# Improving Deep Binary Embedding Networks by Order-aware Reweighting of Triplets

**Jikai Chen**[1], **Hanjiang Lai**[1], **Libing Geng**[1], **Yan Pan**[1],

[1] Schoolf of Data and Computer Science, Sun Yat-sen University

## Abstract

In this paper, we focus on triplet-based deep binary embedding networks for image retrieval task. The triplet loss has been shown to be most effective for the ranking problem. However, most of the previous works treat the triplets equally or select the hard triplets based on the loss. Such strategies do not consider the order relations, which is important for retrieval task. To this end, we propose an order-aware reweighting method to effectively train the triplet-based deep networks, which up-weights the important triplets and down-weights the uninformative triplets. First, we present the order-aware weighting factors to indicate the importance of the triplets, which depend on the rank order of binary codes. Then, we reshape the triplet loss to the squared triplet loss such that the loss function will put more weights on the important triplets. Extensive evaluations on four benchmark datasets show that the proposed method achieves significant performance compared with the state-of-the-art baselines.

## 1 Introduction

With the rapid development of the Internet, the amount of images grows rapidly. The large-scale image retrieval has attracted increasing interest. Hashing methods that encode images into binary codes have been widely studied since the compact binary codes are suitable for fast search and efficient storage. There are a multitude of hashing methods in the literature [Wang *et al.*, 2017; Wang *et al.*, 2016].

Among these methods, the supervised information is given with triplet labels, which have been shown to be most effective since hashing is actually a ranking problem [Zhuang *et al.*, 2016; Lai *et al.*, 2015]. In these works, the triplet ranking loss is defined to learn binary codes that preserve relative similarity relations. In [Lai *et al.*, 2015], an architecture based on deep convolutional neural networks (CNNs) with triplet ranking loss is proposed for image retrieval. In [Zhao *et al.*, 2015], it presents a deep semantic ranking based method to learn hash functions that preserve multi-level semantic similarity between multi-label images. The FaceNet [Schroff *et al.*, 2015] also uses the triplet ranking loss for face recognition and clustering. Due to the huge number of triplets, a collaborative two-stage approach [Zhuang *et al.*, 2016] is employed to reduce the training complexity of the triplet-based deep binary embedding networks.

Not all triplets are of equal importance. In [Hermans *et al.*, 2017], it finds that the triplet loss relatively quickly learns to correctly map most trivial triplets, which makes a large fraction of all triplets uninformative after some point. Thus, the loss decreases quickly at the beginning and slows down drastically after some point [Schroff *et al.*, 2015]. For instance, the triplet $(bird_1, bird_2, dog_3)$ is easier than the triplet $(bird_1, bird_2, bird_3)$ in which three images are from the fine-grained bird database. Intuitively, the hash model which was told over and over again that bird and dog are dissimilar cannot further improve the performance. Hence, up-weighting the informative triplets and down-weighting the uninformative triplets become a crucial problem.

However, most of the existing hashing methods treat the triplets equally [Lai *et al.*, 2015]. Few works select the hard triplets based on the loss [Wu *et al.*, 2017]. For instance, semi-hard negative mining [Schroff *et al.*, 2015] is proposed in FaceNet. It uses all anchor-positive pairs in a mini-batch and selects the negative examplars that are further away from the anchor than the positive examplar. [Wang and Gupta, 2015] investigated to select top $K$ hard negative triplets with highest losses and the other triplets are ignored. *In summary, all existing methods use the loss to select the hard triplets or treat them equally, and totally ignore the order relations in the rank list, which is important in retrieval task.* Since hashing problem is a ranking problem, the losses in the rank lists might not be sufficiently accurate than the order relations.

Inspired by that, we propose an order-aware reweighting method for triplet-based deep binary embedding networks, which up-weights the informative triplets and down-weights the uninformative triplets. We firstly introduce a weighting factor for each triplet. In practice, the weighting factor can be set to the value that indicates how the triplet is misranked by the current hash model. Hence, we use the MAP (*mean average precision*), which is a widely used evaluation measure, to calculate the weights. Specifically, for each mini-batch in the training phase, we encode the images into binary codes via deep CNNs. For an arbitrary triplet with an anchor, a positive code and a negative code, we rank all binary codes, including

the positive code and the negative code, in the mini-batch according to their Hamming distances to the anchor. The weight of this triplet is defined as the change of MAP by only swapping the rank positions of the positive code and the negative code. Besides this order-aware weighting factor, we further use the squared triplet loss instead of the linear form, which up-weights hard triplets and down-weights easy ones from the perspective of the order relation of binary codes in triplets themselves. We conduct extensive evaluations on four benchmark datasets for image retrieval. The empirical results show that the proposed method achieves significant performance over the baseline methods.

## 2 Related Work

Hashing methods [Wang *et al.*, 2017] that learn similarity-preserving hash functions to encode data into binary codes have become popular methods for nearest neighbor search. Many methods have been proposed, which mainly can be divided into three categories: 1) the unsupervised hashing methods [Shen *et al.*, 2018; Liu *et al.*, 2017], 2) the semi-supervised hashing methods [Zhang and Peng, 2017; Wang *et al.*, 2010] and 3) the supervised hashing methods [Gui *et al.*, 2018].

