

# Deep Binaries: Encoding Semantic-Rich Cues for Efficient Textual-Visual Cross Retrieval

Yuming Shen<sup>1</sup>, Li Liu<sup>1,2</sup>, Ling Shao<sup>1</sup>, and Jingkuan Song<sup>3</sup>

<sup>1</sup>School of Computing Sciences, University of East Anglia, Norwich, UK

<sup>2</sup>Malong Technologies Co., Ltd., Shenzhen, China

<sup>3</sup>Future Media Center, University of Electronic Science and Technology of China, Chengdu, China  
yuming.shen@uea.ac.uk, li.liu@malongtech.cn, ling.shao@ieee.org, jingkuan.song@gmail.com

## Abstract

Cross-modal hashing is usually regarded as an effective technique for large-scale textual-visual cross retrieval, where data from different modalities are mapped into a shared Hamming space for matching. Most of the traditional textual-visual binary encoding methods only consider holistic image representations and fail to model descriptive sentences. This renders existing methods inappropriate to handle the rich semantics of informative cross-modal data for quality textual-visual search tasks.

To address the problem of hashing cross-modal data with semantic-rich cues, in this paper, a novel integrated deep architecture is developed to effectively encode the detailed semantics of informative images and long descriptive sentences, named as Textual-Visual Deep Binaries (TVDB). In particular, region-based convolutional networks with long short-term memory units are introduced to fully explore image regional details while semantic cues of sentences are modeled by a text convolutional network. Additionally, we propose a stochastic batch-wise training routine, where high-quality binary codes and deep encoding functions are efficiently optimized in an alternating manner. Experiments are conducted on three multimedia datasets, i.e. Microsoft COCO, IAPR TC-12, and INRIA Web Queries, where the proposed TVDB model significantly outperforms state-of-the-art binary coding methods in the task of cross-modal retrieval.

## 1. Introduction

Data retrieval between image and text modalities has aroused a lot of recent attention, becoming an overwhelming research topic in computer vision. As deep learning technology develops, dramatic progress has been achieved

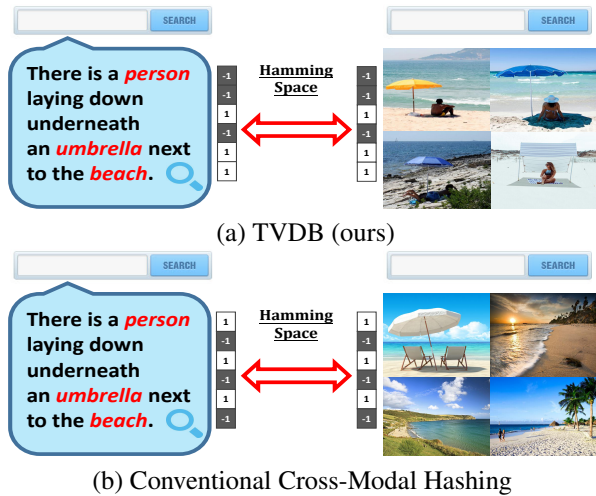


Figure 1. The proposed model (a) aims at encoding all informative words of a sentence and all possible attractive regions in an image. Contrarily, conventional textual-visual binary encoding methods (b) only utilize simple representations of each modality and discard some information (e.g. the instances of *person* and *umbrella* are not well-encoded), resulting in low matching quality.

recently in this area. Handling large-scale image-text retrieval, cross-modal hashing schemes [10, 34, 79, 61, 37, 71, 5, 23, 22, 68, 78] have been proposed to encode heterogeneous data from a high-dimensional feature space into a shared low-dimensional Hamming space where an approximate nearest neighbour of a given query can be found efficiently.

Traditionally, most existing researches in cross-modal hashing only focus on image-tag mapping [10, 34, 79, 61, 3, 37, 71, 5, 23], where holistic image representations and semantic tags feed the shallow binary coding procedure. It is insufficient to simply link images with tags, instead of real

sentences, for multimedia searching problems; the detailed information of images is usually discarded due to the use of holistic representations and the semantics enclosed in descriptive sentences are hardly explored. The recent deep hashing methods [5, 13, 23, 52] also suffer from several drawbacks. First of all, image data are likewise poorly modeled due to the lack of detailed regional information during encoding. We argue that this is not optimal for images with fruitful semantics. Secondly, most of the deep models mentioned above still utilize coarse text representations, which is inappropriate for modeling long sentences. Moreover, the network training efficiency can be further improved by more advanced code learning architectures.

Driven by the drawbacks of the previous works, in this paper, we consider a more challenging task to encode informative multi-modal data, *i.e.*, semantic-rich images and descriptive sentences, into binary codes for cross-modal search, termed as Textual-Visual Deep Binaries (TVDB). Particularly, the popular Region Proposal Network (RPN) [58] and Long Short-Term Memory (LSTM) [77, 63] are introduced to formulate the image binary coding function, so that the regional semantic details in images can be well preserved from dominant to minor. Meanwhile, the latest advances in text Convolutional Neural Networks (CNN) [28, 8, 20] are adopted to build the text binary encoding network, leveraging structural cues between the words in a sentence. The proposed deep architecture produces high-semantic-retentivity binary codes and achieves promising retrieval performance. The intuitive difference between the proposed method and the traditional ones are given in Figure 1. It can be seen that the proposed TVDB encodes as many details as possible from images and sentences, leading to more representative binary codes for matching.

In addition to the novel deep binary encoding networks of TVDB, an efficient stochastic batch-wise code learning procedure is proposed. Inspired by Shen *et al.* [60], the binary codes in TVDB are discretely and alternately optimized during the batch-wise learning procedure. Batching data randomly and iteratively, the proposed training routine guarantees an effective learning objective convergency. The contributions of this work can be summarized as follows:

- The TVDB model is proposed to effectively encode rich regional information of images as well as semantic dependencies and cues between words by exploiting two modal-specific deep binary embedding networks. In this way, the intrinsic semantic correlation between heterogeneous data can be quantitatively measured and captured.
  - A novel stochastic batch-wise training strategy is adopted to optimize TVDB, in which reliable binary codes and the deep encoding functions are optimized in an alternating manner within every single batch.
- The evaluation results of our model on three semantically rich datasets highly surpass those of existing state-of-the-art binary coding methods in cross-modal retrieval.

## 2. Related Work

The cutting-edge studies in vision and language achieve promising results in terms of visual question answering [50, 2, 14], caption generation [75, 12, 24], and real-valued cross-modal retrieval [62, 48, 27, 66, 21, 25, 30, 26, 51, 49, 76]. The best-performing real-valued cross-modal retrieval models typically rely on densely annotated image-region and text pairs for embedding. However, these methods are far from satisfactory for large-scale data retrieval due to the inefficient similarity computation of real-valued embeddings.

