

Central Similarity Multi-View Hashing for Multimedia Retrieval

Jian Zhu^{1†}, Wen Cheng^{1†}, Yu Cui^{1†}, Chang Tang², Yuyang Dai¹, Yong Li¹,
and Lingfang Zeng^{1*}

¹ Zhejiang Lab

{qijian.zhu, chengwen, yu.cui, daiyy, yonglii, zenglf}@zhejianglab.com

² China University of Geosciences

tangchang@cug.edu.cn

Abstract. Hash representation learning of multi-view heterogeneous data is the key to improving the accuracy of multimedia retrieval. However, existing methods utilize local similarity and fall short of deeply fusing the multi-view features, resulting in poor retrieval accuracy. Current methods only use local similarity to train their model. These methods ignore global similarity. Furthermore, most recent works fuse the multi-view features via a weighted sum or concatenation. We contend that these fusion methods are insufficient for capturing the interaction between various views. We present a novel Central Similarity Multi-View Hashing (CSMVH) method to address the mentioned problems. Central similarity learning is used for solving the local similarity problem, which can utilize the global similarity between the hash center and samples. We present copious empirical data demonstrating the superiority of gate-based fusion over conventional approaches. On the MS COCO and NUS-WIDE, the proposed CSMVH performs better than the state-of-the-art methods by a large margin (up to 11.41 mean Average Precision (mAP) improvement).

Keywords: Multi-view Hash · Central Similarity Learning · Multi-modal Hash · Multimedia Retrieval

1 Introduction

Multi-view hashing solves multimedia retrieval problems. The accuracy can be significantly increased with a well-crafted multi-view hashing method. Multi-view hashing, as opposed to single-view hashing, which only searches in a single view, can make use of data from other sources (e.g., image, text, audio, and video). Extraction of heterogeneous features from several views is accomplished via multi-view hashing representation learning. It fuses multi-view features to capture the complementarity of different views.

[†] These authors contributed equally to this work.

* Corresponding author.

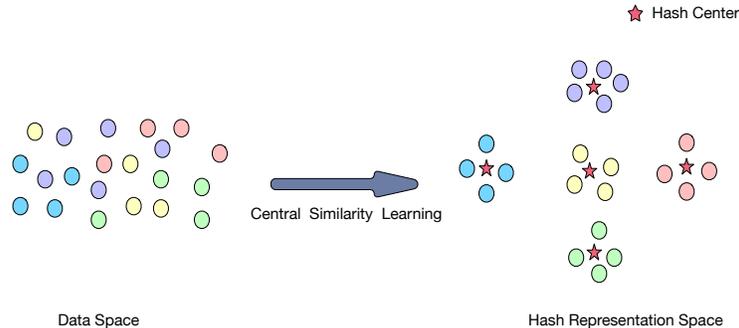


Fig. 1. The inputs are randomly distributed in the data space. Central similarity learning utilizes the hash center to separate samples from different classes while minimizing the class intra-distance.

Retrieval accuracy is currently poor for multi-view hashing methods. The reasons are as follows. First, current methods are fascinated by the information provided by the local similarity. The global similarity is not given enough credit. For example, Flexible Graph Convolutional Multi-modal Hashing (FGCMH) [16] is a GCN-based [22] multi-view hashing method. It constructs the edges of the graph based on similarity. The GCN then aggregates features of adjacent nodes. Hence, the global similarity is not taken into account during this approach. Second, it is not enough to fuse multiple views. To get a global embedding, common multi-view hashing methods (e.g., Deep Collaborative Multi-View Hashing (DCMVH) [29], Flexible Multi-modal Hashing (FMH) [30]) use weighted sum or concatenation to fuse the features. During the fusing process, the relationship between the text and image is disregarded, which results in a lack of expressiveness in the fused features. The aforementioned facts result in low overall retrieval accuracy.

This paper proposes *Central Similarity Multi-View Hashing* method termed CSMVH, which introduces central similarity learning to multi-view hashing. As shown in Fig. 1, samples are distributed randomly in the raw data space. Using central similarity learning, semantically similar samples are close to one another, while semantically different samples are pushed away from each other. This can improve the ability to distinguish sample semantic space. At the same time, because central similarity learning also has the advantage of linear complexity, it is a very effective method.

CSMVH takes advantage of the Gated Multimodal Unit (GMU) [1] to learn the interaction and dependency between the image and text. Different from typical methods, our method fuses multi-view features into a global representation without losing dependency. The semantics-preservation principle of hash representation learning governs the selection of the best embedding space.

We evaluate our method CSMVH on the multi-view hash representation learning tasks of MS COCO and NUS-WIDE datasets. The proposed method

provides up to 11.41% mAP improvement in benchmarks. Our main contributions are as follows:

- We propose a novel multi-view hash method CSMVH. The proposed method achieves state-of-the-art results in multimedia retrieval.
- Central similarity learning is introduced to multi-view hashing for the first time. Our method has lower time complexity than typical pairwise or triplet similarity methods. And it converges faster than previous methods.
- We take advantage of the GMU to learn a better global representation of different views to address the insufficient fusion problem.