Learning the hash codes with deep frameworks, e.g., CNN-based methods [Yang *et al.*, 2018], has been emerged as one of the leading approaches. According to the forms of the supervised information, previous works mainly fall into three categories: 1) the point-wise approaches, 2) the pair-wise approaches and 3) the triplet-based/ranking-wised approaches. The point-wise methods take a single image as input and the loss function is built on individual data [Lin *et al.*, 2015]. The pair-wise hashing methods take the image pairs as input and the loss functions are used to characterize the relationship (i.e., similar or dissimilar) between a pair of two images. For example, DPSH [Li, 2016] and DSH [Liu *et al.*, 2016] learn the hash codes by preserving the similarities among the input pairs of images. The triplet-based methods cast learning-to-hash as a ranking problem. [Lai *et al.*, 2015] proposed a deep triplet-based supervised hashing method. The triplet methods suffer from huge training complexity, thus [Zhuang *et al.*, 2016] further proposed a two-step approach to accelerate the training process of triplet-based hashing network.

Recently, some works [Wu *et al.*, 2017] show that sample selection plays an important role in learning the triplet-based network. The hard or semi-hard triplets are selected to train the network [Wang and Gupta, 2015; Schroff *et al.*, 2015]. The distance weighted sampling [Wu *et al.*, 2017] is proposed to select the informative and stable examples, where the samples are drawn uniformly according to their relative distance from one another. And the focal loss [Lin *et al.*, 2017] reshapes the standard cross entropy loss which down-weights the loss assigned to well-classified examples.

Inspired by these methods, we propose an order-aware method to reweight the triplet loss. The existing methods use the loss to select the hard examples or treat the triplets equally. In contrast, our proposed method introduces the order information [Donmez *et al.*, 2009] to weight the triplets, which is much more effective and accurate.
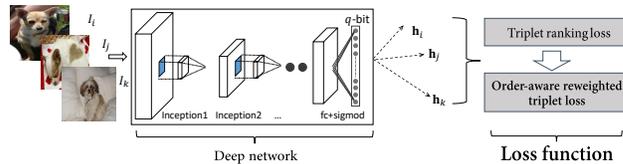


Figure 1: Overview of the triplet-based hashing network. The triplet-based network consists two sequential parts: a deep network and a triplet loss. In this paper, we only focus on the loss function. We reshape the triplet loss to our order-aware reweighted triplet loss.

## 3 Overview of Triplet-based Hashing Networks

In this section, we briefly summarize the triplet-based hashing framework. It takes triplets of images as inputs, i.e., $(I_i, I_j, I_k)$, in which $I_i$ is semantically more similar to $I_j$ than to $I_k$. The triplet hashing network itself can be divided into two sequential parts: a deep network with a stack of convolution, max-pooling and fully-connected layers; and a triplet ranking loss layer as shown in Figure 1.

In deep network, the convolutional layers are applied to produce powerful feature maps, which encode the images into high-level representations. Then the following several fully-connected layers project the feature maps into the desired-length feature vectors, e.g., $q$-dimensional vectors, where $q$ is the length of binary codes. The feature vectors are fed into a sigmoid layer which is smooth and well approximated the threshold function. The outputs of the network are restricted in the range $[0, 1]^q$. We denote the outputs of triplet network as $\mathbf{h}_i = \mathcal{F}(I_i)$, where $I_i$ is the input image and $\mathcal{F}$ is the deep network.

Through the deep network, triplet ranking loss [Lai *et al.*, 2015] is used to preserve the relative similarities of images. Given the input images in the form of $(I_i, I_j, I_k)$, the goal of hash network is to preserve the similarities of the learned binary codes, i.e., the binary code $\mathbf{h}_i$ is closer to $\mathbf{h}_j$ than to $\mathbf{h}_k$. The triplet ranking loss is defined by

$$\begin{aligned}
\ell_{(i,j,k)} &= \ell(\mathbf{h}_i, \mathbf{h}_j, \mathbf{h}_k) \\
&= \max(0, \epsilon - ||\mathbf{h}_i - \mathbf{h}_k||_2^2 + ||\mathbf{h}_i - \mathbf{h}_j||_2^2),
\end{aligned} \tag{1}$$

where $\epsilon$ is a hyper-parameter to control the margin between the two distances.

## 4 Order-aware Reweighting of Triplets

In this section, we only focus on the loss function and propose a simple yet effective order-aware reweighting algorithm for image retrieval. We first introduce the motivation of this work and then elaborate the proposed method.

