

Fast Class-wise Updating for Online Hashing

Mingbao Lin, Rongrong Ji*, *Senior Member, IEEE*, Xiaoshuai Sun, Baochang Zhang, Feiyue Huang, Yonghong Tian, *Senior Member, IEEE*, and Dacheng Tao, *Fellow, IEEE*

Abstract—Online image hashing has received increasing research attention recently, which processes large-scale data in a streaming fashion to update the hash functions on-the-fly. To this end, most existing works exploit this problem under a supervised setting, *i.e.*, using class labels to boost the hashing performance, which suffers from the defects in both adaptivity and efficiency: First, large amounts of training batches are required to learn up-to-date hash functions, which leads to poor online adaptivity. Second, the training is time-consuming, which contradicts with the core need of online learning. In this paper, a novel supervised online hashing scheme, termed **Fast Class-wise Updating for Online Hashing (FCOH)**, is proposed to address the above two challenges by introducing a novel and efficient inner product operation. To achieve fast online adaptivity, a class-wise updating method is developed to decompose the binary code learning and alternatively renew the hash functions in a class-wise fashion, which well addresses the burden on large amounts of training batches. Quantitatively, such a decomposition further leads to at least 75% storage saving. To further achieve online efficiency, we propose a semi-relaxation optimization, which accelerates the online training by treating different binary constraints independently. Without additional constraints and variables, the time complexity is significantly reduced. Such a scheme is also quantitatively shown to well preserve past information during updating hashing functions. We have quantitatively demonstrated that the collective effort of class-wise updating and semi-relaxation optimization provides a superior performance comparing to various state-of-the-art methods, which is verified through extensive experiments on three widely-used datasets.

Index Terms—Image retrieval, similarity preserving, online hashing, binary codes.

1 INTRODUCTION

TRADITIONAL hashing methods [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21], [22], [23] are mostly designed to learn hash functions offline from a fixed collection of training data with/without supervised labels. However, such a setting cannot handle dynamic application scenarios where data are fed into the system in a streaming fashion. Therefore, online hashing has attracted much research attention recently [24], [25], [26], [27], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37]. It aims to online update the hash functions from the sequentially arriving data instances, which merits in its superior adaptivity and scalability for large-scale online retrieval applications.

Ideally, online hashing should efficiently update hash functions based on the streaming data on-the-fly, while preserving the information from the past data stream. Existing works on online hashing can be classified into supervised and unsupervised methods. Representative supervised methods include, but are not limited to, OKH [24], AdaptHash [27], OSH [30], MIHash [33], HCOH [35] and BSODH [36], which use supervised labels to guide the

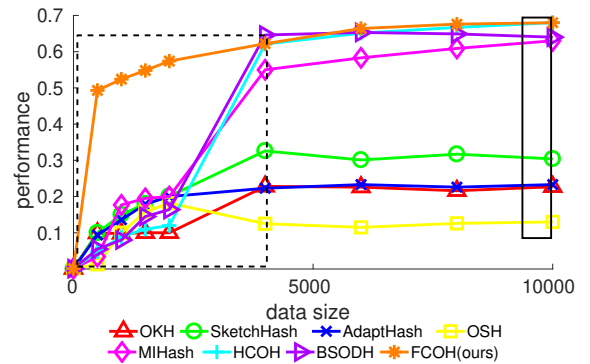


Fig. 1. Comparisons of the existing supervised online hashing methods [24], [26], [27], [30], [33], [34], [35], [36] and ours in terms of adaptivity. Existing methods reach the full performance only at the final training stage (solid rectangle box), which cannot be fast adapted online, and require a large number of training batches to learn up-to-date hash functions (dashed rectangle box). Differently, the proposed FCOH method can achieve superior adaptivity online at the earlier stages with much less training data.

online hashing learning. In contrast, unsupervised online hashing approaches, *e.g.*, SketchHash [26] and FROSH [34] consider a smaller data sketch from a large dataset [38] to update hash functions while preserving the main property of the dataset. Due to the usage of label information, supervised online hashing methods usually yield better results over unsupervised ones, which therefore role as the main trend in the literature and are the focus of this paper.

However, supervised online hashing remains as an open problem due to two issues, as validated in Sec. 4: First, a large number of training batches are required to learn up-to-date hash functions, which suffers poor online adaptivity. Previous works [24], [26], [27], [30], [33], [34], [35], [36] simply focus on designing hash models to obtain the full performance at the final training stage, which thereby lack

- M. Lin, R. Ji (Corresponding author) and X. Sun are with the Media Analytics and Computing Laboratory, Department of Artificial Intelligence, School of Informatics, Xiamen University, 361005, China. E-mail: rrji@xmu.edu.cn.
- R. Ji is also with Institute of Artificial Intelligence, Xiamen University, China.
- B. Zhang is with Beihang University, China.
- F. Huang is with YouTu Laboratory, Tencent, Shanghai, 200233, China.
- Y. Tian is with the Department of Computer Science and Technology, Peking University, Beijing, 100871, China.
- D. Tao is with the School of Computer Science, in the Faculty of Engineering, at The University of Sydney, 6 Cleveland St, Darlingtown, NSW 2008, Australia.

Manuscript received April 19, 2005; revised August 26, 2015.

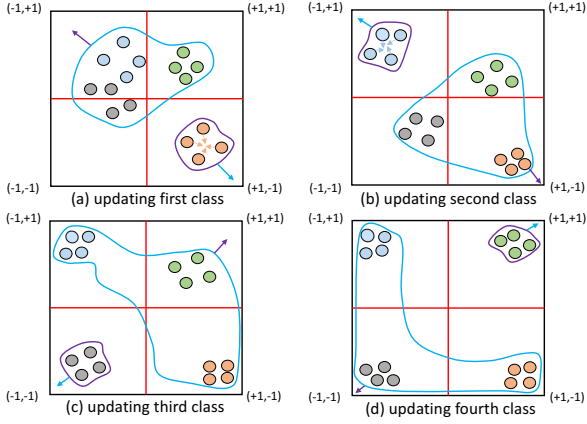


Fig. 2. An illustration of the proposed class-wise updating. For an arriving stream batch, the updating is divided into several independent subprocesses, *i.e.*, (a), (b), (c) and (d). In each subprocess, FCOH aims to discriminate a particular class from the others, which is conducted sequentially ((a)→(b)→(c)→(d)). By this, the neighborhood information is well embedded in the Hamming space.

sufficient adaptivity as shown in Fig. 1. In comparison, such methods perform relatively poorly at the early training stages. Second, the state-of-the-art supervised methods are mostly time-consuming to update each incoming data stream [33], [35], [36]¹. In particular, given a query, MI-Hash [33] has to calculate the Hamming distance between its neighbors and non-neighbors to update the hashing functions, which inevitably increases the burden on time consumption. In [35], a pre-defined ECOC codebook [39] is used to guide the learning of hash functions, and hence eliminates the strong constraints on similarity preserving. However, the quality of short hashing codes still suffers poor performance due to the use of LSH to keep the consistency of code length and the size of ECOC codebook. The work in [36] considers the inner product based scheme, which adopts two equilibrium factors to solve the “data imbalance” problem and enables the usage of discrete optimization [9] in an online setting. However, on one hand, it introduces more variables to be optimized. On the other hand, the usage of discrete optimization faces the convergence problem. These two defects inevitably add more training burden. OSH [30] adopts an Online Boosting algorithm [40] to enhance the online learning ability. However, the usage of boosting inevitably causes training inefficiency. It is worth to note that, the earlier supervised methods like OKH [24] and AdaptHash [27] and unsupervised methods like SketchHash [26] and FROSH [34] are of high training efficiency, while their performance is far from satisfactory, which becomes worse when the training data increases (as quantitatively shown in Fig. 1 and validated later in Sec. 4.2).

To address the above issues, in this paper, we propose a novel supervised online hashing method, dubbed Fast Class-wise Updating for Online Hashing (FCOH). Our key innovation is introducing the usage of inner product to preserve the similarity relationship in the Hamming space, which has been demonstrated to be very effective in learning binary codes offline [12], [36], [41], [42], [43]. Unlike previous inner product based online hashing [36], we address the defects in online adaptivity and training efficiency in a unified framework. In this framework, we first

develop a new updating scheme that alternatively renews the hash functions in a class-wise fashion, the advantages of which are illustrated in Fig. 2. At each round, we only update samples from one class, which are well discriminated from the other classes. Therefore, hash functions can be fast adapted online with much less training data. Such a class-wise updating scheme differs from previous works that require all data involved in the training and therefore avoids large memory consumption. Quantitatively, such an inner product based class-wise updating can dramatically reduce the space complexity by at least 75%. Second, to accelerate the online training, we propose a semi-relaxation optimization which solves the binary constraints based on a divide-and-conquer method, *i.e.*, relaxing the part of the binary constraint as a continuous constraint, while considering the rest as a constant that can be pre-computed by the hash weight learned in the previous stage. In this way, no extra constraints and variables are involved, which highly reduces the training complexity. Thus, the proposed FCOH is more efficient with higher scalability, and can be well applied to large-scale applications with lower complexity than prior works. Besides, we show that the semi-relaxation optimization can well preserve the past information, which further boosts the performance.

The main contributions of our FCOH include:

- We introduce the usage of inner product to preserve the similarity relationship in the Hamming space. Different from previous inner product based online methods, the key innovation of our framework lies in solving the online adaptivity and training efficiency problems in a unified framework.
- We develop a class-wise updating scheme to effectively update the hash functions with much less training data and storage consumption in online hashing. Unlike traditional online hash methods that reach full performance at the final training stage, our method can be fast adapted online, which well addresses the online adaptivity problem.
- We devise a new semi-relaxation optimization, which significantly reduces the training time of the inner product method. Our method solves the binary constraints based on a divide-and-conquer optimization, by relaxing part of the binary constraints to be a continuous constraint. The rest is formulated as a constant that can be pre-computed by the hash weights learned in the last stage, which well preserves the past information.

The collective effort of the proposed class-wise updating and semi-relaxation optimization greatly boosts the performance of online hashing. Extensive experiments on three widely-used benchmarks, *i.e.*, CIFAR-10, Places205 and MNIST, demonstrate that the proposed FCOH obtains superior accuracy and efficiency over the state-of-the-art methods [24], [26], [27], [30], [33], [34], [35], [36].

The rest of this paper is organized as follows: In Sec. 2, we discuss the related work. The proposed FCOH including class-wise updating and semi-relaxation optimization is elaborated in Sec. 3. Sec. 4 reports the experimental results. Finally, we conclude this work in Sec. 5.

1. This statement is quantitatively validated in Sec. 4.4.

2 RELATED WORK

Supervised online hashing leverages label information to learn binary codes. To our best knowledge, Online Kernel Hashing (OKH) [24] is the first of this kind, which requires point pairs to update the hash functions via an online passive-aggressive strategy [44]. Adaptive Hashing (AdaptHash) [27] defines a hinge-like loss, which is approximated by a differentiable Sigmoid function to update the hash functions with SGD. In [30], a more general two-step hashing was introduced, in which binary Error Correcting Output Codes (ECOC) [45], [46], [47] are first assigned to labeled data, and then the hash functions are learned to fit the binary ECOC using online boosting. Cakir *et al.* [33] developed an Online Hashing with Mutual Information (MIHash), which targets at optimizing the mutual information between the neighbors and non-neighbors given a query. Lin *et al.* [35] proposed a Hadamard Codebook based Online Hashing (HCOH), where a more discriminative Hadamard matrix [39] is used as the ECOC codebook to guide the learning of hash functions. Recently, Balanced Similarity for Online Discrete Hashing (BSODH) [36] was developed to investigate the correlation between new data and existing dataset via an inner product fashion. To solve the “data imbalance” problem, BSODH adopts two equilibrium factors to balance the similarity matrix and enables the application of discrete optimization [9] in online learning. Though extensive progress has been made in supervised online hashing [24], [27], [30], [33], [35], [36], there still remain two critical problems: poor online adaptivity and poor training inefficiency, as discussed in Sec. 1.