2 Related Work

Multi-view hashing [28,14,9,20,2,23,12,18,24,10] fuses multi-view features for hash representation learning. These methods use a graph to model the relationships among different views for hash representation learning. Multiple Feature Hashing (MFH) [20] not only preserves the local structure information of each view but also considers global information during the optimization. Multi-view Alignment Hashing (MAH) [12] focuses on the hidden semantic information and captures the joint distribution of the data. Multi-view Discrete Hashing (MvDH) [18] performs spectral clustering to learn cluster labels. Multi-view Latent Hashing (MVLH) [19] learns shared multi-view hash codes from a unified kernel feature space. Compact Kernel Hashing with Multiple Features (MFKH) [14] treats supervised multi-view hash representation learning as a similarity preserving problem. Different from MFKH, Discrete Multi-view Hashing (DMVH) [24] constructs a similarity graph based on Locally Linear Embedding (LLE) [8,17].

Recently, some deep learning-based multi-view hashing methods are proposed. Flexible Discrete Multi-view Hashing (FDMH) [13], Flexible Online Multi-modal Hashing (FOMH) [15], and Supervised Adaptive Partial Multi-view Hashing (SAPMH) [26] seek for a projection from input space to embedding space using nonlinear methods. The learned embeddings are fused into multi-modal embedding for multi-view hashing. Instead of seeking an embedding space, Deep Collaborative Multi-view Hashing (DCMVH) [29] directly learns hash codes using a deep architecture. A discriminative dual-level semantic method is used for their supervised training. FGCMH [16] is based on a graph convolutional network (GCN). It preserves both the modality-individual and modality-fused structural similarity for hash representation learning. To facilitate multi-view hash at the concept aspect, Bit-aware Semantic Transformer Hashing (BSTH) [21] explores bit-wise semantic concepts while aligning disparate modalities.

Compared with previous multi-view hashing methods, we use central similarity learning and gate-based fusion for the first time. First, current multi-view methods underrate the importance of global similarity and only obtain information provided by local similarity. Therefore, we introduce central similarity learning, which can take advantage of global similarity to make our method

learn more helpful information. Second, current hashing methods fuse multi-view features insufficiently, leading to a weak expressiveness of the obtained global representation. To address this issue, we adopt gate-based fusion to learn the interaction and dependency between the image and text features.

3 The Proposed Methodology

The goal of CSMVH is to use central similarity learning to train a deep multi-view hashing network. We first propose the deep multi-view hashing network, which uses gate-based fusion for the multi-view features. Then the loss of central similarity learning is turned to illustrate. Eventually, CSMVH reduces the computational complexity.

3.1 Deep Multi-View Hashing Network

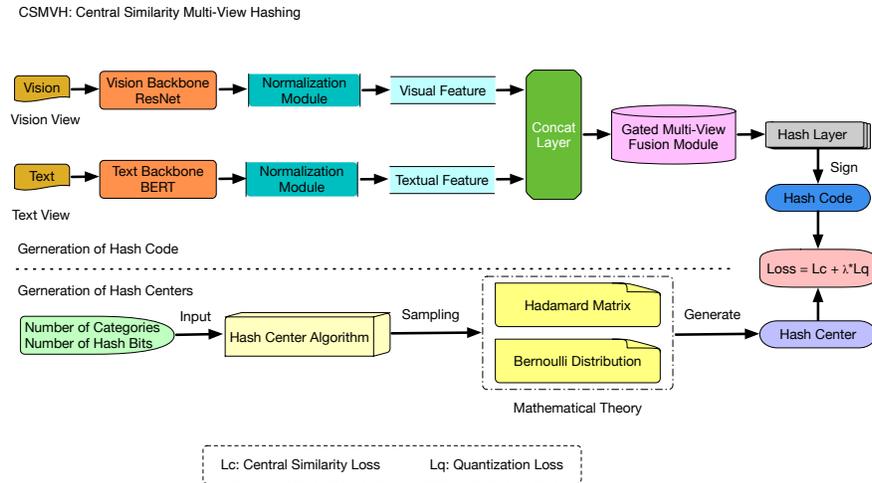


Fig. 2. Architecture Overview of CSMVH. CSMVH has two pipelines, one for the multi-view hash code generation and the other for the hash center generation. In the first pipeline, the image and text features are extracted by the backbones, and then the features are fused by the GMU and hashed by the hash layer. In the second pipeline, the hash center for each class is generated using the Hadamard matrix or Bernoulli distribution. Eventually, CSMVH is optimized by minimizing the distance between the hash code and the corresponding hash center while considering the quantization loss.

The deep multi-view hashing network is made to transform multi-view data into hash code. As shown in Fig. 2, CSMVH comprises a vision backbone, text backbone, normalization modules, multi-view fusion module, and a hash layer.

We utilize deep residual net (ResNet) [6] for image feature extraction and BERT-base [5] for text feature extraction. The extracted features are then passed to the fusion module. These modules are described in detail below. In order to better explain our method, we detailly write the steps of the CSMVH in Algorithm 1.

Algorithm 1 Central Similarity Multi-View Hashing

Input: Training Set $\{\mathbf{x}_i^{(m)}, \mathbf{Y}_i\}_{i=1}^N$, Number of views m , Number of training sets N , Hash Code Length K , Hash Centers $\mathcal{HC} = \{hc_i\}_{i=1}^V \subset \{-1, 1\}^K$, V is the number of categories of Y . Hyper-parameter λ , Training epochs T .

