

# Code Consistent Hashing Based on Information-Theoretic Criterion

Shu Zhang, Jian Liang, Ran He, *Senior Member, IEEE*, and Zhenan Sun, *Member, IEEE*

**Abstract**—Learning based hashing techniques have attracted broad research interests in the Big Media research area. They aim to learn compact binary codes which can preserve semantic similarity in the Hamming embedding. However, the discrete constraints imposed on binary codes typically make hashing optimizations very challenging. In this paper, we present a code consistent hashing (CCH) algorithm to learn discrete binary hash codes. To form a simple yet efficient hashing objective function, we introduce a new code consistency constraint to leverage discriminative information and propose to utilize the Hadamard code which favors an information-theoretic criterion as the class prototype. By keeping the discrete constraint and introducing an orthogonal constraint, our objective function can be minimized efficiently. Experimental results on three benchmark datasets demonstrate that the proposed CCH outperforms state-of-the-art hashing methods in both image retrieval and classification tasks, especially with short binary codes.

**Index Terms**—Supervised hashing, binary codes, code consistent constraint, information-theoretic criterion

## 1 INTRODUCTION

HIGH dimensional big Media data like audios, images and videos are growing rapidly nowadays. Emerging with these increasingly growing volume of data is the need to retrieve relevant contents from such large databases. The fundamental scientific problem behind this need is the nearest neighbor (NN) search problem. However, finding the exact nearest neighbors of a query point from a very large database with  $N$   $d$ -dimensional points is not feasible especially when  $N$  and  $d$  tend to be very large [1]. In order to alleviate this issue, some efforts have been made to reduce the search space, such as *KD-tree* [2], or to reduce the computational cost in calculating the similarities, like hashing.

Hashing methods aim to learn compact binary codes with hamming distance computation. They map the high-dimensional data to a binary embedding while preserving some predefined similarity in the original space (e.g., euclidean space or semantic space). Then bit-wise XOR operations are performed to calculate the individual similarities. Due to the reduced time cost of the Hamming distance computation and the low storage overhead to store the binary codes, hashing has become one of the most popular approximate nearest neighbor (ANN) search techniques for many large scale

computer vision applications, including content-based image retrieval [3], [4], object recognition [5] and image matching [6]. One representative method in the earlier hashing work is Locality-Sensitive Hashing (*LSH*) [7]. Since *LSH* and its variants [8] employ randomly generated projections as their hash functions without exploring data distribution, they often require long bit codes to achieve satisfactory performance.

Recently, learning based hashing methods draw much attention because these methods exploit data distribution and even additional supervisory information. These methods employ statistical learning to directly learn hash functions from the data instead of using randomly generated projections. Various methods have been developed in the literature of learning based hashing. Generally, according to whether and how supervisory information is used, learning based hashing methods can be divided into three main categories, i.e., unsupervised hashing [9], [10], [11], [12], [13], semi-supervised hashing (SSH) [14], [15] and supervised hashing [16], [17], [18], [19], [20]. In unsupervised hashing, iterative quantization (*ITQ*) [9] and hashing with graphs (*AGH*) [10] resort to some information-theoretic criteria as their learning objectives meanwhile taking no account of the external supervisory information. A common criterion is to maximize the information entropy on each bit (which favors a more balanced distribution) and minimize the correlation among them. In [21], deep neural networks are employed to learn nonlinear hash functions with the previously mentioned criteria.

In contrast, semi-supervised and supervised hashing methods explore the external supervisory information for better search of semantically similar neighbors, which is more compatible with the image retrieval tasks in real application scenario. Note that, supervisory information emerges in multiple forms. Hence, how to efficiently utilize the supervisory information effectively becomes very important. Most existing supervised hashing algorithms apply the pairwise similarity matrix as their supervisory information,

- S. Zhang and J. Liang are with the Center for Research on Intelligent Perception and Computing (CRIPAC), National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA), No. 95 ZhongGuanCun East Street, HaiDian District, Beijing 100190, P.R.China. E-mail: {shu.zhang, jian.liang}@nlpr.ia.ac.cn.
- R. He and Z. Sun are with the Center for Research on Intelligent Perception and Computing (CRIPAC), National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA), and the Center for Excellence in Brain Science and Intelligence Technology (CEBSIT), Chinese Academy of Sciences (CAS). E-mail: {rhe, znsun}@nlpr.ia.ac.cn.

Manuscript received 5 July 2015; revised 1 Oct. 2015; accepted 29 Oct. 2015. Date of publication 10 Nov. 2015; date of current version 18 Dec. 2015.

Recommended for acceptance by J. Wang, G.-J. Qi, N. Sebe, and C. Aggarwal. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TBDA.2015.2499191

e.g., the graph Laplacians constrained hashing methods, such as spectral hashing (SH) [22], discrete graph hashing (DGH) [23] and (AGH)[10]. More recently, semantic label information is directly exploited in [9], [24], [25]. Particularly, in supervised discrete hashing (SDH) [24], the learnt binary codes are demanded to be optimal for linear classification. In addition, [26] utilizes the relative ranking information as supervision and optimizes various retrieval evaluation criteria such as mean average precision (mAP) and Normalized Discounted Cumulative Gain (NDCG) [26] directly. Recently, several deep learning based hashing methods [27], [28], [29], [30] are also proposed that exploits multiple forms of supervisory information.

In general, most of the existing hashing methods can be formulated as mix-integer optimization problems. These problems are still NP-hard due to the involvement of discrete constraints on binary codes. The supervised learning framework, which consists of discrete constraints and the empirical loss via supervisory information, is very challenging to optimize and thus satisfactory results are rather hard to obtain. To simplify the optimization, most of the aforementioned approaches usually adopt a spectral relaxation [22] on the thresholding (sign operator) procedure, leading to suboptimal solutions as pointed out by [24]. To alleviate this problem, SDH [24] proposed to directly learn binary codes without relaxations and provide a tractable and scalable solver for the formulated discrete optimization problem, making a breakthrough towards the binary optimization in the hashing literature. However, the way in SDH to leverage the supervisory information might not be optimal. It is still necessary to develop new algorithms that can well capture the supervisory information meanwhile have an efficient optimization process. For a brief review of the literature of hashing techniques, readers can refer to a recent review of hashing methods in [31].

This paper presents a code consistent hashing (CCH) algorithm to learn binary codes for image retrieval. Our basic motivation is that learning hash function in an information-theoretic way may yield high-quality binary codes as well as improve existing supervised hashing schemes. We introduce a new code consistent constraint to leverage the label information and the Hadamard code favored by the information-theoretic criterion. The supervisory label information is mainly exploited with fixed ‘discriminative’ binary prototypes that are induced from information entropy maximization. Specifically, we demand that the data within the same class should be mapped into a unique binary code (termed as the class prototype) in the Hamming embedding. At the same time, from the view of information theory, different prototypes representing particular classes should be as far away from each other as possible. We discover that the Hadamard codes [32] are the perfect candidates for the class prototypes, which are equidistant and potentially balanced and uncorrelated.

To formulate such motivation, we relax the aforementioned ‘unique binary code’ constraints with a Frobenius Norm, so that it is more mathematically sound and can allow for small intra-class variations in the learnt codes. In our formulation, we do not expect the learnt binary codes can be directly used for linear classification in the Hamming space as in [24], but it turns out that the proposed code consistent

term does improve the classification accuracy. Besides, for better generalization performance, we seek linear hashing functions in a kernel space as in [17], [24]. Taking all these factors into consideration, we propose a joint learning framework for supervised hashing with discrete variables. To efficiently optimize the mixed-integer problem, we keep the discrete constraint of binary codes in the optimization process as in [24], meanwhile taking advantage of orthogonal constraints on the transformation matrix. As a result, we can develop a much faster optimization process. Extensive experiments on several benchmark datasets are conducted to verify the effectiveness and efficiency of the proposed method on both retrieval and classification tasks.

There are three major contributions of this work:

1. Based on the information-theoretic criterion (uncorrelated and balanced bits), a new code consistent term with Hadamard code is introduced for supervised hashing. It not only yields high-quality binary codes but also reduces the model complexity of supervised hashing schemes.