### 4.1 Motivation

In retrieval task, the misranked triplets $(\mathbf{h}_i, \mathbf{h}_j, \mathbf{h}_k)$ in which $\mathbf{h}_j$ ranks behind $\mathbf{h}_k$ when $\mathbf{h}_i$ is query, should be further emphasized to boost the model. A simple and intuitive method is to add more weights to these misranked triplets. However, existing methods treat the triplets equally or only use
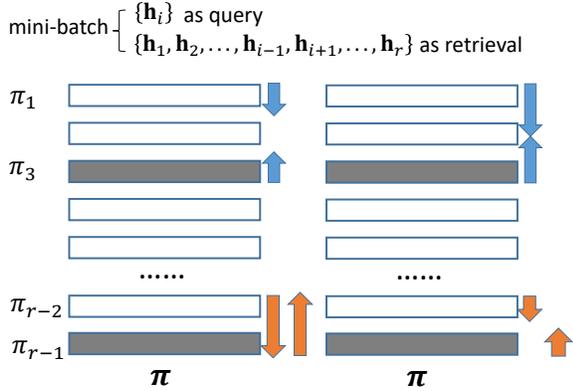
Figure 2: Motivation of our method. Given a query code $\mathbf{h}_i$ in a mini-batch, we return a rank list $\pi$ for $\mathbf{h}_i$. The white bars represent codes that are irrelevant to the query, and black bars represent the relevant codes. The arrows denote the weights and the moving directions of the items. Two example triplets, $\ell(\mathbf{h}_i, \mathbf{h}_{\pi_3}, \mathbf{h}_{\pi_1}) = 1$ and $\ell(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}}, \mathbf{h}_{\pi_{r-2}}) = 4$, are misranked by the current hashing model. Left: the weights of the triplets (see the arrows on the left) based on the loss. However, the $\mathbf{h}_{\pi_3}, \mathbf{h}_{\pi_1}$ are in the top positions and top results are more important in retrieval task. To better quantify the misranked triplet, we adopt the change of MAP by swapping the positions of the positive and negative codes to weight these misranked triplets. Right: the better choice to reweight the two triplets (see the arrows on the right).

the loss to select the hard ones. Image hashing is a ranking problem. Thus, the order relations are desirable. Figure 2 is an example illustration. Given $\mathbf{h}_i$ as the query, we rank the other binary codes according to their Hamming distance to $\mathbf{h}_i$ and $\pi = \{\pi_1, \cdots, \pi_{r-1}\}$ is the returned rank list. Take two misranked triplets $(\mathbf{h}_i, \mathbf{h}_{\pi_3}, \mathbf{h}_{\pi_1})$ and $(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}}, \mathbf{h}_{\pi_{r-2}})$ as examples in which $\ell(\mathbf{h}_i, \mathbf{h}_{\pi_3}, \mathbf{h}_{\pi_1}) = 1$ and $\ell(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}}, \mathbf{h}_{\pi_{r-2}}) = 4$, if we only use the loss to weight the triplets, the triplet $(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}}, \mathbf{h}_{\pi_{r-2}})$ will have larger weight. However, since $\mathbf{h}_{\pi_1}$, $\mathbf{h}_{\pi_3}$ are in the top positions and the top items are more important for the retrieval task, simply swapping the positions of $\mathbf{h}_{\pi_1}$ and $\mathbf{h}_{\pi_3}$ can achieve better performance than simply swapping those of $\mathbf{h}_{\pi_{r-2}}$ and $\mathbf{h}_{\pi_{r-1}}$. Hence, the triplet $(\mathbf{h}_i, \mathbf{h}_{\pi_3}, \mathbf{h}_{\pi_1})$ should be assigned more weight than $(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}}, \mathbf{h}_{\pi_{r-2}})$. The better choice is shown in the right side of Figure 2.

In this paper, we propose a simple algorithm that down-weights the uninformative triplets and up-weights the informative triplets via 1) a order-aware weighting factor and 2) a squared triplet ranking loss.

## 4.2 Order-aware Weighting Factor

Observed by that, we add an order-aware weighting factor to indicate the importance of the triplet. If the triplet is important, the more weight should be assigned to this triplet. We use the following steps to obtain the order-aware weights for triplets.

(1) Triplet generation with order information. Given a mini-batch of $r$ images, these images go through the deep network and are encoded as $\mathbf{h}_i, i = 1, \cdots, r$. Then we construct $r$ rank lists $\pi^{(1)}, \pi^{(2)}, \ldots, \pi^{(r)}$. The $i$-th rank list is

constructed for the $i$-th binary code, in which given the $i$-th code as the query, we rank the other $r - 1$ codes according to their Hamming distance to the $i$-th code, e.g., $\pi^{(i)} = \{\pi_1^{(i)}, \cdots, \pi_{r-1}^{(i)} | D_H(\mathbf{h}_i, \mathbf{h}_{\pi_1^{(i)}}) \leq \cdots \leq D_H(\mathbf{h}_i, \mathbf{h}_{\pi_{r-1}^{(i)}})\}$, where $D_H(\cdot, \cdot)$ denotes Hamming distance function. With these rank lists, we generate the set of triplets for the $i$-th query code: $T^{(i)} = \{(\mathbf{h}_i, \mathbf{h}_{\pi_j^{(i)}}, \mathbf{h}_{\pi_k^{(i)}}) | sim(\mathbf{h}_i, \mathbf{h}_{\pi_j^{(i)}}) > sim(\mathbf{h}_i, \mathbf{h}_{\pi_k^{(i)}})\}$ where $sim(\cdot, \cdot)$ is semantic similarity measure. Note that the first item is always $\mathbf{h}_i$ in the set $T^{(i)}$. Then, the union of all the $r$ sets, $T = \bigcup_{i=1}^{r} T^{(i)}$, is the total set of all triplets in the mini-batch.