Unsupervised online hashing is mainly designed based upon the idea of “data sketch” [38], where a small set of sketch data is used to preserve the main property of a large-scale dataset. To this end, Leng *et al.* [26] proposed an Online Sketching Hashing (SketchHash), which employs an efficient variant of SVD to learn hash functions, with a PCA-based batch learning on the sketch to learn hashing weights. A faster version of Online Sketch Hashing (FROSH) was developed in [34], where the independent Subsampled Randomized Hadamard Transform (SRHT) is employed on different data chunks to make the sketch more compact and accurate, and further accelerate the learning process. In general, unsupervised online hashing suffers low performance due to the lack of supervised labels.

3 THE PROPOSED FRAMEWORK

3.1 Problem Definition

Suppose the dataset is formed by a set of n vectors, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$, accompanied by a set of class labels $\mathbf{L} = [\mathbf{l}_1, \dots, \mathbf{l}_n] \in \mathbb{N}^n$. Besides, we denote $C = |\text{set}(\mathbf{L})|$ as the total categories of \mathbf{L} . We aim to learn a set of r -bit hashing codes $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n] \in \{-1, +1\}^{r \times n}$ such that the desired neighborhood structure is preserved. It is achieved by projecting the dataset \mathbf{X} using a set of r hash functions $H(\mathbf{X}) = \{h_i(\mathbf{X})\}_{i=1}^r$, i.e.,

$$\mathbf{B} = H(\mathbf{X}) = \text{sgn}(\mathbf{W}^T \mathbf{X}), \quad (1)$$

TABLE 1

List of Notations used in this paper.

Notation	Meaning
\mathbf{X}	$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$: set of n training samples; each column of \mathbf{X} corresponds to one sample
\mathbf{X}^t	$\mathbf{X} = [\mathbf{x}_1^t, \dots, \mathbf{x}_{n_t}^t] \in \mathbb{R}^{d \times n_t}$: set of n_t training samples at the t -th stage; each column of \mathbf{X}^t corresponds to one sample
\mathbf{B}	$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n] \in \{-1, +1\}^{r \times n}$: binary code matrix for \mathbf{X}
\mathbf{B}^t	$\mathbf{B}^t = [\mathbf{b}_1^t, \dots, \mathbf{b}_{n_t}^t] \in \{-1, +1\}^{r \times n_t}$: binary code matrix for \mathbf{X}^t
\mathbf{L}	$\mathbf{L} = [\mathbf{l}_1, \dots, \mathbf{l}_n] \in \mathbb{N}^n$: label set for \mathbf{X}
\mathbf{L}^t	$\mathbf{L}^t = [\mathbf{l}_1^t, \dots, \mathbf{l}_{n_t}^t] \in \mathbb{N}^{n_t}$: label set for \mathbf{X}^t
\mathbf{W}	$\mathbf{W} = \{\mathbf{w}_i\}_{i=1}^r \in \mathbb{R}^{d \times r}$: projection matrix for the learned hashing functions
\mathbf{W}^t	$\mathbf{W}^t = \{\mathbf{w}_i^t\}_{i=1}^r \in \mathbb{R}^{d \times r}$: projection matrix updated at the t -th stage
d	feature dimension
r	number of the hashing bits
C	total categories
C_t	total categories at the t -th stage
$\ \cdot\ _F$	Frobenius norm
$\ \cdot\ _1$	l_1 norm

where $\mathbf{W} = \{\mathbf{w}_i\}_{i=1}^r \in \mathbb{R}^{d \times r}$ is the projection matrix and \mathbf{w}_i is the i -th hash function. The sign function is defined as:

$$\text{sgn}(x) = \begin{cases} 1 & x > 1, \\ -1 & \text{otherwise.} \end{cases}$$

In the online setting, \mathbf{X} comes in a streaming fashion. Hence, for streaming data at the t -stage, we denote $\mathbf{X}^t = [\mathbf{x}_1^t, \dots, \mathbf{x}_{n_t}^t] \in \mathbb{R}^{d \times n_t}$ as the input streaming data, denote $\mathbf{B}^t = [\mathbf{b}_1^t, \dots, \mathbf{b}_{n_t}^t] \in \{-1, +1\}^{r \times n_t}$ as the learned binary codes for \mathbf{X}^t , and denote $\mathbf{L}^t = [\mathbf{l}_1^t, \dots, \mathbf{l}_{n_t}^t] \in \mathbb{N}^{n_t}$ as the corresponding label set, where n_t is the size of streaming data at the t -stage. Further, we denote $C^t = |\text{set}(\mathbf{L}^t)|$ as the total categories received at the t -stage. Correspondingly, the parameter \mathbf{W} updated at the t -stage is denoted as \mathbf{W}^t . Noticeably, \mathbf{W}^t can only be deduced via \mathbf{X}^t in an online setting.

In Tab.1, we summarize the notations used in the following contexts.

3.2 The Proposed Method

The proposed framework is built based upon the inner product based formulation [12], [36], [41], [42], [43]. The key ingredient is to map data points into binary codes, the inner product of which can well approximate the similarity matrix $\mathbf{S}^t \in \{-1, +1\}^{n_t \times n_t}$ where $\mathbf{S}_{ij}^t = 1$, if $\mathbf{l}_i^t = \mathbf{l}_j^t$, and -1 otherwise. More precisely, the goal of inner product based methods is to minimize the following objective function:

$$\mathcal{L} = \|(\mathbf{B}^t)^T \mathbf{B}^t - r \mathbf{S}^t\|_F^2 \quad \text{s.t.} \quad \mathbf{B}^t \in \{-1, +1\}^{r \times n_t}, \quad (2)$$

where $\|\cdot\|_F$ is the Frobenius norm.

The above formulation has been shown to be effective in traditional offline hashing methods [12], [41], [42], [43]. However, directly applying this equation to online learning is infeasible due to the “data imbalance” problem, as identified in [36]. Specifically, at the t -th stage, Eq.(2) can be re-written as:

$$\mathcal{L}^t = \underbrace{\sum_{i,j, \mathbf{S}_{ij}^t=1} ((\mathbf{b}_i^t)^T \mathbf{b}_j^t - r)^2}_{\text{term } \mathcal{A}} + \underbrace{\sum_{i,j, \mathbf{S}_{ij}^t=-1} ((\mathbf{b}_i^t)^T \mathbf{b}_j^t + r)^2}_{\text{term } \mathcal{B}} \quad (3)$$

s.t. $\mathbf{b}_i^t \in \{-1, 1\}^r, \mathbf{b}_j^t \in \{-1, 1\}^r$.

In an online setting, the similarity matrix \mathbf{S}^t is highly sparse, *i.e.*, most data pairs are dissimilar. Therefore in practice, we have $\text{term } \mathcal{B} \gg \text{term } \mathcal{A}$, which is referred to as the “data imbalance” problem. Since the learning of binary codes heavily relies on such dissimilar pairs, the retrieval performance cannot be satisfied as expected. To address this problem, the work in [36] introduces two equilibrium factors and enables the usage of discrete optimization in the online setting. However, such a solution needs a large number of training batches to learn up-to-date hash functions, which causes poor online adaptivity at the early training stages, as illustrated in Fig.1. Besides, it also introduces more variables to be optimized, and the discrete optimization further brings about the convergence problem. These two drawbacks unavoidably lead to the training inefficiency.

The novelty of our method lies in two-fold: First, to solve the problem of poor online adaptivity, we propose a “Class-wise Updating” scheme, which decomposes the optimization into several independent subprocesses with each responsible for one category, resulting in high accuracy and at least 75% storage saving as shown latter in Sec.3.2.1. Second, to address the problem of training inefficiency, we propose a “Semi-relaxation Optimization” scheme. In this scheme, one part of the binary constraints in Eq.(2) is relaxed as a continuous variable while the other is considered as a constant variable, through which the time complexity is drastically reduced. We detailedly describe the above two methods as below.

3.2.1 Class-wise Updating

To solve the poor online adaptivity, we develop a class-wise updating scheme as shown in Fig.2. Detailedly, we re-write Eq.(3) as follows:

$$\mathcal{L}^t = \sum_c \left(\underbrace{\sum_{i,j, \mathbf{l}_i^t=\mathbf{l}_j^t=c} ((\mathbf{b}_i^t)^T \mathbf{b}_j^t - r)^2}_{\text{term 1}} + \underbrace{\sum_{i,j, \mathbf{l}_i^t=c, \mathbf{l}_j^t \neq c} ((\mathbf{b}_i^t)^T \mathbf{b}_j^t + r)^2}_{\text{term 2}} \right) \quad (4)$$

s.t. $\mathbf{b}_i^t \in \{-1, 1\}^r, \mathbf{b}_j^t \in \{-1, 1\}^r$,

where c denotes the c -th class at the t -th training stage.

Term 1 stands for the inner product among instances from the c -th class, and *term 2* represents the inner product between instances from the c -th class and instances from classes excluding the c -th class. The basic idea behind class-wise updating is that at the t -th training stage, we decompose the training procedure into C_t independent subprocesses as illustrated in Fig.2. At the c -th subprocess, FCOH focuses on learning discriminative binary codes for the c -th class, which can well push the training instances of the c -th class away from instances from the other classes. Fig.3

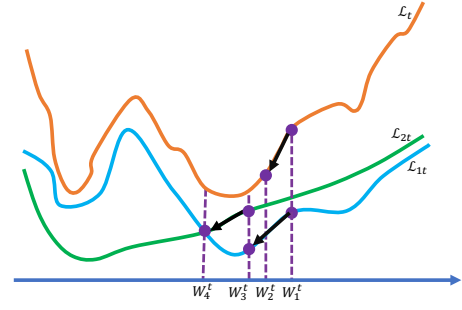


Fig. 3. Analysis on the class-wise updating. The blue and green curves represent the loss for the first class and the second class. The orange curve stands for the overall loss. \mathbf{W}_1^t is the start point for the t -th updating stage. Without class-wise updating, \mathbf{W}_1^t is updated for only once and the renewed weights are \mathbf{W}_5^t , which is not a good choice for both the first class and second class. With class-wise updating, the hash weights are updated gradient-by-gradient. It first updates hash functions using \mathcal{L}_{1t} and obtains \mathbf{W}_3^t , based on which \mathcal{L}_{2t} is used and then \mathbf{W}_4^t is obtained, serving as the final weights for the t -stage and shows an overall better result for class 1 and class 2.

illustrates a toy example that the class-wise updating renews the hash weights via a gradient-by-gradient fashion. After a total of C_t subprocesses, the final hash weights keep an overall better performance for all classes and obssesses a fast adaptivity to the coming data at the t -stage.

