

Cluster-wise Unsupervised Hashing for Cross-Modal Similarity Search

Lu Wang^a, Jie Yang^{a,*}

^a*Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai, China*

Abstract

Large-scale cross-modal hashing similarity retrieval has attracted more and more attention in modern search applications such as search engines and autopilot, showing great superiority in computation and storage. However, current unsupervised cross-modal hashing methods still have some limitations: (1) many methods relax the discrete constraints to solve the optimization objective which may significantly degrade the retrieval performance; (2) most existing hashing model project heterogeneous data into a common latent space, which may always lose sight of diversity in heterogeneous data; (3) transforming real-valued data point to binary codes always results in abundant loss of information, producing the suboptimal continuous latent space. To overcome above problems, in this paper, a novel Cluster-wise Unsupervised Hashing (CUH) method is proposed. Specifically, CUH jointly performs the multi-view clustering that projects the original data points from different modalities into its own low-dimensional latent semantic space and finds the cluster centroid points and the common clustering indicators in its own low-dimensional space, and learns the compact hash codes and the corresponding linear hash functions. A discrete optimization framework is developed to learn the unified binary codes across modalities under the guidance cluster-wise code-prototypes. The reasonableness and effectiveness of CUH is well demonstrated by comprehensive experiments on diverse benchmark datasets.

Keywords: cross-modal similarity retrieval, multi-view clustering, the cluster-wise code-prototypes, cross-modal hashing,

1. Introduction

Due to the explosive growth of big data with multiple modalities in the form of images, text, and videos on social networks, efficient data analysis has gotten an immediate attention to purify the semantic correlations across different heterogeneous modalities. In other word, when we have relevant data in different modalities endowing the semantic correlation structures, it always is desirable to perform cross-modal search, which retrieves the semantically -similar items across the heterogeneous modalities in response to a query. Taking Wikipedia as an example, we can retrieval images of a relevant query tag, or tags of a relevant query image. Nevertheless, as a result of large-scale databases, heterogeneity, diversity and huge semantic gap, it still remains a great challenge for effective and efficient cross-modal retrieval.

Under the circumstances that the searchable database has large volume or that the similarity measure calcu-

lation between query item and database items is expensive, hashing based methods gets great popularity for its low storage cost, fast searching speed and impressive retrieval performance. Moreover, a hash method will search approximate nearest neighbor (ANN) within the reference database for a query item in many tasks such as machine learning [1, 2], data mining [3, 4] and computer vision [5, 6], which could balance retrieval efficiency against retrieval accuracy. The basic principle for hashing is to transform each high-dimensional data point into compact binary code, making close binary codes for the relevant data samples in different modalities.

In recent time, various kinds of attempts have been investigated for cross-modal hashing, which encodes the correlation structures between different heterogeneous modalities when learning hash function and indexing cross-modal data points in the Hamming space [7–9], [10], [11–15]. These existing cross-modal hashing methods always can be induced to a two-step scheme: first, projected multiple heterogeneous data modalities into a continuous common latent space by optimizing inter-modal coherence, and second, quantize the con-

*Corresponding author

Email address: luwang_16@sjtu.edu.cn, jieyang@sjtu.edu.cn (Jie Yang)

tinuous projections into compact binary codes by sign function. While demonstrating successful performance, there are some limitation in the two-step scheme: first, transformation from real-valued data to discrete binary codes always results in abundant loss of information, producing the suboptimal continuous common latent space and the suboptimal compact binary codes [7, 8]; second, solving the optimization objective by relaxing the discrete constraints which may significantly degrade the retrieval performance with great quantization error [11, 12]; third, projecting heterogeneous data into a common latent space can always lose sight of diversity, which could help to learn better binary codes for cross-modal search to some degree. Hence, how to learn compact binary codes with excellent performance is still a great challenge work. Besides, generally speaking, we can roughly classify existing cross-modal hashing methods into unsupervised ones [7–9], [10] and supervised ones [11–15]. The details are represented in the Section 2.

In this paper, we propose Cluster-wise Unsupervised Hashing (CUH), a novel hash model performing effective and efficient cross-modal retrieval. Technically, CUH jointly performs the multi-view clustering that projects the original data points from different modalities into its own low-dimensional latent semantic space and finds the cluster centroid points and the common clustering indicators in its own low-dimensional space, and learns the compact hash codes and the corresponding linear hash functions. The flowcharts of CUH are shown in Fig. 1. To construct a seamless learning framework, we are inspired by the work of class-wise supervised hashing [19] and the work of re-weighted discriminatively embedded K-means for multi-view clustering [20], and create a co-training framework for learning to hash in the unsupervised case, in which we simultaneously realize the multi-view clustering, the learning of hash codes and the learning of hash functions. These above steps are jointly optimized by a unified learning problem, which could keep both inter-modal semantic coherence and intra-modal similarity when minimizing both the multi-view least-absolute clustering residual and the quantization error. The CUH model can generate one extremely compact unified hash code to all observed modalities of any instance for efficient cross-modal search and could scale linearly to the data point size. The reasonableness and effectiveness of CUH is well demonstrated by comprehensive experiments on diverse benchmark datasets.

We summarize the contributions of this paper as follows.

- 1) We propose a cluster-wise unsupervised hashing method, which constructs a co-training framework for learning to hash. In the framework, we simultaneously realize the multi-view clustering and the learning of hash codes.
- 2) We propose an alternately optimization scheme for solving our model. Besides, we develop a discrete optimization method to jointly learn binary codes and the corresponding hash functions for each modality which can improve the performance.

The remainder of this paper is structured as follows. In Section 2, we briefly overview the related works of cross-modal hashing methods. Section 3 elaborates our proposed cluster-wise unsupervised hashing method, along with an efficient discrete optimization algorithm to tackle this problem. In Section 4, we report the experimental results and extensive evaluations on popular benchmark datasets. Finally, we draw a conclusion in Section 5.

2. Related work

As mentioned above, there are two categories cross-modal hashing methods, i.e. unsupervised and supervised ones. The former ones maximize intra-modality and inter-modality relevance of the features of training data for learning hash functions. Meanwhile, the latter ones can better learn the hash functions and acquire superior performance by further utilizing the available supervised information. Actually, for supervised methods, they usually require label information of the entire data, which is difficult when the database is large-scale. Recently, deep learning based cross-modal hashing methods have attracted increasing attention for their significant performance improvements, where an end-to-end deep learning architecture can give binary codes for different modalities, capturing the intrinsic cross-modal relevance [16–18].