On the other hand, there exists several hashing methods [40, 69, 45, 18, 46, 16, 55, 44, 33, 43] aiming at efficient retrieval. For textual-visual hashing, it has been a traditional and common solution to encode images and tags via shallow embedding functions with either unsupervised [34, 61, 10, 56, 80, 67, 65, 42, 39], pairwise based [3, 71, 22, 54, 68, 78, 47] or supervised [72, 6, 37, 79, 52] code learning methods. More recently, deep hashing methods [7, 73, 81, 38, 11, 4, 41] provide promising results in image recognition, which is also adopted in [5, 13, 23, 52, 4] for textual-visual retrieval. Jiang *et al.* proposed DCMH [23] for image-tag retrieval using a set of multi-layer neural networks which simply take deep holistic image features and word count vectors as input. DVSH [5] proposed by Cao *et al.* becomes a more feasible solution for image-sentence hashing as the sequential information of sentence data is better encoded by introducing the Recurrent Neural Networks (RNN) [77]. CDQ [4] combines data representation learning steps with quantization error controlling hash coding methods with deep neural networks, while it still basically addresses image-tag hashing. In general, most of these methods can barely obtain adequate performance in our task due to the coarse image and text representations.

## 3. Deep Encoding Networks for TVDB

This work addresses the problem of data retrieval between informative images and long sentences using deep binary codes. As shown in Figure 2, the proposed TVDB model is composed of two deep neural networks and a batch-wise code learning phase. The two deep neural networks play the role of binary encoding functions for images and sentences denoted as  $f(\cdot)$  and  $g(\cdot)$  respectively. The batch-wise optimization allows using binary codes as supervision of  $f(\cdot)$  and  $g(\cdot)$  in a mini-batch during training.

Some preliminary notation is introduced here. We consider a multi-media data collection  $\mathcal{O} = \{\mathbf{X}, \mathbf{Y}\}$  containing both image data  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and sentence data

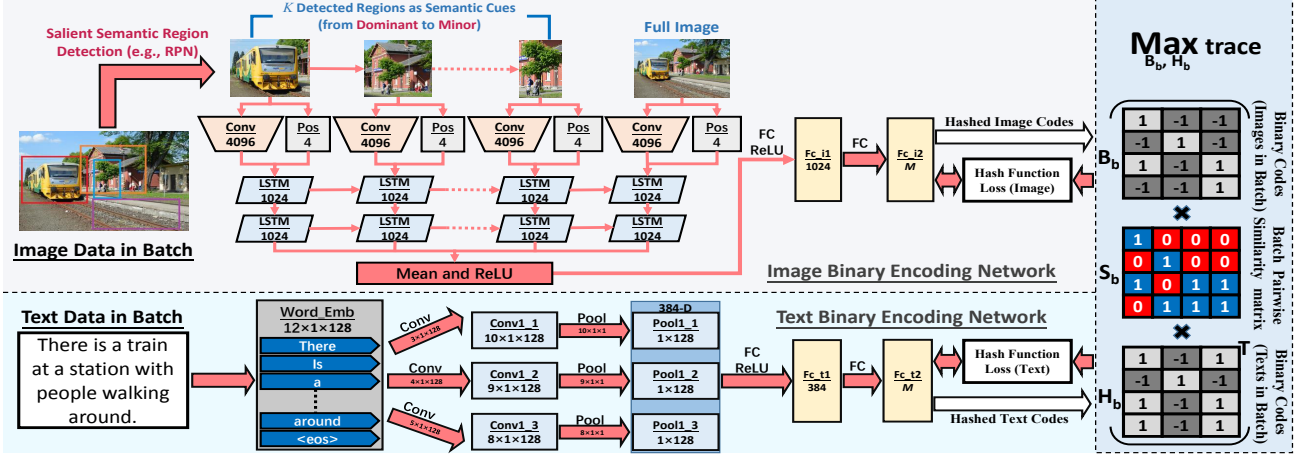


Figure 2. The network architecture of TVDB. RPN, CNN and LSTM are utilized to form the image encoding function  $f(\cdot)$  (the upper network), encoding image regions from dominant to minor. The sentence encoding network  $g(\cdot)$  is built with a text-CNN (the lower network). The rightmost component refers to a **batch-wise** coding procedure. Here the module `Pos` computes the region coordinates discussed in subsection 3.1, while `Word_Emb` refers to the linear word embedding procedure introduced in subsection 3.2.

$\mathbf{Y} = \{y_i\}_{i=1}^N$ , with  $N$  denoting the total number of data points of each modality in the dataset. The two deep binary encoding functions  $f(\mathbf{X}; \Theta)$  and  $g(\mathbf{Y}; \Phi)$  of TVDB are parameterized by  $\Theta$  and  $\Phi$ . The sign function is applied to  $f(\cdot)$  and  $g(\cdot)$  to produce binary representations:

$$\begin{aligned} \mathbf{B} &= \text{sign}(f(\mathbf{X}; \Theta)) \in \{-1, 1\}^{M \times N}, \\ \mathbf{H} &= \text{sign}(g(\mathbf{Y}; \Phi)) \in \{-1, 1\}^{M \times N}, \end{aligned} \quad (1)$$

where  $M$  refers to the target binary encoding length. In the following subsections, we introduce the setups of the deep binary encoding networks.

### 3.1. Image binary encoding network

The architecture of the image encoding network  $f(\cdot)$  is given in this subsection. As discussed previously, encoding the holistic image discards the informative patterns and produces poor coding quality, which is not desirable for our task. We consider a deep neural network architecture that embeds several salient regions of an image into a single binary vector to enrich the encoded semantic information. It is worthwhile to note that we are not directly linking every image region with a certain concept of a sentence as it is not feasible for cross-modal hashing problems and is contrary to the original intention of binary encoding in this work. Instead, regional semantic cues of images are leveraged here to improve the encoding quality.

**Salient Semantic Region Selection.** As shown in Figure 2, for each image in the mini-batch, TVDB firstly detects a number of regional proposals that possibly carry informative parts, *e.g.* recognizable or dominative objects in the image. Recent works in region-based CNN [15, 58] show great potential in detecting semantically meaningful areas

of an image. We adopt the framework of the state-of-the-art RPN [58] as the proposal detection basis of TVDB. A total number of  $K$  semantic regions are sampled for further processing according to a simple heuristic attraction score  $a_k$  in descending order, that is:

$$a_k = (c_k + d_k)/2, \quad (2)$$

where  $c_k \in (0, 1)$  denotes the confidence score determined by the RPN and  $d_k \in (0, 1)$  refers to the normalized proportion of the  $k$ -th detected proposal in an image. We consider those regional proposals with high attraction score  $a_k$  semantically dominating to the whole image as they are usually the most recognizable image parts. This heuristic region selection solution highly fits the task of cross-modal binary encoding since it does not require any additional supervision or fine-grained region-sentence relations.

**Regional Representation and Augmentation.** The selected image regions are fed into CNNs to extract vectorized representations. The benchmark CNN architecture AlexNet [32] is involved here, from which a feature representation of 4096-D is obtained. To make the most of structural information, the feature vector of each region is augmented with four additional digits indicating the normalized height, width, and center coordinates of the corresponding region bounding box, making the whole regional representation a 5000-D vector for each.