Without loss of generalization, we analyze that the decomposed Eq.(4) in the meantime can drastically reduce at least 75% space complexity in Eq.(2). To that effect, we first reformulate Eq.(4) in the following:

$$\mathcal{L}_t = \sum_c \mathcal{L}_{ct}, \quad (5)$$

where $\mathcal{L}_{ct} = \text{term 1} + \text{term 2}$ in Eq.(4), and it represents the loss objective for the c -th class. Further, we rewrite \mathcal{L}_{ct} into a matrix form:

$$\mathcal{L}_{ct} = \underbrace{\|(\mathbf{B}^{ct})^T \mathbf{B}^{ct} - r\|_F^2}_{\text{term 1}} + \underbrace{\|(\mathbf{B}^{ct})^T \mathbf{B}^{\bar{ct}} + r\|_F^2}_{\text{term 2}} \quad (6)$$

s.t. $\mathbf{B}^{ct} \in \{-1, +1\}^{r \times n_{ct}}, \mathbf{B}^{\bar{ct}} \in \{-1, +1\}^{r \times n_{\bar{ct}}}$,

where \mathbf{B}^{ct} is the matrix consisting of binary codes from the c -th class and $\mathbf{B}^{\bar{ct}}$ is the matrix consisting of binary codes excluding the c -th class; n_{ct} is the total number of training instances from the c -th class at the t -stage and $n_{\bar{ct}}$ is the total number of training instances excluding the c -th class at the t -stage. It is easy to derive that the space complexity in Eq.(6) is $\mathcal{O}(n_{ct}n_{\bar{ct}})$ instead of $\mathcal{O}(n_t^2)$ in Eq.(2). The storage saving rate R can be formulated as:

$$R = 1 - \frac{n_{ct}n_{\bar{ct}}}{n_t^2} \quad (7)$$

s.t. $n_{ct} + n_{\bar{ct}} = n_t$.

When $n_{ct} = n_{\bar{ct}} = \frac{1}{2}n_t$, $n_{ct}n_{\bar{ct}} = \frac{1}{4}n_t^2$ requires the largest memory consumption. In this case, the proposed FCOH obtains $R = 75\%$ storage saving rate. When $n_{ct} \neq n_{\bar{ct}}$, $n_{ct}n_{\bar{ct}} < \frac{1}{4}n_t^2$ and $R > 75\%$. Therefore, the proposed FCOH can reduce at least 75% storage consumption.

We note that the early work [48] adopts a similar class-wise updating scheme. However, our work still has distinct difference to [48]: First, [48] solves the offline hashing problem while FCOH solves the online hashing problem. Second,

[48] solves the offline hashing problem based on classifiers while FCOH solves the online hashing problem based on the inner product. Lastly, [48] proposes to class-wise innovate the binary codes, while FCOH proposes to class-wise update the hash weights.

3.2.2 Semi-relaxation Optimization

Due to the discrete constraint in Eq. (6), it is easy to derive that the above problem is highly non-convex and difficult (usually NP-hard) to solve. An alternative and suboptimal solution is to apply the discrete cyclic coordinate descent (DCC) developed in [9], *i.e.*, updating one hashing bit while fixing the other bits. However, as shown in the recent online hashing method [36], the discrete optimization brings more variables to be optimized and the bit-wise scheme is hard to converge. The above two issues inevitably lead to more training burden (as verified in Tab. 7), which should be further addressed in online hashing.

To achieve the efficient optimization, we adopt the relaxation method by partly removing the $\text{sgn}(x)$ function in Eq. (1). However, different from traditional relaxation process where the sign functions are all removed from the binary constraints, we propose to semi-relax the sign function in Eq. (6) based on a divide-and-conquer idea. It removes the sign function in part of the inner product of binary matrix, and regards the other part of binary matrix as a constant.

To that effect, we reformulate Eq. (6) as follows:

$$\begin{aligned} \hat{\mathcal{L}}_{ct} = & \lambda_1 \|((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{c(t-1)} - r\|_F^2 \\ & + \lambda_2 \|((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{\bar{c}(t-1)} + r\|_F^2, \end{aligned} \quad (8)$$

where

$$\mathbf{B}^{c(t-1)} = \text{sgn}((\mathbf{W}^{t-1})^T \mathbf{X}^{ct}) \in \{-1, +1\}^{r \times n_{ct}}, \quad (9)$$

and

$$\mathbf{B}^{\bar{c}(t-1)} = \text{sgn}((\mathbf{W}^{t-1})^T \mathbf{X}^{\bar{c}t}) \in \{-1, +1\}^{r \times n_{\bar{c}t}}, \quad (10)$$

which are *two constant matrices* pre-computed at the $(t-1)$ -th stage. And λ_1 and λ_2 are two hyper-parameters to balance the importance of the two terms.

The advantages of such operations are in two-fold: Firstly, $\mathbf{B}^{c(t-1)}$ and $\mathbf{B}^{\bar{c}(t-1)}$ are pre-computed by the $(t-1)$ -th hash weights. Therefore, not only training time is saved at the t -stage, but also more information is learned from the past stage. Meanwhile, in updating, the term $\mathbf{W}^{t-1} \mathbf{X}^{ct}$ aims to provide new information from the current stage. Therefore, the semi-relaxation can learn information from newly streaming data, and it can also preserve information from the last stage. Secondly, compared with the inner product based method in [36], such a semi-relaxation simplifies the optimization complexity, *i.e.*, only \mathbf{W}^t needs to be optimized, which greatly improves the training efficiency.

Nevertheless, Eq. (8) only considers the clues from the current and the last stage. As the training proceeds, the streaming data from early stages easily suffers large quantization errors. However, quantizing all the past streaming data is time-consuming as the dynamic dataset grows. Besides, it also violates the principle of online hashing since hashing functions should be updated only based on the newly arriving streaming data online. To solve issue,

Algorithm 1 Fast Class-wise Updating for Online Hashing (FCOH)

Input: Dataset \mathbf{X} and its label \mathbf{L} , number of hash bits r , parameter μ , λ_1 and λ_2 .

Output: Hash codes \mathbf{B} for \mathbf{X} and projection coefficient matrix \mathbf{W} .

```

1: Initialize  $\mathbf{W}^0$  with the normal Gaussian distribution.
2: for  $t = 1 \rightarrow T$  do
3:   for  $c = 1 \rightarrow C^t$  do
4:     if  $t = 1$  then
5:        $N_{ct} = n_{ct}$ ;
6:        $\bar{\mathbf{x}}^{ct} = \frac{\sum_{i=1}^{n_{ct}} \mathbf{x}_i^{ct}}{N_{ct}}$ ;
7:     else
8:        $N_{ct} = N_{c(t-1)} + n_{ct}$ ;
9:       Compute  $\bar{\mathbf{x}}^{ct}$  via Eq. (12);
10:    end if
11:    Compute  $\mathbf{B}^{c(t-1)}$  via Eq. (9);
12:    Compute  $\mathbf{B}^{\bar{c}(t-1)}$  via Eq. (10);
13:    Compute  $\frac{\partial \tilde{\mathcal{L}}_{ct}(\mathbf{W}^{t-1})}{\partial \mathbf{W}^{t-1}}$  via Eq. (15);
14:    Update  $\mathbf{W}^t$  via Eq. (14);
15:  end for
16: end for
17: Set  $\mathbf{W} = \mathbf{W}^T$ ;
18: Compute  $\mathbf{B} = \text{sgn}(\mathbf{W}^T \mathbf{X})$ ;
19: Return  $\mathbf{W}$  and  $\mathbf{B}$ .
```

inspired by [49], we further propose to quantize the centre of the c -th class that can be readily tracked, as follows:

$$\| |((\mathbf{W}^{t-1})^T \bar{\mathbf{x}}^{ct})| - 1 \|_1, \quad (11)$$

where $\|\cdot\|_1$ is the ℓ_1 -norm and $|\cdot|$ returns the absolute value, $\mathbf{1}$ is an all-one matrix and

$$\bar{\mathbf{x}}^{ct} = \frac{N_{c(t-1)} \bar{\mathbf{x}}^{c(t-1)} + \sum_{i=1}^{n_{ct}} \mathbf{x}_i^{ct}}{N_{ct}}, \quad (12)$$

where $N_{ct} = N_{c(t-1)} + n_{ct}$ and $N_{c1} = n_{c1}$. N_{ct} denotes the number of the c -th class and $\bar{\mathbf{x}}^{ct}$ is the center of the c -th class.

Combining Eq. (8) and Eq. (11) leads the final objective function as below:

$$\begin{aligned} \tilde{\mathcal{L}}_{ct} = & \| |((\mathbf{W}^{t-1})^T \bar{\mathbf{x}}^{ct})| - 1 \|_1 \\ & + \lambda_1 \|((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{c(t-1)} - r\|_F^2 \\ & + \lambda_2 \|((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{\bar{c}(t-1)} + r\|_F^2, \end{aligned} \quad (13)$$

To optimize Eq. (13), we adopt the SGD optimization to update the hash functions for the c -th class as follows:

$$\mathbf{W}^t \leftarrow \mathbf{W}^{t-1} - \mu \frac{\partial \tilde{\mathcal{L}}_{ct}(\mathbf{W}^{t-1})}{\partial \mathbf{W}^{t-1}}, \quad (14)$$

where μ is the learning rate. The derivative of $\tilde{\mathcal{L}}_{ct}(\mathbf{W}^{t-1})$ w.r.t. \mathbf{W}^{t-1} can be obtained as:

$$\begin{aligned} \frac{\partial \tilde{\mathcal{L}}_{ct}(\mathbf{W}^{t-1})}{\partial \mathbf{W}^{t-1}} = & \bar{\mathbf{x}}^{ct} \sigma((\mathbf{W}^{t-1})^T \bar{\mathbf{x}}^{ct})^T \\ & + 2\lambda_1 \left(\mathbf{X}^{ct} \left(((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{c(t-1)} - r \right) \right) (\mathbf{B}^{c(t-1)})^T \\ & + 2\lambda_2 \left(\mathbf{X}^{ct} \left(((\mathbf{W}^{t-1})^T \mathbf{X}^{ct})^T \mathbf{B}^{\bar{c}(t-1)} + r \right) \right) (\mathbf{B}^{\bar{c}(t-1)})^T, \end{aligned} \quad (15)$$

where the $\sigma(x)$ function is defined as:

$$\sigma(x) = \begin{cases} 1 & x > 1 \quad \text{or} \quad -1 < x < 0, \\ -1 & \text{otherwise.} \end{cases}$$

Without loss of generality, the main procedures of the proposed FCOH are outlined in Alg. 1.