2. Compared to discrete hashing methods [24], [33], we handle the discrete constraints (NP-hard in general) by introducing an orthogonal constraint to the transformation matrix in CCH, which reduces the discrete cyclic coordinate descent (DCC) method [24] to a one-step thresholding operation (sign function) in each iteration.

3. Compared to several state-of-the-art hashing methods, CCH achieves better results on both retrieval and classification tasks. Particularly, its training time is comparable to the simplest CCA-ITQ and is six times faster than SDH.

The rest of this paper is organized as follows: Section 2 introduces some preliminaries related to our method. Then Section 3 describes our formulation and derives the optimization procedure. Next, Section 4 evaluates our method on two datasets. Finally, we end up with some conclusions in Section 5.

## 2 PRELIMINARIES

For the sake of clarity, we first review some preliminaries about the learning based hashing and class prototypes. Meanwhile, the motivation of our proposed work is slightly touched in this section.

### 2.1 Learning Based Hashing

Given a set of training data  $X = \{x_i\}_{i=1}^N \in \mathbb{R}^{d \times N}$ , the task of hashing is to learn a set of hash functions  $[h_1(x), h_2(x), \dots, h_k(x)]$ , each of which can map the  $d$ -dimensional input onto a binary code bit  $\{-1, 1\}$ , where  $k$  is the code length in the binary embedding. Therefore the learnt binary codes for  $X$  are denoted as  $B = \{b_i\}_{i=1}^N \in \{-1, 1\}^{k \times N}$ , with each column  $b_i$  representing the  $k$ -bits binary codes for  $x_i$ .

For learning based hashing methods, the hash functions are learnt from data  $X$  with additional supervisory information. For instance, suppose the supervisory information appears as semantic labels, that is, each training sample  $x_i$  is associated with a label  $y_i \in \{1, 2, \dots, L\}$ . If we treat the learnt binary codes  $b_i$  as a feature for classification, and enforce that the learnt binary codes should be optimal for classification tasks in the binary embedding as in [24], then a joint optimization framework for binary codes classification and hash function learning can be derived as follows:

$$\begin{aligned} \arg \min_{F, W, B} \sum_{i=1}^n \ell\{y_i, f(b_i, W)\} + \lambda \|W\|_F^2 \\ \text{s.t. } b_i = \text{sgn}(F(x_i)), i = 1, \dots, n, \end{aligned} \quad (1)$$

where  $W \in \mathbb{R}^{L \times k}$  is the model parameter for a linear multi-class classifier,  $\ell$  is the loss function for classification,  $\lambda$  is the regularization parameter to prevent overfitting and  $\|\cdot\|_F^2$  denotes the squared Frobenius norm.  $F(x) = P^T X$ , and  $H(x) = \text{sgn}(F(x))$  is the set of hash functions that encodes  $X$  to the binary embedding. Typical choices of loss functions for  $\ell$  include hinge loss, logistic loss and a simple quadratic loss [24]. To more efficiently optimize problem (1), it can be reformulated to the following problem by the regularization methods [34] while keeping the binary constraints of  $b_i$  as in [24]

$$\begin{aligned} \arg \min_{F, W, B} \sum_{i=1}^n \ell\{y_i, f(b_i, W)\} + \lambda \|W\|_F^2 \\ + \nu \sum_{i=1}^n \|b_i - F(x_i)\|_F^2 \quad \text{s.t. } b_i \in \{-1, 1\}^k. \end{aligned} \quad (2)$$

The last term in (2) is derived from the constraints in (1), and serves as the quantization loss from the continuous embedding  $F(x_i)$  to binary codes  $b_i$ ,  $\nu$  is the regularization parameter. The problem in Equation (2) is very informative and easy to optimize using an iterative procedure. Nonetheless, only minimizing the classification error might not be able to generate the optimal hash codes for the image retrieval task as it does not explicitly take the distribution of binary codes into consideration for ANN search. It is quite understandable that the code which is good enough for linear classification might not work well with nearest neighbor classifier. Therefore, extensions should be made to better leverage the label information for optimized retrieval performance. To this end, we explicitly take the distribution of binary codes into consideration by enforcing a class prototype consistent constraint to the learnt binary codes. We refer to this method as code consistent hashing and will elaborate it in the next section.

Now we have briefly reviewed a learning based hashing method called *SDH* which exploits the label information directly in its formulation. In the next part, we will review some basic knowledge of a special kind of binary codes called Hadamard code, which plays an important role in constructing our proposed algorithm.

## 2.2 Hadamard Code

We aim to learn binary codes that can leverage both the label information and the information-theoretic criterion to benefit the retrieval task. Information theory is well exploited in many machine learning tasks and information-theoretic criterions like the maximum correntropy criterion [35], [36], [37] have achieved great success in various applications [38]. Previous observations have shown that binary codes with balanced and uncorrelated bits (which favor a large information entropy) are desired for image retrieval tasks [9]. Various algorithms [22], [33] exploit this observation to formulate the learning of hash function as a constrained optimization problem, i.e., learn binary codes subject to  $B\mathbf{1} = 0$  and  $BB^T = NI_k$ , where  $\mathbf{1} = \{1, 1, \dots, 1\}^T \in$

$\mathbb{R}^N$  and  $\mathbf{0} = \{0, 0, \dots, 0\}^T \in \mathbb{R}^N$ . The constraint  $B\mathbf{1} = 0$  is imposed to maximize the information entropy of each bit, leading to balanced partitioning of the data. Another constraint  $BB^T = NI_k$  forces  $k$  bits to be mutually uncorrelated with each other in order to minimize redundancy among them. However, these constraints are imposed on all the samples (which is often very large) in the training set, making the objective hard to optimize. In this work, we explore another way to satisfy these information-theoretic favored constraints by enforcing the binary codes to be as close as possible to their predefined class prototypes in the binary embedding that satisfy such strict constraints. The class prototypes can be understood as class centers in the binary embedding. We shall address the details of this idea in the next section. Here, we only give a brief introduction of the Hadamard code that serves as the class prototypes in our proposed framework.

The research of the Hadamard code that originates decades ago [32] is a very well-established one in the field of mathematics and has various variants. Here we only consider a specific kind of Hadamard code that has been used as error correcting codes in the multi-class classification problem [39]. As a matter of fact, the code has just the desired properties mentioned above, i.e., balanced and uncorrelated bits.

One can generate a Hadamard code from the Hadamard matrix. A Hadamard matrix is a squared matrix  $H \in \{+1, -1\}^{m \times m}$  that meets  $HH^T = mI$ . This requires that any two rows of the matrix are orthogonal (thus uncorrelated) to each other. The generation of a Hadamard matrix is rather simple using the Sylvester's method [40], [41], where a new Hadamard matrix is produced from the old one by Kronecker product. For instance, given a Hadamard matrix  $H_2 = \begin{bmatrix} ++ \\ +- \end{bmatrix}$ , we can produce  $H_4$  by  $H_4 = H_2 \otimes H_2$  as below, where  $\otimes$  denotes the Kronecker product. Similarly,  $H_8$  is computed by  $H_8 = H_4 \otimes H_2$

$$H_4 = \begin{pmatrix} + & + \\ + & - \end{pmatrix} \otimes \begin{pmatrix} + & + \\ + & - \end{pmatrix} = \begin{pmatrix} + & + & + & + \\ + & - & + & - \\ + & + & - & - \\ + & - & - & + \end{pmatrix}. \quad (3)$$

However, it's easy to see that the size of Hadamard matrix generated with this method is a power of 2, and the first row and column is always 1 for all sizes of Hadamard matrix. Excluding the first row and first column from a Hadamard matrix  $H \in \{+1, -1\}^{m \times m}$ , the Hadamard code is constructed from the Hadamard matrix as  $HC \in \{+1, -1\}^{(m-1) \times (m-1)}$ . It is easy to see the size of Hadamard code is limited to a power of 2 minus 1 (e.g., 15, 31, 63, ...). As inherited from the Hadamard matrix, the orthogonal property still holds for the Hadamard code. With a closer examination, we can see that the Hadamard code have several other properties. First, each row and column of a Hadamard code has  $m/2$  symbols that equals to one; Second, the distance between any two rows (or columns) is  $m/2$ . The first property means that each bit code is balanced, while the second property ensures that when the Hadamard code is used as the predefined class prototypes (i.e., class centers), the distance between any two class prototypes is very large with a constant value of  $m/2$ . It's worth mentioning that the second property also make it very popular as an error



correcting code [39], because it is robust to small errors that may occur in some bits.