IMH [9], SMMH [8], CMFH [10], LSSH [7] are unsupervised cross-modal hashing methods. Song et al. proposed inter-media hashing (IMH), which maximizes the intra-modality and inter-modality consistencies for learning binary codes [9]. Zhen et al. proposed spectral multi-modal hashing (SMMH), which is an extension of spectral analysis of the correlation matrix to obtain binary hash codes [8]. Ding et al. proposed collective matrix factorization hashing (CMFH) that performs collective matrix factorization to learn unified hash codes [10]. Zhou et al. proposed latent semantic sparse hashing (LSSH), which respectively, utilizes sparse coding

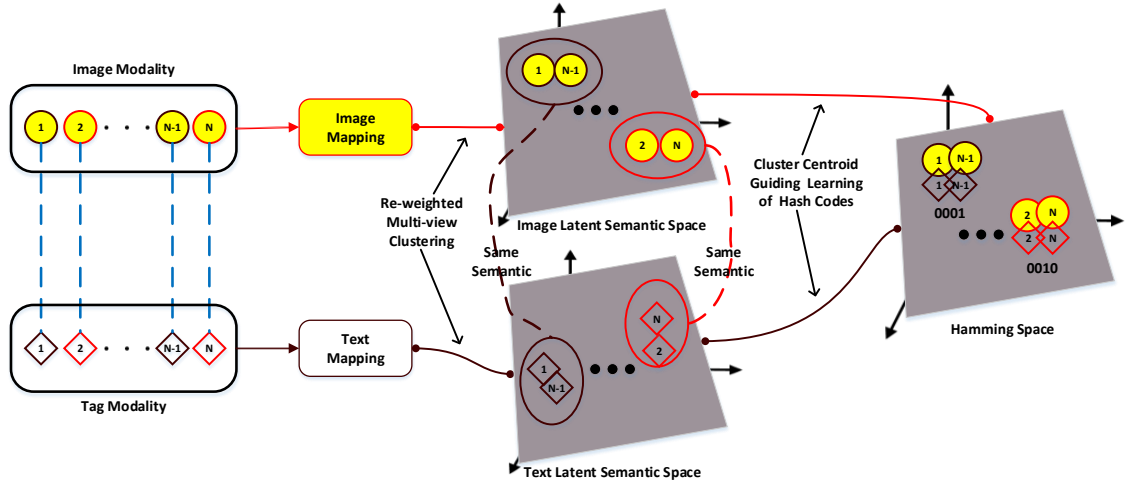


Figure 1: The flowchart of CUH.

for images and matrix factorization for texts to learn their latent semantic features to generate unified hash codes [7].

Differently, CMSSH [15], CVH [13], CRH [14], DCDH [12], SCM [11] are supervised cross-modal hashing methods. Bronatein et al. proposed cross-modality similarity-sensitive hashing (CMSSH) that models the binary classification problems for the projections from features in each modality to hash codes, and utilizes boosting methods to efficiently learn them [15]. Kumar et al. proposed a cross-view hashing (CVH), which is an extension of the single-modal spectral hashing [13]. Zhen et al. proposed co-regularized hashing (CRH) to learn hash function for multi-modal data in a boosted co-regularization framework [14]. Yu et al. proposed discriminative coupled dictionary hashing (DCDH) [12], which learns a coupled dictionary for each modality and unified hash functions. Zhang et al. proposed semantic correlation maximization (SCM) to take semantic labels into consideration for the hash learning procedure in large-scale datasets [11].

In recent years, deep learning methods have acquired great performance improvements on various tasks [16–18]. Inspiring from the advancement of deep learning, many cross-modal hashing methods have developed significant frameworks with deep neural networks, which can bridge the heterogeneous modalities more effectively by overcoming the insufficient character of the hand-crafted features. While these deep models can lead to outstanding performance, there also are some problem such as computational complexity and exhaustive search of learning parameters. Besides, another limitation is that these approaches cannot well reduce

the gap between the Hamming distance and the metric distance on real-valued high-level data representations.

After surveying the existing cross-modal hashing methods, we can clear that well preserving the semantic relevance between instances is the key for reducing the quality loss when retrieving neighbors and achieving better performance. Therefore, it is still desirable to develop a flexible cross-modal retrieval algorithm. Differently, in this paper, the proposed CUH further incorporates the correlations between pairwise Hamming distances to force the to-be-learned hash codes to better preserve the semantic relevance. As will be demonstrated by our experiments, CUH is reasonable and emerges superior performance.

3. Proposed Algorithm

In this section, we will present the detail of the CUH algorithm.

3.1. Notations and Problem Statements

Now, we describe in details the cross-modal retrieval system. Let the two modalities be denoted as $X_1 = [x_1^1, x_1^2, \dots, x_1^N] \in \mathcal{R}^{d_1 \times N}$ and $X_2 = [x_2^1, x_2^2, \dots, x_2^N] \in \mathcal{R}^{d_2 \times N}$, with N being the number of items in either modality and d_1, d_2 being the dimensionality of the data (in general $d_1 \neq d_2$) respectively. Without loss of generality, we assume the input instances in X_1 and X_2 are both zero centered, i.e., $\sum_{i=1}^N (X_v)_i = 0, v = 1, 2$.

Given such data, the goal of CUH is to learn the unified binary codes matrix $B = \{b_i\}_{i=1}^N \in \{+1, -1\}^{r \times N}$ for training instances in both X_1 and X_2 . Besides,

$b_i \in \{+1, -1\}^r$ is the unified r -bits binary codes vector for both instance x_1^i and x_2^i . The modal-specific hash functions aims to map each input instances from corresponding modality to a binary code of r bits through learning r hash functions as follows,

$$\begin{aligned} H_1(x_1^i) &= \text{sgn}(W_1^T x_1^i), \\ H_2(x_2^j) &= \text{sgn}(W_2^T x_2^j), \end{aligned} \quad (1)$$

where $H_1(\cdot)$ and $H_2(\cdot)$ are modal-specific hash functions for image and text modalities, respectively. x_1^i is the i -th input instance from the image modality, x_2^j is the j -th input instance from the text modality. Here, $W_1 \in \mathfrak{R}^{d_1 \times r}$ and $W_2 \in \mathfrak{R}^{d_2 \times r}$ are the linear projection matrices that map the original feature of x_1^i and x_2^j to low-dimensional latent spaces, respectively. The sign function $\text{sgn}(\cdot)$ outputs $+1$ for positive number and -1 otherwise.

3.2. Cluster-wise Unsupervised Hashing

The main framework of CUH is to jointly find the cluster centroid points and the common clustering indicators in its own low-dimensional semantic space and learn the unified hash codes under the guidance of the cluster centroid points, where the cluster centroid points are the cluster-wise code-prototypes to improve the performance of the binary codes.

To realize this mission, we are inspired by the work of re-weighted discriminatively embedded K-means for multi-view (image and text) clustering [20] and induce a robust re-weighted discriminatively embedded K-means for maximizing the inter-modality consistencies to get the cluster centroid points and the common clustering indicators in its own low-dimensional semantic space. In the process for learning of hash codes, we utilize the cluster centroid points as the cluster-wise code-prototypes to guide the learning of the corresponding hash codes of original data for more precise compact codes for semantic information retrieval. We describe how to construct the CUH method under above idea.