- 1: Initialize BERT parameters w_{BERT} and ResNet parameters w_{ResNet} by loading pre-trained model parameters.
- 2: Randomly initialize Normalization Module parameters w_{vnorm} and w_{tnorm} , Multi-View Fusion Module parameters w_{fusion} and Hash Layer parameters w_{hash} .
- 3: Initialize learning rate, mini-batch size, dropout, and λ .
- 4: **for** $epoch = 1$ to T **do**
- 5: **for** $i = 1$ to N **do**
- 6: Backbones extract embedding from view data.
- 7: Normalize embedding to the same dimensions.
- 8: Fuse multi-view embedding to the global representation.
- 9: Hash Layer generates hash code.
- 10: Update L_{total} by Eq.(13).
- 11: Update w_{vnorm} , w_{tnorm} , w_{fusion} , and w_{hash} by the Adam Optimizer.
- 12: **end for**
- 13: **end for**

Output: Multi-View Hash Network Parameters: \mathbf{W}_K^{hash} .

3.1.1 Vision Encoder We use a deep residual network [6] to extract features. Specifically, we take ResNet-50 as the backbone network for acquiring visual features.

3.1.2 Text Encoder We use BERT-base [5] model as the backbone network for extracting text features. We transform the semantic tag of images into descriptive sentences and then input it to the BERT-base model for feature extraction. Finally, the CLS Token embedding of the BERT model output layer is taken as the semantic vector of the image tag.

3.1.3 Normalization Module The normalization module projects multi-view features into the same dimension, and the output range of features is also within a certain threshold. It is implemented through a fully connected layer network.

3.1.4 Gated Multi-View Fusion Module In multimodal hashing, we need to fuse different information sources so that the fused one can provide more information than a single source. Fusion can be performed at the feature level or decision level. Gate mechanism can combine both feature and decision fusion. In this work, we utilize a gate-based fusion module called GMU. GMU is inspired by GRU [4] and LSTM [7]. It can work as an internal unit finding intermediate representation for multi-view data. The data from different sources is modulated. The gate will learn a prediction from input data to fusion weight. The modulated data is combined using weighted summation. Primarily, for image and text input, the equation governing can be represented as follows:

$$h_i = \tanh(W_i \cdot x_i) \quad (1)$$

$$h_t = \tanh(W_t \cdot x_t) \quad (2)$$

$$z = \sigma(W_z \cdot [x_i, x_t]) \quad (3)$$

$$h_f = z * h_i + (1 - z) * h_t \quad (4)$$

Where x_i and x_t are the image and text input, h_i and h_t are the modulated feature. σ represents the sigmoid activation function. z is the predicted weight for summation. We let $\Theta = \{W_i, W_t, W_z\}$, where Θ is the learnable parameters. h_f is the global vector after multi-view fusion.

3.1.5 Hash Layer The hash layer is a linear layer with a *tanh* activation, which is denoted by:

$$h_{k\text{-bit}} = \text{sgn}[\tanh(w_{\text{hash}}X_{\text{fusion}} + b_{\text{hash}})], \quad (5)$$

where *sgn* is the signum function. $w_{\text{hash}} \in \mathbb{R}^{n \times n}$ and $b_{\text{hash}} \in \mathbb{R}^n$ are network parameters. The output has exactly as many dimensions as the hash code has. Each dimension corresponds to a hash bit.

3.2 Central Similarity Learning

Let $\mathcal{X} = \{\{x_i^{(m)}\}, Y_i\}_{i=1}^N$ denotes the train dataset, where each $x_i^{(m)} \in \mathbb{R}^D$ is a multi-view instance, m is the number of views, and N is the total number of dataset samples. And $Y = \{y_1, y_2, \dots, y_N\}$ is set, where y_i represents the label of x_i . We use the symbol $F: x \mapsto he \in \{-1, 1\}^K$ to signify the deep hash function from the input space \mathbb{R}^D to K-bit Hamming space $\{-1, 1\}^K$. Like other deep hashing methods, we investigate a deep multi-view hashing method, which can generate hash code approximation for data points with the same semantic label in Hamming space.

The use of central similarity learning can cause data points from the same class to cluster together around a single hash center and those from other semantic categories to cluster together around several hash centers, respectively. It makes sense that hash centers maintain a suitable mutual separation from one another, which effectively keeps samples from various classes apart in the Hamming space. High-quality hash codes can be produced through central similarity learning in the deep hash function F by learning the overall similarity information between data pairs.

3.2.1 Hash Center Definition We use a similar definition of hash centers as Central Similarity Quantization (CSQ) [25] with a useful modification. We replace the hash code 0 with -1, thus $\mathcal{HC} = \{hc_i\}_{i=1}^V \subset \{-1, 1\}^K$. V is the number of categories of Y . An elegant linear relationship exists between Hamming distance $\text{dist}_H(\cdot, \cdot)$ and inner product $\langle \cdot, \cdot \rangle$:

$$\text{dist}_H(hc_i, hc_j) = \frac{1}{2}(K - \langle hc_i, hc_j \rangle), \quad (6)$$

where K is the length of hash bits. We can easily have the following:

$$\frac{1}{T} \sum_{i \neq j}^V \langle hc_i, hc_j \rangle \leq 0 \quad (7)$$

V is the number of hash centers, and T is the number of combinations of different hc_i and $hc_j \in \mathcal{HC}$.