### 3 CODE CONSISTENT HASHING

As mentioned in Section 2.1, we aim to leverage both the label information and the information-theoretic criterion to learn binary codes in Hamming space. The learnt codes should be optimal for retrieval tasks using nearest neighbor search with the Hamming distance. To achieve this, one intuitive idea is to yield LDA like codes in the binary embedding, i.e., codes from the same class lie close to each other, whereas codes from different classes separate from each other with a large margin. To this end, we propose to enforce the codes lie close to their respective class centers (termed as class prototypes in this work). Those class prototypes are not learnt (so that we do not have to deal with the NP-hard optimization problem with no guaranteed best solution) but are rather predefined with the Hadamard code whose codes are both balanced and uncorrelated.

#### 3.1 Code Consistent Constraint

We formulate the aforementioned idea into the following constrained optimization problem with a code consistent term:

$$\begin{aligned} & \arg \min_{B, P, M} \|CY - M^T B\|_F^2 \\ \text{s.t. } & M^T M = MM^T = I_k, B = \text{sgn}(P^T(X)), \end{aligned} \quad (4)$$

where  $B \in \{-1, 1\}^{k \times N}$  is the binary codes with each column representing one sample.  $P^T \in \mathbb{R}^{k \times h}$  is the projection matrix that transform features  $X \in \mathbb{R}^{d \times N}$  in the original space to the binary embedding with a  $\text{sgn}$  function.  $C \in \{-1, 1\}^{k \times L}$  is composed of the predefined class prototypes corresponding to each class, and here  $Y \in \{0, 1\}^{L \times N}$  is a matrix representing the label ground-truth, with each column  $y_i$  denoting the correct label for that sample.  $y_i$  is a one-hot vector and the position of 1 in each column  $y_i = [0 \cdots 1 \cdots 0]^T$  indicates that the correct class label, e.g., the one-hot vector  $y_i = [1, 0, 0, 0, 0, 0]^T$  denotes the  $i$ th sample belongs to the first class in a classification problem with six classes involved in total. Thus,  $CY$  is a  $k \times N$  matrix, with each column representing the class prototypes of the corresponding samples in  $B$ . The code consistent term  $\|CY - M^T B\|_F^2$  measures the empirical error of the generated codes  $B$ , i.e., the total deviations of  $M^T B$  from their respective class centers  $CY$ . To add another degree of freedom, we introduce a transformation matrix  $M \in \mathbb{R}^{k \times k}$  in this term and here we enforce the transformation matrix  $M$  to be an orthogonal matrix, i.e., a rotation matrix, so that the distances between each class can be preserved during the transformation. For a better understanding of this term, we can rewrite the term  $\|CY - M^T B\|_F^2$  as  $\|M^T CY - B\|_F^2$ , since  $\|U\|_F^2 = \text{Tr}(U^T U) = \text{Tr}(U^T I U) = \text{Tr}(U^T M M^T U) = \|M^T U\|_F^2$ . Therefore, this term tries to close the gap between the learnt binary codes  $B$  and a rotational form of their class prototypes. What is more, as specified in Section 3.2, the introduction of this rotation matrix  $M$  will drastically reduce the computational cost to minimize the overall objective function.

#### Algorithm 1. Code Consistent Hashing (CCH)

**Input:** Data matrix  $X \in \mathbb{R}^{d \times N}$ , the corresponding label matrix  $Y \in \{0, 1\}^{L \times N}$ , class prototypes  $C \in \{-1, 1\}^{m \times L}$ , the kernel width  $\sigma$ , the regularization parameter  $\alpha$  and number of iterations  $IterNum$ .

**Output:** Hash functions in matrix form  $P^T \in \mathbb{R}^{k \times h}$  the learnt binary codes  $B \in \{-1, 1\}^{k \times N}$ .

```

1: Initialization: Normalize all data points:  $x_i \leftarrow x_i / \|x_i\|$ , calculate the kernel features with  $\phi(x) = [\exp(\|x - x_1\|^2 / \sigma), \dots, \exp(\|x - x_h\|^2 / \sigma)]$ , initialize  $B$  with randomly generated binary codes matrix  $B \in \{-1, 1\}^{k \times N}$  and set iteration to 0;
2: while iteration < IterNum do
3:   P-Step:
4:      $P = (\phi(X)\phi(X)^T)^{-1}\phi(X)B^T$ ,
5:   M-Step:
6:      $USV^T = \text{SVD}(CYB^T)$ ,  $M = UV^T$ ,
7:   B-Step:
8:      $B = \text{sgn}(Q)$ , where  $Q = MCY + \alpha P^T \phi(X)$ ,
9:     iteration = iteration + 1
10: end while
11:  $B = \text{sgn}(P^T \phi(X))$ 

```

With regards to the separability of inter-class samples, we specifically make the class prototypes to be the Hadamard code to enforce a large margin between different classes. Next, we will describe the process for generating the class prototype matrix  $C$ . As introduced in Section 2.2, using the Hadamard code as the class prototypes naturally separate all classes with a uniform distance of  $m/2$ , which is very promising for separating different classes with a large margin. Besides, as we mentioned earlier in Section 2.2, the property of having balanced and uncorrelated bits is another reason that makes the Hadamard code an ideal choice for class prototypes. However, in practical applications, the number of classes  $L$  seldom meets the number of the code length  $m$  of  $HC$  ( $m = 2^k - 1, k = 2, 3, \dots$ ). Most of the time, we need to choose  $L$  codewords (columns) from  $HC \in \{+1, -1\}^{m \times m}$  to form a class prototype matrix  $C \in \{+1, -1\}^{k \times L}$ . However, the resultant matrix  $C$  violates the property of having balanced and uncorrelated bits described before. To compensate for this, we employ a greedy search method as in [41] to select  $k$  codewords from  $HC$  to form  $C$  so as to best approximate the balanced bits property. Due to a small search space determined by the number of classes and the code length, this strategy works quite well to achieve class prototype matrix with balanced bits. And in practice, we find that the chosen codes also favor the small correlation property. In general, we can learn hash functions  $P$  with any suitable embedding learning algorithm, linear or nonlinear. Here, for better generalization performance, we learn our hash functions in a kernel embedding. Specifically, we use the simple yet powerful nonlinear form  $\phi(x)$  generated with an RBF kernel mapping process:  $\phi(x) = [\exp(\|x - x_1\|^2 / \sigma), \dots, \exp(\|x - x_h\|^2 / \sigma)]$ , where  $\{x_i\}_{i=1}^h$  are  $h$  chosen anchor points from the training samples and  $\sigma$  is the kernel width. Thus, the learned hash functions  $P$  will project the mapped data  $\phi(X)$  onto the binary embedding with  $\text{sgn}(P^T \phi(X))$ . Similar formulations are widely used as the kernel hash

function in, e.g., *KSH* [17]. Typically, the anchor points can be chosen as the clustering centers with k-means or randomly sampled data points. In this work, we adopt randomly chosen samples as anchor points for simplicity.

### 3.2 Optimization

The constraint optimization problem with both continuous and discrete variables in Equation (4) is NP-hard and very difficult to optimize. Similar to [24], using the regularization methods proposed in [34], we keep the binary constraints of  $b_i$  in the optimization and rewrite (4) as the following problem:

$$\begin{aligned} \arg \min_{P, M, B} & \|CY - M^T B\|_F^2 + \alpha \|B - P^T \phi(X)\|_F^2 \\ \text{s.t. } & B \in \{-1, 1\}^{k \times N}, \quad M^T M = M M^T = I_k. \end{aligned} \quad (5)$$

The joint learning problem in Equation (5) is still non-convex and difficult to solve. Directly optimizing all the variables at the same time is not tractable. Therefore, we employ an iterative optimization procedure, where we minimize the problem with respect to one variable while fixing others at each step and iterate over all the steps. The detailed optimization process is addressed as follows.

