

# Transductive Zero-Shot Hashing For Multilabel Image Retrieval

Qin Zou, Ling Cao, Zheng Zhang, Long Chen, Song Wang

<https://github.com/qinnzou/Zero-Shot-Hashing>

**Abstract**—Hash coding has been widely used in the approximate nearest neighbor search for large-scale image retrieval. Given semantic annotations such as class labels and pairwise similarities of the training data, hashing methods can learn and generate effective and compact binary codes. While some newly introduced images may contain undefined semantic labels, which we call unseen images, zero-shot hashing (ZSH) techniques have been studied for retrieval. However, existing ZSH methods mainly focus on the retrieval of single-label images and cannot handle multilabel ones. In this article, for the first time, a novel transductive ZSH method is proposed for multilabel unseen image retrieval. In order to predict the labels of the unseen/target data, a visual-semantic bridge is built via instance-concept coherence ranking on the seen/source data. Then, pairwise similarity loss and focal quantization loss are constructed for training a hashing model using both the seen/source and unseen/target data. Extensive evaluations on three popular multilabel data sets demonstrate that the proposed hashing method achieves significantly better results than the comparison methods.

**Index Terms**—image retrieval, zero-shot learning, multi-label image, deep hashing, transductive learning.

## I. INTRODUCTION

Hashing methods can transform high dimensional data into compact binary codes while preserving the similarity between them. With high computing efficiency and low storage cost, hashing methods have been widely used for large-scale image retrieval. A number of hashing methods have been proposed in the past decade [1]–[4].

Existing hashing methods can be roughly divided into two categories: supervised [5]–[7] and unsupervised [8]–[12]. The supervised hashing methods incorporate human-annotated information, *e.g.*, semantic labels and pairwise similarities, into the learning process to find an optimal hash function, while the unsupervised methods often learn hash functions by exploiting the intrinsic manifold structure of the unlabeled data. Generally, supervised methods can obtain much higher performance than the unsupervised ones.

In recent years, inspired by the remarkable success of deep neural networks in a broad range of computer-vision applications such as image classification [13]–[15], object

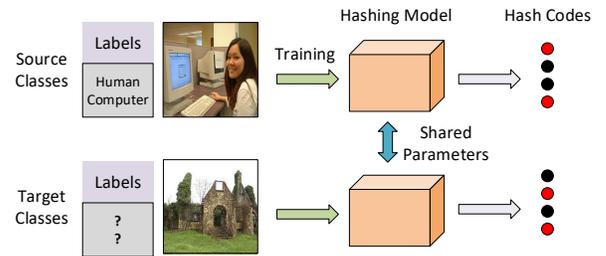


Fig. 1: An illustration of the transductive zero-shot hashing. In the learning procedure, both the source and target data are used for training the hashing model. The categories (labels) of the target data are unknown to the learning system.

detection [16]–[19], semantic segmentation [20], [21], many supervised hashing methods turn to use deep neural networks for hash-code learning [22]–[27]. These deep hashing methods have greatly advanced the retrieval performance on several popular benchmark datasets.

However, with the rapid emerging of new goods and new activities, images may contain concepts (or semantic labels) that are undefined before. For example, various commercial products with different shapes and appearances are released to the market every day, new sports with novel playing scenes are invented from time to time over the world. The images containing these new products or playing scenes are ‘unseen’ as compared to the ‘seen’ images holding pre-defined labels, and are supposed to be annotated with new labels for training the supervised learners. Consequently, the supervised hashing methods may face tremendous challenges due to the lack of timely and reliable annotation of ground-truth labels for the unseen images.

Zero-shot learning (ZSL) [28] is a technique that can potentially solve or alleviate this problem. Zero-shot learning bridges the semantic gap between ‘seen’ and ‘unseen’ categories by transferring supervised knowledge from other modalities or domains, *e.g.*, class-attribute descriptors and word vectors. For instance, the word embeddings of similar words that locate closely in the embedding space can capture the distributional similarity in the textual domain [29] based on a large-scale text corpus such as Wikipedia. Thus, such knowledge transfer can be used to capture the relationship between seen and unseen concepts, and can be helpful to handle unseen images in supervised learning.

For image retrieval under the circumstance of unseen im-

Q. Zou, L. Cao and Z. Zhang are with the School of Computer Science, Wuhan University, Wuhan 430072, P.R. China (E-mails: qzou@whu.edu.cn; lingcao@whu.edu.cn; zhengzhang@whu.edu.cn).

L. Chen is with the School of Data and Computer Science, Sun Yat-Sen University, Guangzhou 518001, P.R. China (E-mail: chenl46@mail.sysu.edu.cn).

S. Wang is with the Department of Computer Science and Engineering, University of South Carolina, Columbia, SC 29201 USA (E-mail: song-wang@cec.sc.edu).

ages, some zero-shot hashing (ZSH) methods [30]–[34] have also been proposed. Nevertheless, these methods focus on single-label image retrieval, in which a one-to-one visual-semantic representation pair can satisfy the training of a hashing model. While in more complicated scenarios, an image often contains multiple object classes, and hence more complex semantics and their relationships. How to represent the complex visual-semantic relationships for multi-label images, in a unified framework, is a difficult problem. To the best of our knowledge, there does not exist any work on multi-label zero-shot hashing.

Another important but easily ignored problem is that, since the underlying distributions of the source data and the target data are different, learning a hash function from a naive knowledge transfer on the source domain without making it adaptive to the target domain may lead to severe domain-shift problems. The knowledge transfer across different domains has become a very important issue in many computer-vision problems [35]–[37], which also needs to be investigated in zero-shot hashing.

Considering the problems discussed above, we propose a novel transductive zero-shot hashing method (T-MLZSH) for multi-label image retrieval. Both the labeled source and the unlabeled target data are used in the training phase, as illustrated in Fig. 1. The labeled source data are used to learn the relationship between visual images and semantic embeddings, and the unlabeled data of target classes are used to alleviate the domain-shift problem. More specifically, we first study a visual-semantic bridge via instance-concept coherence ranking on the source data. In instance-concept coherence ranking, a relatedness score for an image instance and a semantic concept is calculated for each image in the source data, where the score of an instance with a relevant label is larger than that of the same instance with an irrelevant label. Then, we can generate predicted labels for target data, and use these predicted labels as supervised information to guide the learning of hashing models. Moreover, we propose a focal quantization loss for fast and efficient hashing learning.

The contributions of this work lie in three-fold:

- A transductive zero-shot hashing method (T-MLZSH) is proposed to solve the domain-shift problem in multi-label image retrieval. To the best of our knowledge, it is the first work studying the zero-shot hashing for multi-label image retrieval.
- An instance-concept coherence ranking algorithm is proposed for visual-semantic mapping, which can be used to predict the labels for unseen target data and hence improve the performance of zero-shot deep hashing.
- The proposed method obtains very promising results on three popular multi-label datasets, which constructs the benchmark for zero-shot multi-label image retrieval and paves the road for new research in this field.

The rest of this paper is organized as follows. Section II briefly reviews the related work. Section III describes the neural network architecture for zero-shot image retrieval. Section IV demonstrates the effectiveness of the proposed method by experiments. Finally, Section V concludes the paper.

## II. RELATED WORK

**Zero-shot Learning.** Zero-shot learning (ZSL), refers to training on the seen labeled data with seen labels, and testing on the unseen data with unseen labels, where there is no intersection between the seen label set and the unseen label set. The core of ZSL is to obtain the instance labels for the unseen data. According to the ways of obtaining the instances, ZSL can be divided into three categories [38]–[40]: projection-based methods [41]–[43], instance-borrowing methods [44]–[46] and synthesizing methods [47]–[50].

The projection-based methods obtain instances by projecting the elements in the feature space and semantic space into a common space. Generally, the feature space contains the instances of seen classes, and the semantic space contains the prototypes of seen/unseen classes. In [42], a hinge rank loss was constructed to learn a linear transformation from the feature space to the semantic space. In [43], a bilinear mapping function was learned by minimizing the loss of an SVM classifier. In [51], the prediction noise and class bias of label embedding were decreased by deploying multiple classifiers in an ensemble manner. Instance-borrowing methods are based on the similarities of classes, in which the instances belonging to similar classes are taken as positive. In [52], instances in training were borrowed from the seen classes. The borrowed instances have high similarities with the unseen classes. Unlike instance-borrowing methods, the synthesizing methods create the pseudo instances for unseen classes, where the adversarial autoencoder and generative adversarial networks (GANs) are often used. In [53], an optimal latent space was learned to construct a bias-reducing generator network, which can reduce the hubness problem. In [54], an out-of-distribution detector was introduced to reduce the effect of domain shift, and a GAN was employed to synthesize the unseen instances.

To handle multi-label images, the projection is much more complex for the projection-based ZSL. A possible solution is to follow the visual-semantic mapping strategy used in the single-label case. In [?], the meaning of multiple labels for one instance was inferred by summing the word vector representations of individual labels. In [55], an alternative way was proposed to use the direct visual-semantic mapping, which first finds the corresponding area of each semantic label and then extracts the object-level visual presentation for visual-semantic mapping. Some other works try to utilize the co-occurrence among the labels. One typical work is the COSTA [56], which constructs the linear projection matrix between the seen labels and unseen labels by statistic learning on the annotated datasets.

**Hashing-based image retrieval.** Hashing methods for image retrieval can be roughly divided into two categories: the unsupervised and the supervised. The unsupervised hashing methods generate hash codes without any semantic labels. They use clustering techniques or projection strategies to transfer visual information to feature space and generate an optimized hash function to preserve the similarity in Hamming space [10], [12], [57]–[59]. Some classical algorithms, such as SH [11], formulated hash encoding as a spectral graph-partitioning problem and learned a nonlinear mapping to pre-

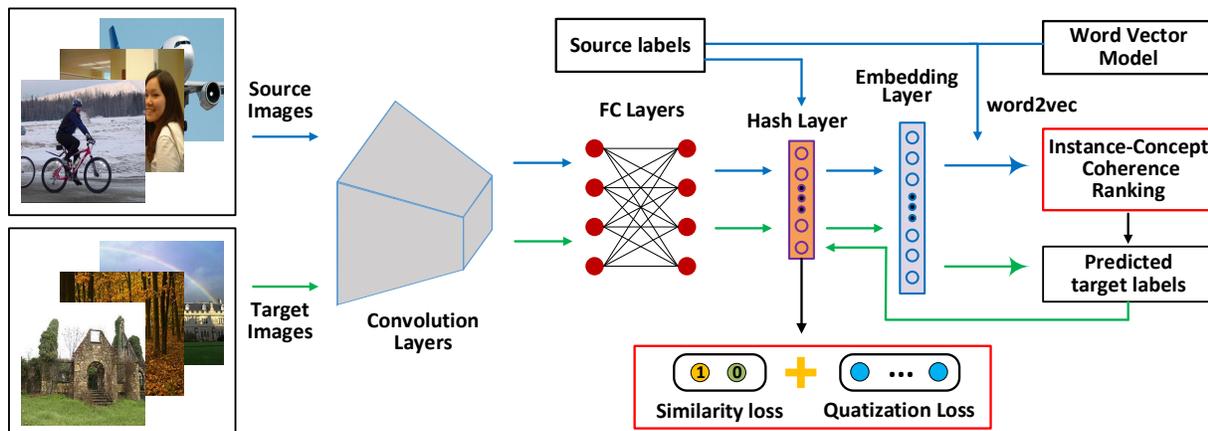


Fig. 2: An overview of the proposed T-MLZSH network. At first, the model constructs a visual-semantic bridge via instance-concept coherence ranking on the source data. It calculates the relatedness scores between visual features and semantic word vectors, under the assumption that related instance-concept pairs should have larger relatedness scores. Given the learned instance-concept coherence relationship, the most relevant concepts for each image instance of the unlabeled target data can be selected to guide the similarity-preserving learning. The whole network can be trained in an end-to-end manner.

serve semantic similarity. Some other methods, e.g., SPQ [60] and multi-k-means [61], decomposed the feature space into a Cartesian product of low-dimensional subspaces and encoded high-dimensional feature vectors into binary codes by clustering-based subspace quantization [62]. ECE [63] treated it as an optimization problem and combined the genetic programming with the boosting-based weight updating. DSTH [64] advocated discrete hash codes and resorted to the semantics augment from auxiliary contextual modalities.