3.3 Time Complexity

In Alg. 1, the majority of the training time is spent on the operations between line 11 to line 13 to update the hash functions on the c -th class at the t -stage. Specifically, the complexity of calculating $\mathbf{B}^{c(t-1)}$ in line 11 is $\mathcal{O}(rdn_{ct})$ and it is $\mathcal{O}(rdn_{\bar{c}t})$ for calculating $\mathbf{B}^{\bar{c}(t-1)}$ in line 12. The time cost for calculating derivative of $\hat{\mathcal{L}}_{ct}(\mathbf{W}^{t-1})$ w.r.t. \mathbf{W}^{t-1} in line 13 is $\mathcal{O}(rdn_{ct} + rn_{ct}^2 + dn_{ct}^2 + rn_{ct}n_{\bar{c}t} + dn_{ct}n_{\bar{c}t} + rdn_{\bar{c}t} + rd)$. Since $r \ll d$ and $n_{ct} \ll n_{\bar{c}t}$, the total complexity is $\mathcal{O}(rdn_{\bar{c}t} + dn_{ct}n_{\bar{c}t})$. We denote $e = \max(r, n_{ct})$ and the complexity can be re-written as $\mathcal{O}(edn_{\bar{c}t})$. Hence, the overall time complexity at the t -th training stage is $\mathcal{O}(C_t edn_{\bar{c}t})$. As can be seen in Sec. 4, C_t , e and $n_{\bar{c}t}$ are small values. Therefore, the time complexity of the proposed method mainly depends on the feature dimension d , which guarantees the efficiency and scalability of the proposed method.

4 EXPERIMENTS

To validate the model accuracy and learning efficiency, we compare our FCOH with several state-of-the-art methods [24], [26], [27], [30], [33], [35], [36] on three widely-used datasets, *i.e.*, CIFAR-10 [50], Places205 [51], and MNIST [52].

4.1 Experimental Settings

4.1.1 Datasets

CIFAR-10 consists of 60,000 images from 10 classes. Each class contains 6,000 instances and each image is represented as a 4,096-dim CNN vector [53]. As in [33], [35], [36], we randomly select 59,000 samples to form a retrieval set. And the remaining 1,000 are used to construct a test set. Besides, we randomly sample 20K images from the retrieval set as the training set to learn hash functions in a streaming fashion.

Places205 is a large-scale dataset [51] that contains 2.5 million images from 205 scene categories. The feature of each image is extracted from the AlexNet [54] and then reduced to a 128-dim feature by PCA. Following [35], 20 instances from each category are randomly sampled to construct the test set, and the others are used to form the retrieval set. Lastly, a random subset of 100K images from the retrieval set is used to learn the hash functions in a streaming fashion.

MNIST contains 70,000 handwritten digit images from 0 to 9 [52]. Each image is represented by $28 \times 28 = 784$ -dim normalized original pixels. According to the experimental setting in [36], we construct the test set by randomly sampling 100 instances from each class. The rest is used to form the retrieval set. Finally, a random subset of 20,000 images from the retrieval set is randomly sampled to form the training set to learn the hash functions in a streaming fashion.

The detailed analysis on the size of data stream, *i.e.*, n_t , is given in Sec. 4.3.1.

TABLE 2
Parameter Configurations on Three Benchmarks.

Method	CIFAR-10	Places205	MNIST
n_t	100	1,000	100
λ_1	1e-2	1e-1	1e-1
λ_2	1e-2	1e-1	1e-2
μ	1e-1	1e-4	1e-2

4.1.2 Evaluation

Following the previous methods [35], [36], we adopt the following protocols to evaluate the performance: mean Average Precision (denoted as *mAP*), Precision within a Hamming ball of radius 2 centered on each query (denoted as Precision@H2), *mAP* vs. different sizes of training instances curves and their corresponding areas under the *mAP* curves (denoted as AUC), Precision of the top K retrieved neighbors curves (denoted as Precision@K) and their corresponding areas under the Precision@K curves (denoted as AUC).

For Precision@K, the values of K are ranged in $\{1, 5, 10, 20, 30, \dots, 90, 100\}$ on all benchmarks, since the values of $K \leq 100$ are common choices and are closer to the user behavior when investigating the search results. Regarding *mAP* vs. different sizes of training instances, on CIFAR-10 and MNIST, the experimental results are sampled when the size of training instances falls in $\{2K, 4K, 6K, \dots, 18K, 20K\}$, and it is $\{5K, 10K, \dots, 95K, 100K\}$ on Places205 due to its large scale. Such settings provide a fair comparison with state-of-the-art method, *i.e.*, BSODH [36], where the best choice of streaming data size is 2K, 2K and 5K on CIFAR-10, MNIST and Places205, respectively.

Notably, when reporting the *mAP* performance on Places205, following the works in [33], [35], [36], we only compute the top 1,000 retrieved items (denoted as *mAP*@1,000) due to its large scale and heavy time consumption. The above metrics are evaluated with hashing bits varying among 8, 16, 32, 48, 64 and 128.

4.1.3 Baselines

To illustrate the effectiveness of the proposed FCOH, we compare our method with several state-of-the-art online hashing algorithms, including Online Kernel Hashing (OKH) [24], Online Sketching Hashing (SketchHash) [26], Adaptive Hashing (AdaptHash) [27], Online Supervised Hashing (OSH) [30], Online Hashing with Mutual Information (MIHash) [33], Hadamard Codebook based Online Hashing (HCOH) and Balanced Similarity for Online Discrete Hashing (BSODH) [36]. The source codes of these methods are publicly available. Our model is implemented with MATLAB. The experiments are conducted on a server with a 3.60GHz Intel Core I7 4790 CPU and 16G RAM.

4.1.4 Settings

To share the same dataset configurations on all the three benchmarks, we directly adopt the parameters as described in [35], [36], which have been carefully validated for each method. The following elaborates the parametric settings.

- **OKH:** The tuple (C, α) is set as (0.001, 0.3), (0.0001, 0.7) and (0.001, 0.3) on CIFAR-10, Places205 and MNIST, respectively.

TABLE 3
mAP And Precision@H2 Comparisons on CIFAR-10 with 8, 16, 32, 48, 64 And 128 Bits. The Best Result Is Labeled with Boldface And The Second Best Is with An Underline.

Method	mAP						Precision@H2					
	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit
OKH	0.100	0.134	0.223	0.252	0.268	0.350	0.100	0.175	0.100	0.452	0.175	0.372
SketchHash	0.248	0.301	0.302	0.327	-	-	0.256	0.431	0.385	0.059	-	-
AdaptHash	0.116	0.138	0.216	0.297	0.305	0.293	0.114	0.254	0.185	0.093	0.166	0.164
OSH	0.123	0.126	0.129	0.131	0.127	0.125	0.120	0.123	0.137	0.117	0.083	0.038
MIHash	0.512	0.640	0.675	0.668	0.667	0.664	0.170	0.673	0.657	0.604	0.500	0.413
HCOH	0.536	0.698	0.688	0.707	0.724	0.734	<u>0.333</u>	<u>0.723</u>	<u>0.731</u>	0.694	0.633	0.471
BSODH	0.564	0.604	<u>0.689</u>	0.656	0.709	0.711	0.305	0.582	0.691	<u>0.697</u>	0.690	<u>0.602</u>
FCOH	0.602	<u>0.674</u>	0.702	0.711	0.739	0.742	0.484	0.738	0.743	0.711	<u>0.648</u>	0.618

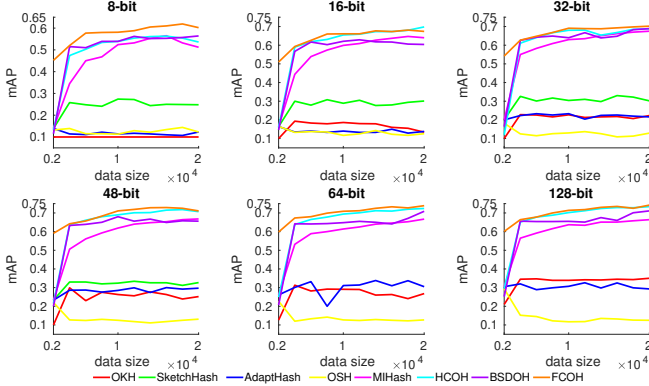


Fig. 4. mAP performance with respect to different sizes of training instances on CIFAR-10.

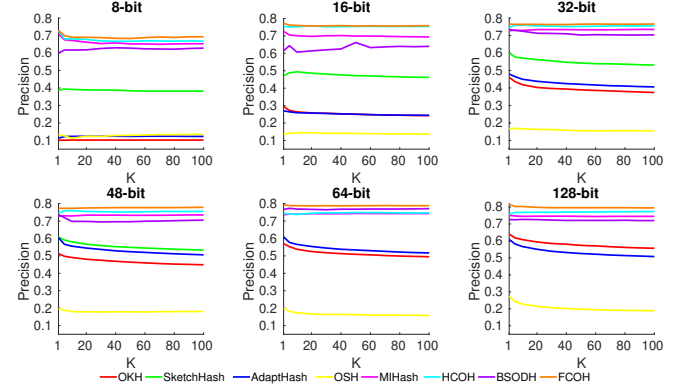


Fig. 5. Precision@K curves of compared algorithms on CIFAR-10.

- **SketchHash:** The tuple $(sketchsize, batchsize)$ is set as (200, 50), (100, 50) and (200, 50) on CIFAR-10, Places205 and MNIST, respectively.
- **AdaptHash:** The tuple (α, λ, η) is set as (0.9, 0.01, 0.1), (0.9, 0.01, 0.1) and (0.8, 0.01, 0.2) on CIFAR-10, Places205 and MNIST, respectively.
- **OSH:** On all datasets, η is set as 0.1 and the ECOC codebook C is populated the same way as in [30].
- **MIHash:** The tuple (θ, \mathcal{R}, A) is set as (0, 1000, 10), (0, 5000, 10) and (0, 1000, 10) on CIFAR-10, Places205 and MNIST, respectively.
- **HCOH:** The tuple (n_t, η) is set as (1, 0.2), (1, 0.1) and (1, 0.2) on CIFAR-10, Places205 and MNIST, respectively.
- **BSODH:** The tuple $(\lambda, \sigma, \eta_s, \eta_d)$ is set as (0.6, 0.5, 1.2, 0.2), (0.3, 0.5, 1.0, 0.0) and (0.9, 0.8, 1.2, 0.2) on CIFAR-10, Places205 and MNIST, respectively.

Specific descriptions of these parameters for each method can be found in [24], [26], [27], [30], [33], [35], [36], respectively. Tab.2 shows the parameter settings of the proposed FCOH, which are carefully tuned and validated. Emphatically, for SketchHash [26], the training size in each data stream has to be larger than the code length. Following previous methods in [35], [36], we only report its experimental results with hashing bit being 8, 16, 32, 48. All the experiments are run over three times and the averaged results are reported.

4.2 Results and Discussions

4.2.1 Results on CIFAR-10

Regarding mAP and Precision@H2, three observations can be found: First, as can be seen in Tab.3, FCOH achieves

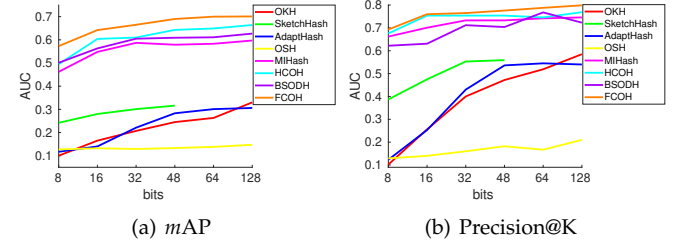


Fig. 6. AUC curves for mAP and Precision@K on CIFAR-10.

the best results in most cases. Quantitatively, FCOH surpasses state-of-the-art methods like HCOH or BSODH by an average of 1.49% mAP and 7.94% Precision@H2, which demonstrates the effectiveness of FCOH. Second, for most methods, the mAP goes up as the code length increases. In the case of 16-bit, BSODH increases a lot, which is even better than the proposed FCOH. Third, different from mAP, Precision@H2 drops a lot for most methods in high bits (64 and 128). To analyze the contrast between the second and the third observations, more information can be encoded in high bits which is beneficial to the linear scan like mAP, while the number of distinct hash buckets to examine grows rapidly with the hash bit r , which is harmful to the retrieval scenario that directly returns data points within Hamming radius 2, i.e., Precision@H2 [18].