We first compute the features in the kernel space, and initialize  $B$  with random binary codes. Our algorithm iterates over the following three steps to minimize the objective function in Equation (5).

**P-Step.** We first fix all the variables except  $P$ , the degenerated problem is a unconstrained regression problem, which can be easily computed as:

$$P = (\phi(X)\phi(X)^T)^{-1} \phi(X) B^T. \quad (6)$$

**M-Step.** By fixing all the variables except  $M$ , this problem degenerates to a regularized least squares problem with a closed-form solution, same as the problem in *ITQ* [9]:

$$USV^T = \text{SVD}(CYB^T), \quad M = VU^T. \quad (7)$$

**B-Step.** In this step, we try to optimize  $B$  with all other variables fixed. The optimization problem in this step takes the following form:

$$\begin{aligned} \arg \min_B & \|CY - M^T B\|_F^2 + \alpha \|B - P^T \phi(X)\|_F^2 \\ \text{s.t. } & B \in \{-1, 1\}^{k \times N}. \end{aligned} \quad (8)$$

By expanding the Frobenius norm in (8) and removing the constants, problem (8) is equivalent to:

$$\begin{aligned} \arg \min_B & \|M^T B\|_F^2 - 2\text{Tr}(B^T Q) \\ \text{s.t. } & B \in \{-1, 1\}^{k \times N}. \end{aligned} \quad (9)$$

Where  $Q = MCY + \alpha P^T \phi(X)$  and  $\text{Tr}(\cdot)$  represents the trace of a matrix. By expanding problem (9) we can derive an equivalent problem as follows:

$$\begin{aligned} \arg \min_B & -\text{Tr}(B^T Q) \\ \text{s.t. } & B \in \{-1, 1\}^{k \times N}. \end{aligned} \quad (10)$$

Due to the orthogonal property of the learned  $M$  ( $MM^T = I_k$ ) from the last step,

$$\|M^T B\|_F^2 = \text{Tr}(B^T M M^T B) = \text{Tr}(B^T B) = \text{const} \quad (11)$$

and it's easy to see that the optimal solution to problem(10) would be

$$B = \text{sgn}(Q). \quad (12)$$

Unlike the *SDH* which iteratively learns  $B$  bit by bit with the discrete cyclic coordinate descent method and needs to iterate over all the bits for two to five times till the procedure converges. Our proposed method updates  $B$  as a whole with a simple thresholding operation and it only needs to be carried out once at each **B-Step**. Compared with *SDH*, the computation cost is greatly reduced in this step.

It is noted that the most time-consuming calculation in the whole procedure is the matrix inversion in the  $M, P$  step. However, it is obvious that the dimensionality of the matrix, whose inverse needs to be calculated, is not decided by the number of training samples. Thus, the number of training samples will not pose any difficulty to the optimization of our proposed method. Therefore, we argue that our proposed method is very efficient and can be scaled to large training set with even millions of training samples.

## 4 EXPERIMENTS AND RESULTS

The proposed approach is evaluated on several benchmark datasets. We present quantitative evaluations in terms of several retrieval metrics and compare our approach Code Consistent Hashing with many popular unsupervised methods: Iterative Quantization (*PCA-ITQ*) [9], Hashing with Graphs (*AGH*) [10], Inductive Hashing on Manifolds (*IMH*) [11] with t-SNE [42] and state-of-the-art supervised methods: FastHash [43], Supervised Discrete Hashing [24], *CCA-ITQ* [9], Binary Reconstructive Embedding (*BRE*) [19], Semi-supervised Hashing [14], Supervised Hashing with Kernels (*KSH*) [17]. For all the compared methods, we used the parameter settings reported in their works for fair comparison. And for the proposed approach, we empirically set  $\alpha$  to 1e-4 and the number of anchor points to 1,000 for evaluation efficiency. In practice, our framework converges in about 3~5 iterations.

In the following sections, we first detail our evaluation protocols concerning both retrieval and classification tasks. And in the subsequent section, brief descriptions of the benchmark datasets are presented, followed by extensive comparative experimental results and detailed analysis on each dataset.

### 4.1 Evaluation Protocols

Our proposed approach is evaluated to validate its effectiveness and efficiency on both retrieval and classification tasks. For retrieval tasks, two commonly [14] used criteria for evaluating hashing methods are adopted:

1. *Hamming ranking*: All the points in the database are ranked according to their Hamming distance from the query and the desired neighbors are returned from the top of the ranked list. Although the complexity of Hamming ranking is linear, it is still very fast in practice due to the simplicity in calculating the Hamming distance.

2. *Hash lookup*: A lookup table is constructed using the database codes, and all the points in the buckets within a

small Hamming radius  $r$  of the query are returned. The complexity of the hash lookups is constant time.

As pointed out by [14], evaluations based on Hamming ranking and hash lookup focus on different characteristics of hashing techniques. In general, hash lookup emphasizes more on the practical search speed and works well with small number of hash bits. However, when the number of bits is large, the Hamming space becomes increasingly sparse; and very few samples fall within the Hamming radius  $r$  ( $r = 2$  in our experiments), resulting in many failed queries. This phenomenon can be clearly seen in our experimental results. In this situation, Hamming ranking provides better quality measurement of the Hamming embedding. Therefore, for comprehensive evaluations, we adopt both criteria in our experiments. All the experiments were conducted with relatively compact codes (from 8 to 36 bits). Since we focus on searching semantically similar points, the search results are evaluated based on whether the returned images and the query sample share the same semantic labels for all the compared methods, supervised or unsupervised. Specifically, we use three metrics, i.e., precision with regards to top  $K$  returned points ( $K = 500$ ), mean average precision and precision @ radius  $r$  with  $r = 2$  to measure the quantitative performance of different methods. The setting of  $K$  and  $r$  are based on standard protocols widely used in previous literature. Moreover, the first two metrics are Hamming ranking based evaluations, whereas the last one is based on the hash lookup criterion. In addition to these quantitative measurements, we also plot the precision-recall curve for an overall evaluation of the system performance.

To evaluate the proposed approach on the classification tasks, the classification accuracy is adopted as the evaluation metric. Specifically, we apply linear SVM (the LIBLINEAR toolbox [44] is used) for classification using the obtained codes with different hashing algorithms. The classification results using binary codes are also compared with the results using the original features for better understanding of the supervised hashing methods.

## 4.2 Data Sets and Results

In our experiments, three image datasets are used, i.e., CIFAR-10,<sup>1</sup> MNIST<sup>2</sup> and NUS-WIDE,<sup>3</sup> with the number of samples ranging from thousands to hundreds of thousands. Since our approach aims at improving searching performance for semantic neighbors, all the datasets we use are fully annotated. The neighbor ground-truth for searching evaluation is defined by the correct label information from the datasets. In addition, the feature vector used in our approach is normalized to have unit norm before computing the kernelized features as described in Section 3.1. And with regards to all the other compared methods, the recommended normalization is also accordingly implemented. In the following part, we provide detailed descriptions of the datasets and extensive evaluations on those datasets.

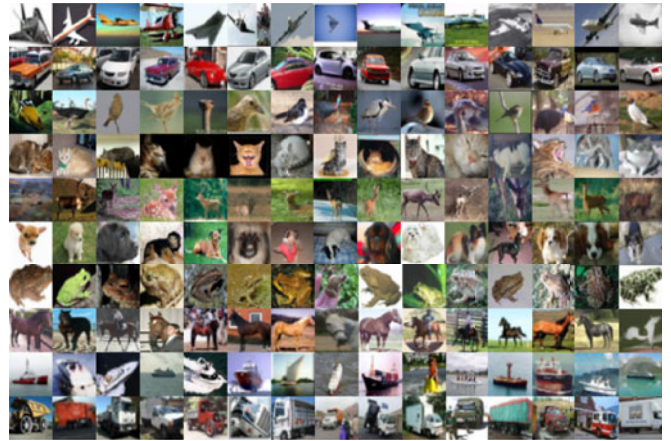


Fig. 1. Sample images of the CIFAR-10 dataset. Each row corresponds to images belonging to a certain class. From top to bottom, the image classes are airplane, automobile, bird, cat, deer, dog, frog, horse, ship and truck.