The supervised hashing methods use the annotation information to learn compact hash codes, which usually perform better than the unsupervised methods. Among these supervised methods, CNN based hashing methods have attracted more and more attention due to the powerful representation ability of deep neural networks [65]. According to the difference of input forms, the existing deep hashing methods can be divided into two kinds. One receives image triplets as the input of the network and generates hash codes by minimizing the triplet ranking loss [5], [6]. These methods consume many computing resources and time to train the hashing model as there are enormous triplet combinations. The other receives the minibatch of images as input and uses the pairwise similarity between images as supervised information to learn the hash codes. Typical works of this form include the DHN [24], DSH [25], and HashNet [26], etc. Considering that existing hashing methods often fall short in concentrating relevant images to be within a small Hamming ball, DCH [66] built a pairwise cross-entropy loss on the Cauchy distribution to improve its capability on this point.

In recent years, hashing-based methods were also developed for multimodal retrieval, which transform high-dimensional data of different modalities into compact binary codes in a common Hamming space. The multimodal hashing supports image retrieval across other domains, e.g., text, video, audio, etc. Some representative work can be found in [67]–[72].

**Zero-shot hashing.** To handle images with unseen categories, some deep learning-based methods formulate the hash-

ing as an unsupervised problem [73], [74]. However, without using reliable supervised information, it is difficult to achieve satisfactory performance. Some other methods [30], [31], [75], from a different perspective, consider it as a zero-shot hashing (ZSH) problem. The goal of ZSH [76], [77] is to transfer the model trained on the seen data to unseen data via other available knowledge, such as word vector representations. Since the underlying data distributions of the seen-categories and the unseen-categories are different, the hashing functions learned by the seen categories without any adaptation to the unseen categories may cause a domain-shift problem. To narrow the domain gap between seen data and unseen data, ZSH-DA [32] first learns a zero-shot hashing model on seen data, and then learns the final hashing model with a domain-adaptation algorithm. In [33], a transductive zero-shot hashing network (TZSH) was proposed, which contains a coarse-to-fine similarity mining to find the most presentative target examples of each unseen labels, and adds these presentative examples and its corresponding predicted labels to the process of supervised hashing learning. [78] uses GCN to learn the zero-shot hashing model for sketch-image retrieval. Although the GCN method is very promising in exploring the relationship between semantic labels, it is not flexible and requires a pre-defined adjacency matrix of nodes and additional training costs.

### III. TRANSDUCTIVE MULTI-LABEL ZERO-SHOT HASHING

#### A. Problem Definition

Suppose  $\mathcal{D}^s = \{\mathcal{I}_i^s, \mathcal{Y}_i^s\}_{i=1}^{N_s}$  is a labeled source dataset including  $N_s$  images, where  $\mathcal{I}_i^s$  is an image and  $\mathcal{Y}_i^s$  is the corresponding label annotated with one or more classes, and  $\mathcal{D}^t = \{\mathcal{I}_i^t\}_{i=1}^{N_t}$  is an unlabeled target dataset, which includes  $N_t$  images of the unseen target classes and has the labels  $\mathcal{Y}^t$  unknown. In the zero-shot setting, the target and source classes are two mutually exclusive label sets, i.e.,  $\mathcal{Y}^t \cap \mathcal{Y}^s = \emptyset$ . For hash-code learning, we construct the similarity matrix

$\mathcal{S}=\{s_{ij}|i,j=1,2,\dots,N_s+N_t\}$ , where  $s_{ij} = 1$  denotes that the pairwise images  $\mathcal{I}_i$  and  $\mathcal{I}_j$  are similar, and  $s_{ij} = 0$  denotes they are dissimilar. The goal of T-MLZSH is to learn a mapping  $\mathcal{F} : \mathcal{I} \mapsto \{-1, +1\}^M$  on the labeled source dataset  $\mathcal{D}^s$  and the unlabeled dataset  $\mathcal{D}^t$  to encode an input image  $\mathcal{I}_i$  into an  $M$ -bit binary code  $\mathcal{F}(\mathcal{I}_i)$ , with the pairwise similarity preserved.

Figure 2 gives a flowchart of the proposed method. The input images firstly go through the deep network with the stacked convolutional and fully-connected layers and are encoded as a high-dimensional feature representation. Then, the outputs of the last fully-connected layer are fed into a hashing layer for compact binary encoding. To transfer the knowledge from seen categories to unseen categories and construct the bridge between visual and semantic modalities, we add a fully-connected layer after hashing layer, which maps features from hamming space to the common embedding space.

### B. Instance-Concept Coherence Ranking

Since there are no label information for target images, we should firstly predict labels for these images by transferring the knowledge from the semantic representations to visual features, before learning a supervised hashing function. Let  $v_i$  be the visual embedding of the  $i$ -th image instance  $\mathcal{I}_i$  and  $u_j$  be the semantic embedding of the  $j$ -th semantic concept, then we can calculate the relatedness score between  $\mathcal{I}_i$  and the  $j$ -th semantic concept in the embedding space:

$$o_{ij} = \langle v_i, u_j \rangle, \quad (1)$$

where  $\langle a, b \rangle = a^T b$  is the inner product operation. The semantic embeddings can be obtained from the existing word vector models, and the visual embeddings are variables that should be learned. During the training process, we can get a score list of source labels  $\{o_{i1}, o_{i2}, \dots, o_{iL_s}\}$ , where  $L_s$  is the number of seen labels. The goal of our embedding model is to learn a mapping function that scores with a relevant label should be higher than that with an irrelevant one, as illustrated by Fig. 3. Inspired by [79], we adopt a RankNet loss function to learn the ranking relationships for instance  $\mathcal{I}_i$ :

$$\mathcal{L}_{rank} = w_i \cdot \left( \sum_{p \in \mathcal{C}_i^+} \sum_{q \in \mathcal{C}_i^-} \log(1 + \exp(o_{iq} - o_{ip})) + \sum_{j \in \mathcal{C}} \log(1 + \exp(-\psi_{ij} o_{ij})) \right), \quad (2)$$

where  $\mathcal{C}_i^+$  and  $\mathcal{C}_i^-$  denote two sets of relevant and irrelevant labels to  $i$ -th instance.  $\psi_{ij}$  is defined as an indicator function, where  $\psi_{ij} = 1$  indicates that  $i$ -th instance is related to  $j$ -th label and  $\psi_{ij} = -1$  indicates that  $i$ -th instance is irrelative to  $j$ -th label.  $w_i = (|\mathcal{C}_i^+| \cdot |\mathcal{C}_i^-| + |\mathcal{C}|)^{-1}$  plays a regularization role. In the bracket of Eq. (2), the first term gives punishment to the situation when the labels irrelevant to  $\mathcal{I}_i$  have higher ranking orders than the relevant ones. The second term is used to enlarge the relatedness scores of the relevant pairs and reduce those of the irrelevant pairs.

Based on the above-trained model, the pairwise relatedness scores for the visual embeddings of target images and the semantic embeddings of target classes can be calculated. We

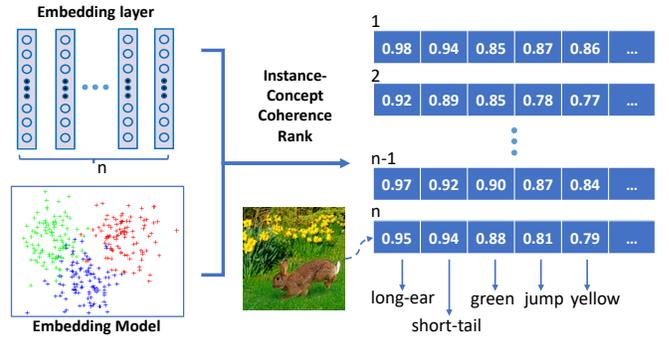


Fig. 3: An illustration of the embedding model. It learns a mapping function where the scores with a relevant label should be higher than that with an irrelevant one.

rank the scores  $\{o_{i1}, o_{i2}, \dots, o_{iL_t}\}$  ( $L_t$  is the number of unseen labels) in descending order, and select the classes of top- $k$  highest scores as the predicted target labels.

### C. Hash Code Learning

For efficient nearest neighbor search, the semantic similarity of original images should be preserved in the Hamming space. Generally, the similarity relationships can be defined with image labels. For a multi-label dataset, if two images share at least one label, they are considered similar, and dissimilar otherwise. Let  $\mathcal{B}$  be a set of hash codes for all images, and  $\mathcal{S}=\{s_{ij}\}$  be the pairwise similarity matrix, then the conditional probability of  $s_{ij}$  can be defined as,

$$p(s_{ij}|\mathcal{B}) = \begin{cases} \sigma(\Omega_{ij}), & s_{ij} = 1, \\ 1 - \sigma(\Omega_{ij}), & s_{ij} = 0, \end{cases} \quad (3)$$

where  $\Omega_{ij} = \langle b_i, b_j \rangle$  is the inner product of hash codes  $b_i$  and  $b_j$ , and  $\sigma(x) = \frac{1}{1+e^{-x}}$  is the sigmoid function, which scales the inner product value within  $[0, 1]$ .

We adopt negative log-likelihood as the cost function to measure the pairwise similarity loss, as formulated by Eq. (4),

$$\begin{aligned} \mathcal{L}_p &= - \sum_{s_{ij} \in \mathcal{S}} \log(p(s_{ij}|\mathcal{B})) \\ &= - \sum_{s_{ij} \in \mathcal{S}} (s_{ij} \cdot \log(\sigma(\Omega_{ij})) + (1 - s_{ij}) \cdot \log(1 - \sigma(\Omega_{ij}))). \end{aligned} \quad (4)$$

As  $\sigma(\Omega_{ij}) = \frac{1}{1+e^{-\Omega_{ij}}}$ , the Eq. (4) can be rewritten as

$$\mathcal{L}_p = \sum_{s_{ij} \in \mathcal{S}} (\log(1 + e^{\Omega_{ij}}) - s_{ij} \cdot \Omega_{ij}). \quad (5)$$

It is very challenging to directly optimize this discrete optimization problem, as the binary constraint  $b_i \in \{-1, 1\}^M$  requires thresholding on the network outputs which may result in the vanishing-gradient problem in backpropagation. We adopt the continuous relaxation strategy [24], [25] to solve this problem. The output of deep hashing layer  $u_i$  is fed to a tanh function  $h_i = \tanh(u_i)$ , which is used as a substitute for binary code  $b_i$ . Thus,  $\Omega_{ij}$  is redefined as  $h_i^T h_j$ .

For more efficient and faster hash learning, we design a focal quantization loss to mitigate the divergence between the discrete binary codes and the continuous output of hashing networks, inspired by [80]. Since the gradient accumulations of a large number of simple samples are not helpful for training, the focal loss attempts to reduce the weights of simple examples to promote the training process. First, we convert the binary code quantization problem into a binary classification problem. We use a sigmoid activation to map the outputs of the hash layer into a probability distribution  $\hat{p}_i = \sigma(u_i)$ . Notice that, tanh and sigmoid are both monotonic increasing functions that hold the same variation trend, *i.e.*, when  $h_i$  asymptotically approaches to -1,  $p_i$  also approaches 0, and vice versa (both approach to 1). Thus the probability of binary classification can reflect the compactness of hash codes effectively.

The focal quantization loss is defined as

$$\mathcal{L}_q = -\frac{1}{N} \sum_{i \in N} \sum_{j \in M} (\hat{y}_{ij} \cdot (1 - \hat{p}_{ij})^\alpha \cdot \log(\hat{p}_{ij}) + (1 - \hat{y}_{ij}) \cdot (\hat{p}_{ij})^\alpha \cdot \log(1 - \hat{p}_{ij})), \quad (6)$$

where  $\hat{y}_i$  is a label indicator that indicates which class (0 or 1) the output of hash layer should be classified as. We adopt a weighted sigmoid function to achieve such effect, *i.e.*,  $\hat{y}_i = \sigma(\beta \cdot u_i)$ ,  $\beta$  is a parameter far greater than 1.

