





# Codebook-Free Compact Descriptor for Scalable Visual Search

Yuwei Wu , Feng Gao, Yicheng Huang, Jie Lin , *Member, IEEE*, Vijay Chandrasekhar, *Member, IEEE*, Junsong Yuan , *Senior Member, IEEE*, and Ling-Yu Duan , *Member, IEEE*

**Abstract**—The MPEG compact descriptors for visual search (CDVS) is a standard toward image matching and retrieval. To achieve high retrieval accuracy over a large scale image/video dataset, recent research efforts have demonstrated that employing extremely high-dimensional descriptors such as the Fisher vector (FV) and the vector of locally aggregated descriptors (VLAD) can yield good performance. Since the FV (or VLAD) possesses high discriminability but small visual vocabulary, it has been adopted by CDVS to construct a global compact descriptor. In this paper, we study the development of global compact descriptors in the completed CDVS standard and the emerging compact descriptors for video analysis (CDVA) standard, in which we formulate the FV (or VLAD) compression as a resource-constrained optimization problem. Accordingly, we propose a codebook-free aggregation method via dual selection to generate a global compact visual descriptor, which supports fast and accurate feature matching free of large visual codebooks, fulfilling the low memory requirement of mobile visual search at significantly reduced latency. Specifically, we investigate both sample-specific Gaussian component redundancy and bit dependency within a binary aggregated descriptor to produce compact binary codes. Our technique contributes to the scalable compressed Fisher vector (SCFV) adopted by the CDVS standard. Moreover, the SCFV descriptor is currently serving as the frame-level hand-crafted video feature, which inspires the inheritance of CDVS descriptors for the emerging CDVA standard. Furthermore, we investigate

the positive complementary effect of our standard compliant compact descriptor and deep learning based features extracted from convolutional neural networks with significant mean average precision gains. Extensive evaluation over benchmark databases shows the significant merits of the codebook-free binary codes for scalable visual search.

**Index Terms**—Visual Search, Compact Descriptor, CDVS, CDVA, Codebook free, Feature Descriptor Aggregation.

## I. INTRODUCTION

WITH the advent of the era of Big Data, huge body of data resources come from the multimedia sharing platforms such as *YouTube*, *Facebook* and *Flickr*. Visual search has attracted considerable attention in multimedia and computer vision [1]–[4]. It refers to the discovery of images/videos contained within a large dataset that describe the same objects/scenes as those depicted by query terms. There exists a wide variety of emerging applications, *e.g.*, location recognition and 3D reconstruction [5], searching logos for estimating brand exposure [6], and locating and tracking criminal suspects from massive surveillance videos [7]. This area involves a relatively new family of visual search methods, namely the Compact Descriptors for Visual Search (CDVS) techniques standardized from the ISO/IEC Moving Pictures Experts Group (MPEG) [8]–[12]. Generally speaking, the MPEG CDVS aims to define the format of compact visual descriptors as well as the feature extraction and visual search process pipeline to enable interoperable design of visual search applications. In this work, we focus on how to generate an extremely low complexity codebook-free<sup>1</sup> global descriptor for CDVS. Moreover, the proposed descriptor contributes to the emerging compact descriptors for video analysis (CDVA) standard [13].

The basic idea of visual search is to extract visual descriptors from images/videos, then perform descriptor matching between the query and dataset to find relevant items. Towards effective and efficient visual search, the visual descriptors need to be discriminative and resource-efficient (*e.g.*, low memory footprint). The discriminability relates to the search accuracy, while the computational resource cost impacts the scalability

<sup>1</sup>The typical compression scheme, Product Quantization (PQ), often requires tens of sub-codebooks while conventional Hashing algorithms involve thousands of hash functions, thereby incurring heavy memory cost. Extensive experiments have shown that our compact descriptor consumes extremely low memory of 0.015 MB with promising search performance. This merit of extremely low memory footprint is referred to as “codebook free”.

