^{\ r}$ , respectively. The total classification loss composes of the  $L_{ori}$ ,  $L_{drop}$  and  $L_{recn}$  as the formula (13) shown.

310 
$$L_{cls} = L_{ori}(Y_{ori}, Y^*) + L_{drop}(Y_{drop}, Y^*) + L_{recn}(Y_{recn}, Y^*)$$
(13)

#### 311 **3.4.2 Loss of SHECM**

Previous deep hashing methods have used the Sigmoid function to define the probability function, however, existing hashing methods usually lack the ability to concentrate the relevant images within a small Hamming sphere, so they may perform poorly for Hamming space retrieval. Therefore, inspired by DCH[35], we use a Bayesian framework to optimize the quantization loss. We use the probability function, as:

318 
$$\sigma\left(d(\boldsymbol{h}_{i},\boldsymbol{h}_{j})\right) = \frac{\gamma}{\gamma + d(\boldsymbol{h}_{i},\boldsymbol{h}_{j})}$$
(14)

where  $\gamma$  is the scale parameter of the Cauchy distribution. When the Hamming distance is small, this function decreases rapidly, resulting in similar points being pulled into a small Hamming radius. For a pair of binary hash codes  $h_i$  and  $h_j$  the Hamming distance is:

323 
$$d(\boldsymbol{h}_i, \boldsymbol{h}_j) = \frac{\kappa}{4} \left\| \frac{\boldsymbol{h}_i}{\|\boldsymbol{h}_i\|} - \frac{\boldsymbol{h}_j}{\|\boldsymbol{h}_j\|} \right\|_2^2 = \frac{\kappa}{2} \left( 1 - \cos(\boldsymbol{h}_i, \boldsymbol{h}_j) \right)$$
(15)

324 Then, the Cauchy quantization loss is deduced as

325 
$$L_q = \sum_{i=1}^N \log\left(1 + \frac{d(|h_i|, 1)}{\gamma}\right)$$
(16)

In addition to reducing the quantization error, our approach takes into account the bitbalancing property, which means that each bit of the hash code has about a 50% chance of being 1 or -1. To produce more discriminative rather than ambiguous hash codes, we add the bit-balancing loss.

330 
$$L_b = \sum_{i=1}^{M} \frac{1}{K} \sum_{j=1}^{K} H_i$$
(17)

where K is the dimensionality of vector H. The purpose of bit balance loss is to 331 generate unbiased informative hash codes. Then, the hash loss  $L_{hash}$  is given by 332 equation (18). 333

$$L_{hash} = L_q + L_b \tag{18}$$

In summary, the total loss includes the loss of GFFAM and SHECM as the formula (19) 335 shown. 336

$$L = L_{ctr} + L_{cls} + L_q + L_b \tag{19}$$

#### **4** Experiments 338

#### **4.1 Datasets and Evaluation Metric** 339

We conducted experiments on three classic fine-grained datasets, CUB200-340 2011[47], Stanford Dog[48] and Stanford Cars[49], and compare them with other fine-341 342 grained retrieval methods.CUB200-2011 contains 200 bird species with 11788 images, in which 5994 images are for training and 5794 images are for testing. Stanford dog 343 contains 120 dog species and 20580 images, in which 12000 images are for training 344 and 8580 images are for testing. Stanford cars contains 196 car species and 16185 345 images, in which 8144 images are for training and 8041 images are for testing. We 346 follow the experimental protocol in [50], using the test images as the query set and the 347 training images as the retrieval database for all images. 348

We use retrieval accuracy and Precision-Recall curve as evaluation metrics for a 349 fair comparison with previous methods. We calculate the retrieval precision as: 350

351 
$$mAP = \frac{1}{n_q} \sum_{i=1}^{n_q} AP, \ AP = \frac{1}{N_+} \sum_{k=1}^n \frac{N_+^k}{k} pos(k)$$
(20)



Where  $n_q$  denotes the number of samples in query set, n denotes the number of top 352

returned samples. The  $N_{+}$  represents the number of positive samples in the *n* returned samples, and the  $N_{+}^{k}$  refers to the number of positive samples in the top *k* returned samples. The value of pos(k) is 1 if the image at position *k* is positive while will be 0 otherwise.

Precision-Recall curve measures the relationship between precision and recall. Specially, for the image retrieval learning, precision represents the proportion of positive samples in all returned samples while recall refers to the proportion of correctly retrieved samples to all positive samples in the database, which can be calculated by the following formula (21) and (22).