3.2.2 Generation of Hash Centers With the definition of hash center, we can naturally notice a boundary scenario:

$$\frac{1}{T} \sum_{i \neq j}^V \langle hc_i, hc_j \rangle = 0 \quad (8)$$

If all hash centers are orthogonal to each other, this equation holds. The most intuitive way to generate orthogonal vectors is by leveraging the Hadamard matrix. The rows of the Hadamard matrix are orthogonal to each other. Thus, we can easily get the hash center with the Hadamard matrix and concatenated transpose. However, we need Bernoulli Distributions for the extra when we need more than $2K$ hash centers.

3.2.3 Loss Function Given the generated centers $\mathcal{HC} = \{hc_i\}_{i=1}^V \subset \{-1, 1\}^K$ in the K -dimensional Hamming space for training data X with V categories, we obtain the semantic hash centers $\mathcal{HC}' = \{hc'_1, hc'_2, \dots, hc'_N\}$ for single-label or multi-label data, where hc'_i denotes the hash center of the data sample x_i . Since hash centers are binary vectors, we use Binary Cross Entropy (BCE) to measure the Hamming distance between the hash code and its center, $\text{dist}_H(hc'_i, hc_i) =$

$BCE(hc'_i, he_i)$. We obtain the optimization objective of the central similarity loss L_c :

$$hc'_{i,k} = \frac{1}{2}(1 + hc'_{i,k}) \quad (9)$$

$$he_{i,k} = \frac{1}{2}(1 + he_{i,k}) \quad (10)$$

$$L_c = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K [hc'_{i,k} \log he_{i,k} + (1 - hc'_{i,k}) \log (1 - he_{i,k})] \quad (11)$$

The output of the neural network is continuous. However, in a Multi-view hashing problem, the output is binary code. Hence, an intuitive method uses a threshold to quantify the network output. Unfortunately, the foundation of the elegant linear relation between Hamming distance and inner product requires the outputs to be exact -1 and 1. Otherwise, it does not hold. Furthermore, quantifying the output using a threshold will introduce an uncontrollable quantization error, which can not pass through the backpropagation. Inspired by DHN [27], we utilize a quantization before constraining the network output. Different from DHN, we use the log cosh function rather than L1-norm for its smoothness, which can be represented as:

$$L_q = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K (\log \cosh(|he_{i,k}| - 1)) \quad (12)$$

Eventually, we have a central similarity optimization problem:

$$L_{total} = L_c + \lambda L_q \quad (13)$$

where L_{total} is the total loss function of our algorithm. λ is the hyper-parameter obtained through grid search in our work.

3.2.4 Computational Complexity Concerning time complexity, CSMVH only has a $O(n)$ value, where n is the number of samples. To achieve similar results, most approaches use pairwise or triplet data similarity, which has a time complexity of $O(n^2)$ or $O(n^3)$. As a result, CSMVH is more effective than the earlier method. For example, Flexible Graph Convolutional Multi-modal Hashing (FGCMH) [16] is a GCN-based multi-view hashing method. It requires an adjacency matrix to represent the relationship between the nodes. Therefore, the complexity of FGCMH is $O(n^2)$. In addition, Eq. (13) can also be concluded that the complexity of our method is $O(n)$.

Through the above analysis, we use central similarity learning to improve the computational complexity, which is that our CSMVH converges faster than the previous methods.

4 Experiments

We evaluate the proposed CSMVH on large-scale multimedia retrieval tasks in experiments. Two genetic datasets are adopted: NUS-WIDE [3] and MS COCO [11]. These datasets have been widely used for evaluating multimedia retrieval performance. We use the mean Average Precision (mAP) as the evaluation metric. The statistics of two datasets used in experiments are summarized in Table 1.

4.1 Evaluation Datasets

Table 1. General statistics of two datasets. The dataset size, number of categories, and feature dimensions are included.

Datasets	Training Size	Retrieval Size	Query Size	Categories	Visual Feature	Text Feature
MS COCO	18000	82783	5981	80	ResNet(768-D)	BERT(768-D)
NUS-WIDE	21000	193749	2085	21	ResNet(768-D)	BERT(768-D)

4.1.1 MS COCO In our experiments, the MS COCO 2014 dataset is adopted. It contains 82,783 training samples and 40,504 validation samples. We randomly select 80 validation samples from each category as the query set to retrieve samples from the training set and use 18000 of them for training.

4.1.2 NUS-WIDE NUS-WIDE contains 269,648 Flickr images with 81 ground-truth semantic concepts. In experiments, we select the 21 most common concepts. We randomly select 100 samples for each concept as the query set. The rest of the images are treated as the retrieval set. We utilize 21,000 samples from the retrieval set for training.