3.2.1. Multi-view Clustering

In order to deal with multi-view and high-dimensional data, re-weighted discriminatively embedded K-means proposes an objective function as follows,

$$\begin{aligned} \min_{W_k, F_k, G} & \sum_{k=1}^2 \|W_k^T X_k - F_k G^T\|_F \\ \text{s.t.} & W_k^T W_k = I_{m_k}, k = 1, 2, \\ & G \in \text{Ind}, \end{aligned} \quad (2)$$

where $W_k \in \mathfrak{R}^{d_k \times m_k}$ represents the projection matrix which reduces the dimensionality from d_k to m_k for each

view, $F_k \in \mathfrak{R}^{m_k \times C}$ is the cluster centroid matrix and each column of G denotes the clustering indicator vector for each sample where $G_{ic} = 1 (i = 1, \dots, N; c = 1, \dots, C)$ if the i -th sample belongs to the c -th class and $G_{ic} = 0$ otherwise. Thus, $G \in \text{Ind}$ can be defined, which denotes a set of matrices with above restrictions. For adaptively learning the weights in a re-weighted manner, the objective function can be defined as:

$$\begin{aligned} \min_{W_k, F_k, G, \alpha_k} & \sum_{k=1}^2 \alpha_k \|W_k^T X_k - F_k G^T\|_F^2 \\ \text{s.t.} & W_k^T W_k = I_{m_k}, k = 1, 2, \\ & G \in \text{Ind}, \end{aligned} \quad (3)$$

where $\alpha_k = (2 \|W_k^T X_k - F_k G^T\|_F)^{-1}$ is the weight for the k -th view and can be calculated by current W_k , F_k and G .

3.2.2. Learning of Hash Codes under Cluster-wise Code-prototypes

We can get hash codes from above multi-view clustering in a co-training framework. In the process for learning of hash codes, the dimension reduced data is used as the approximation for the corresponding hash codes of original data. Besides, the cluster centroid points are the cluster-wise code-prototypes. These cluster-wise code-prototypes can guide the learning of the corresponding hash codes of original data to improve the performance of the binary codes. For this goal, we come up with the following objective function,

$$\begin{aligned} \min_B & \sum_{k=1}^2 \|B - W_k^T X_k\|_F^2 - \beta \text{tr}(B^T F_k G^T) \\ \text{s.t.} & B \in \{+1, -1\}^{r \times N}, \end{aligned} \quad (4)$$

where β is the parameter to balance the reconstruction error and the similarity between cluster-wise code-prototypes and binary codes.

3.2.3. Joint Optimization Framework

To approach CUH, which jointly finds the cluster centroid points and the common clustering indicators in its own low-dimensional semantic space and learns the unified hash codes under the guidance of the cluster centroid points, we combine the aforementioned description. That can bring about the objective function of

CUH is written below:

$$\begin{aligned}
\min_{W_k, F_k, G, \alpha_k, B} \quad & \sum_{k=1}^2 (\alpha_k \|W_k^T X_k - F_k G^T\|_F^2 \\
& + \lambda \|B - W_k^T X_k\|_F^2 - \beta \text{tr}(B^T F_k G^T)) \\
\text{s.t.} \quad & W_k^T W_k = I_r, k = 1, 2, \\
& B \in \{+1, -1\}^{r \times N}, \\
& G \in \text{Ind},
\end{aligned} \tag{5}$$

By minimizing (5), the unified hash binary codes will be obtained directly.

3.3. Optimization

To find a feasible solution for the optimization problem (5), in this section, we present an alternating optimization approach.

1) W_k and F_k -step: fix G , α_k and B , update W_k and F_k

By fixing G , α_k and B , the optimization problem (5) becomes

$$\begin{aligned}
\min_{W_k, F_k} \quad & \sum_{k=1}^2 (\alpha_k \|W_k^T X_k - F_k G^T\|_F^2 \\
& + \lambda \|B - W_k^T X_k\|_F^2 - \beta \text{tr}(B^T F_k G^T)) \\
\text{s.t.} \quad & W_k^T W_k = I_r, k = 1, 2.
\end{aligned} \tag{6}$$

Calculating W_k and F_k is a supervised learning stage when we fix G , α_k and B . Firstly, we rewrite Eq.(6) as following Eq.(7) which is very to implement and can be readily used to solve a general trace minimization problem:

$$\begin{aligned}
\min_{W_k} \quad & \sum_{k=1}^2 (\text{tr}(W_k^T M_k W_k) - 2\text{tr}(W_k^T N_k)) \\
\text{s.t.} \quad & W_k^T W_k = I_r, k = 1, 2,
\end{aligned} \tag{7}$$

where $M_k = (\alpha_k + \lambda)X_k X_k^T - \alpha_k X_k G(G^T G)^{-1} G^T X_k^T$, $N_k = \lambda X_k B^T + \frac{\beta}{2} X_k G(G^T G)^{-1} G^T B^T$. Besides, we obtain $F_k = (\frac{\beta}{2\alpha_k} B + W_k^T X_k)G(G^T G)^{-1}$.

We can use the orthogonal constraint optimization procedure in [21], [22]. Through introducing Lagrangian multipliers, we can rewrite the objective function for optimizing $W_k(k = 1, 2)$ as follows:

$$\begin{aligned}
L(W_k, \Lambda) = & \text{tr}(W_k^T M_k W_k) - 2\text{tr}(W_k^T N_k) \\
& - \text{tr}(\Lambda(W_k^T W_k - I)),
\end{aligned} \tag{8}$$

where Λ consists of Lagrangian multipliers. Since $W_k^T W_k$ is symmetric, Λ is symmetric as well. Setting

the gradient of Eq.(8) with respect to W_k to be zero, we can get

$$\frac{\partial L(W_k, \Lambda)}{\partial W_k} = 2(M_k W_k - N_k - W_k \Lambda) = 0. \tag{9}$$

From Eq.(9), it is clear that we can get $\Lambda = W_k^T M_k W_k - W_k^T N_k$. So $\Lambda = W_k^T M_k W_k - W_k^T N_k = W_k^T M_k W_k - N_k^T W_k$ and $\frac{\partial L(W_k, \Lambda)}{\partial W_k} = 2(M_k W_k - N_k - W_k W_k^T M_k W_k + W_k N_k^T W_k)$. Based on the orthogonal constraint optimization procedure in [21], we can define a skew-symmetric matrix $A = 2(M_k W_k W_k^T - N_k W_k^T - W_k^T M_k W_k + N_k^T W_k)$. Then, we will update W_k by Crank-Nicolsonlike scheme [23]

$$W_k^{(t+1)} = W_k^{(t)} - \frac{\tau}{2} A(W_k^{(t+1)} + W_k^{(t)}), \tag{10}$$

where τ is the step size. By solving (10), we can obtain

$$\begin{aligned}
W_k^{(t+1)} &= Q W_k^{(t)}, \\
Q &= (I + \frac{\tau}{2} A)^{-1} (I - \frac{\tau}{2} A).
\end{aligned} \tag{11}$$

Hereafter, we iteratively update W_k several times based on Eq.(11) with Barzilai-Borwein (BB) method [21]. In addition, please note that when iteratively optimizing W_k , the initial W_k is set to be the one optimized in the last round between B and W_k . For the first round, W_k is randomly initialized.