By integrating the pairwise similarity loss and focal quantization loss, the overall hashing loss can be defined as

$$\mathcal{L}_{hash} = \mathcal{L}_p + \mathcal{L}_q. \quad (7)$$

## IV. EXPERIMENTS

### A. Datasets

To verify the performance of the proposed method, we compare the proposed method with several baselines on three widely used multi-label image datasets, *i.e.*, NUS-WIDE, VOC2012, and COCO.

**NUS-WIDE** [81] is a dataset containing 269,648 public web images. Each image is annotated with one or more class labels from a total of 81 classes. There exists a widely used subset of images associated with the 21 most common labels and each label associated with at least 5,000 images, resulting in a total of 195,834 images.

**VOC2012** [82] is a widely used dataset for object detection and segmentation, which contains 17,125 images. Each image is associated with at least one of the 20 semantic labels.

**COCO** [14] is a dataset for object detection, semantic scene labeling, and indexing, which contains 123,287 images with semantic labels. Each image is associated with one to sixteen labels from a total of 90.

### B. Implementation Details

To construct a zero-shot scenario, we should further split the dataset<sup>1</sup>. Since there are more complex semantic relationships among these multi-label datasets, we use one of these three

image datasets as source data and one as target data. For example, we can set NUS-WIDE as source data and VOC2012 as target data, and vice versa. Before training models based on these datasets, data preprocessing must be done. We set three experiments, including one between NUS-WIDE and VOC2012, one between NUS-WIDE and COCO, and the last one between COCO and VOC2012.

1) *Experiment between NUS-WIDE and VOC2012:* In NUS-WIDE, we remove the common concepts (semantic labels) shared by these two datasets and related images, because there are much more images in NUS-WIDE than in VOC2012. In VOC2012, we remove several ambiguous concepts and related images. Such data-clean operations result in a subset of NUS-WIDE containing 106,389 images and 18 labels, and a subset of VOC2012 containing 16,750 images and 17 labels. For NUS-WIDE, we randomly select 10,000 images as the training set, 2,000 images as the test query set, and the rest as the retrieval database. For VOC2012, we randomly select 4,000 images as the training set, 1,000 images as the test query set, and the rest as the retrieval database.

2) *Experiment between NUS-WIDE and COCO:* We remove the common concepts and relative images from NUS-WIDE and keep COCO unchanged. Finally, a subset of NUS-WIDE containing 100,303 images and 17 labels, and a subset of COCO containing 123,274 images and 80 labels are prepared for the following experiments. For both datasets, we randomly select 10,000 images as the training set, 2,000 images as the test query set, and the rest as the retrieval database.

3) *Experiment between VOC2012 and COCO:* We remove the images related to three ambiguous concepts for VOC2012. After that, all the concepts of VOC2012 are included in COCO. Then, for COCO, we remove the images that have the concepts in VOC2012. Finally, there remains 21,987 images and 60 labels for COCO, and 16,750 images and 17 labels for VOC2012. For both datasets, we randomly select 4,000 images for training, 1,000 images as the test query, and the rest as the retrieval database.

For NUS-WIDE, VOC2012, and COCO, we check the labels and ensure that the training/query set contains all the concepts of the corresponding dataset.

We implement the proposed method (T-MLZSH) using the TensorFlow toolkit. In this paper, we use AlexNet as the backbone CNN. To validate the versatility of the proposed framework, we will also evaluate by replacing the backbone CNN with VGG16 and ResNet50. We use the pre-trained model to initialize the weight parameters and focus on training the hashing layer and embedding layer. Adam method is adopted for stochastic optimization with a mini-batch size of 128, and all input images are resized to  $227 \times 227$ .

We compare our method (T-MLZSH) with nine other methods, including the traditional methods KSH [2], IMH [3], SDH [4], ZSH-DA [32], and ZSH [30], and the deep learning-based methods DHN [24], Hashnet [26], ADSH [83], and TZSH [33]. Among the traditional methods, ZSH-DA and ZSH are two zero-shot hashing methods. Among the deep learning-based methods, TZSH is a transductive zero-shot hashing

<sup>1</sup><https://github.com/qinnzou/Zero-Shot-Hashing>

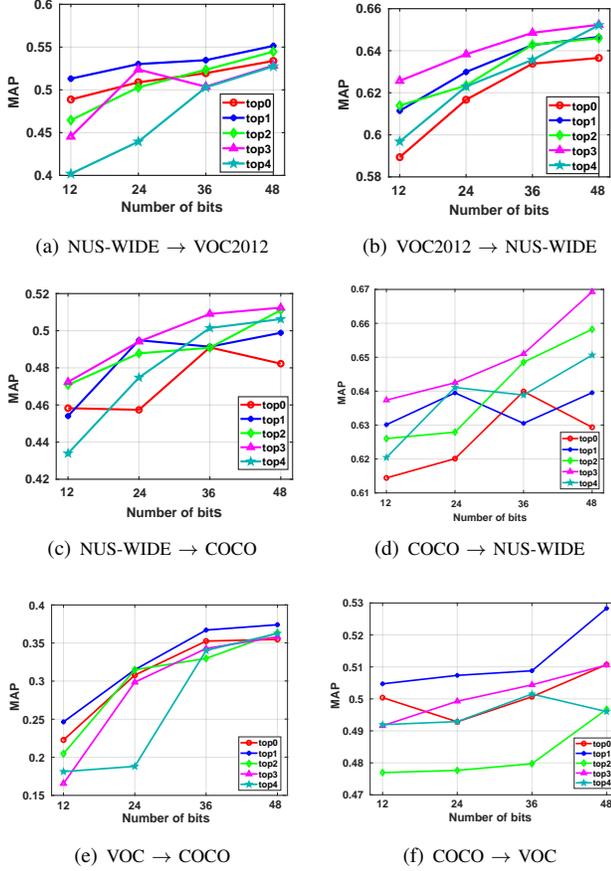


Fig. 4: Performance obtained by using different numbers of predicted labels for target data with hash codes of 12, 24, 36, and 48 bits, respectively. The values are computed based on the top-1000 retrieved images.

method. It is worth noting that, all the methods use the same training and test settings on the newly formed datasets.

In the training, for KSH, IMH, SDH, DHN, Hashnet, and ADSH, the images and labels of the training set are used, while for ZSH-DA, ZSH, TZSH, and the proposed T-MLZSH, in addition to the images and labels in the training set, the raw images (without labels) in the test query set are also used. Note that, ADSH is a supervised method which treats the query points and database points in an asymmetric way, and DIHN [84] successfully extends the asymmetric training strategy in an incremental learning framework. As DIHN uses the labels of the unseen test query data for training, which violates the protocols of our experiments, we do not include it in the comparison.

For deep learning-based methods, we use the raw images as input. For the non-CNN approaches, we use the outputs of the fc7 layer in AlexNet as their visual features. All the experiments are conducted on an NVIDIA TITAN Xp GPU and the Ubuntu 16.04 operating system.

### C. Metrics

The metrics we used to evaluate the image retrieval quality are four widely-used metrics: Average Cumulative Gains

(ACG), Normalized Discounted Cumulative Gains (NDCG), Mean Average Precision (MAP), and Weighted Mean Average Precision (WAP) [7], [23], [85], [86].

MAP is the mean of average precision for each query, which can be calculated by

$$MAP = \frac{1}{Q} \sum_Q AP(q) \quad (8)$$

where

$$AP(q) = \frac{1}{N_{Tr}(q)@n} \sum_n i(Tr(q, i) \frac{N_{Tr}(q)@i}{i}) \quad (9)$$

$Tr(q, i) \in \{0, 1\}$  is an indicator function that if  $I_q$  and  $I_i$  have same class labels,  $Tr(q, i) = 1$ ; otherwise  $Tr(q, i) = 0$ .  $Q$  is the number of query sets and  $N_{Tr}(q)@i$  indicates the number of relevant images w.r.t the query image  $I_q$  within the top  $i$  images. ACG represents the average number of shared labels between the query image and the top  $n$  retrieved images. For a given query image  $I_q$ , the ACG score of the top  $n$  retrieved images is calculated by

$$ACG@n = \frac{1}{n} \sum_n^i C(q, i) \quad (10)$$

where  $n$  denotes the number of top retrieval images and  $C(q, i)$  is the number of shared class labels between  $I_q$  and  $I_i$ . NDCG is a popular evaluation metric in information retrieval. Given a query image  $I_q$ , the DCG score of top  $n$  retrieved images is defined as

$$DCG@n = \sum_n^i \frac{2^{C(q, i)} - 1}{\log(1 + i)} \quad (11)$$

Then, the normalized DCG(NDCG) score at the position  $n$  can be calculated by  $NDCG@n = \frac{DCG@n}{Z_n}$ , where  $Z_n$  is the maximum value of DCG@ $n$ , which constrains the value of NDCG in the range [0,1].

WAP is similar to MAP, the only difference is that WAP is the average value of ACG scores at each top  $n$  retrieval image rather than average precision. WAP can be calculated by

$$WAP = \frac{1}{Q} \sum_Q \left( \frac{1}{N_{Tr}(q)@n} \sum_i^n (Tr(q, i) \times ACG@i) \right) \quad (12)$$

### D. Overall Performance

In this part, we analyze the retrieval results all evaluated on the unseen target data.

1) *The number of predicted unseen categories:* Figure 4 displays the results of using different numbers of predicted labels on the target data. Top- $k$  indicates that the first  $k$  categories in the correlation-score ranking list are used as predicted labels and top-0 means that the labels of target data are set to a vector of all zeros. From Fig. 4(a) and (b), we can see that when setting VOC2012 as target data, using top-1 predicted labels can achieve the best performance. The possible reason is that the average number of objects in each image on VOC2012 is relatively small. With more predicted labels used for supervised hashing learning, it will inevitably

TABLE I: MAP results obtained by using different numbers of bits on (NUS-WIDE, VOC2012). The MAPs are computed based on the top-1000 retrieved images.

Methods	NUS-WIDE $\rightarrow$ VOC2012				VOC2012 $\rightarrow$ NUS-WIDE			
	12-bit	24-bit	36-bit	48-bit	12-bit	24-bit	36-bit	48-bit
KSH [2]	0.4033	0.4079	0.4153	0.4181	0.5476	0.5510	0.5602	0.5670
SDH [4]	0.4097	0.4010	0.4087	0.4095	0.5327	0.5345	0.5368	0.5395
IMH [3]	0.4083	0.4346	0.4320	0.4297	0.5616	0.5723	0.5718	0.5710
DHN [24]	0.4171	0.4282	0.4362	0.4395	0.5664	0.5739	0.5726	0.5688
Hashnet [26]	0.4048	0.4288	0.4349	0.4465	0.5674	0.5940	0.6073	0.6336
ADSH [83]	0.3988	0.4163	0.4060	0.4201	0.5315	0.5473	0.5392	0.5853
ZSH-DA [32]	0.3592	0.3618	0.3770	0.3596	0.5132	0.5166	0.5212	0.5191
ZSH [30]	0.3968	0.4055	0.4111	0.4296	0.5340	0.5566	0.5589	0.5500
TZSH [33]	0.4413	0.4683	0.4644	0.4753	0.5736	0.5805	0.5919	0.5896
<b>T-MLZSH</b>	<b>0.4808</b>	<b>0.4884</b>	<b>0.4894</b>	<b>0.5037</b>	<b>0.6106</b>	<b>0.6131</b>	<b>0.6149</b>	<b>0.6200</b>

TABLE II: MAP results obtained by using different numbers of bits on (NUS-WIDE, COCO). The MAPs are computed based on the top-1000 retrieved images.