We further analyze the mAP curves over different sizes of training instances and the corresponding AUC curves in Fig.4 and Fig.6(a), which reflect our adaptivity for online learning. Two observations can be found in Fig.4: First, FCOH achieves substantially better performance at different training stages. To take an in-depth analysis, the AUC curve for FCOH in Fig.6(a) outperforms the second best (HCOH) by a margin of 8.63%, which demonstrates the robustness of FCOH. Second, FCOH obtains much better performance at

TABLE 4

$mAP@1,000$ And Precision@H2 Comparisons on Places205 with 8, 16, 32, 48, 64 And 128 bits. The Best Result Is Labeled with Boldface And the Second Best Is with An Underline.

Method	$mAP@1,000$						Precision@H2					
	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit
OKH	0.018	0.033	0.122	0.048	0.114	0.258	0.007	0.010	0.026	0.017	0.217	0.075
SketchHash	0.052	0.120	0.202	0.242	-	-	0.017	0.066	0.220	0.176	-	-
AdaptHash	0.028	0.097	0.195	0.223	0.222	0.229	0.009	0.051	0.012	0.185	0.021	0.022
OSH	0.018	0.021	0.022	0.032	0.043	0.164	0.007	0.009	0.012	0.023	0.030	0.059
MIHash	0.094	<u>0.191</u>	0.244	0.288	0.308	0.332	<u>0.022</u>	<u>0.112</u>	0.204	0.242	0.202	0.069
HCOH	0.049	0.173	0.259	0.280	<u>0.321</u>	<u>0.347</u>	0.012	0.082	<u>0.252</u>	0.179	0.114	0.036
BSODH	0.035	0.174	0.250	0.273	0.308	0.337	0.009	0.101	0.241	<u>0.246</u>	<u>0.212</u>	<u>0.101</u>
FCOH	0.099	0.198	0.267	0.302	0.335	0.353	0.024	0.124	0.261	0.263	0.223	0.141

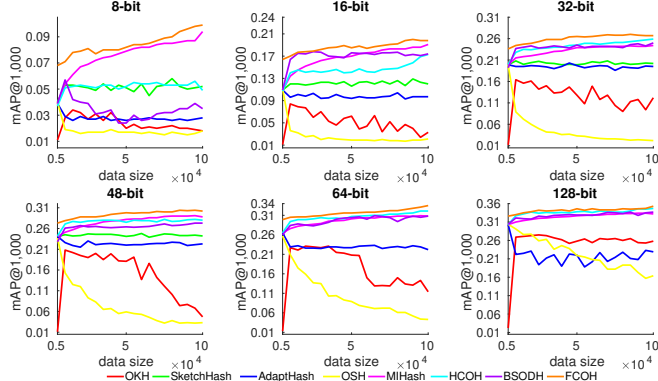


Fig. 7. mAP performance with respect to different sizes of training instances on Places205.

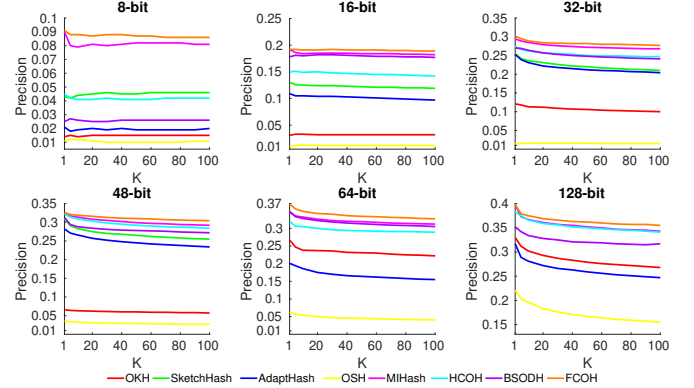


Fig. 8. Precision@K curves of compared algorithms on Places205.

the early training stages (data size is 2,000) in different code lengths, which effectively validates the online adaptivity of the FCOH.

The results of Precision@K and their AUC curves on CIFAR-10 can be seen in Fig.5 and Fig.6(b). Generally, FCOH shows the best performance in all code lengths. When the code length ≥ 48 , the superiority is more evident. Regarding the AUC of Precision@K in Fig.9(b), FCOH transcends the second best, *i.e.*, MIHash, by an average of 6.11%. Notably, combining the Precision@K results in Fig.5 and the mAP results in Tab.3, we can conclude that, given a query, the proposed HCOH not only retrieves more relevant items, but also ranks them in the top places.

4.2.2 Results on Places205

Places205 is a large-scale and challenging benchmark, the results on which can well reflect the performance in real-world applications. We start with an analysis of $mAP@1,000$ and Precision@H2 in Tab.4. As can be observed, FCOH not only ranks first but also obtains a relative increase of 3.84% for $mAP@1,000$ and 12.51% for Precision@H2, comparing to the best baselines, HCOH or BSODH, which demonstrates the scalability of the proposed method for large-scale applications. Moreover, the observations on CIFAR-10 can be also found in Places205, *i.e.*, better mAP and degenerated Precision@H2 in longer code length. The explanations are the same as in Sec.4.2.1, *i.e.*, longer codes are beneficial to the linear scan of mAP while harmful to the Precision@H2 which directly returns data points with Hamming radius 2. Besides, comparing Tab.3 with Tab.4, for all methods, both mAP and Precision@H2 performance on Places205 are much lower than those on CIFAR-10. To analyze, this is due to the large scale of Places205 (millions). It is quite challenging

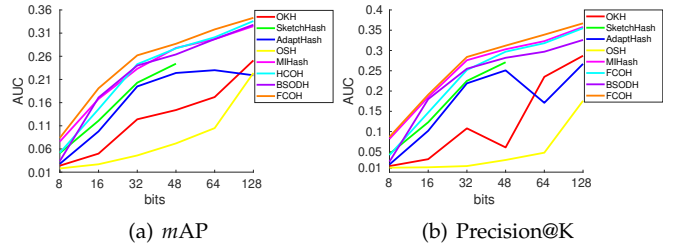


Fig. 9. AUC curves for mAP and Precision@K on Places205.

to obtain high performance on such a benchmark. Nevertheless, with regard to the comparisons between different methods, FCOH yields the best in all the tests, which verifies the feasibility of FCOH in real-world applications.

Given the results of $mAP@1,000$ in Fig.7, we can see that FCOH performs best in all cases. It demonstrates the robustness of FCOH when the training proceeds as reflected in Fig.9(a), where FCOH gains an average of 18.14% AUC improvements compared to HCOH. Besides, it also shows the fast online adaptivity of FCOH since the proposed method grows fast and outperforms other methods by large margins at the early training stages (data size 5,000). Note that when the hashing bits range from 8 to 64, FCOH shows a better margin than other methods. However, in the case of 128-bit, FCOH and BSODH show a similar performance (FCOH is slightly better). To explain, as mentioned in Sec.4.2.1, the mAP increases with the increase of the code length. When the code length is 128, other methods can also obtain relatively high results. Nevertheless, in all code lengths, the proposed method consistently shows the best performance, which further demonstrates the superiority of FCOH.

Lastly, the Precision@K performance in Fig.8 and their AUC results Fig.9(b) demonstrate again the effectiveness of

TABLE 5

mAP And Precision@H2 Comparisons on MNIST with 8, 16, 32, 48, 64 And 128 Bits. The Best Result Is Labeled with Boldface And the Second Best Is with An Underline.

Method	mAP						Precision@H2					
	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit	8-bit	16-bit	32-bit	48-bit	64-bit	128-bit
OKH	0.100	0.155	0.224	0.273	0.301	0.404	0.100	0.220	0.457	0.724	0.522	0.124
SketchHash	0.257	0.312	0.348	0.369	-	-	0.261	0.596	0.691	0.251	-	-
AdaptHash	0.138	0.207	0.319	0.318	0.292	0.208	0.153	0.442	0.535	0.335	0.163	0.168
OSH	0.130	0.144	0.130	0.148	0.146	0.143	0.131	0.146	0.192	0.134	0.109	0.019
MIHash	0.664	0.741	0.744	0.780	0.713	0.681	0.487	0.803	0.814	0.739	0.720	0.471
HCOH	0.536	0.708	0.756	0.772	0.759	<u>0.771</u>	0.350	0.800	0.826	0.766	0.643	0.370
BSODH	0.593	0.700	0.747	0.743	0.766	0.760	0.308	0.709	<u>0.826</u>	<u>0.804</u>	<u>0.814</u>	0.643
FCOH	0.673	<u>0.725</u>	0.786	0.789	0.784	0.801	0.506	0.817	0.849	0.814	0.817	<u>0.620</u>

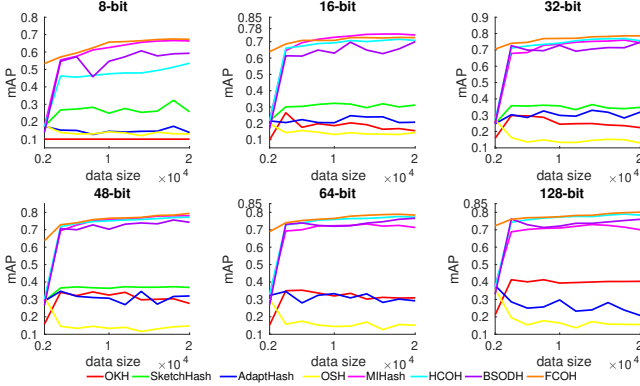


Fig. 10. mAP performance with respect to different sizes of training instances on MNIST.

FCOH. As can be observed, on this large-scale and challenging dataset, FCOH obtains higher precision results in Fig. 8 and transcends the state-of-the-art MIHash with an AUC gain of 3.78%. Together with the $mAP@1,000$ performance in Fig. 7, we can conclude that, given a query, FCOH not only retrieves more relevant items but ranks them in the top position, which highly meets the need of real-world applications.

4.2.3 Results on MNIST

In Tab. 5, the proposed method shows the best results except for the cases of mAP in 16-bit and Precision@H2 in 128-bit, where FCOH ranks the second. The experimental results in Tab. 5 (on MNIST) are similar to those in Tab. 3 (on CIFAR-10). On one hand, the quantitative values on MNIST and CIFAR-10 are much larger than those on Places205. On the other hand, regarding mAP in 16-bit, MIHash obtains the best on both MNIST and CIFAR-10. In terms of Precision@H2 in high bits, BSODH obtains the best on both MNIST (64-bit) and CIFAR-10 (128-bit). This is because MNIST and CIFAR-10 have the similar quantitative scale (60K ~ 70K) and categories (10 for both), which are two relatively simple benchmarks compared with the scale of Places205 (millions of samples and 205 categories). Nevertheless, FCOH still yields an overall better performance. Generally speaking, the proposed FCOH achieves state-of-the-art results with an average increase of 1.76% for mAP and 1.08% for Precision@H2 comparing to the second best. Hence, the superiorities of FCOH are demonstrated.