2) G -step: fix W_k , F_k , α_k and B , update G By fixing W_k , F_k , α_k and B , the optimization problem (5) becomes

$$\begin{aligned}
\min_G \quad & \sum_{k=1}^2 (\alpha_k \|W_k^T X_k - F_k G^T\|_F^2 \\
& - \beta \text{tr}(B^T F_k G^T)) \\
\text{s.t.} \quad & G \in \text{Ind}.
\end{aligned} \tag{12}$$

Obtaining the clustering indicator matrix G via a weighted multi-view K-Means clustering is an unsupervised learning stage. We search the optimal solution of G among multiple low-dimensional discriminative subspaces. By separating X_k and G into independent vectors respectively, Eq.(12) can be re-

placed by the following problem:

$$\begin{aligned}
& \min_G \sum_{k=1}^2 (\alpha_k \|W_k^T X_k - F_k G^T\|_F^2 \\
& \quad - \beta \text{tr}(B^T F_k G^T)) \\
& = \min_G \sum_{i=1}^N \sum_{k=1}^2 (\alpha_k \|W_k^T x_k^i - F_k g_i^T\|_F^2 \\
& \quad - \beta b_i^T F_k g_i^T) \\
& = \min_G \sum_{i=1}^N \sum_{k=1}^2 \alpha_k \left\| W_k^T x_k^i + \frac{\beta}{2\alpha_k} b_i - F_k g_i^T \right\|_F^2 \\
& \quad \text{s.t. } G \in \text{Ind}, g_i \in G, \\
& \quad g_{ic} \in \{0, 1\}, \sum_{c=1}^C g_{ic} = 1, \quad (13)
\end{aligned}$$

where g_i is the i -th row of G which denotes the clustering indicator vector for the i -th sample. Moreover, g_{ic} denotes the c -th element of g_i , and there are C candidates to be g_i and each of them is the c -th row of identity matrix I_C :

$$I_C = [e_1, e_2, \dots, e_C], g_i \in \{e_1^T, e_2^T, \dots, e_C^T\}.$$

The one among C candidates making the objective function reach the minimum value is the solution of Eq.(13). We solve Eq.(13) by separating data matrix X_k along the data points direction and assigning C different e_c^T to the row vector g_i one by one independently. Thus, we can tackle the following problem for the i -th sample:

$$c^* = \arg \min_c \sum_{k=1}^2 \alpha_k \left\| W_k^T x_k^i + \frac{\beta}{2\alpha_k} b_i - F_k e_c \right\|_F^2, \quad (14)$$

where c^* means that the c -th element of g_i is 1 and others are 0. There are only C kinds of candidate clustering indicator vector, so we can easily find out the solution of Eq.(13).

- 3) α_k -step: fix W_k, F_k, G and B , update α_k By fixing W_k, F_k, G and B , Updating the non-negative weight α_k for each view assigns the more discriminative image feature with higher weight. The W_k and G in the t -th iteration are computed from the solution of the current iteration. With the current $W_k^{(t)}$ and $G^{(t)}$, we can derive the closed form solution for $\alpha_k^{(t+1)}$ $\mathbf{1}$

$$\alpha_k^{(t+1)} = (2 \|W_k^{(t)T} X_k - F_k^{(t)} G^{(t)T}\|_F)^{-1}. \quad (15)$$

Note that $W_k^{(t)}$ and $G^{(t)}$ are independent of $\alpha_k^{(t)}$ and can be considered as the constants. We iteratively solve $\alpha_k^{(t+1)}$ based on current $W_k^{(t)}$ and $G^{(t)}$.

- 4) B -step: fix W_k, F_k, G and α_k , update B By fixing W_k, F_k, G and α_k , the optimization problem (5) becomes

$$\begin{aligned}
\min_B Q(B) &= \sum_{k=1}^2 (\lambda \|B - W_k^T X_k\|_F^2 \\
& \quad - \beta \text{tr}(B^T F_k G^T)) \\
& \quad \text{s.t. } B \in \{+1, -1\}^{r \times N}. \quad (16)
\end{aligned}$$

To solve this optimization problem, we further rewrite above problem (16) as follows,

$$Q(B) = \sum_{k=1}^2 (\lambda (\|B\|_F^2 + \|W_k^T X_k\|_F^2 - 2 \text{tr}(B^T W_k^T X_k)) - \beta \text{tr}(B^T F_k G^T)).$$

Since $\|B\|_F^2$ and $\|W_k^T X_k\|_F^2$ are both constant, we have

$$\begin{aligned}
Q(B) &= \sum_{k=1}^2 (-2 \lambda \text{tr}(B^T W_k^T X_k) \\
& \quad - \beta \text{tr}(B^T F_k G^T)) + \text{const} \\
&= -2 \text{tr}(V^T B) + \text{const}, \quad (17)
\end{aligned}$$

where $V = \sum_{k=1}^2 (\lambda W_k^T X_k + \frac{\beta}{2} F_k G^T)$. Minimizing $Q(B)$ is equivalent to maximizing $\text{tr}(V^T B)$. As $B \in \{+1, -1\}^{r \times N}$, the optimal solution for (16) can be obtained by setting

$$B = \text{sgn}(\sum_{k=1}^2 (\lambda W_k^T X_k + \frac{\beta}{2} F_k G^T)). \quad (18)$$

To sum up, by these four steps, we can alternatively update W_k, F_k, G, α_k and B and iterate the procedure above until the objective function get a stable minimum value. The process of CUH can be outlined in Algorithm 1.

3.4. Generating Hash Codes for Queries

Given a new query instance $x_k^q (k = 1, 2)$, generating its binary codes b^q depends on its modality. When $x_k^q (k = 1, 2)$ contains data of only one modality, it is straightforward to predict its unified binary codes via the modality-specific hash function. When $x_k^q (k = 1, 2)$ contains data of both two modalities, its unified binary codes are determined by merging the predicted binary codes from different modalities. Thus, the binary codes generation scheme for $x_k^q (k = 1, 2)$ includes the following two situations:

Algorithm 1 Cluster-wise Unsupervised Hashing.

Require: feature matrices X_1 and X_2 , code length r , parameters λ and β .

Ensure: hash codes B , W_1 and W_2 .

- 1: Initialize W_1, W_2 by identity.
 - 2: Initialize $G \in \text{Ind}$ randomly.
 - 3: Initialize binary codes B randomly, such that +1 and -1 are balanced in the codes.
 - 4: Initialize the weight $\alpha_1 = 0.5, \alpha_2 = 0.5$ for each modality.
 - 5: **repeat**
 - 6: Update $W_k, (t = 1, 2)$, by Eq.(11) and obtain $F_k, (t = 1, 2)$, by:
 - 7: $F_k = (\frac{\beta}{2\alpha_k} B + W_k^T X_k) G (G^T G)^{-1}$.
 - 8: Update G by Eq.(14).
 - 9: Update binary codes B by Eq.(18).
 - 10: Update the weight $\alpha_i, (t = 1, 2)$ by Eq.(15).
 - 11: **until** Objective function of Eqn.(5) converges.
-

Only one modality. In this case, we have x_1^q or x_2^q . For $x_k^q (k = 1, 2)$, we directly compute its binary codes b^q as $b^q = \text{sgn}(W_k^T x_k^q)$.