**Recurrent Network for Encoding.** For our task, it is desirable to use a method which capitalizes on information from the selected ordered regions so that dominating image parts contribute more to the final representation. It has been proved that human eyes sequentially browse image parts from dominant to minor [57]. Simulating this procedure, we sort the  $K$  selected regional proposals according to

their attraction scores in descending order and then the corresponding 5000-D representation for each proposal is sequentially fed into an RNN so that the dominant image parts can be well utilized. Additionally, the CNN feature of the holistic image is also appended to the end of the RNN input sequence, making a total of  $K + 1$  semantic regions for encoding. In particular, a two-level Long Short-Term Memory (LSTM) [19] is implemented as the RNN unit with 1024-D output length, following the popular structure described in [77]. Shown in Figure 2, the outputs of the LSTMs are averaged along the time sequence and appended with a ReLU [53] activation. Two fully-connected layers are applied to the top of the averaged LSTM outputs, with output dimension 1024 and  $M$  respectively. Thus the whole image can be encoded into an  $M$ -bit binary vector using the  $\text{sign}(\cdot)$  function. We choose the identity function as the activation to the fully-connected layer for the convenience of code regression. The convolutional layers for images are built following AlexNet [32].

### 3.2. Sentence binary encoding network

Although recurrent networks are widely adopted in textual-visual tasks [50, 12, 24, 59], it is still argued that RNNs such as LSTMs are not usually a superior choice for specific language tasks due to their non-structural designs [20, 8]. We aim at encoding the structural and contextual cues between the words in a sentence to ensure the produced binary codes have adequate information capacity. To this end, the text-CNN [28] is chosen as the text-side encoding network  $g(\cdot)$ , where each word in a descriptive sentence is firstly embedded into a word vector with a certain dimension and then convolution is performed along the word sequence.

We pre-process text data following the conventional manner where all sentences are appended with an *eos* token and padded or truncated to a certain length with all full stops removed. The sentence length after preprocessing is fixed to 12, which is about the mean length of text data in the datasets used in our experiments. Each word is embedded to a 128-D vector using a linear projection before being fed into the CNN. The text-CNN architecture in TVDB is similar to the one of Kim [28], with more fully-connected layers for coding. The full configuration of our text-CNN is given at the bottom of Figure 2. The `WordEmb` layer in Figure 2 refers to the word embedding, of which the parameters are also involved in the back-propagation (BP) procedure. For the text convolution setups, the first and second digits of the kernel size denote the height and width of the convolutional kernels, with the third digit being the kernel number. Note that, as text convolution is only performed along the word sequence, the second dimension of kernel size of the convolutional layers are always 1. In this work, the number of kernels for all convolutional layers is set to 128. We also

build two fully-connected layers here, where the first one `Fc_t1` takes inputs from all pooling layers followed by a ReLU activation, and for the second one, `Fc_t2`, the identity activation is applied.

## 4. Stochastic Batch-Wise Code Learning

The entire training procedure for  $f(\cdot)$  and  $g(\cdot)$  follows the mini-batch Stochastic Gradient Descent (SGD) since deep neural networks are utilized. We suggest the binary learning solution should provide reliable target codes as network supervision every time a mini-batch feeds.

### 4.1. Batch-wise alternating optimization

Let  $\mathcal{O}_b = \{\mathbf{X}_b, \mathbf{Y}_b\}$  denote a mini-match stochastically taken from the data collection  $\mathcal{O}$ , where  $\mathbf{X}_b = \{\mathbf{x}_i\}_{i=1}^{N_b}$  and  $\mathbf{Y}_b = \{\mathbf{y}_i\}_{i=1}^{N_b}$  are image and sentence data in the mini-batch respectively. As the training process of TVDB is typically batch-based, we introduce an in-batch pair-wise similarity matrix  $\mathbf{S}_b \in \{0, 1\}^{N_b \times N_b}$  for target binary code learning, with  $N_b$  denoting the number of data points within a mini-batch for training. The entry  $s_{pq}$  of  $\mathbf{S}_b$  is defined as follows:

$$s_{pq} = \begin{cases} 1 & \mathbf{x}_p, \mathbf{y}_q \text{ share at least one semantic label,} \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

The aim is to learn a set of target binary codes  $\mathbf{B}_b$  for images and  $\mathbf{H}_b$  for sentences that best describe the in-batch samples. Fully utilizing the pairwise relations  $\mathbf{S}_b$  as supervision, a trace-based prototypic learning objective for hashing is thus built as

$$\mathcal{L}(\mathbf{B}_b, \mathbf{H}_b, \mathbf{S}_b) = -\text{trace}(\mathbf{B}_b \mathbf{S}_b \mathbf{H}_b^\top). \quad (4)$$

A common learning procedure for cross-modal hashing can be formulated by solving the following problem:

$$\begin{aligned} & \min_{\mathbf{B}_b, \mathbf{H}_b, \Theta, \Phi} \mathcal{L}(\mathbf{B}_b, \mathbf{H}_b, \mathbf{S}_b), \\ & \text{s.t. } \mathbf{B}_b = \text{sign}(f(\mathbf{X}_b; \Theta)), \mathbf{H}_b = \text{sign}(g(\mathbf{Y}_b; \Phi)). \end{aligned} \quad (5)$$

Relaxing the binary constraints to be continuous, *i.e.*,  $\mathbf{B}_b = f(\mathbf{X}_b; \Theta)$ ,  $\mathbf{H}_b = g(\mathbf{Y}_b; \Phi)$ , results in a slow and difficult optimization process. Inspired by Shen *et al.* [60], we reformulate the problem of (5) by keeping the binary constraints and regarding  $\mathbf{B}_b$  and  $\mathbf{H}_b$  as auxiliary variables,

$$\begin{aligned} & \min_{\mathbf{B}_b, \mathbf{H}_b, \Theta, \Phi} \mathcal{L}(\mathbf{B}_b, \mathbf{H}_b, \mathbf{S}_b) + \eta (\|\mathbf{B}_b - f(\mathbf{X}_b; \Theta)\|_F^2 \\ & \quad + \|\mathbf{H}_b - g(\mathbf{Y}_b; \Phi)\|_F^2), \\ & \text{s.t. } \mathbf{B}_b \in \{-1, 1\}^{M \times N_b}, \mathbf{H}_b \in \{-1, 1\}^{M \times N_b}, \end{aligned} \quad (6)$$

where  $\eta$  is a penalty hyper parameter and  $\|\cdot\|_F$  refers to the Frobenius norm. The two Frobenius norms here depict

the quantization error between the binary codes  $\mathbf{B}_b$ ,  $\mathbf{H}_b$  and the binary coding function outputs  $f(\cdot)$ ,  $g(\cdot)$ . It has been proved in [60] that with a sufficiently large value of  $\eta$ , Eq. (6) becomes a close approximation to Eq. (5), in which slight disparities between  $\mathbf{B}_b$  and  $f(\mathbf{X}_b; \Theta)$  or  $\mathbf{H}_b$  and  $g(\mathbf{Y}_b; \Phi)$  are tolerant to our binary learning problem.

Therefore, the comprehensive learning objective of TVDB is formulated. We provide optimization schemes below. It is observed that Eq. (6) is a non-convex NP-hard problem due to the binary constraints. To better access it, an alternating solution based on coordinate descent is adopted to sequentially optimize  $\mathbf{B}_b$ ,  $\mathbf{H}_b$ ,  $\Theta$  and  $\Phi$  in every single batch as follows.