Methods	NUS-WIDE $\rightarrow$ COCO				COCO $\rightarrow$ NUS-WIDE			
	12-bit	24-bit	36-bit	48-bit	12-bit	24-bit	36-bit	48-bit
KSH [2]	0.3948	0.4069	0.4113	0.4143	0.5948	0.6167	0.6191	0.6224
SDH [4]	0.3782	0.3917	0.3971	0.4050	0.5681	0.5954	0.6102	0.6051
IMH [3]	0.3905	0.4021	0.4114	0.4188	0.5983	0.5961	0.6098	0.6132
DHN [24]	0.4250	0.4325	0.4529	0.4487	0.6177	0.6421	0.6466	0.6559
Hashnet [26]	0.3908	0.4018	0.4129	0.4439	0.5538	0.5674	0.5830	0.5932
ADSH [83]	0.3783	0.3890	0.3968	0.4056	0.5773	0.5906	0.5871	0.6219
ZSH-DA [32]	0.3597	0.3592	0.3772	0.3744	0.5256	0.5220	0.5230	0.5247
ZSH [30]	0.3832	0.4091	0.4109	0.4286	0.5708	0.5727	0.5753	0.5782
TZSH [33]	0.4436	0.4585	0.4660	0.4800	0.5933	0.6368	0.6070	0.6336
<b>T-MLZSH</b>	<b>0.4724</b>	<b>0.4941</b>	<b>0.5090</b>	<b>0.5124</b>	<b>0.6374</b>	<b>0.6425</b>	<b>0.6510</b>	<b>0.6693</b>

incur misleading information and cause performance degradation. When setting NUS-WIDE as target data, the best results are obtained by using top-3 predicted labels. In the following experiments, we use top-1 and top-3 predicted labels as supervised information for the proposed method in default for VOC2012 and NUS-WIDE, respectively. In Fig. 4(c)-(f), we can see that the best performance can be achieved by using top-3 predicted labels for the experiments on COCO and NUS-WIDE, and using top-1 for experiments on COCO and NUS-WIDE. It may be because that the average number of objects in each image in the (NUS-WIDE, COCO) case is 2.48 and 2.97, and that in the (VOC, COCO) case is 1.30 and 1.48, respectively. Hence, we use top-3 predicted labels as supervised information for the proposed method in default for (COCO, NUS-WIDE) and top-1 for (VOC, COCO) in the following experiments.

2) *Results under different scenarios*: The MAP results of the proposed method and the comparison methods are shown in Table I, II, and III. It can be seen that the non-zero-shot deep hashing methods DHN, Hashnet, and ASDH outperform the traditional hashing methods KSH, IMH, SDH, ZSH-DA, and ZSH. Among the non-zero-shot deep hashing methods, ADSH has lower performance than the other two. The possible reason is that it treats the query points and database points in an asymmetric way which drives the network to pay more

attention to the seen data. Thus, when the dataset is replaced by the unseen dataset, ADSH's performance will be affected. The two traditional zero-shot hashing methods ZSH and ZSH-DA achieve low performance on the multi-label datasets, which indicates that the complex semantics of multi-label images are too hard to be modeled by learning a one-to-one semantic representation. The zero-shot deep hashing methods T-MLZSH and TZSH obtain significantly higher performance than the other methods.

It can also be seen that the proposed T-MLZSH outperforms almost all the comparison methods significantly on all target datasets under different scenarios. On NUS-WIDE and VOC2012, the transductive zero-shot hashing methods, *i.e.*, TZSH and T-MLZSH, achieve higher MAP values than other methods, as shown in Table I. Compared to TZSH, T-MLZSH achieves increments of about 3.1% and 2.8% in the average MAP for different bits on NUS-WIDE and VOC2012, respectively. The possible reason is that TZSH adopts a strategy only utilizing partially-selected target data for hash learning, which limits its performance.

From Table II we can see, compared to TZSH, T-MLZSH achieves increments of about 3.4% or 3.2% in the average MAP for different bits on COCO $\rightarrow$ NUS-WIDE and NUS-WIDE $\rightarrow$ COCO, respectively. In this experiment, the deep supervised hashing method DHN is found to perform better

than TZSH on COCO→NUS-WIDE. The possible reason is that COCO is categorized into more categories and DHN can get more detailed supervised information when setting COCO as the training set. Nevertheless, T-MLZSH still outperforms DHN by about 0.95% in the average MAP.

Table III shows the results on (COCO, VOC2012). Compared to TZSH, T-MLZSH achieves increments of about 6.8% or 1.8% on average MAP for different bits. The proposed method is found to obtain a significantly better performance on VOC2012→COCO, where COCO is unseen. The possible reason is that, although all the instances of VOC2012 are included in COCO, some unlabeled targets in COCO images are similar to that in VOC2012. Features extracted from VOC2012 images may possess high similarity with some COCO images. Since the scale of COCO is much larger than VOC2012, the overall performance is lower when COCO is unseen.

More results in the other three metrics, *i.e.*, ACG, NDCG, and WAP, are presented in Fig. 5. We can see that the overall trends of the performance on the three metrics are consistent. When the code length increases, the performance improves. According to the definition, these three metrics can make a more fair evaluation on multi-label images, as the numbers of shared labels between images are considered. For MAP, the pairwise images that share at least one common object label will be considered as relevant images, and no more comparisons of fine-grained semantic relation between these images are included, which may not stay in step with user demand in multi-label image retrieval. For WAP, the average number of shared class labels among these retrieved similar images is considered. Higher WAP means more high-quality retrieval results having more shared class labels in the nearest search. Although the range of WAP on different datasets would be different, the WAPs of T-MLZSH are stably higher than that of the comparison methods.

Figure 6-10 display more detailed results on ACG, NDCG, and WAP of different numbers of top returned images under different scenarios. Generally, the NDCG curves show a trend of first decreasing and then increasing while ACG and WAP decrease with the increase of retrieved samples. The more samples retrieved, the more low-quality retrieval results having fewer shared class labels will appear in the rear. This will lead to the trends that tend to be gentle. Since most of the compared methods are not specific for zero-shot learning, the results of unseen data may be uncertain. The curve of ZSH-DA is unsmooth. The possible reason is that, while using ZSH-DA for multi-label unseen images, the complex semantics are too hard to be modeled by learning a one-to-one semantic representation, which leads to unstable performances.

3) *Comparison with different backbone CNN blocks:* To justify the versatility of the proposed deep hashing framework, We replace the backbone CNN with VGG16 and ResNet50, both of which achieve more accurate results than AlexNet on the ImageNet competition. We denote these two modifications as ‘T-MLZSH-VGG16’ and ‘T-MLZSH-ResNet50’, respectively. The results are shown in Table IV. We can see from the results that, with more powerful backbone CNN blocks, T-MLZSH generally achieves higher performance on

all these metrics. It simply indicates a good transfer capability and versatility of the proposed deep hashing framework.

### E. Functional Analysis

1) *Influence of the categories of datasets:* In the above two groups of experiments, part of them used NUS-WIDE as the unseen dataset and the results are shown in Table I and Table II. From the left of the tables, we notice that using COCO as seen dataset can achieve a better MAP result, which has an improvement of about 2.6%, 2.9%, 3.6%, and 4.9% in average MAP with different hash bits respectively. These two groups of experiments have the same target domain and the only difference is the source domain. We guess that the possible reason leading to different MAPs is the difference of categories. COCO dataset is more finely divided and more semantic information can be used which make the network much stronger.

2) *Influence of the size of datasets:* Moreover, we explore the influence of the sizes of the source and target datasets by using different amounts of ‘seen’ and ‘unseen’ images to train the model. Three orders of magnitude are considered. Exactly, the number of images from source dataset and target dataset are 10,000 and 10,000, 10,000 and 4,000, 4,000 and 10,000, respectively. The results are shown in Fig. 12. It can be seen that there is a slight difference in the use of different orders of magnitude, but the performance is stable on the whole. It manifests that the proposed model has certain stability even if the number of images from two domains used for training varies from each other.

3) *Necessary of quantization loss:* We also explore the effectiveness of the proposed quantization loss. We compare the proposed method with its variant versions: one adopts the widely used absolute error loss that measures the Euclidean distance between continuous outputs and discrete codes directly, and the other does not use quantization loss. The results are presented in Fig. 13. It can be seen that, without quantization loss, there is a rapid degradation of the performance. The difference in the evaluation index of the MAP is about 0.5%, which illustrates the importance of using quantization loss in deep hashing learning. We can also see that, applying quantified losses has greatly improved the results, the proposed focal quantization loss has a much more advanced performance among all the proposed architecture.

4) *Effectiveness of the visual-semantic bridge:* The visual-semantic bridge is built to predict labels for unseen data. It helps the images belonging to the same category be with higher similarity by transferring the knowledge from the semantic representations to visual features. Since the labels of unseen images are unknown, the predicted labels will be represented as meaningless digital codes for the unseen images.

To validate the effectiveness of the visual-semantic bridge in linking the semantic representations and visual features, we use the proposed model to predict labels for the seen images and compare them to true labels. In the experiment, we use the model trained on the training set of VOC2012 to predict the labels for 10,000 images from the database set. Since the

TABLE III: MAP results obtained by using different numbers of bits on (VOC2012, COCO). The MAPs are computed based on the top-1000 retrieved images.

Methods	VOC2012 $\rightarrow$ COCO				COCO $\rightarrow$ VOC2012			
	12-bit	24-bit	36-bit	48-bit	12-bit	24-bit	36-bit	48-bit
KSH [2]	0.1763	0.2076	0.2346	0.2490	0.4498	0.4781	0.4912	0.5005
SDH [4]	0.1444	0.1762	0.1871	0.1991	0.4656	0.4869	0.5045	0.5060
IMH [3]	0.2005	0.2151	0.2279	0.2317	0.4083	0.4291	0.4347	0.4380
DHN [24]	0.1635	0.1826	0.1954	0.2351	0.4311	0.4370	0.4573	0.4597
Hashnet [26]	0.1521	0.2029	0.2470	0.2585	0.4331	0.4605	0.4766	0.4924
ADSH [83]	0.1368	0.1472	0.1771	0.1884	0.4310	0.4794	0.5018	0.5061
ZSH-DA [32]	0.0853	0.0746	0.0709	0.1104	0.3924	0.4189	0.4247	0.4393
ZSH [30]	0.1602	0.1935	0.2249	0.2227	0.4658	0.4775	0.5003	0.4981
TZSH [33]	0.2336	0.2632	0.2690	0.2631	0.4967	0.5053	0.5128	0.5079
<b>T-MLZSH</b>	<b>0.2465</b>	<b>0.3149</b>	<b>0.3670</b>	<b>0.3740</b>	<b>0.5047</b>	<b>0.5276</b>	<b>0.5296</b>	<b>0.5345</b>

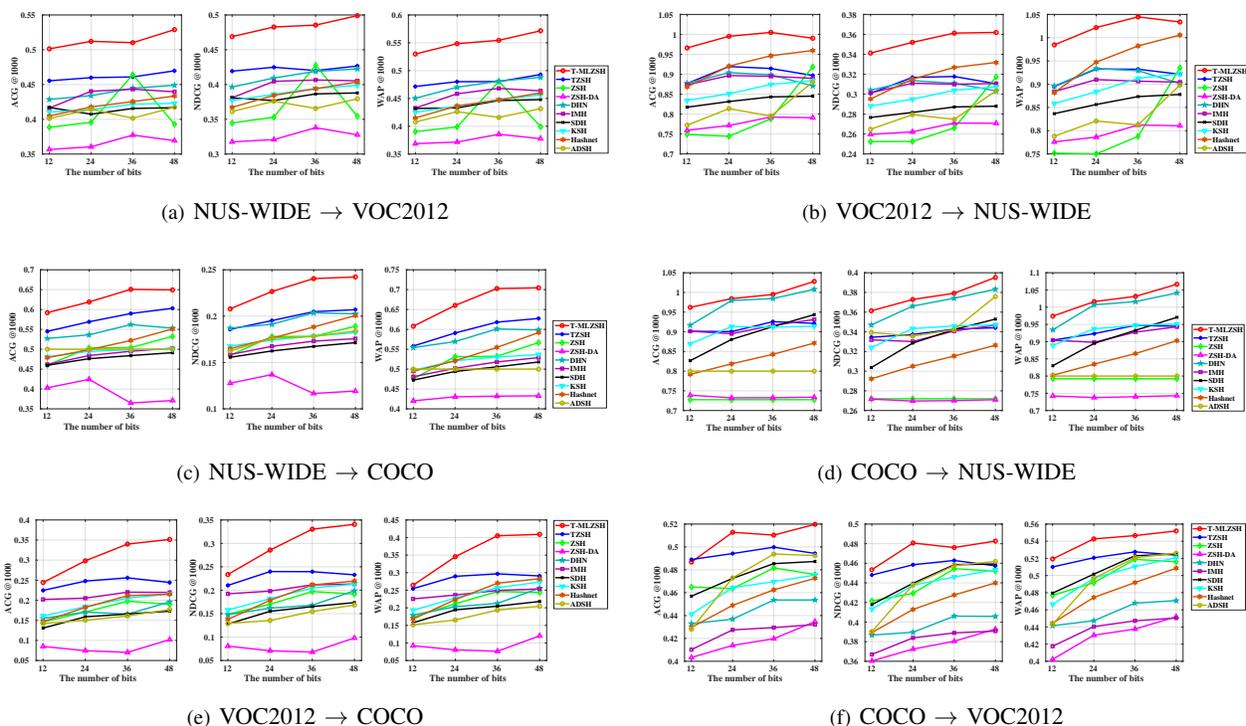


Fig. 5: Comparison of results obtained by the proposed method on ACG, NDCG, and WAP with 12-, 24-, 36-, and 48-bits hash codes, with regarding the top-1000 retrieved images.

average number of objects in each image from VOC2012 is 1.301, we predict one label for each image. The distribution of labels is visualized by t-SNE, as shown in Fig. 14. Comparing Fig. 14(a) and (b), we can find that the overall distributions are similar for the true label and the predicted label. However, because one image may contain multiple labels but is predicted with only one, the area of each predicted label category is observed to be a little smaller than that of the true label category. Meanwhile, we count the correctness of the predicted labels. The predicted labels for 8,283 images fall within their true labels. These experimental results show that the proposed visual-semantic bridge has a high performance.