4.2 Evaluation Metrics

To evaluate the metric of multi-view hashing methods, mean Average Precision (mAP) is utilized. The Hamming ranking measure can be assessed with good steady using the mAP. It is described as follows:

$$AP(q) = \frac{1}{M} \sum_{r=1}^R P_q(r) \text{Ind}(r) \quad (14)$$

$$mAP = \frac{1}{Q} \sum_{i=1}^Q AP(q_i) \quad (15)$$

where Q is the number of queries and P_q is the precision for query q when the top r^{th} similar searching results returned, $Ind(r)$ is an indicator function which is 1 when the r^{th} result has the same label with q and otherwise 0, M is the number of same label samples of query q , and R is the size of the retrieval dataset.

4.3 Implementation Details

Our code is implemented on the PyTorch platform. A ResNet-50 pre-trained on ImageNet is employed as the backbone. We use the BERT-base model as the text backbone. The image and text feature outputs are set to 768-dimensional by normalization modules. The dropout probability is set to 0.1 to improve the generalization capability. The combination coefficient λ of the total loss function is 0.25.

4.4 Baseline

To evaluate the retrieval metric, the proposed CSMVH is compared with twelve comparable multi-view hashing methods, including four unsupervised methods (e.g., Multiple Feature Hashing (MFH) [20], Multi-view Alignment Hashing (MAH) [12], Multi-view Latent Hashing (MVLH) [19], and Multi-view Discrete Hashing (MvDH) [18]) and eight supervised methods (e.g., Multiple Feature Kernel Hashing (MFKH) [14], Discrete Multi-view Hashing (DMVH) [24], Flexible Discrete Multi-view Hashing (FDMH) [13], Flexible Online Multi-modal Hashing (FOMH) [15], Deep Collaborative Multi-View Hashing (DCMVH) [29], Supervised Adaptive Partial Multi-view Hashing (SAPMH) [26], Flexible Graph Convolutional Multi-modal Hashing (FGCMH) [16], and Bit-aware Semantic Transformer Hashing (BSTH) [21]).

4.5 Analysis of Experimental Results

The mAP results are presented in Table 2. The experimental comparisons of all methods are conducted according to the unified conditions of the train set, the retrieval set, and the query set in Table 1. The proposed CSMVH is a top performer in multi-view hashing tasks. On the NUS-WIDE dataset, our method outperforms the prior state-of-the-art method by a large margin (up to 3.70%). However, the more complex dataset is where our method shines. MS COCO contains 80 categories, and the samples are more complex than other datasets. Our method beats the previous state-of-the-art 128-bit hashing result with only a 16-bit hash code. Our method shows a magnificent retrieval accuracy improvement with the number of hash bits increasing. The mAP of BSTH [21] only increased 8.23% from 16-bit to 128-bit, while CSMVH has a 17.67% absolute mAP increase, which indicates a great capability of solving very complex hashing problems. The absolute 128-bit mAP increase over BSTH is 11.47% on the MS COCO dataset.

The main reasons for these superior results come from two aspects:

Table 2. The comparable mAP results on NUS-WIDE and MS COCO. The experimental conditions of all methods are the same. The best results are bolded, and the second-best results are underlined. The * indicates that the results of our method on this dataset are statistical significance.

Methods	Ref.	NUS-WIDE*				MS COCO*			
		16 bits	32 bits	64 bits	128 bits	16 bits	32 bits	64 bits	128 bits
MFH	TMM13	0.3603	0.3611	0.3625	0.3629	0.3948	0.3699	0.3960	0.3980
MAH	TIP15	0.4633	0.4945	0.5381	0.5476	0.3967	0.3943	0.3966	0.3988
MVLH	MM15	0.4182	0.4092	0.3789	0.3897	0.3993	0.4012	0.4065	0.4099
MvDH	TIST18	0.4947	0.5661	0.5789	0.6122	0.3978	0.3966	0.3977	0.3998
MFKH	MM12	0.4768	0.4359	0.4342	0.3956	0.4216	0.4211	0.4230	0.4229
DMVH	ICMR17	0.5676	0.5883	0.6902	0.6279	0.4123	0.4288	0.4355	0.4563
FOMH	MM19	0.6329	0.6456	0.6678	0.6791	0.5008	0.5148	0.5172	0.5294
FDMH	NPL20	0.6575	0.6665	0.6712	0.6823	0.5404	0.5485	0.5600	0.5674
DCMVH	TIP20	0.6509	0.6625	0.6905	0.7023	0.5387	0.5427	0.5490	0.5576
SAPMH	TMM21	0.6503	0.6703	0.6898	0.6901	0.5467	0.5502	0.5563	0.5672
FGCMH	MM21	0.6677	0.6874	0.6936	0.7011	0.5641	0.5273	0.5797	0.5862
BSTH	SIGIR22	<u>0.6990</u>	<u>0.7340</u>	<u>0.7505</u>	<u>0.7704</u>	<u>0.5831</u>	<u>0.6245</u>	<u>0.6459</u>	<u>0.6654</u>
CSMVH	Proposed	0.7360	0.7633	0.7740	0.7819	0.6028	0.7002	0.7485	0.7795

Table 3. Ablation experiments on two datasets. The mAP is shown in the table, indicating the performance impact of different modules.