**Updating  $\mathbf{B}_b$ .** By fixing  $\mathbf{H}_b$ ,  $\Theta$  and  $\Phi$ , the subproblem of (6) w.r.t.  $\mathbf{B}_b$  can be written as

$$\begin{aligned} \min_{\mathbf{B}_b} \eta \|\mathbf{B}_b - f(\mathbf{X}_b; \Theta)\|_{\mathbb{F}}^2 - \text{trace}(\mathbf{B}_b \mathbf{S}_b \mathbf{H}_b^{\top}), \\ \text{s.t. } \mathbf{B}_b \in \{-1, 1\}^{M \times N_b}, \end{aligned} \quad (7)$$

to which the closed-form optimal solution becomes

$$\mathbf{B}_b = \text{sign}(2\eta f(\mathbf{X}_b; \Theta) + \mathbf{S}_b \mathbf{H}_b^{\top}). \quad (8)$$

**Updating  $\mathbf{H}_b$ .** Similar to  $\mathbf{B}_b$ , the solution to the subproblem of (6) w.r.t.  $\mathbf{H}_b$  is

$$\mathbf{H}_b = \text{sign}(2\eta g(\mathbf{Y}_b; \Phi) + \mathbf{B}_b \mathbf{S}_b). \quad (9)$$

**Updating  $\Theta$  and  $\Phi$ .** The subproblem of (6) w.r.t.  $\Theta$  and w.r.t.  $\Phi$  can be respectively written as

$$\begin{aligned} \min_{\Theta} \text{loss}_f &:= \eta \|\mathbf{B}_b - f(\mathbf{X}_b; \Theta)\|_{\mathbb{F}}^2, \\ \min_{\Phi} \text{loss}_g &:= \eta \|\mathbf{H}_b - g(\mathbf{Y}_b; \Phi)\|_{\mathbb{F}}^2, \end{aligned} \quad (10)$$

when all other variables are fixed. These two subproblems are typically in the form of the  $l_2$  norm which are differentiable problems and thus  $\Theta$  or  $\Phi$  can be optimized in the framework of SGD using back-propagation. Here the updated auxiliary binary codes  $\mathbf{B}_b$  and  $\mathbf{H}_b$  act as the supervisions to the binary encoding networks  $f(\cdot)$  and  $g(\cdot)$ .

## 4.2. Stochastic-batched training procedure.

In the above subsection, we have presented the binary code learning algorithm for each mini-batch of TVDB. However, how to apply the batch-wise learning objective to the whole dataset has not yet been discussed. In this subsection, the overall training procedure for mini-batch SGD is introduced. Note that simply keeping data in every mini-batch unaltered in each training epoch usually results in poorly-learned hash functions. This is because the cross-modal in-batch data are not able to interact with the data

---

### Algorithm 1: The Training Process of TVDB

---

**Input:** Image-sentence dataset  $\mathcal{O} = \{\mathbf{X}, \mathbf{Y}\}$ , Max training iteration  $T$

**Output:** Hash function parameters  $\Theta$  and  $\Phi$

Randomly initialize  $\mathbf{B}, \mathbf{H} \in \{-1, 1\}^{M \times N}$

**repeat**

    Get a **stochastic** mini-batch  $\mathcal{O}_b$  from  $\mathcal{O}$

    Get  $\mathbf{B}_b, \mathbf{H}_b$  from  $\mathbf{B}, \mathbf{H}$  with respective indices

    Build  $\mathbf{S}_b$  according to data relations and labels

**Update**  $\mathbf{B}_b \leftarrow \text{Eq.}(8)$

**Update**  $\mathbf{H}_b \leftarrow \text{Eq.}(9)$

$(\text{loss}_f, \text{loss}_g) \leftarrow \text{Eq.}(10)$

**Update**

$(\Theta, \Phi) \leftarrow (\Theta, \Phi) - \Gamma(\nabla_{\Theta} \text{loss}_f, \nabla_{\Phi} \text{loss}_g)$

**until** convergence or max training iter  $T$  is reached;

---

outside of the batch and thus the batch-wise similarity  $\mathbf{S}_b$  skews the statistics of the whole dataset. To be more precise, we consider a data batching scheme to build every  $\mathcal{O}_b$  and  $\mathbf{S}_b$  that well explores the cross-modal semantic relationships across the entire dataset. To this end, a stochastic batching routine is designed. Each data mini-batch is randomly formed before being input into the training procedure. Therefore,  $\mathbf{S}_b$  varies across each batch, which ensures the in-batch data diversity.

Combining the stochastic batching method with the alternating parameter updating schemes, the whole training procedure of TVDB is illustrated in Algorithm 1. The operator  $\Gamma(\cdot)$  in Algorithm 1 indicates the adaptive gradient scaler used for SGD, which is the Adam optimizer [29] in this work. Unlike some existing deep hashing methods [23] which update the target binary codes when a whole epoch is finished, TVDB updates  $\mathbf{B}_b$  and  $\mathbf{H}_b$  instantly when a mini-batch arrives. This training routine proves to achieve fast convergence for effectively learning the encoding networks  $f(\cdot)$  and  $g(\cdot)$  as shown in Figure 5. Once the TVDB model is trained, given an image query  $\mathbf{x}_q$  for example, we compute its binary code by  $\mathbf{b}_q = \text{sign}(f(\mathbf{x}_q; \Theta))$ . While for the retrieval database, the unified binary codes from each sentence is obtained via  $\mathbf{H} = \text{sign}(g(\mathbf{Y}; \Phi))$ . A sentence query can be processed in the similar manner.

## 5. Experiments

Experiments of TVDB on cross-modal retrieval are performed on three semantically fruitful sentence-vision datasets: Microsoft COCO [36], IAPR TC-12 [17] and INRIA Web Queries [31]. The evaluation results are reported according to the following themes: **(a) comparison with state-of-the-arts methods**, **(b) deep encoding network ablation study** and **(c) training routine feasibility**.

For implementation details, we utilize the RPN [58] for

Table 1. Image-sentence cross-modal retrieval MAP results of the proposed model compared with existing methods on the three datasets.