(2) Order-aware reweighting of triplets. Now the problem becomes how to define weighting factors for these triplets. Take the set $T^{(i)}$ as an example, given a triplet $(\mathbf{h}_i, \mathbf{h}_{\pi_j^{(i)}}, \mathbf{h}_{\pi_k^{(i)}})$, we denote $\lambda_{(i,j,k)}$ as the importance weight of this triplet. Since the MAP is a widely used evaluation measure for ranking, we adopt MAP to calculate the weights. More specifically, for the triplet $(\mathbf{h}_i, \mathbf{h}_{\pi_j^{(i)}}, \mathbf{h}_{\pi_k^{(i)}})$, we first calculate the MAP of the rank list $\pi^{(i)}$ for the query $\mathbf{h}_i$. Then we only swap the rank positions of $\pi_j^{(i)}$ and $\pi_k^{(i)}$, and the other rank positions are fixed in rank list $\pi^{(i)}$, through which we can obtain another MAP. The absolute value of the difference between the two MAPs is used as the weight of the triplet. More specifically, let $\pi^{(i)} = \{\cdots, \pi_k^{(i)}, \cdots, \pi_j^{(i)}, \cdots\}$ and $\hat{\pi}^{(i)} = \{\cdots, \pi_j^{(i)}, \cdots, \pi_k^{(i)}, \cdots\}$. Note that other positions in the two rank lists $\pi^{(i)}$ and $\hat{\pi}^{(i)}$ are the same, and only the rank positions of $\pi_j^{(i)}$ and $\pi_k^{(i)}$ are swapped. The order-aware weight for the triplet $(\mathbf{h}_i, \mathbf{h}_{\pi_j^{(i)}}, \mathbf{h}_{\pi_k^{(i)}})$ is defined as

$$\lambda_{(i,j,k)} = |MAP(\pi^{(i)}) - MAP(\hat{\pi}^{(i)})|. \quad (2)$$

## 4.3 Squared Triplet Ranking Loss

Instead of the linear function, we propose a squared triplet ranking loss which aims to down-weight uninformative triplets and focuses on training hard distinguished triplets. The triplet loss function is changed as:

$$\ell_{(i,j,k)} \rightarrow [\ell_{(i,j,k)}]^2. \quad (3)$$

By using the squared triplet loss, the informative triplets will comprise the majority of the loss and dominate the gradient when back propagation is applied to update the weights of the deep networks. Even for the case that there are overwhelming number of uninformative triplets and a small number of informative triplets. Take 10 triplets as an example, there are one harder triplet, e.g., $\ell_1 = 5$ and nine easier triplets, e.g., $\ell_i = 1, i = 2, \cdots, 10$. The overall loss is 14, in which the loss of the easy triplets overwhelms that of the hard triplet. By using the quadratic function, we have $[\ell_1]^2 = 25$ and $[\ell_i]^2 = 1$ for other nine easy triplets. Thus, the loss of hard triplets can dominate the overall loss [1].

---

[1]Note that the conclusion also can be obtained when the losses are less than one, e.g., $\ell_1 = 0.5$ and $\ell_i = 0.1, i = 2, \cdots, 10$.
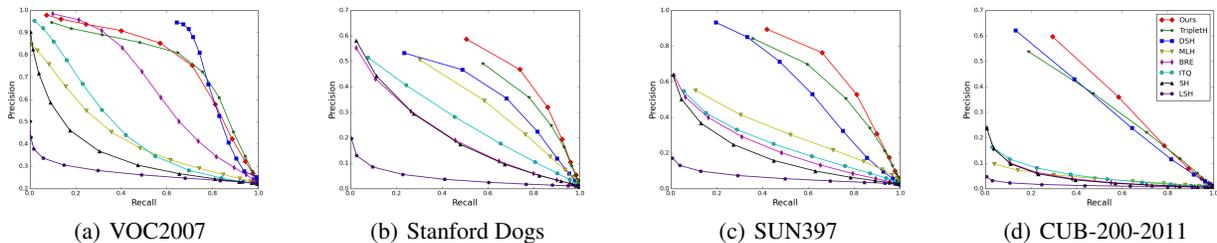
|  |  |  |  |
|---|---|---|---|
| (a) VOC2007 | (b) Stanford Dogs | (c) SUN397 | (d) CUB-200-2011 |

Figure 3: Precision-recall curves



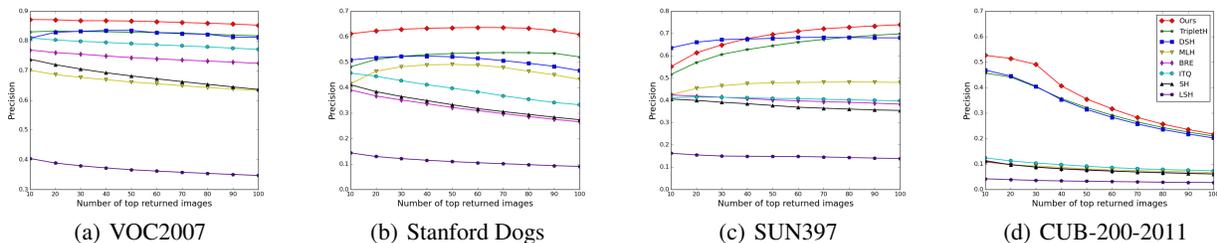|  |  |  |  |
|---|---|---|---|
| (a) VOC2007 | (b) Stanford Dogs | (c) SUN397 | (d) CUB-200-2011 |

Figure 4: Precision curves

## 4.4 Full Objective

The overall objective for order-aware reweighting triplet loss is defined as

$$\min \sum \lambda_{(i,j,k)} \times [\ell_{(i,j,k)}]^2. \qquad (4)$$

Our loss function contains two terms: 1) the order-aware weighting factors, which define the importance of the triplets by considering the whole rank list. It will let the loss function focus on the triplets that have the worst rank positions. And 2) the squared ranking objective is to up-weight the hard triplets by considering the order relation of binary codes in triplets.