362 
$$Precision = \frac{TP}{TP+FP}$$
 (21)

$$Recall = \frac{TP}{TP+FN}$$
(22)

TP, FN, and FP in formula (21) and (22) denote the positive samples of the correct classification, the positive samples of the misclassification and the negative samples of misclassification, respectively.

## 367 **4.2 Implementation Details**

For the GFFAM method, we implement the code based on PyTorch and train the models by preprocessing images to size  $448 \times 448$  with 2 RTX 2080ti GPU. In all our experiments, we use the Resnet-50 as the backbone to extract the features and choose the output of Conv5 layers as feature maps. For the generation of attention maps, the attention maps are obtained by a  $1 \times 1$  convolutional kernel with default *M* of 32, while the dropping threshold  $T_d$  is set to 0.5. For building k-nearest neighbor graph, the number k of nearest neighbors is 3.

Although all data sets are labeled with bounding boxes or part locations, our 375 approach uses only category labels as monitoring information. In the training stage, the 376 377 model is trained for 160 epochs by using the Stochastic Gradient Descent (SGD) with the momentum of 0.9 and batch size of 12. It's worth mentioning that the initial learning 378 rate is set to 0.001, which has an exponential decay of 0.9 after every 2 epochs. In the 379 inference stage, the drop branch is not used, we follow the experimental protocol in 380 [50], using test images as the query set and training images as the retrieval database for 381 all experiments. 382

383 **4.3 State-of-the-arts Comparison** 

Comparisons with Classic Retrieval Algorithms. Experiments in this section, we first compare our approach with other state-of-the-art classic retrieval algorithms. For the sake of fairness, we train all the model using Resnet-50 as feature extractor with different bits of hash code on three fine-grained image datasets.

From Table 1, we can see that our GDF-Net achieves a mAP of 80.43%, 86.01% 388 389 and 86.92% with 16-bit, 32-bit and 48-bit hash codes on CUB 200-2011 data set, outperforming the previous methods by a large margin in all hash bits. Compared with 390 the second-place method of FPH, the mAP of GDF-Net increases by 29.15%, 27.69% 391 and 25.13%, respectively. For the Stanford Dogs and Cars data set, the significant 392 improvement trend is also obvious. The huge gap in results strongly shows the 393 difference among the fine-grained image retrieval and common image retrieval, which 394 395 also proves that specially designed discriminative feature extractor is necessary for the fine-grained image retrieval. In addition, compared with classic retrieval methods 396

which build hash code based on the global image feature, our GDF-Net shows a better stability for the change of hash code bits. For example, when the hash code decreases from 48-bit to 16-bit, the mAP of GDF-Net only decreases by about 6% on the CUB 200-2011 data set, while FPH decreases by more than 10%. This result demonstrates the effectiveness of the proposed method, and can further suggest that it is more efficient by mining the fine-grained local features to build the hash code.

403 Table 1 Mean average precision (mAP) of different hashing retrieval methods on three Fine-grained image datasets

Method	CUB 200-2011			Stanford Dogs			Stanford Cars		
Method	12bit	32bit	48bit	12bit	32bit	48bit	12bit	32bit	48bit
DHN	0.3711	0.4172	0.4602	0.4559	0.5290	0.5736	0.4608	0.5050	0.5574
DQN	0.3789	0.4355	0.4811	0.4676	0.5234	0.5795	0.4897	0.5444	0.5821
HashNet	0.4027	0.4712	0.5103	0.4988	0.5574	0.5981	0.5073	0.5508	0.5832
DCH	0.4602	0.5233	0.5740	0.6081	0.6567	0.6779	0.5488	0.6009	0.6175
FPH	0.5128	0.5832	0.6179	0.6312	0.6909	0.7090	/	/	/
Ours	0.8043	0.8601	0.8692	0.8193	0.8586	0.8664	0.8864	0.9181	0.9268

Comparisons with Fine-grained Retrieval Algorithms. We also conduct 404 405 analytical experiments to evaluate the retrieval performance with other similar finegrained retrieval algorithms, where the mAP results on CUB 200-2011 have been 406 presented in Table 2. From Table 2, we can clearly observe that our approach achieves 407 the best retrieval performance in all cases. Compared with previous works such as 408 409 SCDA, CRL-WSL and DCI-NS without hash module, GDF-Net significantly improved the mAP even if the dimensionality of feature representation is quite smaller. We 410 consider that the possible reason may be the sensitivity of fine-grained tasks to feature 411 412 quality. Though the high dimensional image features contain more objects' information, it also means that the discriminative features can't play a leading role when learning 413 hash code in Hamming space, which makes the hash codes and final retrieval results 414