Manuscript received October 14, 2017; revised May 4, 2018; accepted June 15, 2018. Date of publication July 18, 2018; date of current version January 24, 2019. This work is supported in part by the National Natural Science Foundation of China under Grants 61661146005, U1611461, and 61702037, in part by the National Key Research and Development Program of China under Grant 2016YFB1001501, and in part by the Key Research and Development Program of Beijing Municipal Science & Technology Commission under Grant D171100003517002. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Zhu Li. (*Yuwei Wu and Feng Gao are co-first authors.*) (*Corresponding author: Ling-Yu Duan.*)

Y. Wu is with the Beijing Laboratory of Intelligent Information Technology, School of Computer Science, Beijing Institute of Technology, Beijing 100081, China, and also with PKU-NTU Joint Research Institute, Singapore 639798 (e-mail: wuyuwe@bit.edu.cn).

F. Gao is with the Future Lab, Tsinghua University, Beijing 100084, China, and also with PKU-NTU Joint Research Institute, Singapore 639798 (e-mail: gaofeng2018@tsinghua.edu.cn).

Y. Huang and L.-Y. Duan are with the Institute of Digital Media, School of Electronics Engineering and Computer Science, Peking University, Beijing 100080, China (e-mail: anorange0409@pku.edu.cn; lingyu@pku.edu.cn).

J. Lin and V. Chandrasekhar are with the Institute for Infocomm Research, Singapore 138634, and also with Nanyang Technological University, Singapore 639798 (e-mail: lin-j@i2r.a-star.edu.sg; vijay@i2r.a-star.edu.sg).

J. Yuan is with the Department of Computer Science and Engineering, State University of New York, Buffalo, NY 14260-2500 USA (e-mail: jsyuan@buffalo.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2018.2856628

of a visual search system. Taking mobile visual search as an example, we may extract and compress the visual features of query images at the mobile end and then send compact descriptors to the server end. This is expected to significantly reduce network latency and improve user experience of wearable or mobile devices with limited computational resources or bandwidth. Hereby, we would like to tackle the problem of high performance and low complexity aggregation of local feature descriptors from the perspective of the resource constrained optimization. This is well aligned with the goal of developing an effective and efficient compact feature descriptor representation technique with low memory cost in the “analyze then compress” infrastructure [14].

More recently, various vision tasks are carried out over large scale datasets, and visual search is no exception. To achieve high retrieval accuracy, it is useful to employ high-dimensional hand-crafted descriptors (such as the Fisher Vector (FV) [15], the Vector of Locally Aggregated Descriptors (VLAD) [16]), as well as deep learning based features [17], [18]. In particular, CNNs based deep learning features perform remarkably well on a wide range of visual recognition tasks. Moreover, the concatenation of hand-crafted shallow features and deep features proves to be more distinctive and reliable to characterize object appearances. Recently, Lou *et al.* [19] demonstrated that the combination of deep invariant descriptors and hand-crafted descriptors not only fulfills the lowest bitrate budget of CDVA but also significantly advances the video matching and retrieval performance of CDVA. Therefore, how to effectively and efficiently compress the hand-crafted features is practically useful, no matter from the perspective of CDVS and CDVA standardization, or the emerging applications like memory lightweight augmented reality on wearable or mobile devices. In retrieval and matching tasks, short binary codes can alleviate the issues of limited computational resources or bandwidth. In augmented reality scenarios, the compact descriptors of moderate-scale image database can be stored locally, in which image matching can be performed on local devices.

Hashing [3], [20]–[23] and product quantization (PQ) [24]–[29] are popular compression techniques to obtain binary descriptors. The hashing methods first project the data points into a subspace (or hyperplanes), and then quantize the projection values into binary codes. Learning methods have been involved in projection stage, in which the linear or non-linear projection are applied to convert high dimensional features (*e.g.*, feature  $\mathbf{f} \in \mathbb{R}^N$ ) into binary embedding and then learn a hashing function in the low-dimensional space. However, the computational complexity of the projection is  $\mathcal{O}(N^2)$  [22], [30]. For instance, given that a 426-dimensional dense trajectory feature is extracted from a video and PCA is employed to project each trajectory feature to the 213 dimension, we then obtain a 54,528-dimensional FV when the number of Gaussian components  $K = 128$ , as depicted in Figure 1. In the case of  $N = 54,528$ , a projection matrix alone may take more than 10 GB memory footprint and projecting one vector would cost  $\sim 800$  ms on a single core. On the other hand, the PQ generally divides a data vector into shorter subvectors and uses time-consuming K-means algorithm to train the sub-codebooks. Each subvector

is quantized by its own smaller sub-codebooks. However, as  $K$  (*e.g.*, the number of Gaussian components in FV aggregation) increases, binarizing high-dimensional descriptors using existing hashing or PQ approaches is intractable or infeasible due to the demanding computational cost and memory requirements. Hence, notwithstanding the effectiveness of preserving data similarity structure by the aforementioned hashing or PQ methods, important academic and industrial efforts like CDVS and CDVA have paid more attentions to generate compact binary codes of high-dimensional descriptors subject to the limited computational resources.