5) *Top retrieval results:* Figure 15 shows some retrieval samples of some typical hashing methods according to the ascending Hamming ranking. The query image contains three semantic labels, *i.e.*, building, sky, and water, with the main content of a building. We mark the mismatched image with the red box from the perspective of human vision. The retrieval results of T-MLZSH are better visual plausible while focusing on the main object of the query image, while other compared methods may return some mismatched results like the forest, or return some images related to the less important part of the query image with higher ranking orders.

6) *Running efficiency:* We do the image retrieval by returning the top 1,000 similar images from 16,900 images based

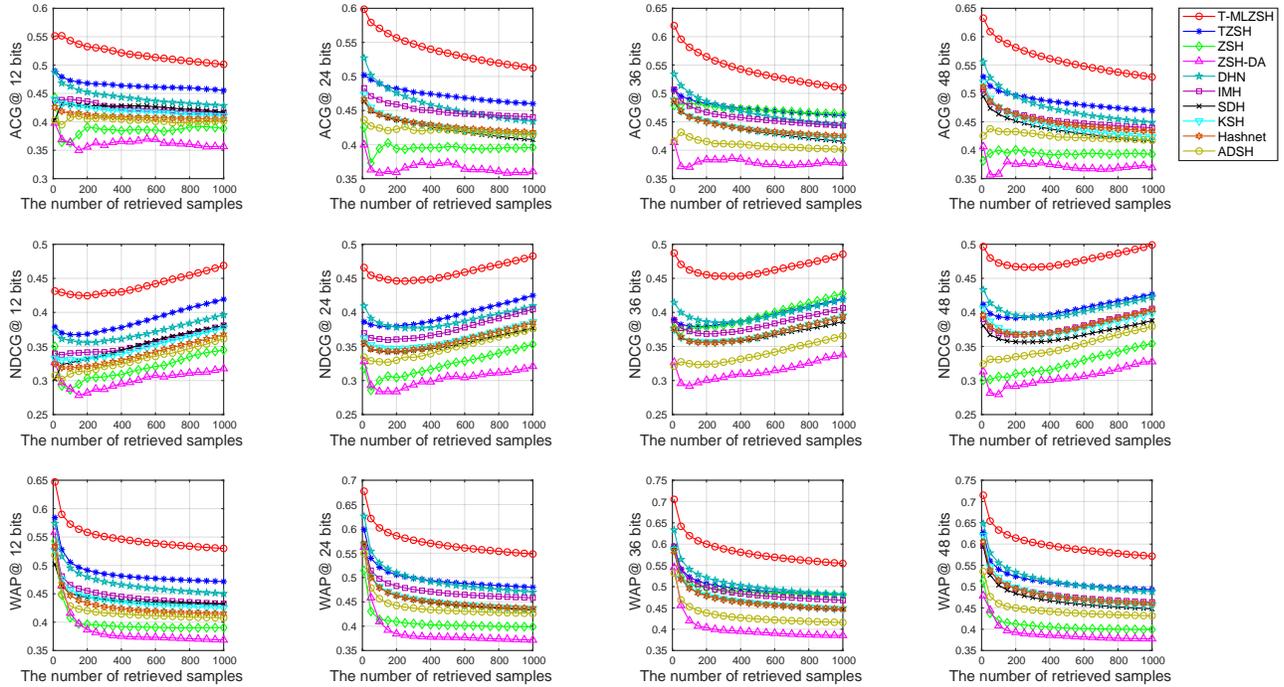


Fig. 6: Performance comparison on NUS-WIDE→VOC2012. The VOC2012 dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

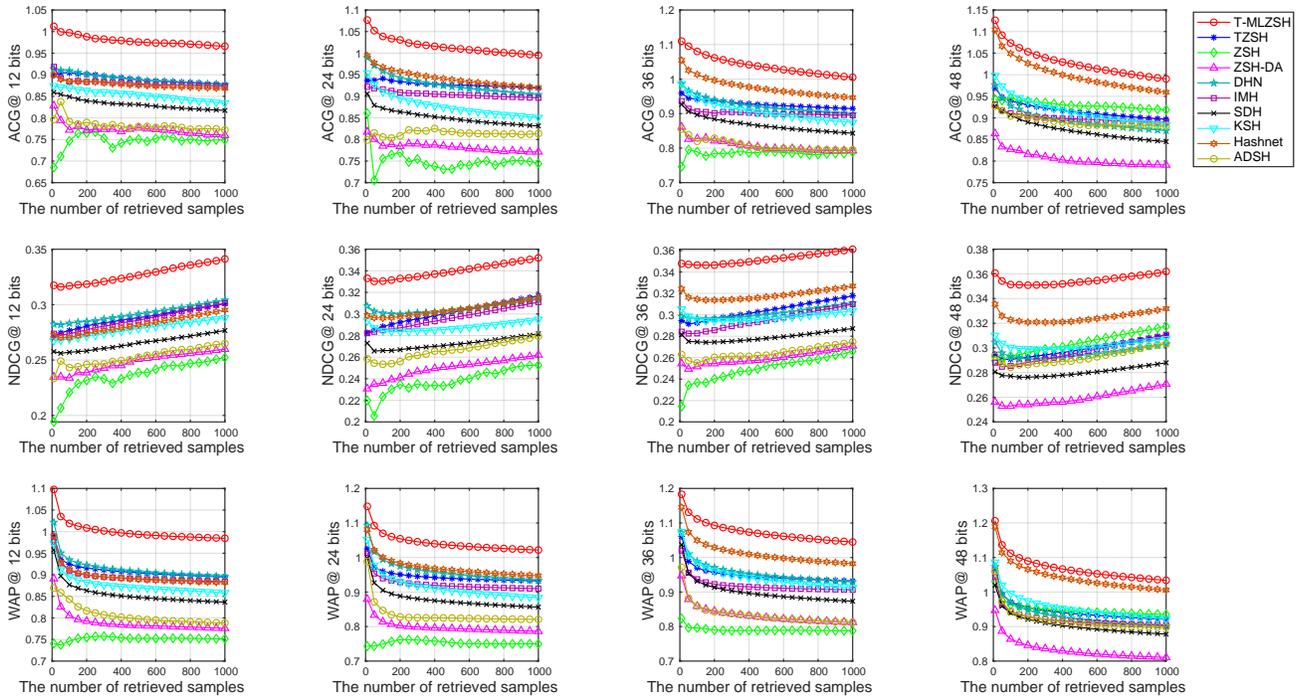


Fig. 7: Performance comparison on VOC2012→NUS-WIDE. The NUS-WIDE dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

on the trained models. The running time includes calculating the hash codes and calculating the hamming distances. In

our experiments, when calculating the 36-bit hash codes of 16,901 images, it takes 17.4616s, 48.2418s and 34.6509s

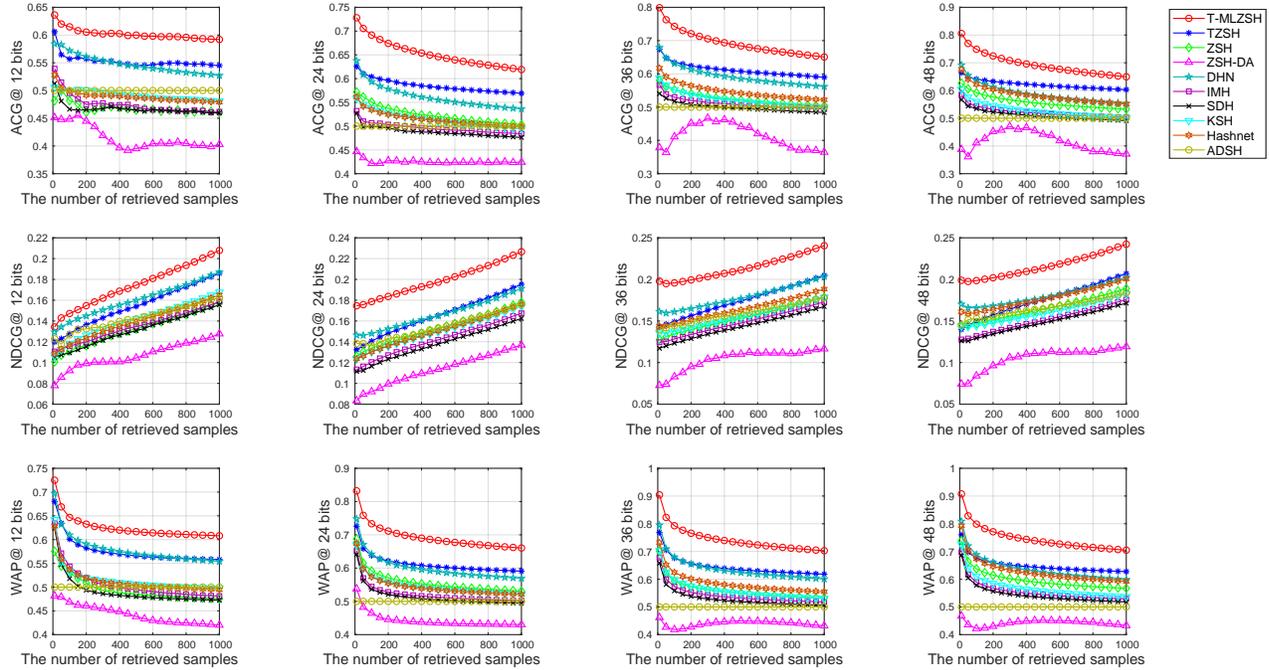


Fig. 8: Performance comparison on NUS-WIDE→COCO. The COCO dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