### 4.2.1 Results on CIFAR-10 Data Set

The CIFAR-10 dataset is a labeled subset of the 80-million tiny images collection [5]. It consists of a total of 60K  $32 \times 32$  color images which are manually categorized into 10 classes, each of which has 6,000 samples. We represent each image with a 512 dimension GIST [45] vector feature in our experiment. Each sample in this dataset is associated with a mutually exclusive class label. A few example images from CIFAR10 dataset are shown in Fig. 1. The entire dataset is partitioned into two parts: a training set with 59,000 samples (5,900 samples each class) and a test set with 1,000 samples (100 samples each class). The training set is used for learning hash functions and constructing the hash lookup tables. For *SSH*, *BRE*, we additionally sample 5,000 random points from the training set for similarity matrix construction based on the image labels. Since methods like *BRE*, *MLH* can not scale well when the number of training samples is very large, we mainly conduct our experiments with 5,000 training samples (500 samples each class and 3,000 labeled samples for constructing similarity matrix in *BRE* and *SSH*) for fair comparisons and evaluation efficiency. Experiments with 59K training samples are also evaluated and reported in Table 1 for a comprehensive comparison.

Firstly, we conduct experiments using 32-bit codes with varied training samples and anchor points and report the result in Table 1. It is clear that *CCH* outperforms all the compared methods on both precision @ radius 2 and mAP whether the number of training samples is 59K or 5K. Same as *CCH*, *SDH* and *KSH* all learn hash functions in the kernel embedding, but they all achieves inferior results with more training time than *CCH*. Thanks to the introduction of the orthogonal matrix  $M$  in the objective function, *CCH* runs approximately five times faster than *SDH*. It is worth noting that the retrieval performance consistently boost with more anchor points, however, the computational cost also increase accordingly. Therefore, for evaluation efficiency, we report results with 5,000 training samples and 1,000 anchor points in the following part. Fig. 2 shows some performance curves illustrating three quantitative results with regards to different number of bits. We only compare our *CCH* with supervised or semi-supervised methods (i.e.,

1. <http://www.cs.toronto.edu/~kriz/cifar.html>

2. <http://yann.lecun.com/exdb/mnist/>

3. <http://lms.comp.nus.edu.sg/research/NUS-WIDE.htm>



TABLE 1  
Results in Precision @ Radius 2, mAP and Training  
Time on CIFAR-10

Method	# training	# anchor	Precision	mAP	Training time
CCH	5,000	300	0.4401	0.3732	0.19
	5,000	1,000	0.4804	0.4925	1.08
	5,000	3,000	0.4937	0.6091	7.20
	59,000	1,000	<b>0.5103</b>	<b>0.4464</b>	9.39
BRE	5,000	-	0.1907	0.2370	1452.9
KSH	5,000	1,000	0.3135	0.4280	2140.6
CCA-ITQ	59,000	-	0.4287	0.3291	6.19
FastHash	59,000	-	0.2621	0.3409	1485.7
SDH	59,000	1,000	0.5078	0.4283	45.07
SSH	59,000	-	0.1924	0.1785	45.20
AGH	59,000	1,000	0.2757	0.1525	6.78
IMH	59,000	1,000	0.2155	0.1679	44.56
PCA-ITQ	59,000	-	0.2425	0.1545	<b>3.65</b>

Results with 32 bits are reported. For our method, the number of anchors varies from 300 to 3,000. The training time is in seconds. The experiments are conducted on a windows server with an Intel Xeon X5660 @ 2.80 GHz and 128 G RAM.

BRE, CCA-ITQ, FSH, KSH, SDH, SSH) in this experiment. Three kernel based methods (CCH, KSH, SDH) achieve comparable results on all three metrics. But our proposed CCH outperforms all the competing methods on all the metrics, even with a large margin when the code length is short (e.g., 8, 12 and 16 bits). Since CCH and SDH adopt the same kind of optimization framework, it is safe to say that the superior performance is due to the incorporation of the code consistent term in our proposed framework. The reason that CCH performs very well with shorter bits is possibly that the generated class prototypes using a greedy search method have better information-theoretic favored property when the code length is short. However, when the code length is big, and the number of classes is small, we need to select a very skinny matrix from the Hadamard matrix as the class prototypes, thus the balanced and uncorrelated property is more easily violated, leading to suboptimal results.

The evaluation of precision @ radius 2 shown in Fig. 2c is a measure of the hash lookup performance, most compared methods (e.g., KSH, CCA-ITQ and FSH) suffer from significant drops with the increasing of used bits. This is because, as stated in Section 4.1, the number of points falling in a bucket decrease exponentially when longer codes are used, leading to many failed queries by not returning any

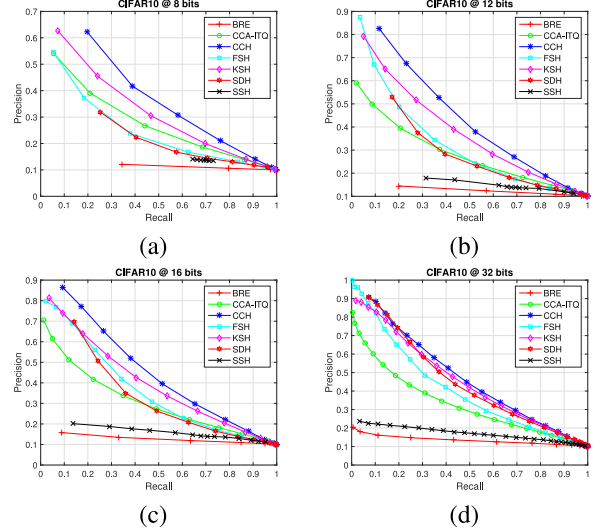


Fig. 3. Precision-recall curves with regards to different number of bits on CIFAR-10. (a) 8 bits. (b) 16 bits. (c) 24 bits. (d) 32 bits.

neighbor even in a Hamming ball of radius 2. Thus, although hash lookup tables have faster searching speed than Hamming ranking, they can't provide sound results when longer bits are used. To evaluate the overall performance of different kinds of hashing methods, we also plot the precision-recall curves for different methods with 8, 16, 24 and 32 bits on Fig. 3. For precision-recall curves on 8 and 16 bits, it is quite obvious that CCH outperforms its competing methods with a large margin, whereas when using longer bits like 24 and 32, three kernel based methods achieve comparable results. But compared with methods that learn hash functions in the original feature space (e.g., CCA-ITQ, SSH and FSH), those three methods achieve far better performance, validating the effectiveness of learning hash function in the kernel embedding.

To further evaluate the learnt binary codes, we conduct classification experiments. Several state-of-the-art supervised hashing methods are adopted to generate the binary codes. We report classification accuracy with code length 8, 16, 32 and 64 in Table 2. The proposed CCH performs best among all the compared methods with all code length except 32, even better than SDH which explicitly incorporates a classification error term in their objective function. And it is interesting to find out that both CCH and SDH

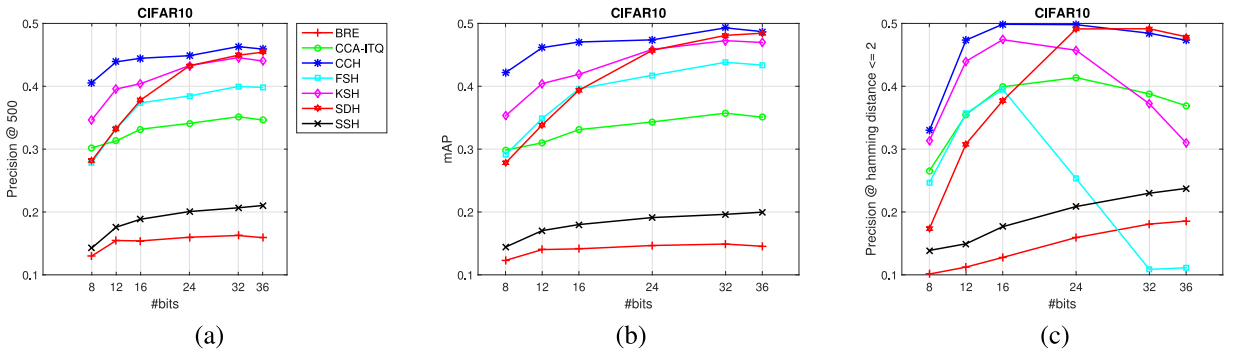


Fig. 2. Quantitative evaluation results on CIFAR-10. (a) Precision @ top 500 returned samples using Hamming ranking with regards to different number of bits. (b) Mean average precision with regards to different number of bits. (c) Precision @ radius 2 with hash lookup with regards to different number of bits.