Methods	NUS-WIDE				MS COCO			
	16 bits	32 bits	64 bits	128 bits	16 bits	32 bits	64 bits	128 bits
CSMVH-central	0.7352	0.7594	0.7683	0.7798	0.5892	0.6893	0.7361	0.7731
CSMVH-quant	0.3085	0.3085	0.3085	0.3085	0.3502	0.3502	0.3502	0.3502
CSMVH-text	0.4873	0.5176	0.5196	0.5242	0.5548	0.6134	0.6407	0.6784
CSMVH-image	0.7208	0.7512	0.7623	0.7679	0.5459	0.6454	0.6833	0.7197
CSMVH-concat	0.7283	0.7566	0.7647	0.7718	0.5844	0.6864	0.7282	0.7580
CSMVH	0.7360	0.7633	0.7740	0.7819	0.6028	0.7002	0.7485	0.7795

- Central similarity learning can take advantage of global similarity to make CSMVH learn more useful information. It enhances the discriminative and semantic capability of hash codes. This can improve the accuracy of multi-media retrieval.
- The multi-view fusion module could deeply fuse the multi-view features into a global representation. It can fully explore the complementarity of multi-view features and improve the ability of feature expression. This can generate high-quality hash codes.

4.6 Ablation Studies

To evaluate the proposed CSMVH component by component, we perform ablation with different experiment settings of our method and report the performance. The experiment settings are as follows:

- *CSMVH-central*: The BCE loss of central similarity learning is used. The quantization loss is removed.
- *CSMVH-quant*: The quantization loss is used. The BCE loss of central similarity learning is removed.
- *CSMVH-image*: Only the visual feature is used. The text embedding is removed.
- *CSMVH-text*: Only the text embedding is used. The visual feature is removed.
- *CSMVH-concat*: The image and text features are fused with concatenation. The multi-view fusion module is removed.
- *CSMVH*: Our full framework.

The comparison results are presented in Table 3. Beginning with the loss function, when we solely use the quantization loss to train the proposed CSMVH, it can not learn anything just as we expected. The quantization loss can not perform any optimization on the embeddings. Because the CSMVH obtains data at random, all of the tasks have relatively poor mAP. On the contrary, the BCE loss of central similarity learning can help the CSMVH learn the embedding well. Due to the absence of a quantization constraint, CSMVH-central performs marginally worse than the entire framework.

Moving on to the multi-view features, CSMVH-text is barely superior to random selection (CSMVH-quant). In every task, CSMVH-image performs significantly better than CSMVH-text. It can be seen that the image features contain more useful information than text. Searching with images is more likely to find a related result. Our method already performs better with concatenated multi-view features than the state-of-the-art methods. But the multi-view fusion module takes it one step further, reaching even higher mAP.

To conclude, the multi-view fusion module can improve the complementarity and representation capability of the fused feature. And the central similarity learning enhances the discriminative capability of hash codes. Both image and text features can provide rich information for the multi-view hashing task.

4.7 Convergence Analysis

We do tests to see how well the CSMVH can generalize and converge. Using the MS COCO dataset, we perform hash benchmarks using various code lengths. The results are shown in Fig. 3. The graphic displays training loss and tests mAP. The loss gradually lessens as the training progresses. The loss stabilized at 40 epochs, indicating that the local minimum is attained. When the experiment starts, the mAP for the test metric quickly increases. The test mAP remains steady after 20 epochs. No worsening on the test MAP is seen after additional

training, which suggests high generalization capacity and no overfitting. With other datasets, we notice a similar convergence outcome. The convergence of our method is promised on the popular datasets.

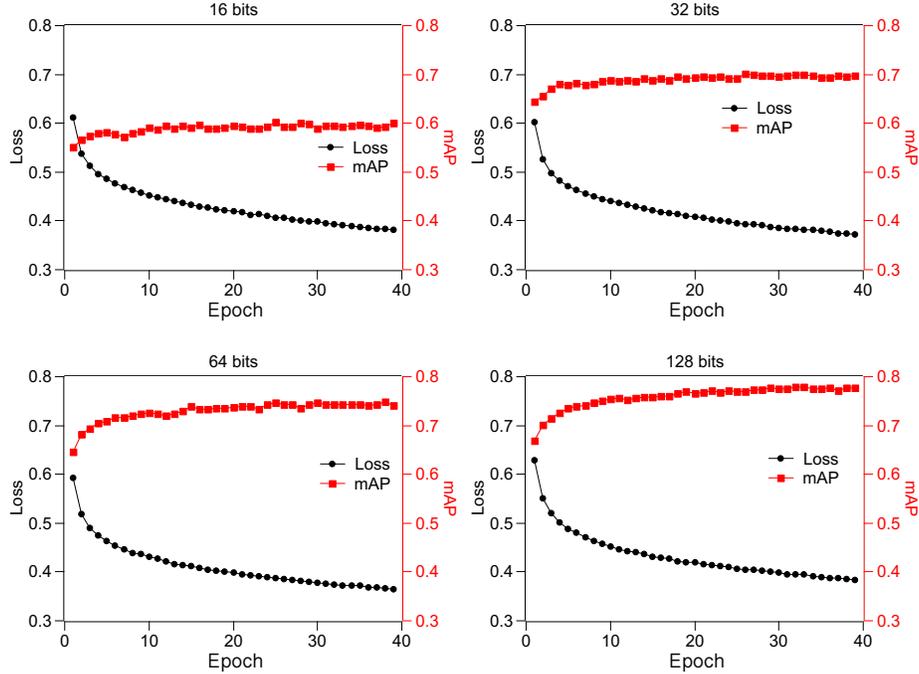


Fig. 3. The training loss curve and test mAP curve on MS COCO dataset. In this figure, the upper curve is the test mAP, and the bottom is the training loss. We notice that after around 40 epochs, our method converges, and the training loss no longer decreases. The test mAP increases rapidly during the first 10 epochs and slowly converges afterward. No overfitting is observed during the experiment.