In this work, we focus on discriminative and compact aggregated descriptor for visual search at very low computational resources. To this end, we consider the following observations on the state-of-the-arts aggregation techniques, *i.e.*, FV and VLAD. (1) As introduced in [9], [15], when local features in an image are aggregated w.r.t Gaussian components to form a residual vector, not all the components are of equal importance to distinguish a sample. (2) In addition to the component level redundancy, there exists the bit-wise dependency within each selected component yet to be removed for better compactness.

Motivated by empirical observations mentioned above and the challenges of developing compact feature descriptors for visual search and analysis (*i.e.*, CDVS and CDVA), we formulate the FV (or VLAD) compression as a resource-constrained optimization problem, and propose a codebook-free binary coding method for scalable visual search. A sample-specific component selection scheme is introduced to remove redundant Gaussian components, and a global structure preserving sparse subspace learning method is presented to suppress the bit-wise dependency within each selected component. Figure 1 takes video search as an example and shows the flow diagram of generating compact global descriptors. Specifically, given a raw Fisher vector with  $K$  Gaussian components, we first binarize the FV by a sign function, leading to a binarized Fisher vector (BFV). Both sample-specific Gaussian component redundancy and bit dependency within a component are removed to produce compact binary codes.

The proposed approach to compressing aggregated descriptor has been adopted by MPEG CDVS standard [9]–[12]. Furthermore, our method serves as the frame level video hand-crafted feature representation towards large-scale video analysis, which leverages the merits of CDVS features for the emerging compact descriptors for video analysis (CDVA) standard. In summary, our contributions are three-fold.

- 1) We formulate a light-weight compression of aggregated descriptors as a resource constrained optimization problem, which investigates how to improve the descriptor performance in terms of both descriptor compactness and computational complexity. To this end, we come up with a codebook-free global compact descriptor via dual selection for visual search. We circumvent the time-consuming codebook training and avoid the heavy memory cost of storing large codebooks, which is particularly beneficial to mobile devices with limited memory usage.
- 2) We investigate both sample-specific Gaussian component redundancy and bit dependency within a binary aggrega-

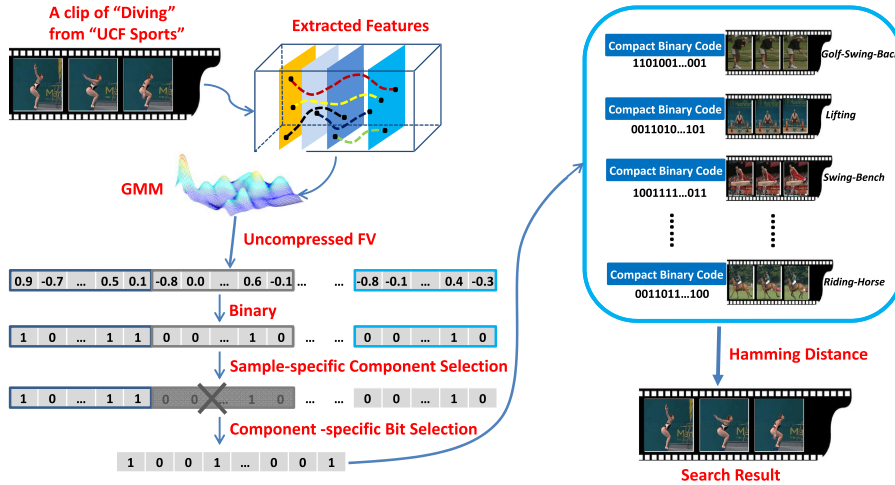


Fig. 1. The flow diagram of visual search by applying the proposed codebook-free compact descriptor. Taking the Fisher vector as an example, we first binarize the FV by a sign function, leading to a binarized Fisher Vector (BFV). Then discriminative bits are selected from the BFV, where both sample-specific Gaussian component redundancy and bit dependency within a binary aggregated descriptor are introduced to produce optimized compact Fisher codes.

gated descriptor to produce compact binary codes. Fast matching with extremely low memory footprint is allowed. The SCFV descriptor derived from our dual selection scheme, has been adopted by the completed MPEG CDVS standard as a normative aggregation technique to produce a compact global feature representation with much lower complexity. Extensive experiments over a wide range of benchmark datasets have shown promising search performance.