Besides, we argue that Fig. 10 and Fig. 12(a) also prove the online adaptivity of FCOH. First, FCOH keeps consistently better mAP in different sizes of training instances and obtains an increase of 9.83% AUC compared with the best

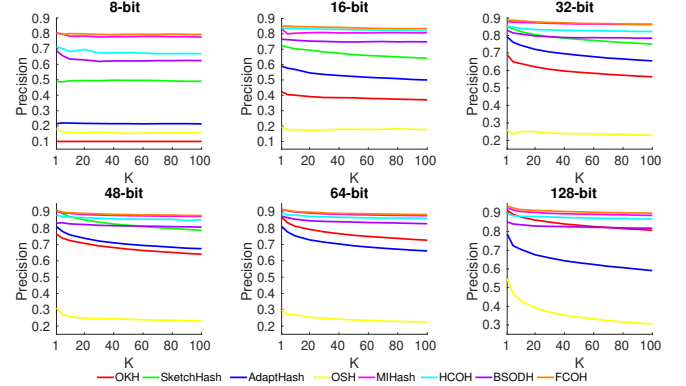


Fig. 11. Precision@K curves of compared algorithms on MNIST.

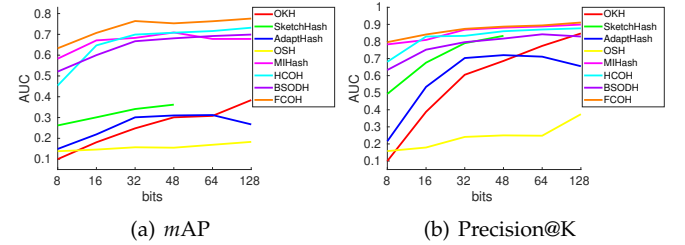


Fig. 12. AUC curves for mAP and Precision@K on MNIST.

baseline HCOH. On the other hand, FCOH transcends other methods in early training stages by a large margin. Taking the hash bit of 64 as an instance, FCOH gets an mAP of 0.689 while it is only 0.323 for MIHash, 0.329 for HCOH and 0.270 for BSODH when the training size is 2,000. Hence, the fast adaptivity ability of FCOH is well demonstrated.

As for Precision@K, FCOH is still competitive in all cases with different code lengths, as reflected in Fig. 11. Quantitatively speaking, FCOH obtains an AUC increase of 1.48% compared with the second best of MIHash. Besides, different from the results on CIFAR-10 and Places205, the increments on MNIST are small. To analyze, MNIST is a simple benchmark and it is easy for existing methods to obtain high performance. When the code length ≤ 16 , existing methods can achieve an accuracy of more than 80%. When the code length ≥ 32 , existing methods can obtain a high accuracy of near 90%. Hence, it is tough to further obtain a significant improvement upon MNIST.

4.2.4 In-depth Analysis on Better Performance

As discussed in Sec. 4.2.1, Sec. 4.2.2 and Sec. 4.2.3, FCOH shows an overall better performance than the SOTAs. To analyze, class-wise updating focuses on minimizing the intra-class distance and pulls other classes away from the

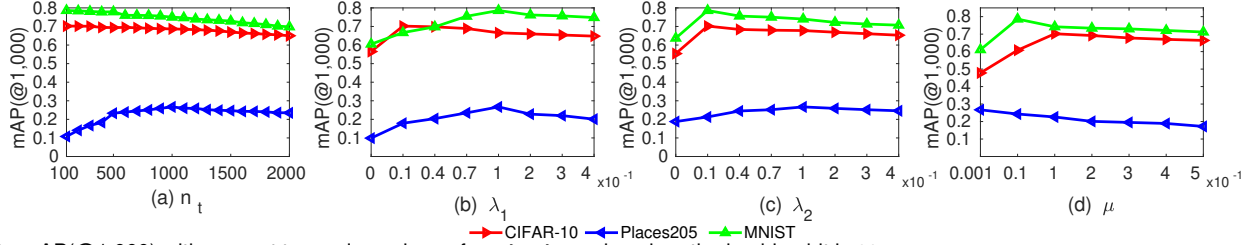


Fig. 13. $mAP(@1,000)$ with respect to varying values of n_t , λ_1 , λ_2 and μ when the hashing bit is 32.

TABLE 6
 $mAP(@1,000)$ Results of 32-Bit for Different Variants of the Proposed Method on the Three Datasets.

	CIFAR-10	Places205	MNIST
BSODH (baseline)	0.689	0.250	0.747
$FCOH_S$	0.691	0.255	0.760
$FCOH_C$	0.696	0.262	0.777
$FCOH$	0.702	0.267	0.786

current optimized class, leading to a fast adaptivity to the new data stream as shown in Sec. 3.2.1. Semi-relaxation optimization treats part of the binary constraints as a constant, which is computed by hash weights learned from the last round. Thus, it can update the hash model by using the new data stream, meanwhile preserving the past information as shown in Sec. 3.2.2. Therefore, the better performance of FCOH is the collective effort of both class-wise updating and semi-relaxation optimization.

By removing the class-wise updating and semi-relaxation, our FCOH is the same with BSODH [36], which can serve as the baseline to show the effectiveness of the proposed class-wise updating and semi-relaxation. To that effect, two variants of FCOH are introduced. The first variant is denoted as $FCOH_S$ where data from all classes are involved in the training at each round and the proposed semi-relaxation is kept. The second variant is denoted as $FCOH_C$ where the class-wise updating is kept and the semi-relaxation is discarded. Instead, we apply the discrete cyclic coordinate descent (DCC) developed in [9], as mentioned in Sec. 3.2.2. DCC is also the standard optimizer of BSODH [36]. Then, we take 32-bit as an example, and show the $mAP(@1,000)$ performance on the three datasets. The experimental results are listed in Tab. 6.

As can be seen, FCOH, which contains both class-wise updating and semi-relaxation optimization, shows the optimal results, and its variants, *i.e.*, $FCOH_S$ and $FCOH_C$ obtain less $mAP(@1,000)$ performance. Among all, BSODH without both class-wise updating and semi-relaxation optimization results in the least performance. Thus, it well verifies the analysis above that the better results of FCOH are the collective effort of the proposed class-wise updating and semi-relaxation optimization.

4.2.5 Convergence of the Class-wise Updating

In Fig. 3, we illustrate a toy example which shows that the class-wise updating can obtain a better solution. In this section, we visualize the updating of Eq. 3 with/without class-wise updating in Fig. 14 for a rigorous demonstration². We implement Fig. 14 with two-class synthetic data and the hash codes are updated using the DCC [9] to exclude

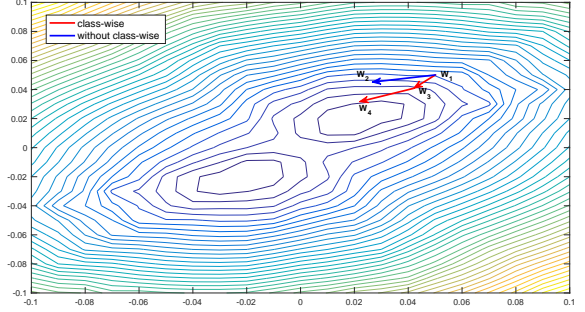


Fig. 14. Visualization with/without the class-wise updating. W_1 is the start point. Without class-wise updating, the renewed weights come at W_2 . With class-wise updating, the hash weights are updated gradient-by-gradient. It first updates hash functions using data from the first class, which returns W_3 , based on which the second-class data is used and then W_4 is obtained. As can be seen, the final W_4 rewards a better solution than W_2 without class-wise updating.

the effect of semi-relaxation optimization. It can be seen that updating with a class-wise fashion reaches a more optimal solution than the one without class-wise updating, which shows that class-wise updating provides a more effective facility for online hash learning.

4.3 Ablation Study

In this subsection, we conduct the ablation studies on the hyper-parameters of FCOH, including the batch size of streaming data, *i.e.*, n_t , as defined in Sec. 3.1; the balance parameters, *i.e.*, λ_1 and λ_2 , as defined in Eq. (8); and the learning rate μ , as defined in Eq. (14). Without loss of generality, we conduct experiments with varying values of these hyper-parameters with respect to mAP ($mAP@1,000$) in the case of 32-bit in Fig. 13 (Detailed values used in this paper is outlined in Tab. 2.). Similar experimental results can be observed in other hashing bits.

4.3.1 The Influence of n_t

We first start with the analysis on the varying values of batch size n_t at the t -stage. The experiments are sampled with n_t ranged in $\{100, 200, \dots, 2,000\}$ and the results are shown in Fig. 13(a). As can be observed, on CIFAR-10 and MNIST, the $mAP(mAP@1,000)$ degenerates along with the increase of batch size n_t , while on Places205, the $mAP(mAP@1,000)$ achieves the optimal when $n_t \approx 1,000$. In the experiments, we empirically set the values of n_t as 100 on CIFAR-10 and MNIST, and 1,000 on Places205 as stated in Tab. 2. Such settings conform with the requirements of online learning since $n_t \ll n$ in our experiments ($n = 20K$ on CIFAR-10 and MNIST, and $n = 100K$ on Places205).

We observe that for CIFAR-10 and MNIST, they have similar data size and share the same batch size. As for

² The figure is plotted using contour function with MATLAB: <https://www.mathworks.com/help/matlab/ref/contour.html>.

Places205, it has a much larger data size than CIFAR-10 and MNIST, and the required batch size is also larger than the others. To explain, Places205 is a large-scale retrieval set. Setting a smaller value makes it hard to capture and adapt to the data property. The three datasets represent different scenes in real life. These different scenes typically require different configurations to deploy models. Hence, it's intuitive that different datasets have different optimal n_t values.

4.3.2 The Influence of λ_1

As shown in Eq. (13), λ_1 is used to reflect the importance of similar items in the inner product based learning scheme. Fig. 13(b) plots the influence of different values of λ_1 to the performance. Generally speaking, when $\lambda_1 = 0.01$ on CIFAR-10, $\lambda_1 = 0.1$ on Places205 and MNIST, FCOH obtains the best mAP ($mAP@1,000$) on the three benchmarks (0.702 on CIFAR-10, 0.267 on Places205 and 0.786 on MNIST). Besides, when $\lambda_1 = 0$, FCOH suffers a great performance loss as can be seen in Fig. 13(b). More specifically, in this case, the mAP ($mAP@1,000$) scores are 0.565, 0.099 and 0.605 on CIFAR-10, Places205 and MNIST, respectively. To analyze, when $\lambda_1 = 0$, FCOH only relies on dissimilar items to learn hash functions. The to-be-learned binary codes will be separable from each other. Therefore, the performance degenerates a lot. In the experiments, we empirically set the values of λ_1 as 0.01 on CIFAR-10, and 0.1 on Places205 and MNIST.

4.3.3 The Influence of λ_2

As shown in Eq. (13), λ_2 is used to reflect the importance of dissimilar items in the inner product based learning scheme. From Fig. 13(c), we can observe that when $\lambda_2 = 0.01$ on CIFAR-10 and MNIST, $\lambda_2 = 0.1$ on Places205, FCOH obtains the best mAP ($mAP@1,000$) of 0.702, 0.267 and 0.786, respectively. Similar to λ_1 , when $\lambda_2 = 0$, FCOH also suffers a great performance loss (0.554 on CIFAR-10, 0.188 on Places205 and 0.637 on MNIST). However, the explanation behind this phenomenon is different from the case of $\lambda_1 = 0$. Specifically, when $\lambda_2 = 0$, the proposed FCOH learns hash functions simply through similar items. Under this situation, the learned binary codes will be as close as possible, which inevitably causes the performance degeneration. In the experiments, we empirically set the values of λ_2 as 0.01 on CIFAR-10 and MNIST, and 0.1 on Places205.