Two modalities. In this case, we both have x_1^q and x_2^q . For CUH, we add up the results computed by the hash functions of two modalities and generate b^q as $b^q = \text{sgn}(W_1^T x_1^q + W_2^T x_2^q)$.

3.5. Complexity Analysis

We discuss the computational complexity of the proposed CUH. In the training phase, the time consuming of each iteration including updating the projection matrices W_1 and W_2 , the cluster centroid matrix F_1 and F_2 , the clustering indicator matrix G , the binary codes B and the weight α_1 and α_2 . Typically, solving Eq.(11), calculating F_1 and F_2 , solving Eq.(14), solving Eq.(18) and solving Eq.(15) require $O(d_k^2 r + d_k^3)$, $O(rN + rCN + C^2 N + C^3 + rC^2)$, $O(rd_k N + rCN + rN + CN)$, $O(rd_k N + rCN + rN)$ and $O(rd_k N + rCN + rN)$. Therefore, the time complexity of each iteration is $O(f_1 d_k^2 + f_2 N + f_3 C^2)$, where $f_1 = \max(r, d_k)$, $f_2 = \max(r, rC, C^2, rd_k, r, C)$ and $f_3 = \max(r, C^2)$. The time complexity of all iterations is $O((f_1 d_k^2 + f_2 N + f_3 C^2)T)$, where T is the number of iterations. It can be observed that the training time is linear to the training set size N . Besides, in the experiments part, we will show that CUH usually only needs few iterations (T is very small) to achieve the best modal parameters. Once the training stage is done, the time and space complexities for generating binary codes for a new query are both $O(d_k r)$ in the query stage, which is extremely efficient. In general, CUH is scal-

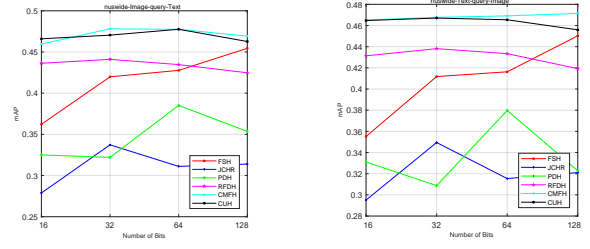


Figure 2: mAP values versus bits on Nuswide.

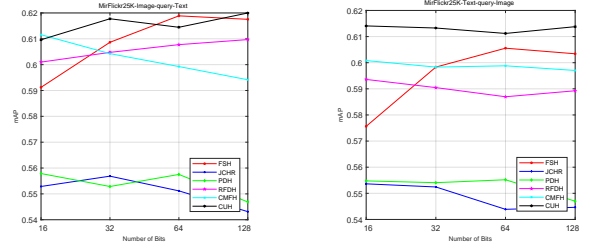


Figure 3: mAP values versus bits on MirFlickr25K.

able for large-scale data sets with most existing cross-modal hashing methods and efficient for encoding new query.

4. Experiments and evaluations

In this section, we execute comprehensive retrieval performance evaluation of CUH on three multimodal benchmark data sets against several state-of-the-art unsupervised cross-modal hashing methods. We present the details of concrete content in data sets, evaluation criteria, comparison methods, and implementation details for the first time. Next, We investigate the experimental results and discussions in terms of fair comparisons. Finally, the convergence and parameter sensitivity of CUH are further reported.

4.1. Data Sets

The effectiveness and efficiency of the proposed CUH model are conducted on three multimodal benchmark data sets: Wiki [27], MIRFlickr25K [28], NUS-WIDE [29]. Specifically, some statistical characteristic of all data sets are depicted in the following.

Wiki [27] consists of 2, 866 image-text pairs collected from Wikipedias articles. It is grouped into 10 semantic categories, where each image-text pair is belong to one of the 10 semantic concepts in the categories. It makes a 128-dimensional bag-of-visual-words vector

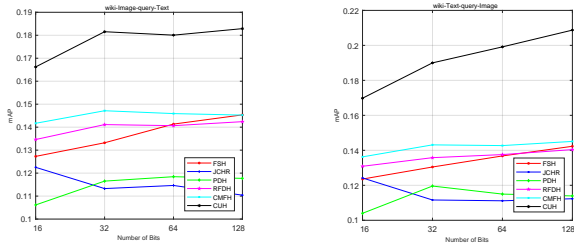


Figure 4: mAP values versus bits on Wiki.

constructed from the SIFT feature to represent every image, and a 10-dimensional topics vector learned by a latent Dirichlet allocation (LDA) model to represent every text. In Wiki, a training set contains 2,173 image-text pairs randomly selected from the whole data set, and a query set contains 693 image-text pairs with the remaining pairs. Besides, training set also is used as the database for retrieval evaluation. In the evaluation, we define the true semantic neighbors for a query through the associated labels.

MIRFlickr25K [28] comprises 25,000 images related to their tags from the Flickr website, in which all image-tag pairs are itemized into the 24 semantic categories. In addition, one pair may have multiple labels from some of the 24 semantic categories above. In the evaluation, we discard image-tag pairs that do not have tags or manually annotated labels and only select tags that appear at least 20 times. By this pretreatment, we can employ a multimodal benchmark of 20,015 image-tag pairs in experiment. It represents an image by a 150-dimensional edge histogram vector, and a text by a 500-dimensional vector extracted by PCA transforming the tag index vector in each pair. In MIRFlickr25K, the query set of 2000 image-tag pairs are randomly taken from the whole data set, and the left pairs are used as the training set, which also serve as the database. In the evaluation, we define the true semantic neighbors for a query as those having at least one label with it at the same time.

NUS-WIDE [29] contains 269,648 images, which are downloaded from real-world website Flickr, it also collecting over 5,000 tags from user. There are 81 concepts fully labeled in the entire data set for performance evaluation. Following [29], we only keep the image-tag pairs belonged to one of the 10 most frequent concepts, and the whole data set is pruned as a new data set comprising 186,577 image-tag pairs. In the experiments, 500-dimensional bag-of-visual-words feature vector is choose to represent image, and an index feature vector of the most common 1,000 tags is choose to represent text. In NUS-WIDE, the query set contains 2,000

image-tag pairs randomly taken from 186,577 image-tag pairs, and the remaining 184,577 pairs are treated as the database. Besides, we randomly selected 5000 image-text pairs as the training set, which is used to learn the hash model. In the evaluation, we define the true semantic neighbors for a query as those having at least one label with it at the same time.

4.2. Evaluation Criteria

To perform a fair evaluation, three widely metrics mAP, topN-precision, and precision-recall are adopt in the evaluation of the retrieval performance for the proposed method and comparison methods. The definitions of these three metrics are as follows:

- (1) mAP. Given a query and a list of R retrieved documents, the value of its average precision (AP) is defined as

$$AP = \frac{1}{N} \sum_{k=1}^R P(k)\delta(r), \quad (19)$$

where N is the number of relevant documents in retrieved set, $P(k)$ denotes the precision of the top k retrieved documents, and $\delta(r) = 1$ if the k -th retrieved document is a true neighbor of the query, and otherwise $\delta(r) = 0$. Then the APs of all queries are averaged to obtain the mAP. R is set to 1000 in the following experiments.