415	more susceptible to be disturbed by similar objects of the same or different categories.
416	Compared with the hash-based FGIR methods such as ExchNet and FCAENet, the best
417	mAP of our GDF-Net outperforms by more than 15% when the hash code is 48bit. In
418	addition, we also notice that the performance of ExchNet and FCAENet changes greatly
419	with the hash bits, which is improved by more than 40% when the hash code increases
420	from 12-bit to 32-bit. While for our GDF-Net, the mAP only increases by 6% and even
421	increases less than 1% when hash code changes from 32-bit to 48-bit. The more stable
422	results not only show the compactness and efficiency of our hash code but also verify
423	the effectiveness of our module GFFAM from the side.

To the best of our knowledge, [21] is the closest work related to our GDF-Net. As 424 shown in Table 2, we achieve an mAP of 86.92% with the 48bit hash code, which is 425 426 1.43% higher than [21]. And when the hash codes reduce to 32-bit, we also outperform 1.75%. We suppose the main reason for the improved results comes from the acquisition 427 of higher quality fine-grained image features through mining the context correlation 428 among different discriminative features, which we will analyze more detailed in the 429 next section. Furthermore, for the sake of verifying the superiority of our method, we 430 also evaluate our method with precision curves of 48-bit hash codes w.r.t different recall 431 rates as shown in Fig. 4. And the larger the area enclosed by the PR curve and axes, the 432 better retrieval performance. Experimental results on three datasets show that our GDF-433 Net outperforms all the methods. 434



Table 2 mAP of GDF-Net compared with Fine-grained image retrieval methods

Method	Model	#Dim	mAP
SCDA	VGG16	4096	0.5957

CRL-WSL VGG16		1024	0.6590
DCI -NS	Resnet- 50	1024	0.6790
	Description	12	0.2514
ExchNet	Resnet-	32	0.6774
	50	48	0.7105
	D (	12	0.3476
FCAENet	Resnet-	32	0.7385
	50	48	0.8014
	Descret	12	0.7901
[21]	Resnet-	32	0.8426
	50	48	0.8549
	Description	12	0.7628
	Resnet-	32	0.8165
0	18	48	0.8244
Ours	Descret	12	0.8043
	Kesnet-	32	0.8601
	50	48	0.8692

436



Fig. 4 Precision-recall curve with topN@48 bits on three Fine-grained image datasets

## 437 **4.4 Ablation Studies**

Quantitative Evaluations of GCN. As we have discussed in Section 3.2, GDFNet constructs a KNN graph convolutional neural network to capture the correlation
between different local features, which aims at reconstructing highly discriminative and
contextual features for hash encoding. To evaluate the effectiveness of our KNN graph,
we further conduct a set of experiments to compare the mAP between GCN and other
correlation mining strategies quantitatively.

444 As shown in Table 3, our baseline is directly feeding M original local features  $f_k$ 

445	into the semantic hash encoding module, which achieves a mAP of 76.84% and 82.02%
446	for 12-bit and 48-bit on the most difficult CUB data set. To aggregate local features, we
447	also report the results of some additional strategies, containing the application of a
448	global average pooling(GAP)and a global max pooling(GMP) operation on all local
449	features $f_k$ , which achieves a significant improvement in retrieval performance. For
450	example, GMP obtains a mAP of 84.96% with a 48-bit hash code on CUB, which is
451	over 2% higher than the baseline. All the results verify our hypothesis that aggregating
452	independent discriminative features are conducive to improving the feature quality and
453	performance of fine-grained retrieval task. For our GCN, we achieve a mAP of 86.27%
454	and 86.92% on CUB200-2011 when reconstructing the higher quality features with
455	global average pooling and global max pooling, respectively. Compared with simple
456	GAP/GMP, we achieve a ~2% improvement. These results show the effectiveness of
457	our GCN which reconstruct feature under the guidance of discriminative features' k-
458	nearest neighbor relationship.