## 5 Experiments

### 5.1 Datasets and Evaluation Measures

In this section, we conduct extensive evaluations of the proposed method on four benchmark datasets:

- **VOC2007** [Everingham *et al.*, 2010]: It consists of 9,963 annotated consumer photographs collected from the Flickr [2] photo-sharing website. There are 20 object classes in this dataset, and each image is annotated with 1.5 labels on average.

- **Stanford Dogs** [Khosla *et al.*, 2011]: It contains 20,580 images of 120 breeds of dogs from around the world, which has been built using images and annotation from ImageNet for the task of fine-grained image categorization.

- **SUN397** [Xiao *et al.*, 2010]: It contains 397 scene categories. The number of images varies across categories, but there are at least 100 images per category and 108,754 images in total.

- **CUB-200-2011** [Wah *et al.*, 2011]: It is a challenging dataset of 200 bird species. All 11,788 images and annotations were filtered by multiple users of Mechanical Turk.

In VOC2007, Stanford Dogs and CUB-200-2011 three datasets, we utilize the official train/test partitions to construct the query sets and the retrieval databases. The testing samples are used as the query set, and the rest images, i.e., the training samples, are used as the retrieval database. The training samples are also used to train the hash models. Note that the validation set of VOC2007 is included in the retrieval database but not used in training.

In SUN397 dataset, we follow the setting of [Do *et al.*, 2016] and use the subset of images associated with the 42 categories with each containing more than 500 images [3]. The randomly sampled 4,200 images (100 images per class) are constructed as the query set. The rest images are used as the database for retrieval. We randomly select 400 samples per class from the retrieval database to form the training set.

To evaluate the quality of hashing, we use the following evaluation metrics: mean average precision (MAP), precision-recall curves, and precision curves w.r.t. different numbers of top returned samples.

Table 1: MAP of Hamming ranking w.r.t. different numbers of bits on VOC2007, Standford Dogs and SUN397 datasets.

| Methods | VOC2007(MAP) | | | | Stanford Dogs(MAP) | | | | SUN397 (MAP) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 16 bits | 32 bits | 48 bits | 64 bits | 16 bits | 32 bits | 48 bits | 64 bits | 16 bits | 32 bits | 48 bits | 64 bits |
| Ours | **0.7672** | **0.8041** | **0.8147** | **0.8227** | **0.6745** | **0.7101** | **0.7252** | **0.7293** | **0.7505** | **0.8068** | **0.8152** | **0.8301** |
| TripletH | 0.7434 | 0.7745 | 0.7858 | 0.7897 | 0.5889 | 0.6478 | 0.6744 | 0.6904 | 0.6893 | 0.7660 | 0.7918 | 0.8027 |
| DSH | 0.7061 | 0.7234 | 0.7143 | 0.7116 | 0.4994 | 0.6134 | 0.6420 | 0.6453 | 0.5999 | 0.7031 | 0.7545 | 0.7725 |
| MLH | 0.4990 | 0.5044 | 0.4566 | 0.4786 | 0.4053 | 0.5031 | 0.6135 | 0.6284 | 0.3806 | 0.4632 | 0.4988 | 0.5212 |
| BRE | 0.6111 | 0.6323 | 0.6453 | 0.6517 | 0.2429 | 0.3158 | 0.3683 | 0.3936 | 0.2546 | 0.3185 | 0.3574 | 0.3675 |
| ITQ | 0.5793 | 0.5808 | 0.5796 | 0.5723 | 0.3275 | 0.4256 | 0.4682 | 0.4977 | 0.2789 | 0.3684 | 0.3873 | 0.4013 |
| SH | 0.4418 | 0.4175 | 0.4004 | 0.3835 | 0.2461 | 0.3076 | 0.3523 | 0.3714 | 0.2070 | 0.2436 | 0.2449 | 0.2404 |
| LSH | 0.2720 | 0.3098 | 0.3273 | 0.3551 | 0.0629 | 0.1293 | 0.1690 | 0.2245 | 0.0794 | 0.0981 | 0.1211 | 0.1486 |

Table 2: MAP of Hamming ranking w.r.t. different numbers of bits on CUB-200-2011 dataset.