Task	Method	Binary Code	Microsoft COCO				IAPR TC-12				INRIA Web Queries			
			16 bits	32 bits	64 bits	128 bits	16 bits	32 bits	64 bits	128 bits	16 bits	32 bits	64 bits	128 bits
Image Query Sentence	CVFL [74]	✗	0.488 (4096-D)				0.499 (4096-D)				0.137 (4096-D)			
	PLSR [70]	✗	0.528 (4096-D)				0.485 (4096-D)				0.258 (4096-D)			
	CCA [64]	✗	0.544	0.536	0.508	0.469	0.468	0.434	0.406	0.389	0.204	0.217	0.189	0.152
	CMFH [10]	✓	0.488	0.506	0.508	0.383	0.442	0.437	0.361	0.362	0.197	0.236	0.245	0.238
	CVH [34]	✓	0.489	0.474	0.436	0.402	0.422	0.401	0.384	0.374	0.204	0.200	0.187	0.163
	SCM [79]	✓	0.529	0.563	0.587	0.601	0.486	0.505	0.515	0.522	0.115	0.215	0.271	0.308
	IMH [61]	✓	0.615	0.650	0.657	0.677	0.463	0.490	0.510	0.521	0.233	0.250	0.277	0.295
	QCH [71]	✓	0.572	0.595	0.613	0.634	0.526	0.555	0.579	0.605	-	-	-	-
	CMSSH [3]	✓	0.405	0.489	0.441	0.448	0.345	0.337	0.348	0.374	0.151	0.162	0.155	0.151
	SePH [37]	✓	0.581	0.613	0.625	0.634	0.507	0.513	0.515	0.530	0.126	0.137	0.134	0.142
	CorrAE [13]	✓	0.550	0.556	0.570	0.580	0.495	0.525	0.558	0.589	-	-	-	-
	CM-NN [52]	✓	0.556	0.560	0.585	0.594	0.516	0.542	0.577	0.600	-	-	-	-
	DNH-C [35]	✓	0.535	0.556	0.569	0.582	0.480	0.509	0.526	0.535	-	-	-	-
	DCMH [23]	✓	0.562	0.597	0.609	0.646	0.443	0.491	0.559	0.556	0.241	0.268	0.301	0.384
DVSH [5]	✓	0.587	0.713	0.739	0.755	0.570	0.632	0.696	0.724	-	-	-	-	
<b>TVDB</b>	✓	<b>0.702</b>	<b>0.781</b>	<b>0.797</b>	<b>0.818</b>	<b>0.629</b>	<b>0.697</b>	<b>0.731</b>	<b>0.772</b>	<b>0.368</b>	<b>0.405</b>	<b>0.419</b>	<b>0.446</b>	
Sentence Query Image	CVFL [74]	✗	0.561 (4096-D)				0.495 (4096-D)				0.291 (4096-D)			
	PLSR [70]	✗	0.538 (4096-D)				0.492 (4096-D)				0.257 (4096-D)			
	CCA [64]	✗	0.545	0.542	0.513	0.468	0.471	0.438	0.413	0.395	0.212	0.228	0.188	0.150
	CMFH [10]	✓	0.574	0.507	0.510	0.472	0.447	0.445	0.438	0.365	0.178	0.231	0.248	0.131
	CVH [34]	✓	0.486	0.470	0.434	0.401	0.424	0.403	0.386	0.376	0.204	0.200	0.179	0.164
	SCM [79]	✓	0.523	0.544	0.569	0.585	0.495	0.514	0.523	0.529	0.123	0.226	0.297	0.349
	IMH [61]	✓	0.610	0.679	0.728	0.740	0.516	0.526	0.534	0.527	0.251	0.275	0.306	0.284
	QCH [71]	✓	0.574	0.606	0.638	0.667	0.500	0.536	0.565	0.589	-	-	-	-
	CMSSH [3]	✓	0.375	0.384	0.340	0.360	0.363	0.377	0.365	0.348	0.154	0.153	0.156	0.147
	SePH [37]	✓	0.613	0.649	0.672	0.693	0.471	0.480	0.481	0.495	0.226	0.256	0.291	0.319
	CorrAE [13]	✓	0.559	0.581	0.611	0.626	0.498	0.520	0.533	0.550	-	-	-	-
	CM-NN [52]	✓	0.579	0.598	0.620	0.645	0.512	0.539	0.549	0.565	-	-	-	-
	DNH-C [35]	✓	0.525	0.559	0.590	0.634	0.469	0.484	0.491	0.505	-	-	-	-
	DCMH [23]	✓	0.595	0.601	0.633	0.658	0.486	0.487	0.499	0.541	0.227	0.305	0.322	0.380
DVSH [5]	✓	0.591	0.737	0.758	0.767	0.604	0.640	0.680	0.675	-	-	-	-	
<b>TVDB</b>	✓	<b>0.713</b>	<b>0.779</b>	<b>0.787</b>	<b>0.810</b>	<b>0.674</b>	<b>0.678</b>	<b>0.704</b>	<b>0.721</b>	<b>0.353</b>	<b>0.462</b>	<b>0.464</b>	<b>0.470</b>	

image proposal detection to pick  $K = 20$  informative regions for further LSTM-based encoding, and the value of  $\eta$  in problem (6) is set to  $10^{-4}$  via cross-validation. For the image-side CNNs, AlexNet [32] without its `fc_8` layer is adopted with the pre-trained parameters from ImageNet classification [9]. The TVDB framework is implemented using Tensorflow [1].

### 5.1. Experimental settings

The experiments of sentence-vision retrieval are taken on three multimedia datasets. Following the conventional textual-visual retrieval measures [37, 5], the relevant instances for a query are defined by sharing at least one label.

**Microsoft COCO [36].** The COCO dataset contains a training image set of 80,000 samples with about 40,000 validation images. Each image is assigned five sentence descriptions and labeled with 80 semantic topics. To be consistent with [5], we randomly select 5,000 images from the validation set and thus the retrieval gallery becomes around 85,000 images, from which we explicitly take 5,000 pairs as the query set and 50,000 images for training.

**IAPR TC-12 [17].** It consists of 20,000 images. Each image is provided with 1.7 descriptive sentences on average. In addition, category annotations are given on all images with 275 concepts. Following the setting in [5], we use 18,000 image-sentence pairs that belong to the most frequent 22 topics as the retrieval gallery, from which we take

2,200 pairs as the query data and 5,000 as the training set.

**INRIA Web Queries [31].** This dataset contains about 70,000 images categorized into 353 conceptual labels. Sentence descriptions are provided to most of the images. We select images belonging to the 100 most frequent concepts, making the gallery 25,015 image-sentence pairs. For the query set, 10 images and sentences are randomly selected from each category and the rest are used as training data.

### 5.2. Comparison with existing methods

The overall retrieval performance of TVDB is analysed and compared using Mean Average Precision (MAP) and Precision-Recall curves.

**Baselines.** Several baselines of traditional cross-modal hashing methods are adopted for comparison, including CMFH [10], CVH [34], SCM[79], IMH [61], QCH [71], CMSSH [3], SePH [37], while CVFL [74], CCA [64] and PLSR [70] are considered as real-valued methods. The deep-learning-based cross-modal hashing methods, *i.e.*, CorrAE [13], CM-NN [52], DNH [35], DCMH [23] and DVSH [5], are also included here. To make fair comparisons, we utilize deep features for all traditional baseline methods mentioned above if the codes are available. For image features, we directly use the 4096-D AlexNet [32] pre-trained representations. For text features, a multi-label classification text-CNN, which shares most of its structure with our text encoding network  $g(\cdot)$  excluding the last layer, is

Table 2. Ablation study: cross-modal retrieval MAP comparison of TVDB with several deep network baselines.