4.8 mAP@K and Recall@K

The matching mAP@K and Recall@K curves for the MS COCO dataset with increasing numbers of retrieved samples are shown in Fig. 4 for various code lengths. The recall curve exhibits quick linear growth, whereas the mAP of the four graphs somewhat declines as the $TOPK$ grows. The tendency demonstrates how well our method does in the retrieval tasks. The majority of people focus their attention on the first few results of the received data. In this case, the precision of our method is considerably higher. Compared to normal users, experts typically browse more results. With respect to the expansion of

the retrieved data, our method can deliver a recall that grows linearly. When conducting their search, professionals may anticipate dependable, high-quality outcomes. The CSMVH can purposely give satisfying retrieval results for various user groups.

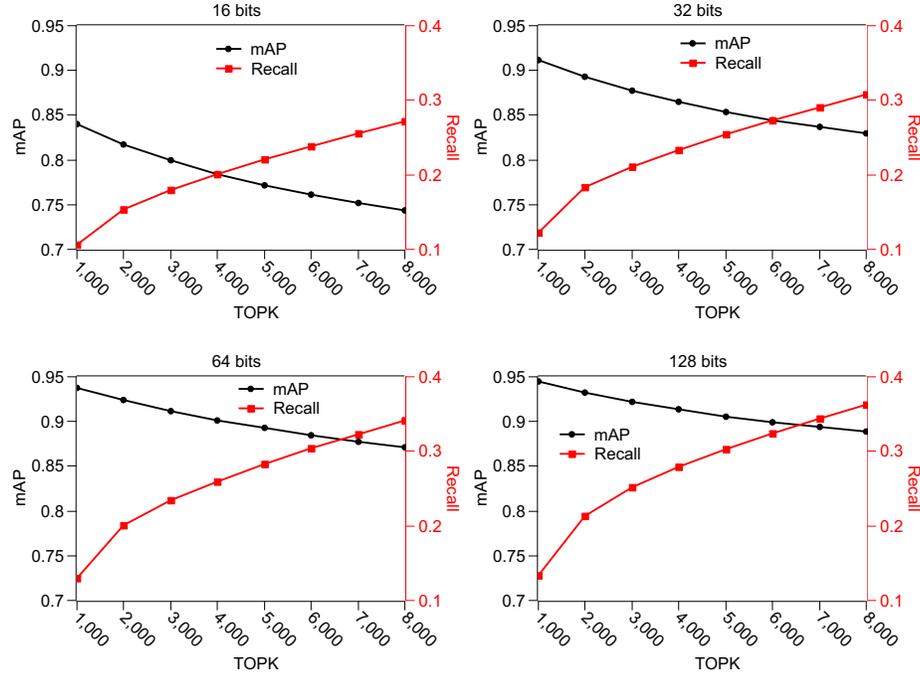


Fig. 4. The mAP@K and Recall@K curves on the MS COCO dataset. The recall curves show a nice linear increase as K increases, while the mAP only slightly decreases. We notice some sweet spots during the experiment. For example, in the 128-bit group, when $TOPK = 5000$, the mAP keeps the same as $TOPK = 4000$, but the recall keeps linearly increasing.

5 Conclusion and Future Work

We propose a new deep multi-view hashing method termed CSMVH. To address multi-view hashing issues, it makes use of central similarity learning. The CSMVH is proposed to conquer two main challenges of the multi-view hashing problem. In our experiment, CSMVH tends to yield consistent improvement in mAP with the growing length of hash code without any signs of performance degradation or overfitting. CSMVH is less computationally intensive than cur-

rent methods. Under multiple experiment settings, it delivers an impressive performance increase over the state-of-the-art methods. Impressively, our method performs even better in more complex hashing tasks. We also find some problems during the process of our experiments. For instance, the performance improvement becomes less and less noticeable as the hash code length increases for some datasets. To further enhance the proposed method, we will work on these problems.

6 Acknowledgment

This work is supported in part by the Zhejiang provincial “Ten Thousand Talents Program” (2021R52007), the National Key R&D Program of China (2022YFB4500405), and the Science and Technology Innovation 2030-Major Project (2021ZD0114300).