| Methods | CUB-200-2011(MAP) | | | |
|---|---|---|---|---|
| | 16 bits | 32 bits | 48 bits | 64 bits |
| Ours | **0.5137** | **0.6519** | **0.6807** | **0.6949** |
| TripletH | 0.4508 | 0.5503 | 0.5963 | 0.6312 |
| DSH | 0.4374 | 0.4933 | 0.5553 | 0.6073 |
| MLH | 0.1069 | 0.1510 | 0.1768 | 0.2091 |
| BRE | 0.0716 | 0.0912 | 0.1127 | 0.1244 |
| ITQ | 0.0899 | 0.1266 | 0.1427 | 0.1599 |
| SH | 0.0703 | 0.0875 | 0.1017 | 0.1111 |
| LSH | 0.0272 | 0.0422 | 0.0573 | 0.0625 |

## 5.2 Experimental Setting

All deep CNN-based methods, including ours and previous baselines, are based on the same CNN architecture, i.e., GoogLeNet [Szegedy *et al.*, 2014]. We make the following modifications for hashing problem: 1) the last fully-connected layer is removed since it is for 1,000 classifications, and 2) another fully-connected layer with $q$ dimensional output is added to generate the binary codes. The weights are initialized with the pre-trained GoogleNet model [4] that learns from the ImageNet dataset. These experiments are implemented by using the open source *Caffe* framework. All networks are trained by stochastic gradient descent with 0.9 momentum and 0.0005 weight decay. The base learning rate is 0.001 and it is changed to one tenth of the current value after every 50 epochs. The total epoch is 150 and the batch size is 100. We implement both our methods and comparison ones for varied hash code lengths, e.g., 16 bits, 32 bits, 48 bits and 64 bits. For fair comparison, the hyper-parameters for all deep-network-based methods are the same, including training iterations, batch sizes and etc. For other non-deep-network-based methods, the input features are also extracted by the same pre-trained GoogleNet model, i.e., the last layer's output 1024 dimensional vector.

The source code of the proposed method will be made publicly available at the first author's homepage.

## 5.3 Experimental Results

### Comparison with State-of-the-art Methods

In this set of experiments, we evaluate and compare the performance of the proposed method with several state-of-the-art algorithms.

---

[4]http://dl.caffe.berkeleyvision.org/bvlc_googlenet.caffemodel

LSH [Gionis *et al.*, 1999], SH [Weiss *et al.*, 2008], ITQ [Gong and Lazebnik, 2011], MLH [Norouzi and Blei, 2011], BRE [Kulis and Darrell, 2009], triplet hashing (TripletH) [Lai *et al.*, 2015] and DSH [Liu *et al.*, 2016] are selected as the baselines. TripletH is one of the representative triplet-based methods and DSH is one of the representative pairwise-based methods. The results of these comparison methods are carefully obtained by the implementations provided by their authors, respectively. In TripletH, we replace the original divide-and-encode structure with the fully connected layer to generate hash codes.

Table 1 and Table 2 show the comparison results of MAP on the four datasets. It can be observed that the proposed method performs significantly better than all previous methods. Specifically, on VOC2017, our method obtains a MAP of 0.8227 on 64 bits, compared with 0.7897 of the existing triplet based method. On Stanford Dogs, our method shows an increase of 2% in comparison with the TripletH. Figure 3 and Figure 4 show the precision-recall and precision curves on 16 bits. Again, for most levels, our method yields the better accuracy. The results show that our proposed method can achieve better performance than the existing state-of-the-art methods.

### Effects of the Order-aware Weight and Squared Triplet Loss

In the second set of our experiments, we do ablation study to clarify the impact of each part of our method on the final performance.

In the first baseline, we only explore the effect of the order-aware weighting factors. The loss function is formulated as

$$\min \sum \lambda_{(i,j,k)} \ell_{(i,j,k)}. \tag{5}$$

In the second baseline, we set the order-aware weighting factors to be one for all triplets, e.g., $\lambda_{i,j,k} = 1$, and only explore the effects of the squared ranking loss. The objective is defined as

$$\min \sum [\ell_{(i,j,k)}]^2. \tag{6}$$

The last baseline is the existing triplet hashing (TripletH), which the objective is formulated as

$$\min \sum \ell_{(i,j,k)}. \tag{7}$$

Note that all baselines and our method use the same network and the only difference is the loss function, these comparisons can show us whether the proposed order-aware

Table 3: Quantitative ablation study on four databases.

| Methods | 16 bits | 32 bits | 48 bits | 64 bits |
|---|---|---|---|---|
| VOC2007 | | | | |
| Ours | 0.7672 | 0.8041 | 0.8147 | 0.8227 |
| Order-aware Weight | 0.7540 | 0.7851 | 0.7975 | 0.8052 |
| Squared Loss | 0.7552 | 0.7997 | 0.8049 | 0.8160 |
| TripletH | 0.7434 | 0.7745 | 0.7858 | 0.7897 |
| Stanford Dogs | | | | |
| Ours | 0.6745 | 0.7101 | 0.7252 | 0.7293 |
| Order-aware Weight | 0.6548 | 0.7058 | 0.7291 | 0.7368 |
| Squared Loss | 0.6261 | 0.6891 | 0.6899 | 0.7090 |
| TripletH | 0.5889 | 0.6478 | 0.6744 | 0.6904 |
| SUN397 | | | | |
| Ours | 0.7505 | 0.8068 | 0.8152 | 0.8301 |
| Order-aware Weight | 0.7117 | 0.7768 | 0.8065 | 0.8170 |
| Squared Loss | 0.7293 | 0.7732 | 0.8042 | 0.8191 |
| TripletH | 0.6893 | 0.7660 | 0.7918 | 0.8027 |
| CUB-200-2011 | | | | |
| Ours | 0.5137 | 0.6519 | 0.6807 | 0.6949 |
| Order-aware Weight | 0.5169 | 0.6158 | 0.6507 | 0.6779 |
| Squared Loss | 0.4698 | 0.5723 | 0.6371 | 0.6754 |
| TripletH | 0.4508 | 0.5503 | 0.5963 | 0.6312 |