TABLE 2  
Classification Results with Learnt Binary Codes on CIFAR and MNIST

	#bits	BRE	CCA-ITQ	CCH	FastHash	KSH	SDH	SSH	Original Feature
CIFAR	8	22.4	43.7	<b>55.6</b>	42.4	46.4	42.8	18.7	54.3
	16	30.3	48.8	<b>58.6</b>	47.4	53.2	57.1	19.3	54.3
	32	33.4	52.4	59.0	51.3	56.1	<b>59.3</b>	35.5	54.3
	64	36.7	51.8	<b>58.9</b>	51.9	58.5	58.7	39.1	54.3
MNIST	8	9.6	79.1	<b>92.0</b>	84.8	85.5	90.0	39.8	88.5
	16	68.2	84.3	<b>93.1</b>	90.7	89.6	92.4	47.5	88.5
	32	73.6	86.4	93.3	89.5	92.3	<b>93.6</b>	57.8	88.5
	64	79.0	85.7	<b>93.0</b>	89.3	92.3	92.5	72.9	88.5

encounter a performance saturation at bit length 32, increasing the code length from this point will not yield better performance. The reason for this phenomenon might be that the redundancy and noise introduced by longer bits is not helpful to increase the discriminability of the learnt codes but will instead deteriorate the classification performance. Despite this interesting phenomenon, the classification results still further validate the effectiveness of the proposed code consistent term. And it is noted that the binary codes obtained by *CCH* and *SDH* are more discriminative than the original GIST features, even with a short bit length. This is expected because the learning process of hash codes can be seen as a nonlinear dimensionality reduction process from the original space. Redundancy is removed and discriminative information is incorporated in this process to get a more discriminative representation.

#### 4.2.2 Results on MNIST Data Set

The MNIST dataset consists of 10 handwritten digits ranging from '0' to '9', each with 7K  $28 \times 28$  grayscale images. Each image is represented with its pixel values in a 784 dimension vector. Same as CIFAR-10, each sample in this dataset is also associated with a mutually exclusive class label. A query set consists of 1,000 samples is sampled uniformly from the whole dataset and 5,000 (500 each class) images are used as training set. For *BRE*, *SSH*, we additionally sample 3,000 random points from the training set for similarity matrix construction based on the image labels. Because this is a relatively simple dataset, 5,000 training samples are enough for most of the compared methods to produce satisfactory results. Thus, we mainly conduct our experiments with 5,000 training samples (500 samples each

class and 1,000 labeled samples for constructing similarity matrix in *BRE* and *SSH*) for fair comparisons and evaluation efficiency. We also conduct experiments with 69K training samples, the training time and quantitative results show the same trend as we analyzed in the CIFAR-10 dataset, so we do not report them here.

Fig. 4 shows the quantitative evaluation results on this dataset. Because this is a relatively easy dataset, *CCH*, *KSH* and *SDH* give very similar performance on all those metrics. Although somewhat negligible, *CCH* still outperforms all its competitors. The same phenomenon on the precision @ Hamming radius 2 is also observed on Fig. 4c. Fig. 5 gives the precision-recall curves with regards to 12, 16 and 24 bits respectively. It is clear that our *CCH* outperforms all the competing methods especially when the code length is short. Classification results in shown in Table 2, it's not surprising to see that the binary codes generated by *CCH* boost the classification accuracy by a large margin compared to the original features and perform best among all the compared hashing methods.

#### 4.2.3 Results on NUS-WIDE Data Set

The NUS-WIDE dataset [46] is a set of Flickr consumer images collected by NUS lab. This dataset contains around 270,000 images associated with 81 ground truth concept tags, with each image assigned to multiple semantic labels. Although searching with multi-label images is gaining traction in recent studies [47], it is a relatively more advanced topic that needs special treatment to acquire a reasonable performance. And all the compared methods, including our own, mainly focus on the traditional visual search problem that aims at searching images with mutually exclusive labels.

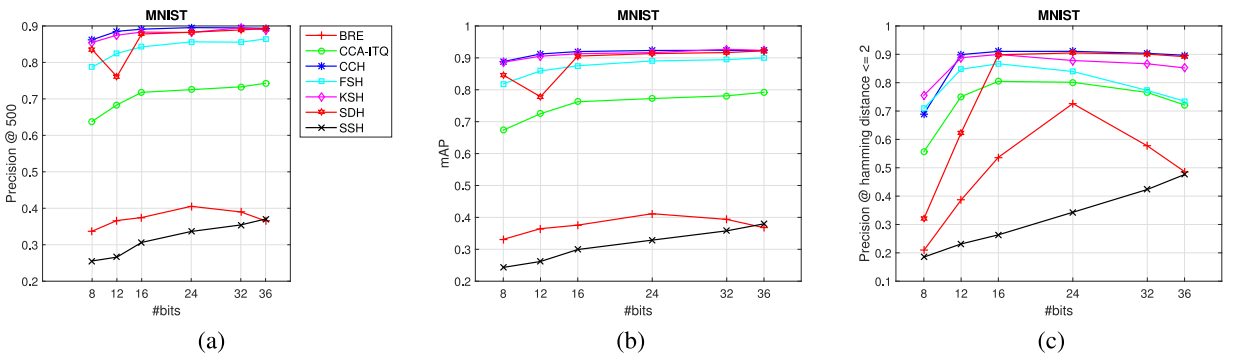


Fig. 4. Quantitative evaluation results on MNIST. (a) Precision @ top 500 returned samples using Hamming ranking with regards to different number of bits. (b) Mean average precision with regards to different number of bits. (c) Precision @ radius 2 with hash lookup with regards to different number of bits.



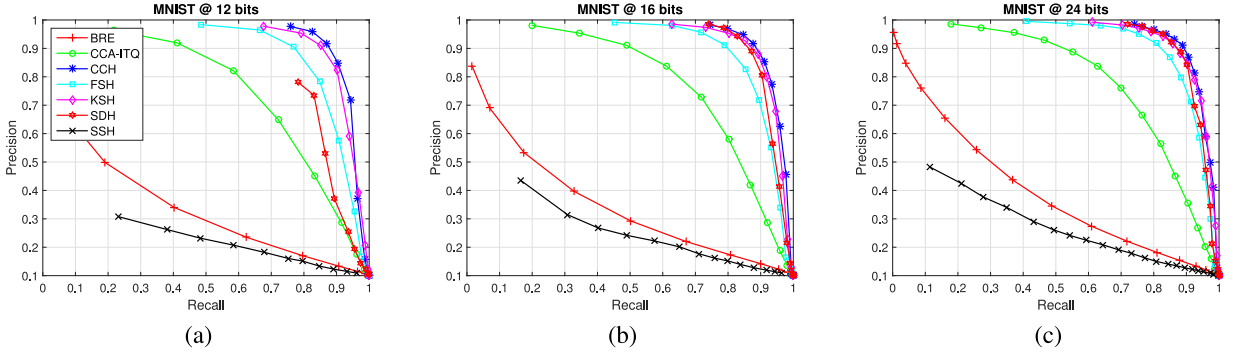


Fig. 5. Precision-recall curves with regards to different number of bits on MNIST. (a) 12 bits. (b) 16 bits. (c) 24 bits.