Table 3 Comparison of mAP between GCN and other correlation mining strategies

Completion mining strategy	CUB 200-2011		Stan	ford Dogs	Stanford Cars		
Correlation mining strategy	12bit	48bit	12bit	48bit	12bit	48bit	
baseline	0.7684	0.8202	0.7721	0.8234	0.8401	0.8876	
GAP	0.7743	0.8325	0.7836	0.8367	0.8495	0.8941	
GMP	0.7862	0.8496	0.7879	0.8482	0.8603	0.9004	
GCN with GAP	0.8012	0.8627	0.8122	0.8601	0.8796	0.9203	
GCN with GMP	0.8043	0.8692	0.8193	0.8664	0.8864	0.9268	

Effectiveness to GCN structure. When constructing GCN, we use local

discriminative features as the input of graph nodes. The initial number of nodes in the

464	graph network will significantly affect model performance. For this reason, we analyze
465	the influence of nodes numbers by ablation studies on three Fine-grained benchmarks.
466	As shown in Table 4, more initial nodes usually contribute to better performance. On
467	CUB 200-2011 data set, GDF-Net achieve a mAP of 85.42% with 4 initial graph nodes,
468	while the mAP increases to 86.75% and 86.92% when the initial graph nodes are set to
469	16 and 32. The similar rising trends is also shown on other two data sets. Considering
470	the balance between the computational overhead and retrieval performance, we finally
471	choose to set the number of graph nodes as 32 in our GDF-Net.

Table 4 Comparison of mAP with different numbers of initial nodes on three Fine-grained benchmarks

Numbers of initial nodes	CUB 200-2011	Stanford Dogs	Stanford Cars
4	0.8542	0.8487	0.9124
8	0.8603	0.8571	0.9196
16	0.8675	0.8652	0.9241
32	0.8692	0.8664	0.9268
64	0.8695	0.8664	0.9269

Effectiveness of Hash loss. Another important module in our GDF-Net 473 framework is the semantic hash encoding, which designs a semantic hash coding 474 module before the classification layer with the supervision of  $L_q$  and  $L_b$ . In order to 475 verify the effectiveness of different losses for hash coding, we conduct a set of 476 experiments to quantitatively evaluate the retrieval results when combining different 477 loss functions, as shown in Table 5. For CUB, the simple baseline trains hash model 478 without any loss, whose mAP can increase to 82.34% and 85.65% with 48-bits hash 479 480 code after introducing the  $L_q$  and  $L_b$ . For our SHECM with 48-bits hash code, it can finally obtain 86.92% on CUB, 86.64% on Stanford Dogs and 92.68% on Stanford Cars. 481 All the results show that Cauchy quantization loss and bit-balancing loss are beneficial 482

484

Table 5 Comparison of mAP with different losses on three Fine-grained benchmarks

T	L <sub>b</sub>	CUB 200-2011		Stanfor	d Dogs	Stanford Cars		
ц		12bit	48bit	12bit	48bit	12bit	48bit	
/	/	0.7123	0.8028	0.7583	0.7911	0.7867	0.8259	
/	$\checkmark$	0.7795	0.8234	0.7689	0.8018	0.8336	0.8617	
$\checkmark$	/	0.7998	0.8565	0.8006	0.8412	0.8719	0.9013	
$\checkmark$	$\checkmark$	0.8043	0.8692	0.8193	0.8664	0.8864	0.9268	

Influence of  $T_d$  and K.  $T_d$  is the dropping threshold for masking image 485 generation, while the K is the number of neighbor nodes aggregated for graph 486 487 propagation. To quantitatively analyze the impact of the above two hyperparameters on retrieval performance, we conduct additional ablation experiments on the CUB dataset 488 and Resnet50 with 48-bits hash codes, whose results are shown in the Table 6 and Table 489 7, respectively. When  $T_d$  changes from 0.1 to 0.9 with a stride of 0.1, the mAP 490 491 increases from 79.65% to the peaks of 86.92% and then gradually decreases to 80.31%, which show the results are sensitive to the dropping threshold. Too small threshold will 492 discard most of the area on the original image and only remain background pixels, 493 which is not conducive to finding other sub-discriminator features. While too large 494 threshold makes the mined discriminative features almost unchanged, which will 495 inhibit the diversity of local features. 496

497

Table 6 Influence analysis of  $T_d$  with Resnet50 and 48-bits hash codes on CUB

T <sub>d</sub>	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
mAP	0.7965	0.4689	0.7250	0.8468	0.8692	0.8620	0.8384	0.8252	0.8031

498 Similar to the changing trend of the dropping threshold, it turns out that considering
499 too much neighbor node information will drastically increase the difficulty of
500 propagation and convergence of the graph convolutional module, whose mAP will

decrease from 86.92% to 84.87% when *K* changes from 3 to 6. While set a smaller K = 2, the results will decrease about 1%, which shows the effectiveness of our GFFAM considering the spatial context information of different nodes.