- 3) Furthermore, our work has provided a groundwork of compact hand-crafted features for the development of emerging MPEG CDVA standard, which is an important extension for large-scale video matching and retrieval. We have shown that the proposed global descriptors and the state-of-the-art CNN descriptors are positive complementary to each other, and their combination can significantly improve the performance for video retrieval. We argue that, especially towards real-time (mobile) visual search and augmented reality applications, how to harmoniously leverage the merits of highly efficient and low complexity handcrafted descriptors, and the cutting edge CNNs based descriptors, is a competitive and promising topic, which has been validated by CDVA.

The rest of the paper is organized as follows. Section II reviews the related work. In Section III, we introduce the problem statement of our work. Section IV presents the details of our compact binary aggregated descriptor via the dual selection scheme. We report and discuss the experimental results in Section VI, and conclude the paper in Section VII.

## II. RELATED WORK

Our work focuses on high performance and low complexity compact binary aggregated hand-crafted descriptors for CDVS and CDVA. Recent study has shown that there is great complementarity between CNNs and hand-crafted descriptors for better performance in visual search [19], which has been leveraged in

the emerging CDVA standard. In this section, we review typical techniques of compact representation. Readers are referred to [8], [13] for more comprehensive review.

### A. Compact Representation of Hand-Crafted Features

Due to the low storage overhead and fast retrieval speed, hashing-based and quantization-based techniques are popular compression techniques to obtain binary descriptors, and widely applied in approximate nearest neighbor (ANN) search. In general, hashing based methods often exploit the Hamming distance as the dissimilarity between the codes to get faster retrieval speed compared with quantization-based approaches, while quantization-based approaches can achieve superior performance over hashing codes with a little higher query time cost.

Hashing methods aim to map data points into Hamming space, which benefits fast Hamming distance computation and compact representation. In general, there are two categories of mainstream hashing approaches, *i.e.*, data-independent methods versus data-dependent methods. Data-independent methods do not rely on any training data, which is flexible, but often require long codes for satisfactory performance. Locality Sensitive Hashing [31] (LSH) adopts random projections to generate binary codes. It is independent of training data, and aims to make the similar items map to the same buckets as much as possible. Shift Invariant Kernel Hashing [32] (SIKH) chooses random projections and applies the shifted cosine function to generate binary codes.

Different from data-independent methods, data-dependent methods learn hashing functions from training data and significantly outperform data-independent methods. Weiss *et al.* [33] proposed Spectral Hashing (SH) by spectral partitioning for the graph constructed from the data similarity relationships. Gong and Lazebnik [20] proposed Iterative Quantization (ITQ) which figures out an orthogonal rotation matrix to refine the initial projection matrix. Wang *et al.* [34] proposed Minimal

Reconstruction Bias Hashing (MRH) to learn similarity preserving binary codes that jointly optimize both projection and quantization stages. Liu *et al.* [35] introduced a novel unsupervised hashing approach called min-cost ranking for large-scale data retrieval by learning and selecting salient bits. Many other representative data-dependent hashing methods include k-means hashing [36], and graph hashing [37].

The other alternative approaches of learning binary codes include clustering-based quantization which quantizes the space into cells via k-means. Product quantization [24] divides the data space into multiple disjoint subspaces, where a number of clusters are obtained by k-means over each subspace. Then a data point is approximated by concatenating the nearest cluster center of each subspace, yielding a short code comprising the indices of the nearest cluster centers. The distance between two vectors is computed by a pre-computed look-up table. Optimized Product Quantization [38] and its equivalent Cartesian k-means [29] improve Product Quantization by finding out an optimal feature space rotation and then performing product quantization over the rotated space. Composite Quantization [39] further improves the quantization method by approximating a data point using the summation of dictionary words selected from different dictionaries. It is demonstrated that quantization methods usually yield better search performance with comparable code length than hashing algorithms. However, the retrieval of quantization methods is often less efficient as Hashing methods apply the fast Hamming matching. Besides, quantization methods require a large codebook of k-means centroids and a look-up table for distance computation. In contrast, our codebook-free method does not involve memory cost for quantization which can be readily applied in those resource limited scenarios, such as mobile visual search.