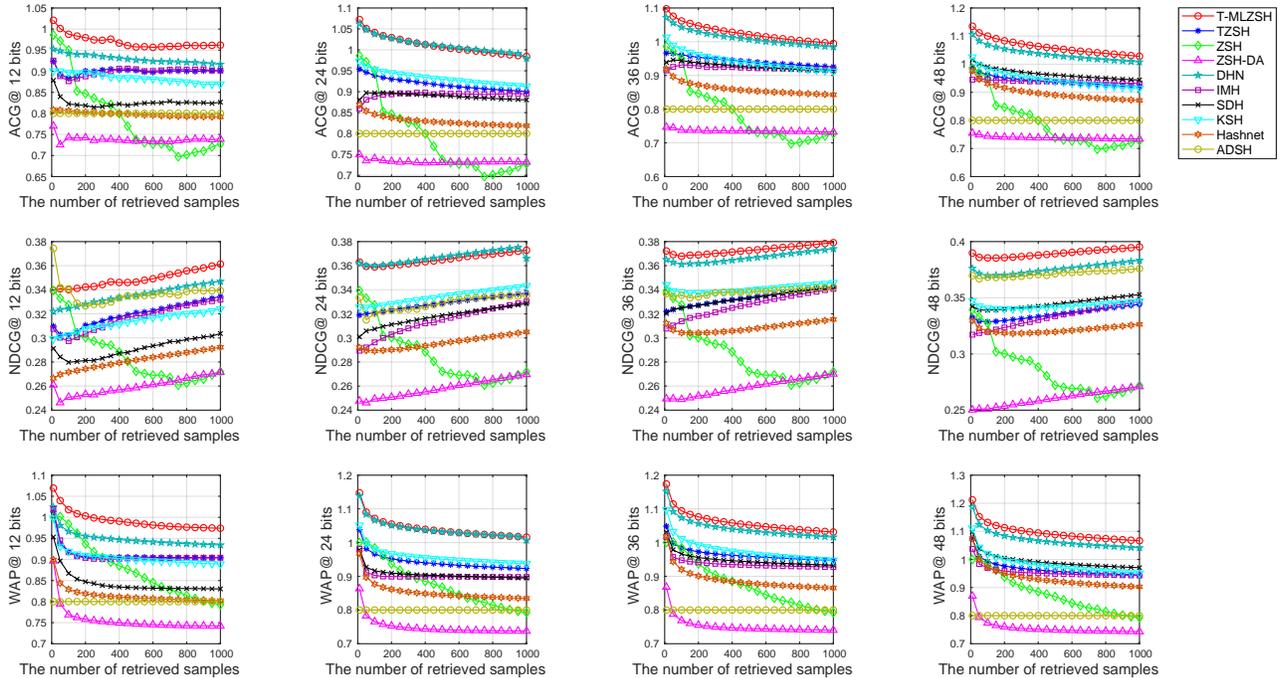


Fig. 9: Performance comparison on COCO→NUS-WIDE. The NUS-WIDE dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

for T-MLZSH-AlexNet, T-MLZSH-VGG16 and T-MLZSH-ResNet50, respectively. That is, it will take about 1.03ms, 2.85ms, and 2.05ms for the three methods to calculate the

hash codes of a new image. The hash codes are computed on an NVIDIA TITAN Xp GPU. The GPU memory usage is about 2.5GB, 4.3GB, and 1.2GB for the three models in

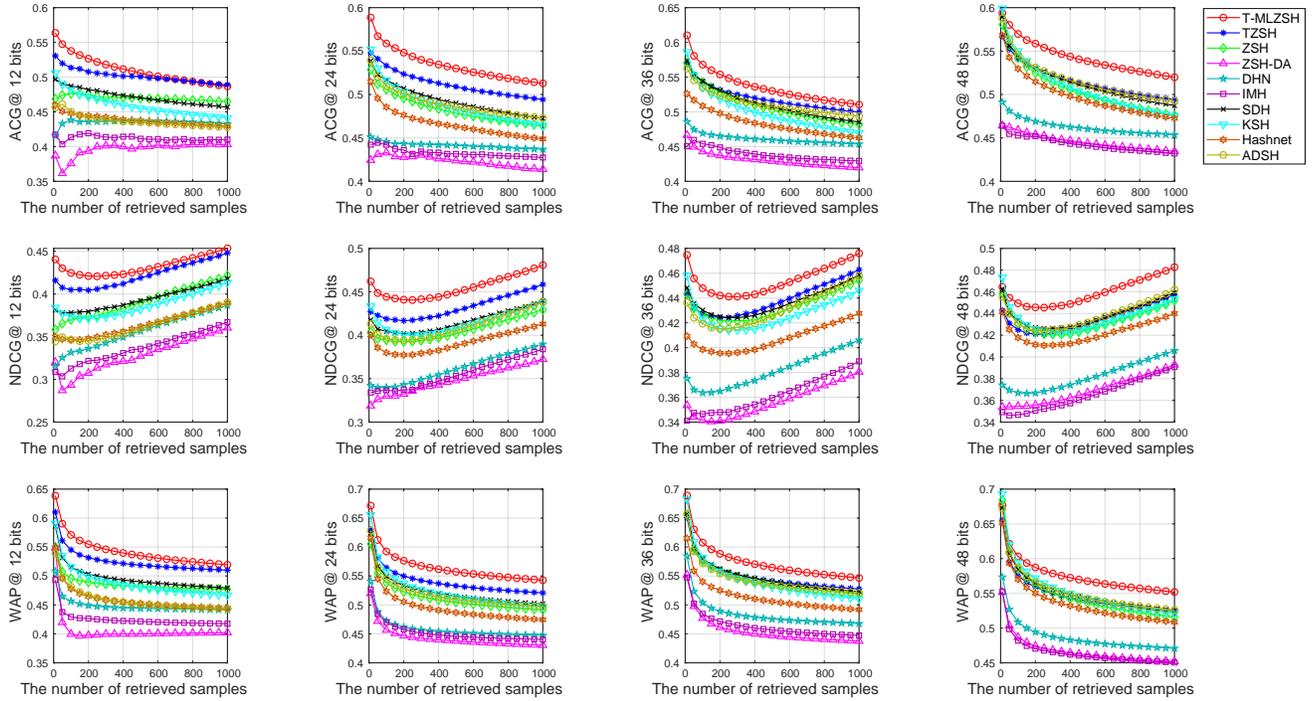


Fig. 10: Performance comparison on COCO→VOC2012. The VOC2012 dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

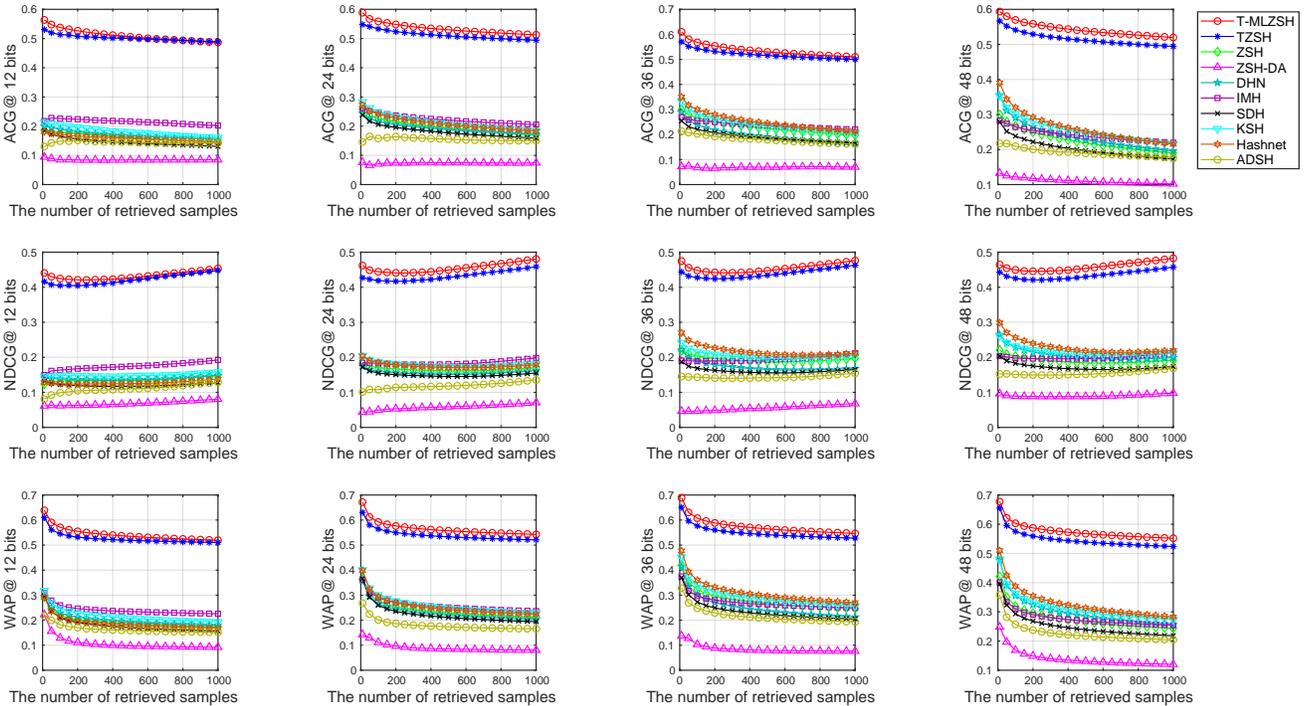


Fig. 11: Performance comparison on VOC2012→COCO. The COCO dataset is unseen. From top to bottom, there are ACG, NDCG, and WAP w.r.t. different top returned samples with hash codes of 12, 24, 36, and 48 bits, respectively.

the retrieval test. For calculating the hamming distance of 36-bit hash codes, the running time on the 16,900 images is 0.006742s, which is about  $4 \times 10^{-7}$ s for one calculation.

The hamming distance is computed by a 2.0GHz core of an Intel(R) Xeon(R) E5-2620 CPU.

TABLE IV: Performance under different backbone CNN blocks. The values are computed on the top-1000 retrieved images.