- (2) topN-precision. It expresses the variation of precision with respect to the number of retrieved instances.
- (3) precision-recall. It conveys the precision at different recall level, which can be gotten by changing the Hamming radius of retrieval and evaluating the precision and recall at the same time.

In general, the larger the values of three popular metrics are, the better the performance. Detailed description of the above evaluation criteria can be referred to [26].

4.3. Baseline Methods

The proposed CUH model is compared with the following five state-of-the-art unsupervised multimodal hashing methods: PDH [30], CMFH [10], RFDH [29], FSH [31], JCHR [32]. The parameters in above methods are set according to the corresponding papers.

4.4. Implementation Details

Initialization. Following [25], we will use a stable method to initialize G . Hence, we initialize G as follows:

$$G = \mathbf{1} \otimes Z_C, \quad (20)$$

where $I_C \in \mathbb{R}^{C \times C}$ is a identity matrix and $Z_C \in \mathbb{R}^{C \times C}$ is a binary matrix by randomly sorting the rows of I_C , and $\mathbf{1} \in \mathbb{R}^{[N \div C] \times 1}$ is a column vector with all elements being 1. This method uses direct product of vector $\mathbf{1}$ and matrix Z_C to initialize G . If N cannot be divisible by C , we need to extra select $r = N - C \times \lfloor N \div C \rfloor$ rows from Z_C randomly to fill the indivisible part. This initialization method can make the mapping relationships between different labels of different categories nearly invariable, which can lead to a more stable initialization. In our experiment, for all datasets, we applied this new initialization on all methods.

Parameter setting. The CUH model is related to three model parameters: the quantization error hyper-parameter λ , the cluster-wise code-prototypes regularization hyper-parameter β and the number of cluster centroid points hyper-parameter numCluster. For CUH, the quantization error hyper-parameter λ is set to 10^{-1} , the cluster-wise code-prototypes regularization hyper-parameter β is set to 10^{-4} , and the number of cluster centroid points hyper-parameter numCluster is set to 40 throughout the comparative study. We will study parameter sensitivity in Section 4.7 to validate that CUH can consistently outperform the state of the arts with a wide range of parameter configurations.

4.5. Results and Discussions

In Fig. 2, 3, and 4, the mAP evaluation results are exhibited on all three data sets, i.e. Wiki, MIR-Flickr25K, and NUS-WIDE respectively. From these figures, for all cross-modal tasks (i.e. image-query-text and text-query-image), CUH achieves significantly better result than all comparison methods On Wiki and MIRFlickr25K. Besides, CUH also shows comparable performance with CMFH, outperforming other remaining comparison methods on NUS-WIDE. Superiority of CUH can be attributed to their capability to reduce the effect of information loss, adjust the weights adaptively, as well as avoid the large quantization error. The above observations show the effectiveness of the proposed CUH.

The topN-precision curves with code length 32 bits on all three data sets are demonstrated in Figs. 5, 6, and 7 respectively. From the experimental results, the topN-precision results are in accordance with mAP evaluation values. CUH have better performance than others

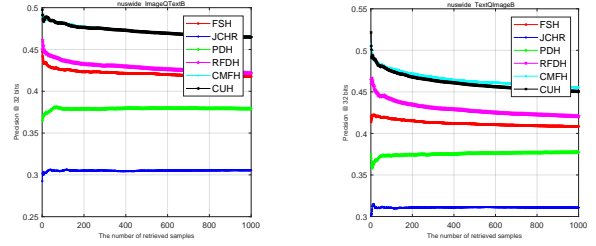


Figure 5: TopN-precision Curves @ 32 bits on Nuswide.

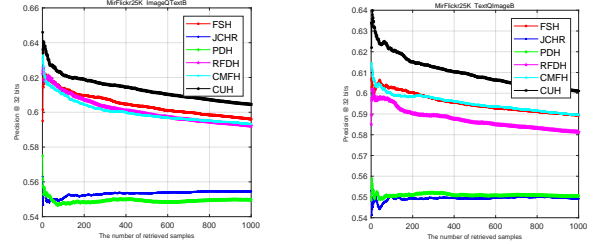


Figure 6: TopN-precision Curves @ 32 bits on Mir-Flickr25K.

comparison methods about cross-modal hashing search tasks on Wiki and MIRFlickr25K. Further more, CUH demonstrates comparable performance with CMFH, outperforming other remaining comparison methods on NUS-WIDE. In retrieval system, we focus more on the front items in the retrieved list returned by search algorithm. Hence, CUH achieves better performance on all retrieval tasks in some sense.

From Figs. 2-7, CUH usually demonstrates large margins on performance when compared with other methods about cross-modal hashing search tasks on Wiki and MIRFlickr25K. At the same time, CUH also exhibits comparable performance with CMFH, better than other remaining comparison methods on NUS-WIDE. We consider two possible reasons, explaining this phenomenon. First, CUH utilizes least-absolute clustering residual in multi-view clustering for learning binary codes, which can be robust to data (i.e. image and text data) outliers and noises. Thus, CUH can achieve improvement on performance. Second, CUH keeps the inter-modal semantic coherence by multi-view clustering, which can extract the high-level hidden semantic features in the image and text. Therefore, CUH could find the common clustering indicators, that reflect the semantic properties more precise. On the consequences, under the guidance of the cluster-wise code-prototypes, CUH can achieve better performance

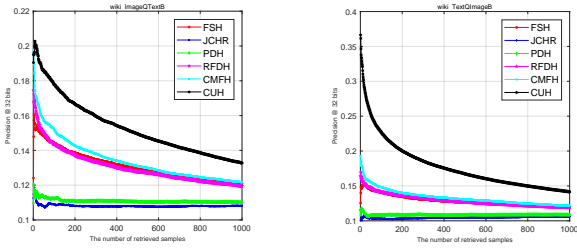


Figure 7: TopN-precision Curves @ 32 bits on Wiki.

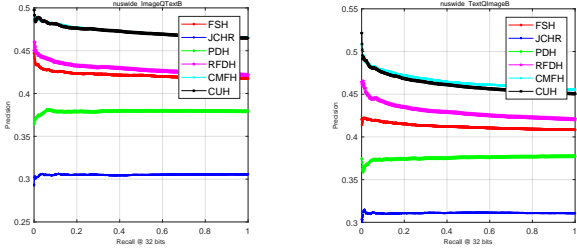


Figure 8: Precision-Recall Curves @ 32 bits on Nuswide.

on cross-modal retrieval tasks.

The precision-recall curves with code length of 32 bits are also demonstrated in Fig. 8, 9, and 10. By calculation of the area under precision-recall curves, we can discover that CUH outperforms comparison methods about cross-modal hashing search tasks on Wiki and MIRFlickr25K. In addition, CUH has comparable performance with CMFH, better than other remaining comparison methods on NUS-WIDE.