Task	Method	Microsoft COCO				IAPR TC-12			
		16 bits	32 bits	64 bits	128 bits	16 bits	32 bits	64 bits	128 bits
Image Query Sentence	TVDB-I1 (full image only)	0.565	0.597	0.621	0.679	0.468	0.526	0.597	0.589
	TVDB-I2 (ave-pooled regions)	0.677	0.692	0.704	0.726	0.533	0.572	0.647	0.656
	TVDB-T1 (text bag-of-words)	0.553	0.565	0.591	0.631	0.436	0.508	0.503	0.548
	TVDB-T2 (text LSTM)	0.654	0.717	0.742	0.775	0.551	0.619	0.679	0.726
	<b>TVDB (full model)</b>	<b>0.702</b>	<b>0.781</b>	<b>0.797</b>	<b>0.818</b>	<b>0.629</b>	<b>0.678</b>	<b>0.731</b>	<b>0.772</b>
Sentence Query Image	TVDB-I1 (full image only)	0.573	0.634	0.653	0.657	0.462	0.514	0.516	0.568
	TVDB-I2 (ave-pooled regions)	0.633	0.645	0.687	0.717	0.541	0.568	0.615	0.629
	TVDB-T1 (text bag-of-words)	0.554	0.609	0.613	0.654	0.508	0.516	0.553	0.590
	TVDB-T2 (text LSTM)	0.642	0.734	0.761	0.782	0.613	0.606	0.635	0.696
	<b>TVDB (full model)</b>	<b>0.713</b>	<b>0.779</b>	<b>0.787</b>	<b>0.810</b>	<b>0.674</b>	<b>0.697</b>	<b>0.704</b>	<b>0.721</b>

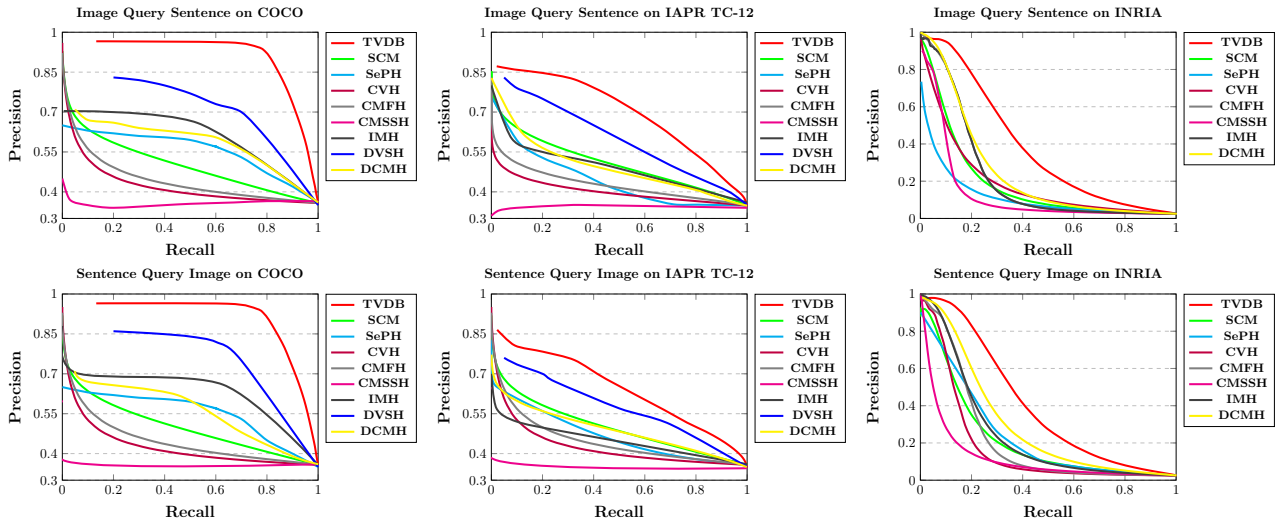


Figure 3. Cross-modal retrieval Precision-Recall curves of TVDB and some existing methods with code length  $M = 32$ .

pre-trained on each dataset. Then a 384-D feature is extracted for each sentence from the pooling layers. In implementing DCMH [23], we build an identical text coding network to ours, enabling it to handle sentence data. For DNH [35], we cite the performances of its variants DNH-C provided in [5] since the same settings are used.

**Results and Analysis.** The retrieval MAP performances on the three datasets are reported in Table 1. In general, the proposed TVDB model outperforms existing methods on the three datasets with large margins. Most traditional cross-modal hashing techniques are not usually designed specifically for images and sentences, which limits their performances on the three datasets. This suggests that image-tag hashing and retrieval are unrepresentative for vision-language tasks. The recent deep hashing model DVSH [5] hits the closest overall figures to ours as its modal-specific deep networks are able to explore the intrinsic semantics of images and sentences. TVDB provides even superior performance since both regional image information and relations between the words are well encoded, making the output binary codes more discriminative. Some results are left blank because the corresponding

baseline codes are not available and the performances are never reported. Note that the overall MAP scores on INRIA Web Queries [31] are relatively lower than those on the other two. This is probably due to the relatively low image quality. The corresponding precision-recall curves with 32-bit code length are given in Figure 3. A *Sentence Query Image* example is given in Figure 4 with top-5 closest retrieved candidates on 128-bit COCO. TVDB provides well-matched results with detailed information preserved (*e.g.* people, bike, street), while the compared methods do not.

### 5.3. Ablation study of deep encoding network

We demonstrate the impact of the deep encoding networks for TVDB in this subsection.

**Baselines.** Four variants of TVDB are built as baselines by modifying the deep encoding networks with some other architectures: **(a) TVDB-I1** is built by replacing the region-based image encoding network of TVDB by a holistic AlexNet CNN. **(b) TVDB-I2** mixes the image region features using average pooling instead of rendering them to the LSTM units. **(c) TVDB-T1** takes text bag-of-words as sentence features with the original text-CNN removed. **(d)**

Query: A group of people on bicycles wait along a street at nighttime.



Figure 4. Intuitive *Sentence Query Image* retrieval top-5 results on Microsoft COCO 128 bits. TVDB carries the best matching candidates, where most of the objects mentioned in the query are included (e.g. people, bike, street).

**TVDB-T2** is a variant of TVDB where the text-CNN is replaced by a two-layer LSTM structure.

**Results and Analysis.** Self-comparison MAP results on cross-modal retrieval are shown in Table 2. As we expected, the MAP scores drop dramatically with simple image CNN (TVDB-I1, TVDB-I2) and bag-of-words features (TVDB-T1), but are still acceptable compared to some existing methods. Image binary encoding without regional information is still far from satisfactory. TVDB-T2 with text-LSTM obtains reasonable performance and is in general superior to the state-of-the-art DVSH [5], since LSTM is also capable of modeling sentences. However, TVDB-T2 performs poorer than the original TVDB, suggesting that the proposed text-CNN architecture is a suitable choice for the cross-modal hashing task. To this end, it can be seen that our proposed binary encoding network is successfully designed and all components are reasonably implemented.

#### 5.4. Training efficiency

The training efficiency and convergency of the proposed *stochastic batch-wise learning routine* are illustrated in this subsection, by comparing it with several baselines.