Metrics	MAP				NDCG				ACG				WAP			
	12b	24b	36b	48b												
NUS-WIDE->VOC2012																
T-MLZSH-AlexNet	0.4808	0.4884	0.4894	0.5037	0.4689	0.4827	0.4855	0.4990	0.5298	0.5484	0.5545	0.5716	0.5298	0.5484	0.5545	0.5716
T-MLZSH-VGG16	0.5340	0.5346	0.5843	0.6149	0.6363	0.6270	0.6520	0.6782	0.5299	0.5124	0.5546	0.5798	0.5504	0.5521	0.6065	0.6379
T-MLZSH-ResNet50	<b>0.5711</b>	<b>0.6054</b>	<b>0.6518</b>	<b>0.6541</b>	<b>0.6471</b>	<b>0.6820</b>	<b>0.6915</b>	<b>0.7377</b>	<b>0.5505</b>	<b>0.5761</b>	<b>0.5850</b>	<b>0.6171</b>	<b>0.5902</b>	<b>0.6274</b>	<b>0.6739</b>	<b>0.6791</b>
VOC2012->NUS-WIDE																
T-MLZSH-AlexNet	0.6106	0.6131	0.6149	0.6200	0.3412	0.3520	0.3613	0.3619	0.9660	1.0207	1.0448	1.0338	0.9845	1.0217	1.0448	1.0338
T-MLZSH-VGG16	0.6049	0.6503	0.6635	0.6799	0.3280	0.3670	0.3657	0.3885	0.9175	1.0104	1.0238	1.0587	0.9315	1.0361	1.0597	1.1040
T-MLZSH-ResNet50	<b>0.6542</b>	<b>0.6715</b>	<b>0.6957</b>	<b>0.6977</b>	<b>0.3677</b>	<b>0.3836</b>	<b>0.4056</b>	<b>0.4037</b>	<b>1.0267</b>	<b>1.0599</b>	<b>1.1146</b>	<b>1.0882</b>	<b>1.0485</b>	<b>1.0858</b>	<b>1.1373</b>	<b>1.1282</b>
NUS-WIDE->COCO																
T-MLZSH-AlexNet	0.4724	0.4941	0.5291	0.5124	0.2097	0.2296	0.2405	0.2423	0.5921	0.6191	0.6508	0.7046	0.6080	0.6604	0.6508	0.7046
T-MLZSH-VGG16	0.4970	<b>0.5705</b>	0.6031	0.6075	0.3014	<b>0.3683</b>	0.3978	0.3935	0.6504	<b>0.6724</b>	0.6889	0.7093	0.6374	<b>0.6609</b>	0.6657	0.6700
T-MLZSH-ResNet50	<b>0.4994</b>	0.5519	<b>0.6030</b>	<b>0.6209</b>	<b>0.3022</b>	0.3639	<b>0.3983</b>	<b>0.4112</b>	<b>0.6543</b>	0.6698	<b>0.6934</b>	<b>0.7102</b>	<b>0.6309</b>	0.6532	<b>0.6694</b>	<b>0.6741</b>
COCO->NUS-WIDE																
T-MLZSH-AlexNet	0.6374	0.6425	0.6510	0.6693	0.3614	0.3728	0.3791	0.3953	0.9621	0.9014	0.9948	1.0280	0.9741	1.0162	1.0314	1.0665
T-MLZSH-VGG16	0.6333	0.6686	0.6709	0.6858	0.4380	0.4808	0.4798	0.4985	0.9718	0.9796	1.0297	1.1145	0.9932	1.0345	1.1136	1.1203
T-MLZSH-ResNet50	<b>0.6588</b>	<b>0.6843</b>	<b>0.6883</b>	<b>0.6914</b>	<b>0.4713</b>	<b>0.4925</b>	<b>0.5044</b>	<b>0.5020</b>	<b>0.9821</b>	<b>0.9902</b>	<b>1.1342</b>	<b>1.1201</b>	<b>1.0023</b>	<b>1.0962</b>	<b>1.1361</b>	<b>1.1903</b>
VOC2012->COCO																
T-MLZSH-AlexNet	0.2465	0.3149	0.3670	0.3740	0.2336	0.2857	0.3304	0.3406	0.2443	0.2983	0.3399	0.3515	0.2642	0.3455	0.4050	0.4097
T-MLZSH-VGG16	0.3707	0.4183	0.4757	0.4901	0.3631	0.3890	0.4430	0.4651	0.4033	0.4504	0.5151	0.4295	0.3968	0.4234	0.4596	0.4775
T-MLZSH-ResNet50	<b>0.4161</b>	<b>0.4695</b>	<b>0.5184</b>	<b>0.5320</b>	<b>0.3992</b>	<b>0.4372</b>	<b>0.4766</b>	<b>0.4921</b>	<b>0.4558</b>	<b>0.5052</b>	<b>0.5673</b>	<b>0.5804</b>	<b>0.4033</b>	<b>0.4504</b>	<b>0.5151</b>	<b>0.5295</b>
COCO->VOC1012																
T-MLZSH-AlexNet	0.5047	0.5074	0.5088	0.5283	0.4502	0.4567	0.4665	0.4860	0.4868	0.4791	0.4930	0.5134	0.5193	0.5244	0.5338	0.5555
T-MLZSH-VGG16	0.5583	0.6538	0.6652	<b>0.6844</b>	0.5131	0.6084	0.6100	<b>0.6314</b>	0.5761	<b>0.6804</b>	<b>0.6905</b>	<b>0.7146</b>	0.5506	0.6149	0.6224	0.6298
T-MLZSH-ResNet50	<b>0.6402</b>	<b>0.6672</b>	<b>0.6796</b>	0.6761	<b>0.6056</b>	<b>0.6148</b>	<b>0.6330</b>	0.6272	<b>0.6210</b>	0.6223	0.6327	0.6274	<b>0.6614</b>	<b>0.6939</b>	<b>0.7071</b>	<b>0.7027</b>

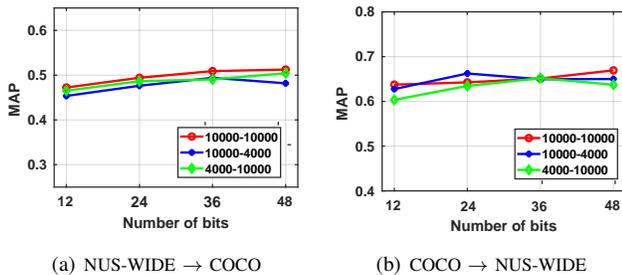


Fig. 12: Performance obtained by using different amounts of training data. The numbers indicate the training data in the form of ‘source - target’.

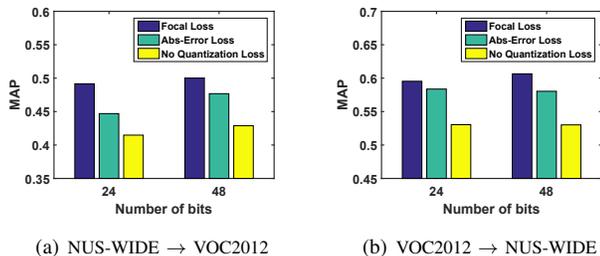


Fig. 13: Performance obtained by different quantization losses.

## V. CONCLUSION

In this paper, a novel transductive zero-shot hashing method was proposed for multi-label image retrieval. It utilized the instance-concept coherence to construct a bridge for connecting the seen and unseen labels. Based on these connections, several categories with the highest relatedness scores were selected as the predicted labels for target data. Then, hashing learning was performed on both the source data and target data in a supervised manner. Experimental results on three widely-used multi-label datasets showed that the proposed method outperformed state-of-the-art methods with a significant margin. Moreover, the superiority of the proposed focal

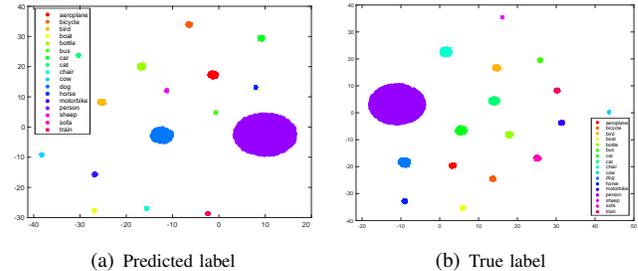


Fig. 14: Visualization of the predicted labels and the true labels using t-SNE.

loss was verified by ablation studies, the effectiveness of the visual-semantic bridge was demonstrated through feature visualization, and the high performance on image retrieval was illustrated with visual comparisons.

## ACKNOWLEDGEMENTS

The authors would like to thank Dr. Song Bai from the University of Oxford for helpful suggestions.

## REFERENCES

- [1] B. Kulis and K. Grauman, “Kernelized locality-sensitive hashing,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 6, pp. 1092–1104, 2011.
- [2] W. Liu, J. Wang, R. Ji, Y.-G. Jiang, and S.-F. Chang, “Supervised hashing with kernels,” in *CVPR*, 2012, pp. 2074–2081.
- [3] F. Shen, C. Shen, Q. Shi, A. Van Den Hengel, and Z. Tang, “Inductive hashing on manifolds,” in *CVPR*, 2013, pp. 1562–1569.
- [4] F. Shen, C. Shen, W. Liu, and H. Tao Shen, “Supervised discrete hashing,” in *CVPR*, 2015, pp. 37–45.
- [5] H. Lai, Y. Pan, Y. Liu, and S. Yan, “Simultaneous feature learning and hash coding with deep neural networks,” in *CVPR*, 2015, pp. 3270–3278.
- [6] R. Zhang, L. Lin, R. Zhang, W. Zuo, and L. Zhang, “Bit-scalable deep hashing with regularized similarity learning for image retrieval and person re-identification,” *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 4766–4779, 2015.
- [7] Z. Zhang, Q. Zou, Y. Lin, L. Chen, and S. Wang, “Improved deep hashing with soft pairwise similarity for multi-label image retrieval,” *IEEE Trans. Multimedia*, pp. 1–13, 2019.

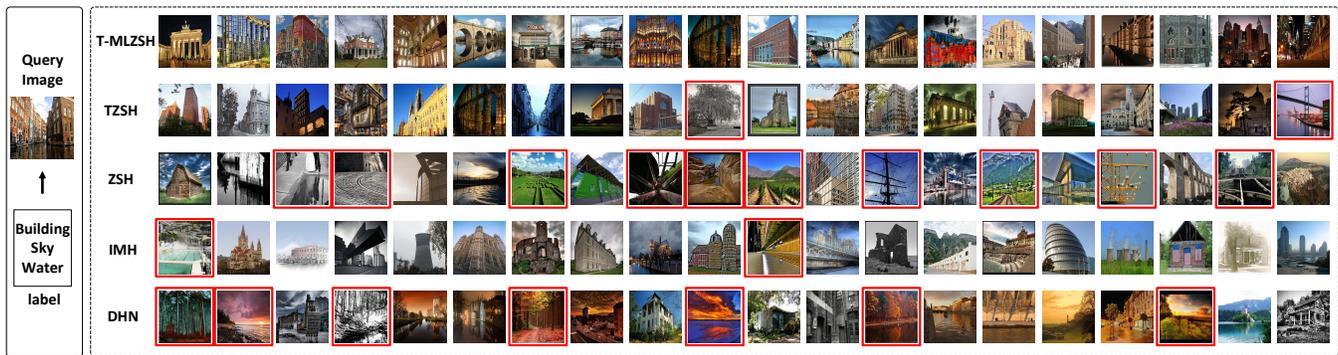


Fig. 15: Top-20 images retrieved by the proposed method and the comparison methods using the Hamming ranking on 48-bit hash codes.