4.6. Convergence Analysis

Since CUH is solved in iterative steps, we empirically analyse its convergence property. Fig. 11 demonstrates that the value of the objective (the value is averaged by the number of training data) can fall steadily with the number of iterations. From Fig. 11, we can realize the value of the objective will converge with 15 iterations on all three datasets at 32 bit. This result verifies the effectiveness of Algorithm 1.

4.7. Computational Complexity Analysis

In this section, the train and test time about different cross-modal hashing methods measured in the study. The test time refer to the time that implements an out-of-sample binary code extension for all query and database instances. Our comparison is performed on a PC, which has configuration of 2.20GHz i7-8750H CPU

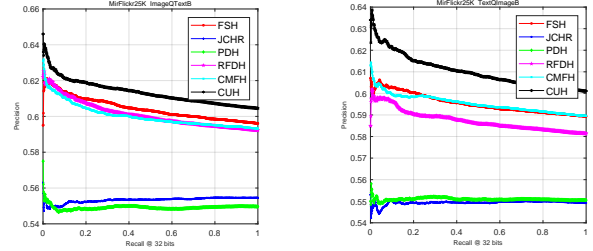


Figure 9: Precision-Recall Curves @ 32 bits on Mir-Flickr25K.

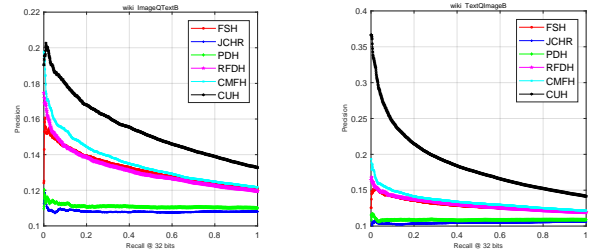


Figure 10: Precision-Recall Curves @ 32 bits on Wiki.

and 16.0GB RAM. We evaluate the time cost of training and testing on the Wiki data set containing 2, 173 training pairs and on the NUS-WIDE data set consisting of 5,000 training pairs. The comparison of training and

Table 1: Training time (s) of different hashing methods on Wiki and NUS-WIDE at 32 bits.

Method	Wiki	NUS-WIDE
	Training time	Training time
FSH	9.1582	21.0636
JCHR	4.4338	9.6238
PDH	25.8539	137.0271
RFDH	33.0569	395.9939
CMFH	4.3356	8.9293
CUH	9.8471	20.8080

testing time complexity at 32 bits is shown in Table. 1 and Table. 2. As shown in Table. 1, RFDH needs most time for learning the model, since it is a two-step learning scheme despite its good performance, in which it first learns binary codes, then trains hash functions. On the other hand, our CUH can learn the hash codes and hash functions in an acceptable speed, which also has better retrieval performance than existing cross-modal hashing methods. Besides, the testing time of comparison in Table. 2 is nearly identical for the compared cross-modal hashing model.

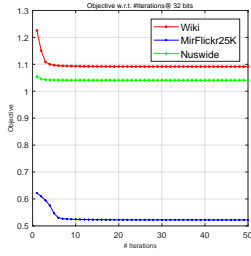


Figure 11: Convergence Analysis.

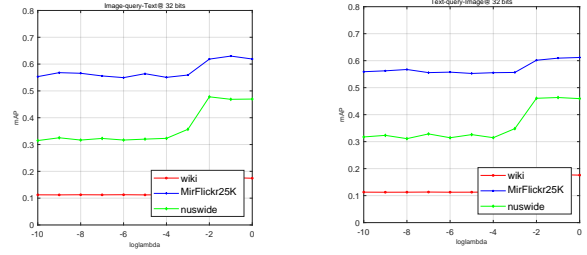


Figure 13: mAP values versus parameter λ .

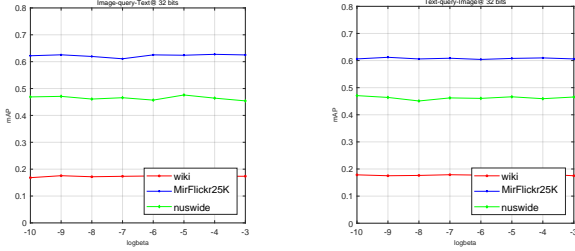


Figure 12: mAP values versus parameter β .

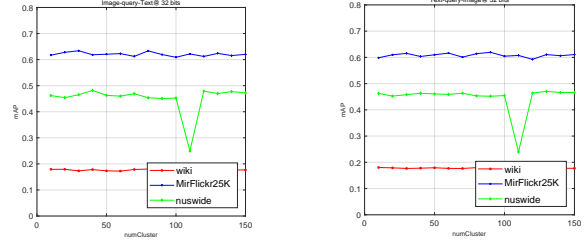


Figure 14: mAP values versus parameter $numCluster$.

4.8. Parameter Sensitivity Analysis

In this section, we conduct the parameter sensitivity to verify that the proposed CUH can achieve stable and superior performance under a large range of parameter values. We test the performance effects about the algorithm in different settings on all datasets. Here, we utilize mAP performance at 32 bits for reporting the variation of performance with respect to parameter values. Our CUH has three hyper-parameters, which include the quantization error hyper-parameter λ , the cluster-wise code-prototypes regularization hyper-parameter β and the number of cluster centroid points hyper-parameter

Table 2: Testing time (s) of different hashing methods on Wiki and NUS-WIDE at 32 bits.

Tasks	Method	Wiki	NUS-WIDE
		Testing time	Testing time
I→T	FSH	0.0506	9.1641
	JCHR	0.0501	9.2376
	PDH	0.0502	9.1492
	RFDH	0.0508	9.1982
	CMFH	0.0495	9.3734
	CUH	0.0497	9.6722
T→I	FSH	0.0529	9.1838
	JCHR	0.0506	9.1974
	PDH	0.0501	9.1726
	RFDH	0.0526	9.2026
	CMFH	0.0510	9.3257
	CUH	0.0508	9.2291

$numCluster$.

The parameter λ balances the reconstruction quantization error and clustering error in the CUH model. It can be observed from Fig. 13 that the performance of CUH goes down slightly when λ increasing. We find CUH can achieve best performance around $\lambda = 10^{-1}$ on all three datasets. Fortunately, when we select λ form the range $[10^{-2}, 1]$, the robust performance of the proposed CUH can be guaranteed.

The parameter β is a hyper-parameter, which balances the cluster-wise code-prototypes regularization and clustering error in the CUH model. From Fig. 12, we can see that Wiki, MirFlickr25K and NUS-WIDE achieve the best around $\beta = 10^{-4}$. Besides, we can observe that CUH achieves stable and superior performance under a large range of β .

The parameter $numCluster$ is a hyper-parameter, which controls the number of cluster centroid points in the CUH model. From Fig. 13, we can see that Wiki, MirFlickr25K and NUS-WIDE achieve the best around $numCluster = 40$. Besides, we can observe that CUH achieves stable and superior performance under a large range of $numCluster$.