504

Table 7 Influence analysis of K with Resnet50 and 48-bits hash codes on CUB

K	2	3	4	5	6
mAP	0.8506	0.8692	0.8681	0.8553	0.8487

## 505 **4.5 Visualization Results**

Visualization of Top 10 Results. Fig.5 shows the top 10 samples returned by our GDF-Net and other retrieval algorithms, including DCH[35] and Fine-grained hash retrieval algorithm [21] on the CUB Birds, Stanford dogs and Stanford Cars. Thanks to our powerful fine-grained feature extraction module, GDF-Net yields much more relevant and user-desired retrieval results than the state-of-the-art method.



Fig.5 Examples of top 10 retrieved images and precision @ 10 on three datasets.

511 **Visualization of hash code compactness.** In order to analyze the quality of hash 512 code, we visualize the hash code generated by GDF-Net and DCH using the t-513 distributed stochastic neighbor embedding (t-SNE) on Stanford Cars. As the Fig.6 shown, we sample 10 categories and find that the results of our method are better, because the results returned by DCH focus more on the corresponding global information, while the GEFAM module in GDF-Net designs a specialized feature extractor for fine-grained images.





(b) Our

Fig.6 The t-SNE of hash codes learned by GDF-Net and DCH on Stanford Cars

# 518 **5 Conclusion**

In this paper, we propose a novel deep hash fine-grained retrieval method based on 519 graph convolutional neural network, which consists of Global Fine-grained Feature 520 Aggregation Module (GFFAM) and Semantic Hash Encoding Module (SHECM). 521 GFFAM constructs a k-nearest neighbor graph to capture the correlation among 522 different independent discriminative features, which can guide the fusion of these local 523 features adaptively to enhance the feature quality. To reduce the storage overhead, 524 GFFAM designs a hash encoding layer before the final classification layer, which can 525 generate semantic and compact hash codes with the guidance of Cauchy quantization 526 loss and bit-balancing loss. Extensive experimental results on CUB Birds, Stanford 527 Dogs and Stanford Cars demonstrate that our proposed method is comparable to the 528

529	state-of-art methods for fine-grained retrieval task. In the future, as we notice that fine-				
530	grained images tend to have high intra-class variances and low inter-class variances, we				
531	1 will further strengthen learning tolerance to improve the retrieval accuracy.				
532	<b>Reference</b>				
533	[1]	L. Xie, J. Wang, B. Zhang, and Q. Tian, "Fine-grained image search," IEEE Transactions on			
534		Multimedia, vol. 17, no. 5, pp. 636-647, 2015.			
535	[2]	K. Lin, F. Yang, Q. Wang, and R. Piramuthu, "Adversarial learning for fine-grained image			
536		search," in 2019 IEEE International Conference on Multimedia and Expo (ICME), 2019: IEEE,			
537		рр. 490-495.			
538	[3]	K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image			
539		recognition," arXiv preprint arXiv:1409.1556, 2014.			
540	[4]	K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in			
541		Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770-			
542		778.			
543	[5]	C. Szegedy et al., "Going deeper with convolutions," in Proceedings of the IEEE conference on			
544		computer vision and pattern recognition, 2015, pp. 1-9.			
545	[6]	S. Jin, H. Yao, X. Sun, S. Zhou, L. Zhang, and X. Hua, "Deep saliency hashing for fine-grained			
546		retrieval," IEEE Transactions on Image Processing, vol. 29, pp. 5336-5351, 2020.			
547	[7]	J. Wang, T. Zhang, N. Sebe, and H. T. Shen, "A survey on learning to hash," IEEE transactions			
548		on pattern analysis and machine intelligence, vol. 40, no. 4, pp. 769-790, 2017.			
549	[8]	TY. Lin, A. RoyChowdhury, and S. Maji, "Bilinear cnn models for fine-grained visual			
550		recognition," in Proceedings of the IEEE international conference on computer vision, 2015,			