### B. Compact Representation of Deep Features

While the aforementioned methods have certainly achieved great success, most of them work on hand-crafted features, which do not capture the semantic information and thus limit the retrieval accuracy of binary codes. The great success in deep neural network for representation learning has inspired deep feature coding algorithms. Lin *et al.* [40] developed a deep neural network to learn binary descriptors in an unsupervised manner with three criterions on binary codes, *i.e.*, minimal loss quantization, evenly distributed codes and uncorrelated bits. Liong *et al.* [41] defined a loss to penalize the difference between binary hash codes and real values, and introduced the bit balance and bit uncorrelation conditions.

Undoubtedly, deep learning based representations are becoming the dominant approach to generate compact descriptors for instance retrieval. The major concern is that deep neural network models may require a large storage of up to hundreds of megabytes, which makes it inconvenient to deploy in mobile applications or memory lightweight hardware. Gong *et al.* [42] proposed to leverage deep learning technique to compress FV features, while our method considers the statistics of the FV or VLAD feature to derive a high performance compact aggregated descriptor at extremely low resources. Different from the deep learning based hashing methods that use a weakly-supervised

fine-tuning scheme to update network parameters, our unsupervised dual selection scheme possesses superior computational efficiency. In addition, we have shown that the combination of the proposed hand-crafted compact descriptor and CNNs based descriptor can significantly improve the retrieval performance in the emerging CDVA standard [13].

### C. Compression of High Dimensional versus Low Dimensional Descriptors

In the aforementioned methods, an intermediate embedding is generated in certain ways, then the projected values are quantized into binary codes. However, embedding matrices incur considerable memory and extra computational cost. In addition, most of them work on relatively low-dimensional descriptors (*e.g.*, SIFT and GIST). Few work is dedicated to high-dimensional descriptors such as FV and VLAD.

Perronnin *et al.* [43] proposed a ternary encoding method to compress FV. Ternary encoding is fast and memory free, however, it will result in long codesize. Sivic *et al.* [44] binarized the BoW histogram entries over large vocabularies. Chen *et al.* [45] presented a discriminative projection to reduce the dimension of uncompressed descriptors before binarization, leading to comparable accuracy with the original descriptor at smaller codes, but causing additional memory cost because of the projection matrices. Gong *et al.* [30] employed compact bilinear projections instead of a single large projection matrix to get similarity-preserving binary codes. Parkhi *et al.* [46] considered face tracks as the unit of face recognition in videos and developed a compact and discriminative Fisher vector. For example, in the scenario of cross-modal retrieval, high-dimensional modality-specific features contain more abundant information that helps to bridge the modality gap. Therefore, it is indispensable to compress the high-dimensional semantic relevance to achieve satisfactory performance.

The above-mentioned methods have shown that a high-dimensional binary code is often necessary to preserve the discriminative power of aggregated descriptors. Hence, this paper concerns a compact binary aggregated descriptor coding approach for fast ANN search. The proposed method significantly reduces the codesize of raw aggregated descriptors, without degrading the search accuracy or introducing additional memory footprint.

## III. PROBLEM STATEMENT

To compress FV (or VLAD) into a compact descriptor for the considerable reduction on storage space and the low complexity of Hamming distance computation, we formulate the FV (or VLAD) compression as a resource-constrained optimization problem. Let  $A(\cdot)$  denotes search accuracy,  $R(\cdot)$  denotes descriptor size and  $C(\cdot)$  denotes complexity of compression. Our goal is to design a quantizer  $q(\cdot)$  to quantize Fisher vector  $g$  that is able to maximize search performance  $A(q(g))$  subject to the constraints of descriptor compactness  $R_{budget}$ , compression complexity  $C_{budget}$  in terms of memory and time:

$$\max_q A(q(g)) \quad s.t. \quad R(q) \leq R_{budget} \quad and \quad C(q) \leq C_{budget}. \quad (1)$$