References

1. Arevalo, J., Solorio, T., Montes-y Gómez, M., González, F.A.: Gated multimodal units for information fusion. arXiv preprint arXiv:1702.01992 (2017)
2. Chen, Y., Zhang, N., Yan, J., Zhu, G., Min, G.: Optimization of maintenance personnel dispatching strategy in smart grid. *World Wide Web* **26**(1), 139–162 (2023)
3. Chua, T.S., Tang, J., Hong, R., Li, H., Luo, Z., Zheng, Y.: Nus-wide: a real-world web image database from national university of singapore. In: *Proceedings of the ACM international conference on image and video retrieval*. pp. 1–9 (2009)
4. Chung, J., Gulcehre, C., Cho, K., Bengio, Y.: Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555 (2014)
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
7. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
8. Hou, Y., Zhang, P., Xu, X., Zhang, X., Li, W.: Nonlinear dimensionality reduction by locally linear inlaying. *IEEE transactions on neural networks* **20**(2), 300–315 (2009)
9. Kang, Y., Kim, S., Choi, S.: Deep learning to hash with multiple representations. In: *2012 IEEE 12th International Conference on Data Mining*. pp. 930–935. IEEE (2012)
10. Kim, S., Choi, S.: Multi-view anchor graph hashing. In: *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*. pp. 3123–3127. IEEE (2013)
11. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *European conference on computer vision*. pp. 740–755. Springer (2014)

12. Liu, L., Yu, M., Shao, L.: Multiview alignment hashing for efficient image search. *IEEE Transactions on image processing* **24**(3), 956–966 (2015)
13. Liu, L., Zhang, Z., Huang, Z.: Flexible discrete multi-view hashing with collective latent feature learning. *Neural Processing Letters* **52**(3), 1765–1791 (2020)
14. Liu, X., He, J., Liu, D., Lang, B.: Compact kernel hashing with multiple features. In: *Proceedings of the 20th ACM international conference on multimedia*. pp. 881–884 (2012)
15. Lu, X., Zhu, L., Cheng, Z., Li, J., Nie, X., Zhang, H.: Flexible online multi-modal hashing for large-scale multimedia retrieval. In: *Proceedings of the 27th ACM international conference on multimedia*. pp. 1129–1137 (2019)
16. Lu, X., Zhu, L., Liu, L., Nie, L., Zhang, H.: Graph convolutional multi-modal hashing for flexible multimedia retrieval. In: *Proceedings of the 29th ACM International Conference on Multimedia*. pp. 1414–1422 (2021)
17. Saul, L.K., Roweis, S.T.: Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of machine learning research* **4**(Jun), 119–155 (2003)
18. Shen, X., Shen, F., Liu, L., Yuan, Y.H., Liu, W., Sun, Q.S.: Multiview discrete hashing for scalable multimedia search. *ACM Transactions on Intelligent Systems and Technology (TIST)* **9**(5), 1–21 (2018)
19. Shen, X., Shen, F., Sun, Q.S., Yuan, Y.H.: Multi-view latent hashing for efficient multimedia search. In: *Proceedings of the 23rd ACM international conference on Multimedia*. pp. 831–834 (2015)
20. Song, J., Yang, Y., Huang, Z., Shen, H.T., Luo, J.: Effective multiple feature hashing for large-scale near-duplicate video retrieval. *IEEE Transactions on Multimedia* **15**(8), 1997–2008 (2013)
21. Tan, W., Zhu, L., Guan, W., Li, J., Cheng, Z.: Bit-aware semantic transformer hashing for multi-modal retrieval. In: *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. pp. 982–991 (2022)
22. Welling, M., Kipf, T.N.: Semi-supervised classification with graph convolutional networks. In: *J. International Conference on Learning Representations (ICLR 2017)* (2016)
23. Xu, D., Chen, Y., Cui, N., Li, J.: Towards multi-dimensional knowledge-aware approach for effective community detection in lbsn. *World Wide Web* pp. 1–24 (2022)
24. Yang, R., Shi, Y., Xu, X.S.: Discrete multi-view hashing for effective image retrieval. In: *Proceedings of the 2017 ACM on international conference on multimedia retrieval*. pp. 175–183 (2017)
25. Yuan, L., Wang, T., Zhang, X., Tay, F.E., Jie, Z., Liu, W., Feng, J.: Central similarity quantization for efficient image and video retrieval. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3083–3092 (2020)
26. Zheng, C., Zhu, L., Cheng, Z., Li, J., Liu, A.A.: Adaptive partial multi-view hashing for efficient social image retrieval. *IEEE Transactions on Multimedia* **23**, 4079–4092 (2020)
27. Zhu, H., Long, M., Wang, J., Cao, Y.: Deep hashing network for efficient similarity retrieval. In: *Proceedings of the AAAI conference on Artificial Intelligence*. vol. 30 (2016)
28. Zhu, J., Huang, Z., Ruan, X., Cui, Y., Cheng, Y., Zeng, L.: Deep metric multi-view hashing for multimedia retrieval. *arXiv preprint arXiv:2304.06358* (2023)

29. Zhu, L., Lu, X., Cheng, Z., Li, J., Zhang, H.: Deep collaborative multi-view hashing for large-scale image search. *IEEE Transactions on Image Processing* **29**, 4643–4655 (2020)
30. Zhu, L., Lu, X., Cheng, Z., Li, J., Zhang, H.: Flexible multi-modal hashing for scalable multimedia retrieval. *ACM Transactions on Intelligent Systems and Technology (TIST)* **11**(2), 1–20 (2020)