Table 4: Comparison results of our method against hard triplet selection methods on four datasets.

| Methods | 16 bits | 32 bits | 48 bits | 64 bits |
|---|---|---|---|---|
| VOC2007 | | | | |
| Ours | 0.7672 | 0.8041 | 0.8147 | 0.8227 |
| HNM | 0.7545 | 0.7768 | 0.8010 | 0.8030 |
| Semi-hard | 0.7347 | 0.7442 | 0.7548 | 0.7655 |
| TripletH | 0.7434 | 0.7745 | 0.7858 | 0.7897 |
| Stanford Dogs | | | | |
| Ours | 0.6745 | 0.7101 | 0.7252 | 0.7293 |
| HNM | 0.5905 | 0.6763 | 0.7057 | 0.7187 |
| Semi-hard | 0.5933 | 0.6662 | 0.6967 | 0.7054 |
| TripletH | 0.5889 | 0.6478 | 0.6744 | 0.6904 |
| SUN397 | | | | |
| Ours | 0.7505 | 0.8068 | 0.8152 | 0.8301 |
| HNM | 0.7124 | 0.7801 | 0.7992 | 0.8142 |
| Semi-hard | 0.7167 | 0.7809 | 0.8002 | 0.8097 |
| TripletH | 0.6893 | 0.7660 | 0.7918 | 0.8027 |
| CUB-200-2011 | | | | |
| Ours | 0.5137 | 0.6519 | 0.6807 | 0.6949 |
| HNM | 0.4834 | 0.5793 | 0.6217 | 0.6641 |
| Semi-hard | 0.4859 | 0.6017 | 0.6346 | 0.6539 |
| TripletH | 0.4508 | 0.5503 | 0.5963 | 0.6312 |

weights and the squared triplet loss can contribute to the accuracy or not.

Table 3 shows the comparison results. We can observe that using both the order-aware weights and squared ranking loss performs best. And the proposed order-aware weight and squared loss perform better than the TripletH. It is desirable to reweight the triplets for triplet-based hashing networks.

**Comparison with Hard Triplet Selection Methods**

Our method is an order-aware method for reweighting the triplets. To show the advantages of the proposed method, we compare it to the hard triplet selection methods.

The first baseline is hard negative mining (HNM) [Wang and Gupta, 2015] for triplet sampling. It includes two steps: 1) random selection. They firstly randomly sample the triplets. After 10 epochs of training using data selected randomly, they do 2) hard negative mining, where it selects the top 4 negative triplets with highest losses for each anchor-positive pair. Similar to that, we first use the TripletH to train a hash model (random selection), then we fine-tune the network by using hard negative mining. In each mini-batch, we select the top 4 negative triplets as the suggestion by HNM. The second baseline is semi-hard triplet selection [Schroff *et al.*, 2015]. It uses all anchor-positive pairs in a mini-batch and selects the negative examplars that are further away from the anchor than the positive examplar.
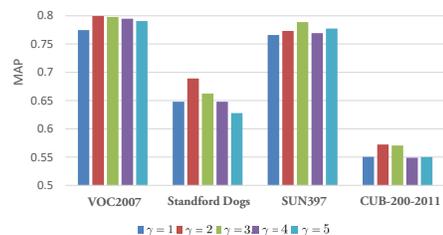
Table 4 shows the comparison results w.r.t. the MAP on the four datasets. We can see that the proposed method performs better than the loss-based hard triplet mining methods. The results show that order relations can further improve the performance.

**Effects of Different Functions for Triplet Loss**

In this paper, we use the squared triplet loss function to replace the linear form. In this set of experiments, we explore the effects of different functions. In general, the triplet loss



Figure 5: MAP of Hamming ranking w.r.t. different $\gamma$

can be written into more general form:

$$\min \sum [\ell_{(i,j,k)}]^{\gamma}. \tag{8}$$

When $\gamma = 1$, it equals to the traditional triplet ranking loss, when $\gamma = 2$, it is the squared triplet loss, etc.

Figure 5 shows the comparison results, which is implemented on 32 bits, on different functions: $\gamma = 1, \cdots, 5$. We can observe that the best results are obtained when $\gamma = 2$ or $\gamma = 3$. Hence, we use the squared triplet ranking loss in the paper.

## 6 Conclusion

In this paper, we proposed an order-aware reweighted method for training the triplet-based deep binary embedding networks. In the proposed deep architecture, images go through the deep network with stacked layers and are encoded into the binary codes. Then, we proposed to up-weight the informative triplets and down-weight the easy triplets by considering the order relations. One is the order-aware weighting factor, which is used to calculate the importance of the triplets. Another is the squared triplet ranking loss, which is used to put more weight on the triplets in which the codes are misranked.