5. Conclusions

In this paper, we have formally found a novel way out of cross-modal similarity retrieval task through the pro-

posed cluster-wise unsupervised hashing (CUH) in the unsupervised case. It integrates multi-view clustering and learning of hash codes under the help of cluster-wise code-prototypes, i.e. cluster centroid points in multi-view clustering into an unified binary optimization framework, which generates better compact binary codes that sufficiently contain enough both inter-modal semantic coherence and intra-modal similarity. The binary codes across modalities are learnt under the guidance cluster-wise code-prototypes in its own latent semantic space, which is use for the purpose of key to the efficacy for the proposed CUH method. The reasonableness and effectiveness of CUH is well demonstrated by comprehensive experiments on diverse benchmark datasets. In the future, developing more non-linear mapping models such as boosting or a deep neural network seems an interesting work.

6. References

References

[1] W. Liu, J. Wang, S. Kumar, S.-F. Chang, Hashing with graphs, in: Proceedings of the 28th International Conference on International Conference on Machine Learning, Omnipress, 2011, pp. 1–8.

[2] G. Lin, C. Shen, D. Suter, A. Van Den Hengel, A general two-step approach to learning-based hashing, in: Proceedings of the IEEE international conference on computer vision, 2013, pp. 2552–2559.

[3] W. Xiao, G. Shi, B. Li, J. Xu, F. Wu, Fast hash-based inter-block matching for screen content coding, *IEEE Transactions on Circuits and Systems for Video Technology* 28 (5) (2016) 1169–1182.

[4] C. Wu, J. Zhu, D. Cai, C. Chen, J. Bu, Semi-supervised nonlinear hashing using bootstrap sequential projection learning, *IEEE Transactions on Knowledge and Data Engineering* 25 (6) (2012) 1380–1393.

[5] D. C. Ngo, A. B. Teoh, A. Goh, Biometric hash: high-confidence face recognition, *IEEE transactions on circuits and systems for video technology* 16 (6) (2006) 771–775.

[6] P. Xu, L. Zhang, K. Yang, H. Yao, Nested-sift for efficient image matching and retrieval, *IEEE MultiMedia* 20 (3) (2013) 34–46.

[7] J. Zhou, G. Ding, Y. Guo, Latent semantic sparse hashing for cross-modal similarity search, in: Proceedings of the 37th international ACM SIGIR conference on Research development in information retrieval, ACM, 2014, pp. 415–424.

[8] Y. Zhen, Y. Gao, D.-Y. Yeung, H. Zha, X. Li, Spectral multi-modal hashing and its application to multimedia retrieval, *IEEE Transactions on cybernetics* 46 (1) (2015) 27–38.

[9] J. Song, Y. Yang, Y. Yang, Z. Huang, H. T. Shen, Inter-media hashing for large-scale retrieval from heterogeneous data sources, in: Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, ACM, 2013, pp. 785–796.

[10] G. Ding, Y. Guo, J. Zhou, Collective matrix factorization hashing for multimodal data, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 2075–2082.

[11] D. Zhang, W.-J. Li, Large-scale supervised multimodal hashing with semantic correlation maximization, in: Twenty-Eighth AAAI Conference on Artificial Intelligence, 2014, pp. 2177–2183.

[12] Z. Yu, F. Wu, Y. Yang, Q. Tian, J. Luo, Y. Zhuang, Discriminative coupled dictionary hashing for fast cross-media retrieval, in: Proceedings of the 37th international ACM SIGIR conference on Research development in information retrieval, ACM, 2014, pp. 395–404.

[13] S. Kumar, R. Udupa, Learning hash functions for cross-view similarity search, in: Twenty-Second International Joint Conference on Artificial Intelligence, 2011, pp. 1360–1365.

[14] Y. Zhen, D.-Y. Yeung, Co-regularized hashing for multimodal data, in: Advances in neural information processing systems, 2012, pp. 1376–1384.

[15] M. M. Bronstein, A. M. Bronstein, F. Michel, N. Paragios, Data fusion through cross-modality metric learning using similarity-sensitive hashing, in: 2010 IEEE computer society conference on computer vision and pattern recognition, IEEE, 2010, pp. 3594–3601.

[16] Y. Cao, M. Long, J. Wang, Q. Yang, P. S. Yu, Deep visual-semantic hashing for cross-modal retrieval, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2016, pp. 1445–1454.

[17] Q.-Y. Jiang, W.-J. Li, Deep cross-modal hashing, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 3232–3240.

[18] Y. Cao, M. Long, J. Wang, H. Zhu, Correlation autoencoder hashing for supervised cross-modal search, in: Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval, ACM, 2016, pp. 197–204.

[19] L.-K. Huang, S. J. Pan, Class-wise supervised hashing with label embedding and active bits., in: IJCAI, 2016, pp. 1585–1591.

[20] J. Xu, J. Han, F. Nie, X. Li, Re-weighted discriminatively embedded k -means for multi-view clustering, *IEEE Transactions on Image Processing* 26 (6) (2017) 3016–3027.

[21] Z. Wen, W. Yin, A feasible method for optimization with orthogonality constraints, *Mathematical Programming* 142 (1-2) (2013) 397–434.

[22] B. Wu, Q. Yang, W.-S. Zheng, Y. Wang, J. Wang, Quantized correlation hashing for fast cross-modal search, in: Twenty-Fourth International Joint Conference on Artificial Intelligence, 2015, pp. 3946–3952.

[23] P. Lax, Numerical solution of partial differential equations, *The American Mathematical Monthly* 72 (sup2) (1965) 74–84.

[24] X. Cai, F. Nie, H. Huang, Multi-view k -means clustering on big data, in: Twenty-Third International Joint conference on artificial intelligence, 2013, pp. 2598–2604.

[25] J. Xu, J. Han, F. Nie, Discriminatively embedded k -means for multi-view clustering, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 5356–5364.

[26] H. Schutze, C. D. Manning, P. Raghavan, Introduction to information retrieval, in: Proceedings of the international communication of association for computing machinery conference, 2008, p. 260.

[27] R. He, M. Zhang, L. Wang, Y. Ji, Q. Yin, Cross-modal subspace learning via pairwise constraints, *IEEE Transactions on Image Processing* 24 (12) (2015) 5543–5556.

[28] X. Xu, F. Shen, Y. Yang, H. T. Shen, X. Li, Learning discriminative binary codes for large-scale cross-modal retrieval, *IEEE Transactions on Image Processing* 26 (9) (2017) 1–1.

[29] D. Wang, Q. Wang, X. Gao, Robust and flexible discrete hashing for cross-modal similarity search, *IEEE Transactions on Circuits and Systems for Video Technology* PP (99) (2017) 1–1.

- [30] M. Rastegari, J. Choi, S. Fakhraei, D. Hal, L. Davis, Predictable dual-view hashing, in: International Conference on Machine Learning, 2013, pp. 1328–1336.
- [31] H. Liu, R. Ji, Y. Wu, F. Huang, B. Zhang, Cross-modality binary code learning via fusion similarity hashing, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 7380–7388.
- [32] Y. Liu, Z. Chen, C. Deng, X. Gao, Joint coupled-hashing representation for cross-modal retrieval, in: Proceedings of the International Conference on Internet Multimedia Computing and Service, ACM, 2016, pp. 35–38.