Combining Sec. 4.3.2 and Sec. 4.3.3 together, the optimal values for the tuple (λ_1, λ_2) are (0.01, 0.01) on CIFAR-10, (0.1, 0.1) on Places205, and (0.1, 0.01) on MNIST, respectively. Note that, the tuple set above shows the equal (or approximate) importance between similar items and dissimilar items since λ_1 and λ_2 share the same (or similar) values on the three benchmarks. That is to say, the proposed class-wise updating strategy can well solve the problem of “data imbalance” since both similar information and dissimilar information are equally learned in our framework³.

3. To analyze, the proposed class-wise updating method reduces the quantitative difference between similar items and dissimilar items each time when updating the hash functions. For example, there have 10 similar items and 100 dissimilar items. Through class-wise updating, there may have 2 similar items and 10 dissimilar items when updating the c -th category.

4.3.4 The Influence of μ

Lastly, the influence of learning rate μ is shown in Fig. 13(d). Generally, $mAP(mAP@1,000)$ is sensitive to the varying μ . To obtain the best $mAP(mAP@1,000)$, in the experiments, we empirically set the values of μ as 0.01, 0.0001 and 0.1 on CIFAR-10, Places205 and MNIST, respectively.

Besides the observation that $\lambda_1 \approx \lambda_2$ as mentioned in Sec. 4.3.3, we empirically find that the suitable value of λ_1 , λ_2 and μ fall in 10^{-m} , where $m = 1, 2, 3, 4, 5$. Hence, in practical applications, the users can try these different values and choose the best one.

4.4 Training Efficiency

To quantitatively evaluate the training efficiency of FCOH, including the updating of hash functions and hash table, we conduct experiments given hashing bit $r = 32$ in Tab. 7. Similar observations can also be found in other code lengths.

Hash Function Updating. Generally, OKH and Sketch-Hash are the most efficient ones, which however suffer poor mAP ($mAP@1,000$). Compared with state-of-the-art methods (MIHash and BSODH), FCOH shows consistently better efficiency by a margin. As analyzed in Sec. 1, MIHash has to calculate the Hamming distance between the neighbors and non-neighbors for each instance, while the discrete optimization adopted in BSODH brings about more variables and has the convergence problem. Compared with HCOH, the proposed method still costs less training time on CIFAR-10 and MNIST, while on Places205 HCOH is more efficient. To analyze, FCOH asks for a larger batch size on Places205 ($n_t = 2,000$ in Tab. 2.) while $n_t = 1$ for HCOH. Nevertheless, HCOH suffers poor performance especially in low hashing bits in the experiments. It can be observed that FCOH costs more training time on CIFAR-10 and Places205 than on MNIST. We argue that this is owing to the high feature dimension of CIFAR-10 (4,096-dim) and large scale of Places205 (100,000).

Hash Table Updating. We further compare the time consumption on hash table updating. The hash table needs to be renewed every time when updating the hash functions. It is easy to know that the total times spent on the hash table updating depends on the times of updating hash functions, i.e., $\frac{n}{n_t}$, where n is the size of the dataset and n_t is the size of image batch at each updating stage. From Tab. 7, it can be seen that OKH and HCOH are most inefficient in hash table updating. This is because the batch size for them is 1, which means more updating of hash table is required. Note that, MIHash also requires $n_t = 1$, however, it consumes less time in hash table updating. To analyze, MIHash does not update hash table each time the hash functions are updated. Instead, it devises a “gating” mechanism where the hash table is updated only when the mutual information of updated hash model is better than the currently best hash model. Hence, MIHash requires fewer times of hash table updating. Besides, we observe that BSODH requires the least time of hash table updating. Taking a deeper analysis, BSODH has the largest batch size ($n_t = 2,000$ for CIFAR-10 and MNIST, and 5,000 for Places205). Lastly, FCOH also merits in its efficient hash table updating as reflected in Tab. 7, especially better than the gating-designed MIHash. Though BSODH is more efficient in hash table updating

TABLE 7
Time Consumption on Hash Function Updating and Hash Table Updating on the Three Benchmarks under 32-Bit Hashing Codes.

Method	CIFAR-10 (s)			Places205 (s)			MNIST (s)		
	Hash Function	Hash Table	Total	Hash Function	Hash Table	Total	Hash Function	Hash Table	Total
OKH	4.53	1,600	1,604.53	15.66	12,500	12,515.66	4.58	400	404.58
SketchHash	4.98	64	68.98	3.52	500	503.52	1.27	16	17.27
AdaptHash	20.73	1,600	1,620.73	14.49	12,500	12,514.49	6.26	400	406.26
OSH	93.45	3,200	3,293.45	56.68	25,000	25,056.68	24.07	800	824.07
MIHash	120.10	103.20	223.30	468.77	437.25	906.02	97.59	24.52	122.11
HCOH	12.34	3,200	3,212.34	10.54	25,000	25,010.54	4.01	800	804.01
BSODH	36.12	1.60	37.72	69.73	5	74.73	4.83	0.40	5.23
FCOH	7.81	32	39.81	15.27	25	40.27	1.64	8	9.64

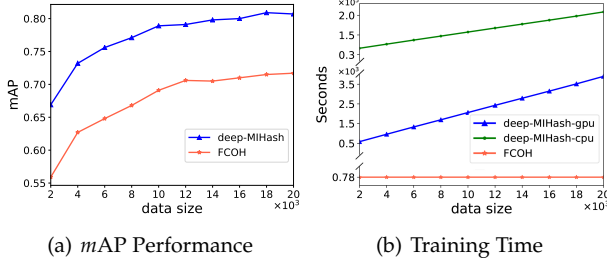


Fig. 15. mAP performance and time consumption on hash function updating at each training stage on CIFAR-10.

than FCOH, the large batch size makes it inefficient in hash function updating.

Above all, the proposed method is efficient in both hash function updating and hash table updating. Notably, the main focus of this paper is to design an efficient and effective method for hash function updating. Nevertheless, some existing works [55], [56] which are particularly designed for hash table updating can be integrated into the above online hashing methods to further improve the efficiency of hash table updating.

4.5 Comparisons with offline Hashing Method

The main focus of online hashing is to efficiently deal with the streaming data. Traditional offline hashing could also do this job by re-training the model with the accumulated data, which however has extremely heavy training consumption, but merits in better accuracy. In this subsection, we further conduct experiments to compare with the state-of-the-art deep-MIHash [23], which is an offline extension of MIHash [33]. To that effect, we conduct experiments in 32-bit on CIFAR-10, and analyze the mAP performance and their training time, the results of which are shown in Fig. 15. For deep-MIHash, we display its training time with both GPU and CPU. As can be seen, deep-MIHash obtains better mAP performance, *i.e.*, around 10% improvements over FCOH at each updating stage. However, FCOH gains a magnitude-order reduction of time consumption in updating hash functions, compared to deep-MIHash with either GPU or CPU. To analyze, deep-MIHash has to accumulate all the available data to re-train the hash functions. Thus, it avoids the information loss of past streaming data. However, the training time also linearly grows with the increase of streaming data.

4.6 Limitation and Future Work

Compared with another inner product based method, BSODH, FCOH is an inner product based scheme with

two innovations, *i.e.*, *class-wise updating* and *semi-relaxation optimization*. It merits in three aspects, *i.e.*, *at least 75% storage saving* (see Sec.3.2.1), *better accuracy and fast online adaptivity* (see Sec.4.2) and *training efficiency* (see Sec.4.4). Nevertheless, it can only deal with the single-label data due to the design of the class-wise updating scheme, which requires mutually exclusive class labels. Such a strict requirement might not be satisfied in many cases, *e.g.*, multi-label datasets, which could limit its application in some real-world applications. More efforts will be made to solve this issue in our future work.

5 CONCLUSION

In this paper, we present a novel supervised online hashing method, termed FCOH, which is an inner product based scheme, which addresses the poor online adaptivity and training inefficiency problems that prevail in existing online hashing methods. To this end, we first develop a class-wise updating scheme that iteratively renews the hash functions at each stage, through which FCOH achieves better performance with much less training data and storage consumption. Second, we propose a semi-relaxation optimization that relaxes a part of the binary constraints while treating the other part as a constant binary matrix by pre-computing it, through which no extra variables are introduced with less time consumption. Extensive experiments on three widely-used benchmarks show that FCOH achieves state-of-the-art accuracy and efficiency in the online retrieval tasks.

ACKNOWLEDGMENTS

This work is supported by the National Natural Science Foundation of China (No. 62025603, No.U1705262, No.61772443, No.61572410, No.61802324 and No.61702136). It is also sponsored by CCF-Baidu Open Fund and Australian Research Council Project FL-170100117.

REFERENCES

- [1] Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in *Proceedings of the NeurIPS*, 2009.
- [2] W. Liu, J. Wang, S. Kumar, and S. Chang, "Hashing with graphs," in *Proceedings of the ICML*, 2011.
- [3] J. Gui and P. Li, "R2sdh: Robust rotated supervised discrete hashing," in *Proceedings of the ACM SIGKDD*, 2018, pp. 1485–1493.
- [4] Y. Gong and S. Lazebnik, "Iterative quantization: A procrustean approach to learning binary codes," in *Proceedings of the CVPR*, 2011.
- [5] J. Wang, S. Kumar, and S.-F. Chang, "Semi-supervised hashing for large-scale search," *IEEE TPAMI*, 2012.
- [6] X. Zhu, Z. Huang, H. T. Shen, and X. Zhao, "Linear cross-modal hashing for efficient multimedia search," in *Proceedings of the ACM MM*, 2013.