- **551** pp. 1449-1457.
- P. Li, J. Xie, Q. Wang, and Z. Gao, "Towards faster training of global covariance pooling
  networks by iterative matrix square root normalization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 947-955.
- 554 On computer vision and patient recognition, 2010, pp. 947-955.
- [10] Z. Yang, T. Luo, D. Wang, Z. Hu, J. Gao, and L. Wang, "Learning to navigate for fine-grained
  classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018,
  pp. 420-435.
- 558 [11] H. Zheng, J. Fu, T. Mei, and J. Luo, "Learning multi-attention convolutional neural network for
  559 fine-grained image recognition," in *Proceedings of the IEEE international conference on*
- *computer vision*, 2017, pp. 5209-5217.
- 561 [12] J. Fu, H. Zheng, and T. Mei, "Look closer to see better: Recurrent attention convolutional neural
- network for fine-grained image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4438-4446.
- 564 [13] X.-S. Wei, J.-H. Luo, J. Wu, and Z.-H. Zhou, "Selective convolutional descriptor aggregation
- for fine-grained image retrieval," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp.
  2868-2881, 2017.
- 567 [14] Q. Cui, Q.-Y. Jiang, X.-S. Wei, W.-J. Li, and O. Yoshie, "ExchNet: A Unified Hashing Network
- for Large-Scale Fine-Grained Image Retrieval," in *European Conference on Computer Vision*,
  2020: Springer, pp. 189-205.
- 570 [15] T. Hu, H. Qi, Q. Huang, and Y. Lu, "See better before looking closer: Weakly supervised data
- 571 augmentation network for fine-grained visual classification," *arXiv preprint arXiv:1901.09891*,
- **572** 2019.

- 573 [16] W. Liu, A. Rabinovich, and A. C. Berg, "Parsenet: Looking wider to see better," *arXiv preprint*574 *arXiv:1506.04579*, 2015.
- 575 [17] M. Zhou, Y. Bai, W. Zhang, T. Zhao, and T. Mei, "Look-into-object: Self-supervised structure
- 576 modeling for object recognition," in *Proceedings of the IEEE/CVF Conference on Computer*

577 *Vision and Pattern Recognition*, 2020, pp. 11774-11783.

- 578 [18] D. Zoran, M. Chrzanowski, P.-S. Huang, S. Gowal, A. Mott, and P. Kohli, "Towards robust
- 579 image classification using sequential attention models," in *Proceedings of the IEEE/CVF*580 *conference on computer vision and pattern recognition*, 2020, pp. 9483-9492.
- 581 [19] X. Zheng, R. Ji, X. Sun, B. Zhang, Y. Wu, and F. Huang, "Towards optimal fine grained retrieval
  582 via decorrelated centralized loss with normalize-scale layer," in *Proceedings of the AAAI*
- 583 *Conference on Artificial Intelligence*, 2019, vol. 33, no. 01, pp. 9291-9298.
- 584 [20] X. Zheng, R. Ji, X. Sun, Y. Wu, F. Huang, and Y. Yang, "Centralized Ranking Loss with Weakly
  585 Supervised Localization for Fine-Grained Object Retrieval," in *IJCAI*, 2018, pp. 1226-1233.
- 586 [21] H. Sun, Y. Fan, J. Shen, N. Liu, D. Liang, and H. Zhou, "A Novel Semantics-Preserving Hashing
- 587 for Fine-Grained Image Retrieval," *IEEE Access*, vol. 8, pp. 26199-26209, 2020, doi:
  588 10.1109/access.2020.2970223.
- 589 [22] M. Gori, G. Monfardini, and F. Scarselli, "A new model for learning in graph domains," in
   590 *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, 2005, vol.
- **591** 2: IEEE, pp. 729-734.
- 592 [23] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena, "Structural-rnn: Deep learning on spatio593 temporal graphs," in *Proceedings of the ieee conference on computer vision and pattern*594 *recognition*, 2016, pp. 5308-5317.

- 595 [24] C.-W. Lee, W. Fang, C.-K. Yeh, and Y.-C. F. Wang, "Multi-label zero-shot learning with
- structured knowledge graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1576-1585.
- 598 [25] X. Liang, X. Shen, J. Feng, L. Lin, and S. Yan, "Semantic object parsing with graph lstm," in
- 599 *European Conference on Computer Vision*, 2016: Springer, pp. 125-143.
- 600 [26] M. Henaff, J. Bruna, and Y. LeCun, "Deep convolutional networks on graph-structured data,"
  601 *arXiv preprint arXiv:1506.05163*, 2015.
- 602 [27] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph
- cnn for learning on point clouds," *Acm Transactions On Graphics (tog)*, vol. 38, no. 5, pp. 1-12,
  2019.
- 605 [28] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, "Locality-sensitive hashing scheme based
  606 on p-stable distributions," in *Proceedings of the twentieth annual symposium on Computational*
- 607 *geometry*, 2004, pp. 253-262.