Empirical evaluations on four datasets show that the proposed method achieves better performance than the state-of-the-art baselines.

# References

[Do *et al.*, 2016] Thanh-Toan Do, Anh-Dzung Doan, Duc-Thanh Nguyen, and Ngai-Man Cheung. Binary hashing with semidefinite relaxation and augmented lagrangian. In *ECCV*, pages 802–817, 2016.

[Donmez *et al.*, 2009] Pinar Donmez, Krysta M Svore, and Christopher JC Burges. On the local optimality of lambdarank. In *SIGIR*, pages 460–467, 2009.

[Everingham *et al.*, 2010] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010.

[Gionis *et al.*, 1999] Aristides Gionis, Piotr Indyk, and Rajeev Motwani. Similarity search in high dimensions via hashing. In *VLDB*, pages 518–529, 1999.

[Gong and Lazebnik, 2011] Yunchao Gong and Svetlana Lazebnik. Iterative quantization: A procrustean approach to learning binary codes. In *CVPR*, pages 817–824, 2011.

[Gui *et al.*, 2018] Jie Gui, Tongliang Liu, Zhenan Sun, Dacheng Tao, and Tieniu Tan. Fast supervised discrete hashing. *TPAMI*, 40(2):490–496, 2018.

[Hermans *et al.*, 2017] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.

[Khosla *et al.*, 2011] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *CVPRW*, Colorado Springs, CO, June 2011.

[Kulis and Darrell, 2009] Brian Kulis and Trevor Darrell. Learning to hash with binary reconstructive embeddings. In *NIPS*, pages 1042–1050, 2009.

[Lai *et al.*, 2015] Hanjiang Lai, Yan Pan, Ye Liu, and Shuicheng Yan. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*, pages 3270–3278, 2015.

[Li, 2016] WuJun Li. Feature learning based deep supervised hashing with pairwise labels. In *IJCAI*, pages 3485–3492, 2016.

[Lin *et al.*, 2015] Kevin Lin, Huei-Fang Yang, Jen-Hao Hsiao, and Chu-Song Chen. Deep learning of binary hash codes for fast image retrieval. In *CVPR*, pages 27–35, 2015.

[Lin *et al.*, 2017] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *arXiv preprint arXiv:1708.02002*, 2017.

[Liu *et al.*, 2016] Haomiao Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Deep supervised hashing for fast image retrieval. In *CVPR*, pages 2064–2072, 2016.

[Liu *et al.*, 2017] Qingshan Liu, Guangcan Liu, Lai Li, Xiao-Tong Yuan, Meng Wang, and Wei Liu. Reversed spectral hashing. *TNNLS*, 2017.

[Norouzi and Blei, 2011] Mohammad Norouzi and David M Blei. Minimal loss hashing for compact binary codes. In *ICML*, pages 353–360, 2011.

[Schroff *et al.*, 2015] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015.

[Shen *et al.*, 2018] Fumin Shen, Yan Xu, Li Liu, Yang Yang, Zi Huang, and Heng Tao Shen. Unsupervised deep hashing with similarity-adaptive and discrete optimization. *TPAMI*, 2018.

[Szegedy *et al.*, 2014] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.

[Wah *et al.*, 2011] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.

[Wang and Gupta, 2015] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, pages 2794–2802, 2015.

[Wang *et al.*, 2010] Jun Wang, Sanjiv Kumar, and Shih-Fu Chang. Semi-supervised hashing for scalable image retrieval. In *CVPR*, pages 3424–3431, 2010.

[Wang *et al.*, 2016] Jun Wang, Wei Liu, Sanjiv Kumar, and Shih-Fu Chang. Learning to hash for indexing big data—a survey. *P IEEE*, 104(1):34–57, 2016.

[Wang *et al.*, 2017] Jingdong Wang, Ting Zhang, Nicu Sebe, Heng Tao Shen, et al. A survey on learning to hash. *TPAMI*, 2017.

[Weiss *et al.*, 2008] Yair Weiss, Antonio Torralba, and Rob Fergus. Spectral hashing. In *NIPS*, pages 1753–1760, 2008.

[Wu *et al.*, 2017] Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *CVPR*, pages 2840–2848, 2017.

[Xiao *et al.*, 2010] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492, 2010.

[Yang *et al.*, 2018] Huei-Fang Yang, Kevin Lin, and Chu-Song Chen. Supervised learning of semantics-preserving hash via deep convolutional neural networks. *TPAMI*, 40(2):437–451, 2018.

[Zhang and Peng, 2017] Jian Zhang and Yuxin Peng. Ssdh: semi-supervised deep hashing for large scale image retrieval. *TCSVT*, 2017.

[Zhao *et al.*, 2015] Fang Zhao, Yongzhen Huang, Liang Wang, and Tieniu Tan. Deep semantic ranking based hashing for multi-label image retrieval. In *CVPR*, 2015.

[Zhuang *et al.*, 2016] Bohan Zhuang, Guosheng Lin, Chunhua Shen, and Ian Reid. Fast training of triplet-based deep binary embedding networks. In *CVPR*, pages 5955–5964, 2016.