Here, a problem arises: how to solve the abstract problem defined in Eq. (1) in a concrete implementation. To this end, we develop a codebook-free compact binary coding method via dual selection to compress the raw FV (or VLAD). Suppose that a raw FV is represented as  $\mathbf{g} = [\mathbf{g}(1), \dots, \mathbf{g}(K)]$  with  $K$  Gaussian components, where  $\mathbf{g}(i)$  ( $1 \leq i \leq K$ ) denotes the  $i$ -th sub-vector<sup>2</sup>. Firstly, to fulfill the constraint of compression complexity, we binarize the FV by a sign function and get a binary aggregated descriptor  $\mathbf{b} = \{\mathbf{b}(1), \dots, \mathbf{b}(K)\}$ , where each binary sub-vector code  $\mathbf{b}(i) = \text{sgn}(\mathbf{g}(i))$ . Secondly, we propose to select discriminative bits from the binary aggregated descriptor  $\mathbf{b}$  to maximize search performance, subject to the constraint of descriptor compactness. A dual selection scheme, *i.e.*, the sample-specific component selection and the component-specific bit selection, is introduced to obtain the discriminative bits towards a low complexity and high performance FV binary representation. In our work, the sample-specific component selection is derived by the variance of each sub-vector  $\mathbf{g}(i)$  of the raw FV, and the component-specific bit selection is applied to further reduce the dimensionality of components, while maintaining search performance as well. In the next sections, we will present the key ideas mentioned above in detail.

**Remark:** This work is tightly related to the MPEG CDVS standard [9]–[12]. CDVS adopts the scalable compressed Fisher Vector (SCFV) representation, in which the quantization procedure was briefly reported in [9]. To compress the high dimensional FV, a subset of Gaussian components in the Gaussian Mixture Model (GMM) are selected by ranking the standard deviation of each sub-vector. The number of selected Gaussian components relies on the bit budget, and then one-bit scalar quantizer is applied to generate binary codes. SCFV derived from the proposed dual selection scheme concentrates on the scalability of compact descriptors. Beyond [9] and [12], this work will systematically investigate how to leverage both sample-specific Gaussian component redundancy and bit dependency within a binary aggregated descriptor to produce the codebook-free compact descriptor. A component-specific bit selection scheme is introduced by global structure preserving sparse subspace learning. In addition, we study how to further improve the performance by combining the proposed compact descriptors with CNNs descriptors, which has been adopted by the emerging CDVA standard. Our technique benefits the aggregated descriptor in improving compactness, and removing bits redundancy to ameliorate discriminative power of the descriptor.

#### IV. COMPACT BINARY AGGREGATED DESCRIPTOR VIA DUAL SELECTION

Our goal is to compress aggregated descriptors, without incurring considerable loss of search accuracy. The coarse binary quantization would degrade search performance. We propose a dual selection model which utilizes both sample-specific component selection and component-specific bit selection to collect

the informative bits. In the stage of sample-specific component selection, a subset of Gaussian components is adaptively determined, which can significantly boost search performance. Furthermore, the component-specific bit selection focuses on the compactness to reduce the dimensionality of components, while maintaining search performance.

##### A. Sample-Specific Component Selection

For FV aggregation, a Gaussian Mixture model (GMM)  $U_\lambda(x) = \sum_{i=1}^K \omega_i p_i(x)$  with  $K$  Gaussians (visual words) is to estimate the distribution of local features over a training set. We denote the set of Gaussian parameters as  $\lambda_i = \{\omega_i, \mu_i, \sigma_i^2, i = 1, \dots, K\}$ , where  $\omega_i$ ,  $\mu_i$  and  $\sigma_i^2$  are the weight, mean vector and variance vector of the  $i$ -th Gaussian  $p_i$ , respectively. Let  $\mathbf{F} = \{\mathbf{f}_1, \dots, \mathbf{f}_T\}$  denote a collection of  $T$  SIFT local features extracted from an image. Projecting the dimensionality of each local SIFT feature to dimension  $D$  is beneficial to the overall performance [9], [16]. The gradient vectors of all local features w.r.t. the mean  $\mu_i$  are aggregated into a  $D$ -dimensional vector

$$\mathbf{g}(i) = \frac{1}{T\sqrt{\omega_i}} \sum_{t=1}^T \gamma(\mathbf{f}_t, i) \sigma_i^{-1} (\mathbf{f}_t - \mu_i), \quad (2)$$