**Baselines.** (a) **TVDB-S** varies TVDB by keeping in-batch images and sentences unaltered with each epoch. (b) **TVDB-N** is a variant of TVDB where  $\mathbf{B}$  and  $\mathbf{H}$  are initialized by  $\max_{\mathbf{B}, \mathbf{H}} \text{trace}(\mathbf{B}\mathbf{S}\mathbf{H}^T)$  and are not updated during SDG training. Here,  $\mathbf{S}$  refers to the similarity matrix on the whole training set as in Eq. (3). (c) **TVDB-E1** is similar to the optimization of DCMH [23], performing epoch-wise binary code learning instead of batch-wise, i.e., updating both  $\mathbf{B}$ ,  $\mathbf{H}$  in a similar manner to Eq. (8) and (9) on the whole training set after each epoch. (d) **TVDB-E2** is similar to TVDB-E1, but code learning and updating is performed after every five epochs of training instead of each one.

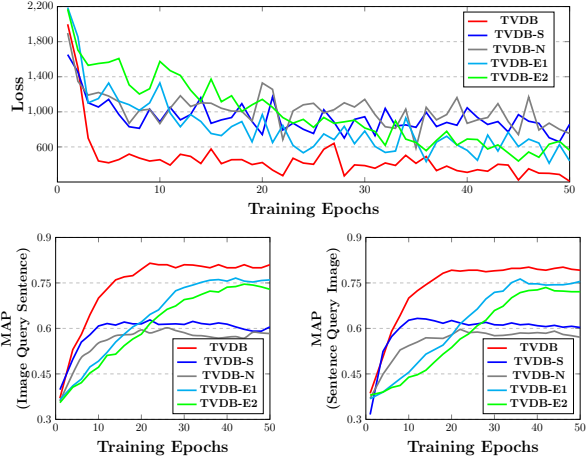


Figure 5. The 128-bit retrieval MAP of *Image Query Sentence* (bottom left) and *Sentence Query Image* (bottom right) w.r.t. training epochs on Microsoft COCO dataset are shown here. The corresponding losses are also given on the top.

**Results and Analysis.** Experiments are conducted on the Microsoft COCO [36] dataset with code length  $M = 128$ . Cross-modal retrieval MAP scores and corresponding learning losses w.r.t. training epochs are shown in Figure 5. It is obvious that TVDB converges quickly to an acceptable MAP score and then gradually hits the best performance at about the 20<sup>th</sup> epoch. TVDB is generally superior to the compared baselines both in terms of peak performance and training efficiency. It can be seen that TVDB-N obtains a similar rate of convergency to TVDB for the first five training epochs, but ends up with significantly lower retrieval performance since the network parameters  $\Theta$  and  $\Phi$  are disjointly optimized with  $\mathbf{B}$  and  $\mathbf{H}$ . TVDB-S follows a close path to TVDB-N with a slightly higher performance. It is clear that our code learning strategy with unaltered in-batch data is not appealing in exploring the generalized optima of  $\mathbf{B}$  and  $\mathbf{H}$  on the whole dataset. TVDB-E1 and TVDB-E2 carry out acceptable retrieval performances but are still outperformed by TVDB. This demonstrates that each aspect of TVDB is necessary to obtain optimal performance.

#### 6. Conclusion

In this paper, we proposed a deep binary encoding method termed as Textual-Visual Deep Binaries (TVDB) which is able to encode information-rich images and descriptive sentences. Two modal-specific binary encoding networks were built using LSTM and text-CNN, leveraging image regional information and semantics between the words to obtain high-quality binary representations. In addition, we proposed a *stochastic batch-wise code learning routine* that performs effective and efficient training. Our experiments justified that both the proposed deep encoding networks and the training routine contribute greatly to the final outstanding cross-modal retrieval performance.



## References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh. Vqa: Visual question answering. In *ICCV*, 2015.
- [3] M. M. Bronstein, A. M. Bronstein, F. Michel, and N. Paragios. Data fusion through cross-modality metric learning using similarity-sensitive hashing. In *CVPR*, 2010.
- [4] Y. Cao, M. Long, J. Wang, and S. Liu. Collective deep quantization for efficient cross-modal retrieval. In *AAAI*, 2017.
- [5] Y. Cao, M. Long, J. Wang, Q. Yang, and P. S. Yu. Deep visual-semantic hashing for cross-modal retrieval. In *ACM SIGKDD*, 2016.
- [6] Y. Cao, M. Long, J. Wang, and H. Zhu. Correlation autoencoder hashing for supervised cross-modal search. In *ICMR*, 2016.
- [7] Y. Cao, M. Long, J. Wang, H. Zhu, and Q. Wen. Deep quantization network for efficient image retrieval. In *AAAI*, 2016.
- [8] A. Conneau, H. Schwenk, L. Barrault, and Y. Lecun. Very deep convolutional networks for natural language processing. *arXiv preprint arXiv:1606.01781*, 2016.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [10] G. Ding, Y. Guo, and J. Zhou. Collective matrix factorization hashing for multimodal data. In *CVPR*, 2014.
- [11] T.-T. Do, A.-D. Doan, and N.-M. Cheung. Learning to hash with binary deep neural network. In *ECCV*, 2016.
- [12] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015.
- [13] F. Feng, X. Wang, and R. Li. Cross-modal retrieval with correspondence autoencoder. In *ACM Multimedia*, pages 7–16, 2014.
- [14] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question. In *NIPS*, 2015.
- [15] R. Girshick. Fast r-cnn. In *ICCV*, 2015.
- [16] Y. Gong and S. Lazebnik. Iterative quantization: A proustian approach to learning binary codes. In *CVPR*, 2011.
- [17] M. Grubinger, P. Clough, H. Müller, and T. Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International Workshop OntoImage*, 2006.
- [18] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, and S.-E. Yoon. Spherical hashing. In *CVPR*, 2012.
- [19] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [20] B. Hu, Z. Lu, H. Li, and Q. Chen. Convolutional neural network architectures for matching natural language sentences. In *NIPS*. 2014.
- [21] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell. Natural language object retrieval. In *CVPR*, 2016.
- [22] Y. Hu, Z. Jin, H. Ren, D. Cai, and X. He. Iterative multi-view hashing for cross media indexing. In *ACM Multimedia*, 2014.
- [23] Q.-Y. Jiang and W.-J. Li. Deep cross-modal hashing. In *CVPR*, 2017.
- [24] J. Johnson, A. Karpathy, and L. Fei-Fei. Denscap: Fully convolutional localization networks for dense captioning. In *CVPR*, 2016.
- [25] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei. Image retrieval using scene graphs. In *CVPR*, 2015.
- [26] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015.
- [27] A. Karpathy, A. Joulin, and F. F. F. Li. Deep fragment embeddings for bidirectional image sentence mapping. In *NIPS*, 2014.
- [28] Y. Kim. Convolutional neural networks for sentence classification. In *EMNLP*, 2014.
- [29] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [30] B. Klein, G. Lev, G. Sadeh, and L. Wolf. Fisher vectors derived from hybrid gaussian-laplacian mixture models for image annotation. In *CVPR*, 2015.
- [31] J. Krapac, M. Allan, J. Verbeek, and F. Juried. Improving web image search results using query-relative classifiers. In *CVPR*, 2010.
- [32] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [33] B. Kulis and K. Grauman. Kernelized locality-sensitive hashing for scalable image search. In *ICCV*, 2009.
- [34] S. Kumar and R. Udupa. Learning hash functions for cross-view similarity search. In *IJCAI*, 2011.
- [35] H. Lai, Y. Pan, Y. Liu, and S. Yan. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*, 2015.
- [36] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [37] Z. Lin, G. Ding, M. Hu, and J. Wang. Semantics-preserving hashing for cross-view retrieval. In *CVPR*, 2015.
- [38] H. Liu, R. Wang, S. Shan, and I. Reid. deep supervised hashing for fast image retrieval. In *CVPR*, 2016.
- [39] L. Liu, Z. Lin, L. Shao, F. Shen, G. Ding, and J. Han. Sequential discrete hashing for scalable cross-modality similarity retrieval. *IEEE Transactions on Image Processing*, 26(1):107–118, 2017.
- [40] L. Liu and L. Shao. Sequential compact code learning for unsupervised image hashing. *IEEE transactions on neural networks and learning systems*, 27(12):2526–2536, 2016.
- [41] L. Liu, F. Shen, Y. Shen, X. Liu, and L. Shao. Deep sketch hashing: Fast free-hand sketch-based image retrieval. In *CVPR*, 2017.
- [42] L. Liu, M. Yu, and L. Shao. Latent structure preserving hashing. *International Journal of Computer Vision*, 122(3):439–457, 2017.