- [8] W. Liu, J. Wang, S. Kumar, and S.-F. Chang, "Hashing with graphs," in *ICML*, 2011, pp. 1–8.
- [9] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, "Locality-sensitive hashing scheme based on p-stable distributions," in *Annual Symposium on Computational Geometry*, 2004, pp. 253–262.
- [10] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, and S.-E. Yoon, "Spherical hashing: binary code embedding with hyperspheres," in *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 11, 2015, pp. 2304–2316.
- [11] Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in *NIPS*, 2009, pp. 1753–1760.
- [12] L. V. Erin, J. Lu, G. Wang, P. Moulin, and J. Zhou, "Deep hashing for compact binary codes learning," in *CVPR*, 2015, pp. 2475–2483.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
- [14] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014, pp. 740–755.
- [15] L. Chen, W. Zhan, W. Tian, Y. He, and Q. Zou, "Deep integration: A multi-label architecture for road scene recognition," *IEEE Transactions on Image Processing*, vol. 28, no. 10, pp. 4883–4898, 2019.
- [16] Y. Sun, Y. Chen, X. Wang, and X. Tang, "Deep learning face representation by joint identification-verification," in *NIPS*, 2014, pp. 1988–1996.
- [17] L. Zhou, Z. Wang, Y. Luo, and Z. Xiong, "Separability and compactness network for image recognition and superresolution," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3275–3286, 2019.
- [18] Z. Shao, L. Wang, Z. Wang, W. Du, and W. Wu, "Saliency-aware convolution neural network for ship detection in surveillance video," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 3, pp. 781–794, 2020.
- [19] L. Chen, Q. Zou, Z. Pan, D. Lai, and L. Zhu, "Surrounding vehicle detection using an fpga panoramic camera and deep cnns," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [20] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015, pp. 3431–3440.
- [21] Q. Zou, Z. Zhang, Q. Li, X. Qi, Q. Wang, and S. Wang, "Deepcrack: Learning hierarchical convolutional features for crack detection," *IEEE Transactions on Image Processing*, vol. 28, no. 3, pp. 1498–1512, 2019.
- [22] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan, "Supervised hashing for image retrieval via image representation learning," in *AAAI*, 2014.
- [23] F. Zhao, Y. Huang, L. Wang, and T. Tan, "Deep semantic ranking based hashing for multi-label image retrieval," in *CVPR*, 2015, pp. 1556–1564.
- [24] H. Zhu, M. Long, J. Wang, and Y. Cao, "Deep hashing network for efficient similarity retrieval," in *AAAI*, 2016, pp. 2415–2421.
- [25] H. Liu, R. Wang, S. Shan, and X. Chen, "Deep supervised hashing for fast image retrieval," in *CVPR*, 2016, pp. 2064–2072.
- [26] Z. Cao, M. Long, J. Wang, and P. S. Yu, "Hashnet: Deep learning to hash by continuation," in *ICCV*, 2017, pp. 5609–5618.
- [27] C. Yan, B. Gong, Y. Wei, and Y. Gao, "Deep multi-view enhancement hashing for image retrieval," *T-PAMI*, 2020.
- [28] C. H. Lampert, H. Nickisch, and S. Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," in *CVPR*, 2009, pp. 951–958.
- [29] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.
- [30] Y. Yang, Y. Luo, W. Chen, F. Shen, J. Shao, and H. T. Shen, "Zero-shot hashing via transferring supervised knowledge," in *ACM MM*, 2016.
- [31] Y. Guo, G. Ding, J. Han, and Y. Gao, "Sitnet: discrete similarity transfer network for zero-shot hashing," in *IJCAI*, 2017, pp. 1767–1773.
- [32] S. Pachori, A. Deshpande, and S. Raman, "Hashing in the zero shot framework with domain adaptation," *Neurocomp.*, vol. 275, pp. 2137–2149, 2018.
- [33] H. Lai, "Transductive zero-shot hashing via coarse-to-fine similarity mining," in *ICMR*, 2018, pp. 196–203.
- [34] H. Zhang, Y. Long, and L. Shao, "Zero-shot hashing with orthogonal projection for image retrieval," *Pattern Recognition Letters*, vol. 117, pp. 201–209, 2019.
- [35] C. Yan, L. Li, C. Zhang, B. Liu, Y. Zhang, and Q. Dai, "Cross-modality bridging and knowledge transferring for image understanding," *IEEE T-MM*, pp. 2675–2685, 2019.
- [36] C. Yan, H. Xie, J. Chen, Z. Zha, X. Hao, Y. Zhang, and Q. Dai, "A fast uyghur text detector for complex background images," *IEEE T-MM*, pp. 3389–3398, 2018.
- [37] C. Yan, B. Shao, H. Zhao, R. Ning, Y. Zhang, and F. Xu, "3d room layout estimation from a single rgb image," *IEEE T-MM*, 2020.
- [38] W. Wang, V. W. Zheng, H. Yu, and C. Miao, "A survey of zero-shot learning: Settings, methods, and applications," *TIST*, vol. 10, no. 2, pp. 1–37, 2019.
- [39] X. Xu, T. Hospedales, and S. Gong, "Transductive zero-shot action recognition by word-vector embedding," *IJCV*, vol. 123, no. 3, pp. 309–333, 2017.
- [40] A. Lazaridou, E. Bruni, and M. Baroni, "Is this a wampimuk? cross-modal mapping between distributional semantics and the visual world," in *ACL*, 2014, pp. 1403–1414.
- [41] T. Mikolov, Q. V. Le, and I. Sutskever, "Exploiting similarities among languages for machine translation," *arXiv preprint arXiv:1309.4168*, 2013.
- [42] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *NIPS*, 2013, pp. 2121–2129.
- [43] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele, "Evaluation of output embeddings for fine-grained image classification," in *CVPR*, 2015, pp. 2927–2936.
- [44] C. Gan, Y. Yang, L. Zhu, D. Zhao, and Y. Zhuang, "Recognizing an action using its name: A knowledge-based approach," *IJCV*, vol. 120, no. 1, pp. 61–77, 2016.
- [45] Y. Guo, G. Ding, J. Han, and Y. Gao, "Zero-shot learning with transferred samples," *TIP*, vol. 26, no. 7, pp. 3277–3290, 2017.
- [46] —, "Zero-shot recognition via direct classifier learning with transferred samples and pseudo labels," in *AAAI*, 2017.
- [47] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata, "Feature generating networks for zero-shot learning," in *CVPR*, 2018, pp. 5542–5551.
- [48] Y. Zhu, M. Elhoseiny, B. Liu, X. Peng, and A. Elgammal, "A generative adversarial approach for zero-shot learning from noisy texts," in *CVPR*, 2018, pp. 1004–1013.
- [49] Y. Guo, G. Ding, J. Han, and Y. Gao, "Synthesizing samples for zero-shot learning," in *IJCAI*, 2017, pp. 1774–1780.
- [50] Y. Zhu, J. Xie, B. Liu, and A. Elgammal, "Learning feature-to-feature translator by alternating back-propagation for generative zero-shot learning," in *CVPR*, 2019, pp. 9844–9854.

- [51] M. Ye and Y. Guo, "Progressive ensemble networks for zero-shot recognition," in *CVPR*, June 2019.
- [52] C. Gan, Y. Yang, L. Zhu, D. Zhao, and Y. Zhuang, "Recognizing an action using its name: A knowledge-based approach," *IJCV*, vol. 120, no. 1, pp. 61–77, 2016.
- [53] A. Paul, N. C. Krishnan, and P. Munjal, "Semantically aligned bias reducing zero shot learning," in *CVPR*, 2019, pp. 7056–7065.
- [54] D. Mandal, S. Narayan, S. K. Dwivedi, V. Gupta, S. Ahmed, F. S. Khan, and L. Shao, "Out-of-distribution detection for generalized zero-shot action recognition," in *CVPR*, June 2019.
- [55] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille, "Multi-instance visual-semantic embedding," *arXiv preprint arXiv:1512.06963*, 2015.
- [56] T. Mensink, E. Gavves, and C. G. Snoek, "Costa: Co-occurrence statistics for zero-shot classification," in *CVPR*, 2014, pp. 2441–2448.
- [57] G. Irie, L. Z., W. X.-M., and C. S.-F., "Locally linear hashing for extracting non-linear manifolds," in *CVPR*, 2014, pp. 2115–2122.
- [58] A. Torralba, R. Fergus, and Y. Weiss, "Small codes and large image databases for recognition," in *CVPR*, 2008, pp. 1–8.
- [59] Y. Weiss, R. Fergus, and A. Torralba, "Multidimensional spectral hashing," in *ECCV*, 2012, pp. 340–353.
- [60] Q. Ning, J. Zhu, Z. Zhong, S. C. H. Hoi, and C. Chen, "Scalable image retrieval by sparse product quantization," *IEEE Transactions on Multimedia*, vol. 19, no. 3, pp. 586–597, 2017.
- [61] S. Ercoli, M. Bertini, and A. D. Bimbo, "Compact hash codes for efficient visual descriptors retrieval in large scale databases," *IEEE Transactions on Multimedia*, vol. 19, no. 11, pp. 2521–2532, 2017.
- [62] S. Xia, D. Peng, and D. Meng, "A fast adaptive k-means with no bounds," *IEEE T-PAMI*, 2020.
- [63] L. Liu and L. Shao, "Sequential compact code learning for unsupervised image hashing," *IEEE T-NNLS*, vol. 27, no. 12, pp. 2526–2536, 2015.
- [64] L. Zhu, Z. Huang, Z. Li, L. Xie, and H. T. Shen, "Exploring auxiliary context: discrete semantic transfer hashing for scalable image retrieval," *IEEE T-NNLS*, vol. 29, no. 11, pp. 5264–5276, 2018.
- [65] J. Gui, T. Liu, Z. Sun, D. Tao, and T. Tan, "Supervised discrete hashing with relaxation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 3, pp. 608–617, 2016.
- [66] Y. Cao, M. Long, B. Liu, and J. Wang, "Deep cauchy hashing for hamming space retrieval," in *CVPR*, June 2018.
- [67] D. Wang, Q. Wang, and X. Gao, "Robust and flexible discrete hashing for cross-modal similarity search," *TCSVT*, vol. 28, no. 10, pp. 2703–2715, 2018.
- [68] D. Wang, X. Gao, X. Wang, and L. He, "Label consistent matrix factorization hashing for large-scale cross-modal similarity search," *TPAMI*, vol. 41, no. 10, pp. 2466–2479, 2019.
- [69] D. Tian, Y. Wei, and D. Zhou, "Learning decorrelated hashing codes with label relaxation for multimodal retrieval," *IEEE Access*, vol. 8, pp. 79 260–79 272, 2020.
- [70] S. Su, Z. Zhong, and C. Zhang, "Deep joint-semantics reconstructing hashing for large-scale unsupervised cross-modal retrieval," in *ICCV*, 2019, pp. 3027–3035.
- [71] D. Wang, X. Gao, X. Wang, and L. He, "Semantic topic multimodal hashing for cross-media retrieval," in *IJCAI*, 2015, pp. 3890–3896.
- [72] D. Wang, X. Gao, X. Wang, L. He, and B. Yuan, "Multimodal discriminative binary embedding for large-scale cross-modal retrieval," *IEEE Transactions on Image Processing*, pp. 4540–4554, 2016.
- [73] H. Zhang, L. Liu, Y. Long, and L. Shao, "Unsupervised deep hashing with pseudo labels for scalable image retrieval," *IEEE Trans. Image Process.*, vol. 27, no. 4, pp. 1626–1638, 2018.
- [74] F. Shen, Y. Xu, L. Liu, Y. Yang, Z. Huang, and H. T. Shen, "Unsupervised deep hashing with similarity-adaptive and discrete optimization," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.
- [75] S. J. Pan and Q. Yang, "survey on transfer learning," *IEEE Trans. Know. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [76] Y. Fu, T. M. Hospedales, T. Xiang, and S. Gong, "Transductive multi-view zero-shot learning," *IEEE T-PAMI*, vol. 37, no. 11, pp. 2332–2345, 2015.
- [77] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Pan-athanathan, "Deep hashing network for unsupervised domain adaptation," in *CVPR*, 2017.
- [78] Y. Shen, L. Liu, F. Shen, and L. Shao, "Zero-shot sketch-image hashing," in *CVPR*, 2018, pp. 3598–3607.
- [79] Q. Wang and K. Chen, "Multi-label zero-shot human action recognition via joint latent embedding," *arXiv preprint arXiv:1709.05107*, 2017.
- [80] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.
- [81] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "Nus-wide: a real-world web image database from national university of singapore," in *ICIVR*, 2009, p. 48.
- [82] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *IJCV*, vol. 88, no. 2, pp. 303–338, 2010.
- [83] Q. Jiang and W. Li, "Asymmetric deep supervised hashing," in *AAAI*, 2018, pp. 3342–3349.
- [84] D. Wu, Q. Dai, J. Liu, B. Li, and W. Wang, "Deep incremental hashing network for efficient image retrieval," in *CVPR*, June 2019.
- [85] Y. Cao, L. Ju, Q. Zou, C. Qu, and S. Wang, "A multichannel edge-weighted centroidal voronoi tessellation algorithm for 3d super-alloy image segmentation," in *CVPR*, 2011, pp. 17–24.
- [86] K. Jarvelin and J. Kekalainen, "Cumulated gain-based evaluation of IR techniques," *ACM Trans. on Inform. Systems*, vol. 20, no. 4, pp. 422–446, 2002.

**Qin Zou** received the B.E. degree in information engineering in 2004 and the Ph.D. degree in computer vision in 2012, from Wuhan University, China. From 2010 to 2011, he was a visiting PhD student at the Computer Vision Lab, University of South Carolina, USA. Currently, he is an associate professor with the School of Computer Science, Wuhan University. He is a co-recipient of the National Technology Invention Award of China in 2015. His research activities involve computer vision, pattern recognition, and machine learning. He is a senior member of the IEEE and a member of the ACM.

**Ling Cao** received the B.S. degree in computer science from Yunnan University in 2018, and is now working towards her M.S. degree in computer application at the School of Computer Science, Wuhan University, China. Her research interests include deep learning and image/video retrieval.

**Zheng Zhang** received the B.S. degree in computer science and M.S. degree in computer application from the School of Computer Science, Wuhan University, in 2016 and 2019, respectively. He won the first prize in 'China Undergraduate Contest in Internet of Things' in 2015. His research interests include deep learning and its application in image classification and retrieval.

**Long Chen** received the B.Sc. degree in communication engineering and the Ph.D. degree in signal and information processing from Wuhan University, Wuhan, China, in 2007 and in 2013, respectively. From October 2010 to November 2012, he was co-trained PhD Student at National University of Singapore. He is currently an Associate Professor with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. His areas of interest include deep learning and cognitive perception.

**Song Wang** received the PhD degree in electrical and computer engineering from the University of Illinois at Urbana-Champaign (UIUC) in 2002. From 1998 to 2002, he also worked as a research assistant in the Image Formation and Processing Group at the Beckman Institute of UIUC. In 2002, he joined the Department of Computer Science and Engineering at the University of South Carolina, where he is currently a professor. His research interests include computer vision, medical image processing, and machine learning. He is currently serving as an Associate Editor of IEEE Transactions on Pattern Analysis and Machine Intelligence and an Associate Editor of Pattern Recognition Letters. He is a senior member of the IEEE and the IEEE Computer Society.