- [7] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin, "Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval," *IEEE TPAMI*, 2013.
- [8] L. Liu and L. Shao, "Sequential compact code learning for unsupervised image hashing," *IEEE TNNLS*, 2015.
- [9] F. Shen, C. Shen, W. Liu, and H. T. Shen, "Supervised discrete hashing," in *Proceedings of the CVPR*, 2015.
- [10] M. Yu, L. Liu, and L. Shao, "Binary set embedding for cross-modal retrieval," *IEEE TNNLS*, 2016.
- [11] Y. Long, L. Liu, F. Shen, L. Shao, and X. Li, "Zero-shot learning using synthesised unseen visual data with diffusion regularisation," *IEEE TPAMI*, 2017.
- [12] P.-F. Zhang, C.-X. Li, M.-Y. Liu, L. Nie, and X.-S. Xu, "Semi-relaxation supervised hashing for cross-modal retrieval," in *Proceedings of the ACM MM*, 2017.
- [13] H. Liu, M. Lin, S. Zhang, Y. Wu, F. Huang, and R. Ji, "Dense auto-encoder hashing for robust cross-modality retrieval," in *Proceedings of the ACM MM*, 2018.
- [14] M. Long, Y. Cao, Z. Cao, J. Wang, and M. I. Jordan, "Transferable representation learning with deep adaptation networks," *IEEE TPAMI*, 2018.
- [15] J. Song, H. Zhang, X. Li, L. Gao, M. Wang, and R. Hong, "Self-supervised video hashing with hierarchical binary auto-encoder," *IEEE TIP*, 2018.
- [16] F. Shen, Y. Xu, L. Liu, Y. Yang, Z. Huang, and H. T. Shen, "Unsupervised deep hashing with similarity-adaptive and discrete optimization," *IEEE TPAMI*, 2018.
- [17] J. Wang, T. Zhang, N. Sebe, and H. T. Shen, "A survey on learning to hash," *IEEE TPAMI*, 2018.
- [18] Y. Cao, M. Long, B. Liu, and J. Wang, "Deep cauchy hashing for hamming space retrieval," in *Proceedings of the CVPR*, 2018.
- [19] H. Liu, R. Ji, J. Wang, and C. Shen, "Ordinal constraint binary coding for approximate nearest neighbor search," *IEEE TPAMI*, 2018.
- [20] K. He, F. Cakir, S. Adel Bargal, and S. Sclaroff, "Hashing as tie-aware learning to rank," in *Proceedings of the CVPR*, 2018.
- [21] L. Jin, K. Li, Z. Li, F. Xiao, G.-J. Qi, and J. Tang, "Deep semantic-preserving ordinal hashing for cross-modal similarity search," *IEEE TNNLS*, 2018.
- [22] L. Wu, H. Ling, P. Li, J. Chen, Y. Fang, and F. Zou, "Deep supervised hashing based on stable distribution," *IEEE Access*, 2019.
- [23] C. Fatih, K. He, S. A. Bargal, and S. Sclaroff, "Hashing with mutual information," *IEEE TPAMI*, 2019.
- [24] L. Huang, Q. Yang, and W. Zheng, "Online hashing," in *Proceedings of the IJCAI*, 2013.
- [25] M. Lin, R. Ji, H. Liu, X. Sun, S. Chen, and Q. Tian, "Hadamard matrix guided online hashing," *IJCV*, 2020.
- [26] C. Leng, J. Wu, J. Cheng, X. Bai, and H. Lu, "Online sketching hashing," in *Proceedings of the CVPR*, 2015.
- [27] C. Fatih and S. Sclaroff, "Adaptive hashing for fast similarity search," in *Proceedings of the ICCV*, 2015.
- [28] L. Xie, J. Shen, and L. Zhu, "Online cross-modal hashing for web image retrieval," in *Proceedings of the AAAI*, 2016.
- [29] M. Lin, R. Ji, S. Chen, X. Sun, and C.-W. Lin, "Similarity-preserving linkage hashing for online image retrieval," *IEEE TIP*, 2020.
- [30] C. Fatih, S. A. Bargal, and S. Sclaroff, "Online supervised hashing," *CVIU*, 2017.
- [31] M. Qi, Y. Wang, and A. Li, "Online cross-modal scene retrieval by binary representation and semantic graph," in *Proceedings of the ACM MM*, 2017.
- [32] L. Xie, J. Shen, J. Han, L. Zhu, and L. Shao, "Dynamic multi-view hashing for online image retrieval," in *Proceedings of the IJCAI*, 2017.
- [33] C. Fatih, K. He, S. A. Bargal, and S. Sclaroff, "Mihash: Online hashing with mutual information," in *Proceedings of the ICCV*, 2017.
- [34] X. Chen, I. King, and M. R. Lyu, "Frosh: Faster online sketching hashing," in *Proceedings of the UAI*, 2017.
- [35] M. Lin, R. Ji, H. Liu, and Y. Wu, "Supervised online hashing via hadamard codebook learning," in *Proceedings of the ACM MM*, 2018.
- [36] M. Lin, R. Ji, H. Liu, X. Sun, Y. Wu, and Y. Wu, "Towards optimal discrete online hashing with balanced similarity," in *Proceedings of the AAAI*, 2019.
- [37] T. Yao, G. Wang, L. Yan, X. Kong, Q. Su, C. Zhang, and Q. Tian, "Online latent semantic hashing for cross-media retrieval," *PR*, 2019.
- [38] E. Liberty, "Simple and deterministic matrix sketching," in *Proceedings of the ACM SIGKDD*, 2013.
- [39] K. J. Horadam, *Hadamard matrices and their applications*. Princeton university press, 2012.
- [40] B. Babenko, M.-H. Yang, and S. Belongie, "A family of online boosting algorithms," in *Proceedings of the ICCV (Workshops)*, 2009.
- [41] W. Liu, J. Wang, R. Ji, Y. Jiang, and S.-F. Chang, "Supervised hashing with kernels," in *Proceedings of the CVPR*, 2012.
- [42] F. Shen, W. Liu, S. Zhang, Y. Yang, and H. Tao Shen, "Learning binary codes for maximum inner product search," in *Proceedings of the ICCV*, 2015.
- [43] Q. Jiang and W. Li, "Asymmetric deep supervised hashing," in *Proceedings of the AAAI*, 2018.
- [44] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *JMLR*, 2006.
- [45] J. Jiang and Z. Tu, "Efficient scale space auto-context for image segmentation and labeling," in *Proceedings of the CVPR*, 2009.
- [46] J. Kittler, R. Ghaderi, T. Windeatt, and J. Matas, "Face verification using error correcting output codes," in *Proceedings of the CVPR*, 2001.
- [47] R. E. Schapire, "Using output codes to boost multiclass learning problems," in *Proceedings of the ICML*, 1997.
- [48] M. Rastegari, A. Farhadi, and D. Forsyth, "Attribute discovery via predictable discriminative binary codes," in *Proceedings of the ECCV*, 2012.
- [49] H. Liu, R. Wang, S. Shan, and X. Chen, "Deep supervised hashing for fast image retrieval," in *Proceedings of the CVPR*, 2016.
- [50] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Technical report, University of Toronto*, 2009.
- [51] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *Proceedings of the NeurIPS*, 2014.
- [52] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, 1998.
- [53] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proceedings of the ICLR*, 2015.
- [54] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proceedings of the NeurIPS*, 2012.
- [55] C. Ma, I. W. Tsang, F. Peng, and C. Liu, "Partial hash update via hamming subspace learning," *IEEE TIP*, 2017.
- [56] Z. Weng and Y. Zhu, "Online hashing with efficient updating of binary codes," in *Proceedings of the AAAI*, 2020.



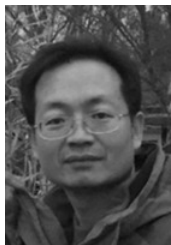
Mingbao Lin finished his M.S. study in Xiamen University, China, in 2018. He is currently pursuing the Ph.D degree with Xiamen University, China. His research interests include information retrieval and computer vision.



Rongrong Ji (SM'14) is currently a Professor and the Director of the Intelligent Multimedia Technology Laboratory, and the Dean Assistant with the School of Information Science and Engineering, Xiamen University, Xiamen, China. His work mainly focuses on innovative technologies for multimedia signal processing, computer vision, and pattern recognition, with over 100 papers published in international journals and conferences. He is a member of the ACM. He was a recipient of the ACM Multimedia Best Paper Award and the Best Thesis Award of Harbin Institute of Technology. He serves as an Associate/Guest Editor for international journals and magazines such as Neurocomputing, Signal Processing, Multimedia Tools and Applications, the IEEE Multimedia Magazine, and the Multimedia Systems. He also serves as program committee member for several Tier-1 international conferences.



Xiaoshuai Sun received the B.S. degree in computer science from Harbin Engineering University, Harbin, China, in 2007, and the M.S and Ph.D. degrees in computer science and technology from the Harbin Institute of Technology, Harbin, in 2009 and 2015, respectively. He was a Post-Doctoral Research Fellow with the University of Queensland from 2015 to 2016 and served as a Lecturer with the Harbin Institute of Technology from 2016 to 2018. He is currently an Associate Professor with Xiamen University, China. He was a recipient of the Microsoft Research Asia Fellowship in 2011.



and wavelets.

Baochang Zhang received the B.S., M.S., and Ph.D. degrees in computer science from the Harbin Institute of Technology, Harbin, China, in 1999, 2001, and 2006, respectively. From 2006 to 2008, he was a Research Fellow with The Chinese University of Hong Kong, Hong Kong, and also with Griffith University, Brisbane, Australia. He is currently a Full Professor with Beihang University, Beijing, China. His current research interests include pattern recognition, machine learning, face recognition,



Yonghong Tian (S'00-M'06-SM'10) is currently a Boya Distinguished Professor with the Department of Computer Science and Technology, Peking University, China, and is also the deputy director of Artificial Intelligence Research Center, PengCheng Laboratory, Shenzhen, China. His research interests include neuromorphic vision, brain-inspired computation and multimedia big data. He is the author or coauthor of over 200 technical articles in refereed journals such as IEEE TPAMI/TNNLS/TIP/TMM/TCSVT/TKDE/TPDS, ACM CSUR/TOIS/TOMM and conferences such as NeurIPS/CVPR/ICCV/AAAI/ACMMM/WWW. Prof. Tian was/is an Associate Editor of IEEE TCSVT (2018.1-), IEEE TMM (2014.8-2018.8), IEEE Multimedia Mag. (2018.1-), and IEEE Access (2017.1-). He co-initiated IEEE Int'l Conf. on Multimedia Big Data (BigMM) and served as the TPC Co-chair of BigMM 2015, and also served as the Technical Program Co-chair of IEEE ICME 2015, IEEE ISM 2015 and IEEE MIPR 2018/2019, and General Co-chair of IEEE MIPR 2020 and ICME2021. He is the steering member of IEEE ICME (2018-) and IEEE BigMM (2015-), and is a TPC Member of more than ten conferences such as CVPR, ICCV, ACM KDD, AAAI, ACM MM and ECCV. He was the recipient of the Chinese National Science Foundation for Distinguished Young Scholars in 2018, two National Science and Technology Awards and three ministerial-level awards in China, and obtained the 2015 EURASIP Best Paper Award for Journal on Image and Video Processing, and the best paper award of IEEE BigMM 2018. He is a senior member of IEEE, CIE and CCF, a member of ACM.



Dacheng Tao (F'15) is currently a Professor of Computer Science and an ARC Laureate Fellow in the School of Computer Science and the Faculty of Engineering at The University of Sydney. He mainly applies statistics and mathematics to artificial intelligence and data science. His research is detailed in one monograph and over 200 publications in prestigious journals and proceedings at prominent conferences such as IEEE TPAMI, TIP, TNNLS, IJCV, JMLR, NIPS, ICML, CVPR, ICCV, ECCV, AAAI, IJCAI, ICDM and ACM SIGKDD, with several best paper awards, such as the Best Theory/Algorithm Paper Runner Up Award at IEEE ICDM'07, the Distinguished Paper Award at 2018 IJCAI, the 2014 ICDM 10-year Highest-Impact Paper Award, and the 2017 IEEE Signal Processing Society Best Paper Award. He received the 2015 Australian Scopus-Eureka Prize and the 2018 IEEE ICDM Research Contributions Award. He is a fellow of the Australian Academy of Science, AAAS, ACM and IEEE.



Feiyue Huang received the Ph.D degree from Tsinghua University, China, in 2008. He is currently the vice general manager of Tencent YouTu Lab. His research interests include computer vision and pattern recognition.