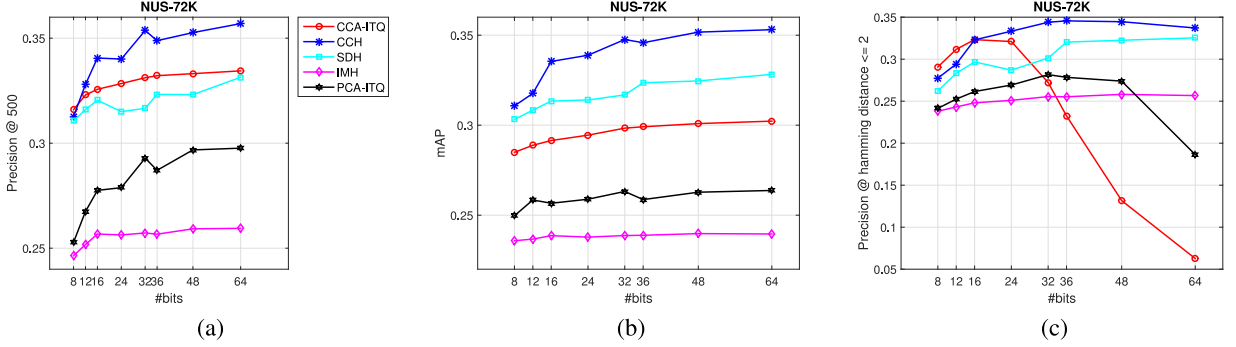


Fig. 6. Quantitative evaluation results on NUS-WIDE. (a) Precision @ top 500 returned samples using Hamming ranking with regards to different number of bits. (b) Mean average precision with regards to different number of bits. (c) Precision @ radius 2 with hash lookup with regards to different number of bits.

Therefore, we select a subset that belongs to the 21 largest classes with each image exclusively belonging to one of the 21 classes, which results in a subset with 72,219 images. It is noted that the class labels do not have a uniform distribution, i.e., different classes have different number of images. The images in the dataset are represented with a Bag-of-Visual-Word model with local SIFT [48] features. Particularly, a visual vocabulary with 500-length code book and a soft assignment strategy are used for deriving the image features, as described in [45]. For each class, 1/10 of the images are sampled as the query set and the remaining images are used as the training and gallery set. As the NUS-WIDE dataset is relatively larger, longer bits might be needed for better performance, therefore, the performance is evaluated with different code lengths varying from 8- to 64-bit.

Since this dataset is relatively large, some methods like *BRE* and *FSH*, whose training process take several hours, do not scale well to this kind of magnitude of data. And previous experimental results have shown that the performance of *SDH* is a representative one among all the compared methods. Therefore, for evaluation efficiency, we only compare *CCH* with *SDH* and the very efficient *CCA-ITQ*. Performances of some efficient unsupervised methods (*PCA-ITQ*, *IMH*) are also reported to yield a comprehensive comparison. Like we did on CIFAR-10 and MNIST, three quantitative metrics and the precision-recall curves are presented to evaluate the performance. Fig. 6 gives the quantitative evaluation results with regards to different number of bits. Our proposed *CCH* achieves the best performance in most cases. And it is surprising to see from Fig. 6c that both *SDH* and *CCH* continue to perform better as the number of

bits increases, while all the other compared methods reveal significant drops with longer bits. This is possibly because the classification error term in *SDH* and the code consistent term in *CCH* can both penalize large intra-class variations, thus the learnt binary codes are more consolidated to their class centers than other methods. From both Figs. 6 and 7, it is clear that supervised methods have great advantages over unsupervised ones when it comes to searching semantically similar neighbors.

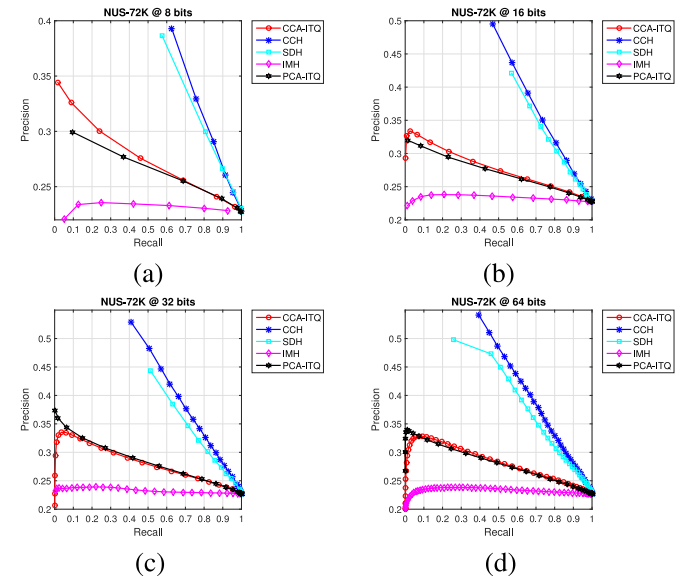


Fig. 7. Precision-recall curves with regards to different number of bits on NUS-WIDE. (a) 8 bits. (b) 16 bits. (c) 32 bits. (d) 64 bits.

## 5 CONCLUSION

This paper introduces a new code consistent constraint to leverage label information and Hadamard code in information theory, which results in a simple yet efficient hashing method, named code consistent hashing. Specifically, we choose the Hadamard code as the class prototypes to implicitly leverage the information theory criteria so as to generate codes with balanced and uncorrelated bits. An efficient optimization process for discrete codes has been developed to solve the resultant objective function. Experimental results show that CCH outperforms state-of-the-art hashing methods on various types of visual benchmarks, demonstrating its remarkable effectiveness, scalability, and efficiency for large-scale retrieval and image classification tasks.

## ACKNOWLEDGMENTS

The authors would like to thank the associate editor and the reviewers for their valuable comments and advice. This work is supported by the National Basic Research Program of China (Grant No. 2012CB316300) and the Strategic Priority Research Program of the Chinese Academy of Sciences (Grant No. XDB02000000). S. Zhang and J. Liang contributed equally to this work and should be considered co-first authors.