609

608 [29] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin, "Iterative Quantization: A Procrustean

Approach to Learning Binary Codes for Large-Scale Image Retrieval," IEEE Transactions on

- 610 Pattern Analysis & Machine Intelligence, vol. 35, no. 12, pp. 2916-2929, 2013.
- 611 [30] H. Liu, R. Ji, Y. Wu, and W. Liu, "Towards optimal binary code learning via ordinal embedding,"
- 612 in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016, vol. 30, no. 1.
- 613 [31] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan, "Supervised hashing for image retrieval via image
- 614 representation learning," in *Proceedings of the AAAI conference on artificial intelligence*, 2014,
  615 vol. 28, no. 1.
- 616 [32] H. Lai, Y. Pan, Y. Liu, and S. Yan, "Simultaneous feature learning and hash coding with deep

- 617 neural networks," in *Proceedings of the IEEE conference on computer vision and pattern*618 *recognition*, 2015, pp. 3270-3278.
- 619 [33] Z. Cao, M. Long, J. Wang, and P. S. Yu, "Hashnet: Deep learning to hash by continuation," in
- 620 *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5608-5617.
- 621 [34] H. Zhu, M. Long, J. Wang, and Y. Cao, "Deep hashing network for efficient similarity retrieval,"

622 in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016, vol. 30, no. 1.

623 [35] C. Yue, M. Long, B. Liu, and J. Wang, "Deep Cauchy Hashing for Hamming Space Retrieval,"

624 in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018.

- 625 [36] Y. Cao, M. Long, J. Wang, H. Zhu, and Q. Wen, "Deep quantization network for efficient image
- 626 retrieval," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016, vol. 30, no.
- 627

1.

- 628 [37] W.-J. Li, S. Wang, and W.-C. Kang, "Feature learning based deep supervised hashing with
  629 pairwise labels," *arXiv preprint arXiv:1511.03855*, 2015.
- 630 [38] S. Su, C. Zhang, K. Han, and Y. Tian, "Greedy hash: Towards fast optimization for accurate
- hash coding in cnn," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 806-815.
- 633 [39] Q. Y. Jiang and W. J. Li, "Asymmetric Deep Supervised Hashing," 2017.
- 634 [40] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected
  635 networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- 636 [41] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs
- 637 with fast localized spectral filtering," Advances in neural information processing systems, vol.
- 638 29, pp. 3844-3852, 2016.

- 639 [42] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks,"
  640 *arXiv preprint arXiv:1609.02907*, 2016.
- 641 [43] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated graph sequence neural networks,"
  642 *arXiv preprint arXiv:1511.05493*, 2015.
- [44] T.-Y. Lin and S. Maji, "Improved bilinear pooling with cnns," *arXiv preprint arXiv:1707.06772*,
  2017.
- 645 [45] G. Zhong, H. Xu, P. Yang, S. Wang, and J. Dong, "Deep hashing learning networks," in 2016
  646 *International Joint Conference on Neural Networks (IJCNN)*, 2016: IEEE, pp. 2236-2243.
- [46] Y. Wen, K. Zhang, Z. Li, and Q. Yu, "A Discriminative Feature Learning Approach for Deep
  Face Recognition," in *European Conference on Computer Vision*, 2016.
- 649 [47] WelinderP, BransonS, WahC, SchroffF, BelongieS, and PeronaP, "Caltech-UCSD Birds 200,"
- 650 *California Institute of Technology*, 2010.
- 651 [48] A. Khosla, N. Jayadevaprakash, B. Yao, and F. Li, "Novel dataset for fine-grained image
  652 categorization," 2013.
- 653 [49] J. Krause, M. Stark, J. Deng, and F. F. Li, "3D Object Representations for Fine-Grained
- 654 Categorization," in *IEEE International Conference on Computer Vision Workshops*, 2014.
- 655 [50] Y. Yang, L. Geng, H. Lai, Y. Pan, and J. Yin, "Feature pyramid hashing," in *Proceedings of the*
- 656 2019 on International Conference on Multimedia Retrieval, 2019, pp. 114-122.

657