- [43] L. Liu, M. Yu, F. Shen, and L. Shao. Discretely coding semantic rank orders for image hashing. In *CVPR*, 2017.
- [44] W. Liu, J. Wang, R. Ji, Y.-G. Jiang, and S.-F. Chang. Supervised hashing with kernels. In *CVPR*, 2012.
- [45] W. Liu, J. Wang, S. Kumar, and S.-F. Chang. Hashing with graphs. In *ICML*, 2011.
- [46] X. Liu, X. Fan, C. Deng, Z. Li, H. Su, and D. Tao. Multilinear hyperplane hashing. In *CVPR*, 2016.
- [47] M. Long, Y. Cao, J. Wang, and P. S. Yu. Composite correlation quantization for efficient multimodal retrieval. In *ACM SIGIR*, 2016.
- [48] X. Lu, F. Wu, X. Li, Y. Zhang, W. Lu, D. Wang, and Y. Zhuang. Learning multimodal neural network with ranking examples. In *ACM Multimedia*, 2014.
- [49] L. Ma, Z. Lu, L. Shang, and H. Li. Multimodal convolutional neural networks for matching image and sentence. In *ICCV*, 2015.
- [50] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 2015.
- [51] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille. Deep captioning with multimodal recurrent neural networks (m-rnn). In *ICLR*, 2015.
- [52] J. Masci, M. M. Bronstein, A. M. Bronstein, and J. Schmidhuber. Multimodal similarity-preserving hashing. *IEEE transactions on pattern analysis and machine intelligence*, 36(4):824–830, 2014.
- [53] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [54] M. Ou, P. Cui, F. Wang, J. Wang, W. Zhu, and S. Yang. Comparing apples to oranges: a scalable solution with heterogeneous hashing. In *ACM SIGKDD*, 2013.
- [55] M. Raginsky and S. Lazechnik. Locality-sensitive binary codes from shift-invariant kernels. In *NIPS*, 2009.
- [56] M. Rastegari, J. Choi, S. Fakhraei, H. Daumé III, and L. S. Davis. Predictable dual-view hashing. In *ICML*, 2013.
- [57] K. Rayner. Eye movements and attention in reading, scene perception, and visual search. *The quarterly journal of experimental psychology*, 62(8):1457–1506, 2009.
- [58] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.
- [59] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele. Grounding of textual phrases in images by reconstruction. In *ECCV*, 2016.
- [60] F. Shen, C. Shen, W. Liu, and H. Tao Shen. Supervised discrete hashing. In *CVPR*, 2015.
- [61] J. Song, Y. Yang, Y. Yang, Z. Huang, and H. T. Shen. Inter-media hashing for large-scale retrieval from heterogeneous data sources. In *ACM SIGMOD*, 2013.
- [62] N. Srivastava and R. Salakhutdinov. Multimodal learning with deep boltzmann machines. *Journal of Machine Learning Research*, 15:2949–2980, 2014.
- [63] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014.
- [64] B. Thompson. Canonical correlation analysis. *Encyclopedia of statistics in behavioral science*, 2005.
- [65] D. Wang, P. Cui, M. Ou, and W. Zhu. Learning compact hash codes for multimodal representations using orthogonal deep structure. *IEEE Transactions on Multimedia*, 17(9):1404–1416, 2015.
- [66] L. Wang, Y. Li, and S. Lazechnik. Learning deep structure-preserving image-text embeddings. In *CVPR*, 2016.
- [67] W. Wang, B. C. Ooi, X. Yang, D. Zhang, and Y. Zhuang. Effective multi-modal retrieval based on stacked auto-encoders. *Proceedings of the VLDB Endowment*, 7(8):649–660, 2014.
- [68] Y. Wei, Y. Song, Y. Zhen, B. Liu, and Q. Yang. Scalable heterogeneous translated hashing. In *ACM SIGKDD*, 2014.
- [69] Y. Weiss, A. Torralba, and R. Fergus. Spectral hashing. In *NIPS*, 2009.
- [70] H. Wold. Partial least squares. *Encyclopedia of statistical sciences*, 1985.
- [71] B. Wu, Q. Yang, W.-S. Zheng, Y. Wang, and J. Wang. Quantized correlation hashing for fast cross-modal search. In *IJCAI*, 2015.
- [72] F. Wu, Z. Yu, Y. Yang, S. Tang, Y. Zhang, and Y. Zhuang. Sparse multi-modal hashing. *IEEE Transactions on Multimedia*, 16(2):427–439, 2014.
- [73] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan. Supervised hashing for image retrieval via image representation learning. In *AAAI*, 2014.
- [74] W. Xie, Y. Peng, and J. Xiao. Cross-view feature learning for scalable social image analysis. In *AAAI*, 2014.
- [75] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- [76] F. Yan and K. Mikolajczyk. Deep correlation for matching images and text. In *CVPR*, 2015.
- [77] W. Zaremba and I. Sutskever. Learning to execute. *arXiv preprint arXiv:1410.4615*, 2014.
- [78] D. Zhai, H. Chang, Y. Zhen, X. Liu, X. Chen, and W. Gao. Parametric local multimodal hashing for cross-view similarity search. In *IJCAI*, 2013.
- [79] D. Zhang and W.-J. Li. Large-scale supervised multimodal hashing with semantic correlation maximization. In *AAAI*, 2014.
- [80] X. Zhu, Z. Huang, H. T. Shen, and X. Zhao. Linear cross-modal hashing for efficient multimedia search. In *ACM Multimedia*, 2013.
- [81] B. Zhuang, G. Lin, C. Shen, and I. Reid. Fast training of triplet-based deep binary embedding networks. In *CVPR*, 2016.