## REFERENCES

- [1] P. Indyk and R. Motwani, "Approximate nearest neighbors: Towards removing the curse of dimensionality," in *Proc. Annu. ACM Symp. Theory Comput.*, 1998, pp. 604–613.
- [2] J. L. Bentley, "K-d trees for semidynamic point sets," in *Proc. Annu. Symp. Comput. Geometry*, 1990, pp. 187–197.
- [3] B. Kulis, P. Jain, and K. Grauman, "Fast similarity search for learned metrics," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 12, pp. 2143–2157, Dec. 2009.
- [4] B. Wu, Q. Yang, W.-S. Zheng, Y. Wang, and J. Wang, "Quantized correlation hashing for fast cross-modal search," in *Proc. Int. Joint Conf. Artif. Intell.*, 2015, pp. 3946–3952.
- [5] A. Torralba, R. Fergus, and W. T. Freeman, "80 million tiny images: A large data set for nonparametric object and scene recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 11, pp. 1958–1970, Nov. 2008.
- [6] C. Strecha, A. M. Bronstein, M. M. Bronstein, and P. Fua, "Ldhash: Improved matching with smaller descriptors," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 1, pp. 66–78, Jan. 2012.
- [7] A. Gionis, P. Indyk, R. Motwani, et al., "Similarity search in high dimensions via hashing," in *Proc. Int. Conf. Very Large Data Bases*, 1999, vol. 99, pp. 518–529.
- [8] B. Kulis and K. Grauman, "Kernelized locality-sensitive hashing for scalable image search," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2009, pp. 2130–2137.
- [9] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin, "Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 12, pp. 2916–2929, Dec. 2013.
- [10] W. Liu, J. Wang, S. Kumar, and S.-F. Chang, "Hashing with graphs," in *Proc. Int. Conf. Mach. Learn.*, 2011, pp. 1–8.
- [11] F. Shen, C. Shen, Q. Shi, A. Van Den Hengel, and Z. Tang, "Inductive hashing on manifolds," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2013, pp. 1562–1569.
- [12] L. Zhang, Y. Zhang, X. Gu, J. Tang, and Q. Tian, "Scalable similarity search with topology preserving hashing," *IEEE Trans. Image Process.*, vol. 23, no. 7, pp. 3025–3039, Jul. 2014.
- [13] A. Krizhevsky and G. E. Hinton, "Using very deep autoencoders for content-based image retrieval," in *Proc. Eur. Symp. Artif. Neural Netw.*, 2011, pp. 489–494.
- [14] J. Wang, S. Kumar, and S.-F. Chang, "Semi-supervised hashing for large-scale search," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 12, pp. 2393–2406, Dec. 2012.
- [15] Q. Yang, L.-K. Huang, W.-S. Zheng, and Y. Ling, "Smart hashing update for fast response," in *Proc. Int. Joint Conf. Artif. Intell.*, 2013, pp. 1855–1861.
- [16] G. Lin, C. Shen, D. Suter, and A. van den Hengel, "A general two-step approach to learning-based hashing," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2013, pp. 2552–2559.
- [17] W. Liu, J. Wang, R. Ji, Y.-G. Jiang, and S.-F. Chang, "Supervised hashing with kernels," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2012, pp. 2074–2081.
- [18] L.-K. Huang, Q. Yang, and W.-S. Zheng, "Online hashing," in *Proc. Int. Joint Conf. Artif. Intell.*, 2013, pp. 1422–1428.
- [19] B. Kulis and T. Darrell, "Learning to hash with binary reconstructive embeddings," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, pp. 1042–1050.
- [20] R. Salakhutdinov and G. Hinton, "Semantic hashing," *Int. J. Approximate Reasoning*, vol. 50, no. 7, pp. 969–978, 2009.
- [21] V. Erin Liong, J. Lu, G. Wang, P. Moulin, and J. Zhou, "Deep hashing for compact binary codes learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 2475–2483.
- [22] Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, pp. 1753–1760.
- [23] T. Ge, K. He, and J. Sun, "Graph cuts for supervised binary coding," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 250–264.
- [24] F. Shen, C. Shen, W. Liu, and H. T. Shen, "Supervised discrete hashing," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 37–45.
- [25] R. He, Y. Cai, T. Tan, and L. Davis, "Learning predictable binary codes for face indexing," *Pattern Recog.*, vol. 48, pp. 3160–3168, 2015.
- [26] G. Lin, C. Shen, and J. Wu, "Optimizing ranking measures for compact binary code learning," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 613–627.
- [27] H. Lai, Y. Pan, Y. Liu, and S. Yan, "Simultaneous feature learning and hash coding with deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 3270–3278.
- [28] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan, "Supervised hashing for image retrieval via image representation learning," in *Proc. 28th AAAI Conf. Artif. Intell.*, 2014, pp. 2156–2162.
- [29] G. Zhong, P. Yang, S. Wang, and J. Dong, "A deep hashing learning network," *arXiv preprint arXiv:1507.04437*, 2015.
- [30] H.-F. Yang, K. Lin, and C.-S. Chen, "Supervised learning of semantics-preserving hashing via deep neural networks for large-scale image search," *arXiv preprint arXiv:1507.00101*, 2015.
- [31] J. Wang, H. T. Shen, J. Song, and J. Ji, "Hashing for similarity search: A survey," *arXiv preprint arXiv:1408.2927*, 2014.
- [32] J. Hadamard, "Résolution d'une question relative aux déterminants," *Bull. Sci. Math.*, vol. 17, no. 1, pp. 240–246, 1893.
- [33] W. Liu, C. Mu, S. Kumar, and S.-F. Chang, "Discrete graph hashing," in *Proc. Adv. Neural Inf. Process. Syst.*, 2014, pp. 3419–3427.
- [34] J. Mallick, J. Povh, F. Rendl, and A. Wiegele, "Regularization methods for semidefinite programming," *SIAM J. Optimization*, vol. 20, no. 1, pp. 336–356, 2009.
- [35] B. Chen and J. C. Principe, "Maximum correntropy estimation is a smoothed MAP estimation," *IEEE Signal Process. Lett.*, vol. 19, no. 8, pp. 491–494, Aug. 2012.
- [36] B. Chen, L. Xing, J. Liang, N. Zheng, and J. C. Principe, "Steady-state mean-square error analysis for adaptive filtering under the maximum correntropy criterion," *IEEE Signal Process. Lett.*, vol. 21, no. 7, pp. 880–884, Jul. 2014.
- [37] B. Chen, J. Wang, H. Zhao, N. Zheng, and J. Principe, "Convergence of a fixed-point algorithm under maximum correntropy criterion," *IEEE Signal Process. Lett.*, vol. 22, no. 10, pp. 1723–1727, Oct. 2015.
- [38] R. He, Y. Zhang, Z. Sun, and Q. Yin, "Robust subspace clustering with complex noise," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 4001–4013, Nov. 2015.
- [39] J. Langford and A. Beygelzimer, "Sensitive error correcting output codes," in *Proc. 18th Annu. Conf. Learn. Theory*, 2005, pp. 158–172.
- [40] A. Hedayat, W. Wallis, et al., "Hadamard matrices and their applications," *Ann. Statist.*, vol. 6, no. 6, pp. 1184–1238, 1978.
- [41] S. Yang, P. Luo, C. C. Loy, K. W. Shum, and X. Tang, "Deep representation learning with target coding," in *Proc. AAAI Conf. Artif. Intell.*, 2015, pp. 3848–3854.
- [42] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, no. 2579–2605, p. 85, 2008.

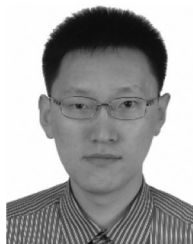
- [43] G. Lin, C. Shen, Q. Shi, A. van den Hengel, and D. Suter, "Fast supervised hashing with decision trees for high-dimensional data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2014, pp. 1971–1978.
- [44] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "Liblinear: A library for large linear classification," *J. Mach. Learn. Res.*, vol. 9, pp. 1871–1874, 2008.
- [45] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3, pp. 145–175, 2001.
- [46] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "NUS-WIDE: A real-world web image database from National University of Singapore," in *Proc. Int. Conf. Image Video Retrieval*, 2009, p. 48.
- [47] F. Zhao, Y. Huang, L. Wang, and T. Tan, "Deep semantic ranking based hashing for multi-label image retrieval," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 1556–1564.
- [48] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.



**Shu Zhang** received the BE degree in automation from Northwestern Polytechnical University in July 2012. He is currently working toward the PhD degree in computer application technology in NLPR. His research interests focus on biometrics, pattern recognition, and computer vision.



**Jian Liang** received the BE degree in electronic information and technology from Xi'an Jiaotong University in July 2013. He is currently working toward the PhD degree in pattern recognition and intelligent systems in NLPR. His research interests focus on machine learning, pattern recognition, and computer vision.



**Ran He** received the BE degree in computer science from Dalian University of Technology, the MS degree in computer science from Dalian University of Technology, and the PhD degree in pattern recognition and intelligent systems from Institute of Automation, Chinese Academy of Sciences in 2001, 2004, and 2009, respectively. Since September 2010, he has joined NLPR, where he is currently a full professor. He currently serves as an associate editor of *Neurocomputing* (Elsevier) and serves on the program committee of several conferences. His research interests focus on information theoretic learning, pattern recognition, and computer vision. He is a senior member of the IEEE.



**Zhenan Sun** received the BE degree in industrial automation from Dalian University of Technology in 1999, the MS degree in system engineering from Hua zhong University of Science and Technology in 2002, and the PhD degree in pattern recognition and intelligent systems from CASIA in 2006. He is a full professor in the Institute of Automation, Chinese Academy of Sciences (CASIA). In March 2006, he joined the Center of Biometrics and Security Research (CBSR) in the National Laboratory of Pattern Recognition (NLPR) of CASIA as a faculty member. His research focuses on biometrics, pattern recognition, and computer vision. He is a member of the IEEE and the IEEE Computer Society.

▷ **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).**