where  $\gamma(\mathbf{f}_t, i) = \omega_i p_i(\mathbf{f}_t) / \sum_{j=1}^k \omega_j p_j(\mathbf{f}_t)$  denotes the posterior probability of local feature  $\mathbf{f}_t$  being assigned to the  $i$ -th Gaussian. By concatenating the sub-vector  $\mathbf{g}(i)$  of all the  $K$  components, we form the FV  $\mathbf{g} = [\mathbf{g}(1), \dots, \mathbf{g}(K)] \in \mathbb{R}^{KD}$ , where  $\mathbf{g}(i) \in \mathbb{R}^D$ . Note that the gradient vectors can be extended to the deviates w.r.t. variance as well, which was adopted in CDVS standard to improve the performance at high bitrates [8]. Similarly, the VLAD can be derived from FV by replacing the GMM soft clustering with k-means clustering, *i.e.*,

$$\mathbf{g}(i) = \sum_{t: NN(\mathbf{f}_t) = \mu_i} (\mathbf{f}_t - \mu_i), \quad (3)$$

where  $NN(\mathbf{f}_t)$  indicates  $\mathbf{f}_t$ 's nearest neighbor centroids.

Furthermore, we apply a one-bit quantizer to binarize the high dimensional  $\mathbf{g} \in \mathbb{R}^{KD}$ . Specifically, we generate binary aggregated descriptors by quantizing each dimension of FV or VLAD into a single bit 0/1 by a sign function. Formally speaking,  $\text{sgn}(\mathbf{g})$  is used to map each element  $g_j$  of the descriptor  $\mathbf{g}$  to 1 if  $g_j > 0$ ,  $j = 1, 2, \dots, KD$ ; otherwise, 0, yielding a binary aggregated descriptor  $\mathbf{b} = \{\mathbf{b}(1), \dots, \mathbf{b}(K)\}$  with  $KD$  bits, where  $\mathbf{b}(i) \in \mathbb{R}^D$ .

FV exhibits a natural 2-D structure, as shown in Figure 2. Referring to Eq. (2), the aggregated FV is formed by concatenating residual vectors of all Gaussian components, while each residual vector is aggregated from local features being assigned to the corresponding Gaussian component. In other words, not all Gaussian components equally contribute to the discriminative power. The role of a Gaussian component relates to the number of local features quantized to that component. The occurrence of different Gaussian components may vary for different image samples. Those Gaussian components with low occurrence are supposed to be less discriminative. If none of local features is assigned to a Gaussian component  $i$ , then all the elements of the

<sup>2</sup>Since the VLAD is a simplified non-probabilistic version of FV, in what follows, we take the FV as an example to elaborate how to obtain the compact binary codes for convenient discussion.

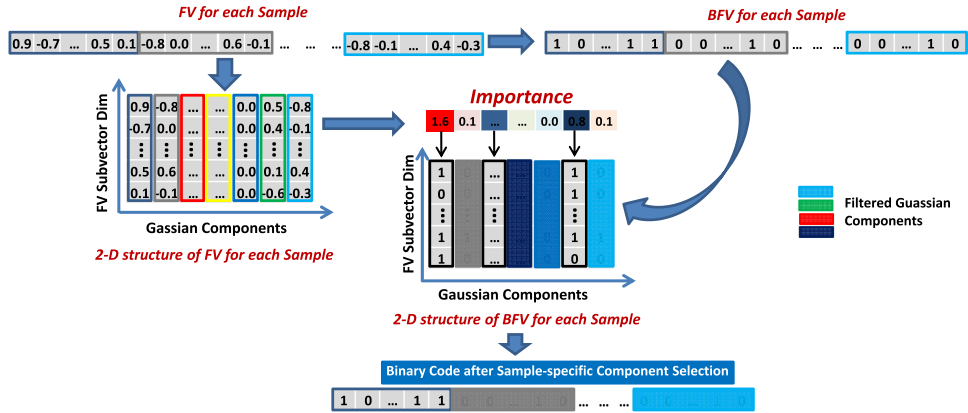


Fig. 2. Sample-specific component selection. The aggregated FV is partitioned into disjoint components by Gaussian components. We select all the bits of Gaussian components with high *importance* and the rest of components are discarded. In this work, the *importance* is defined as the sum of posterior probabilities of local features being assigned to the  $i$ -th Gaussian component.

corresponding subvector  $\mathbf{g}(i)$  in Eq. (2) are zero. Accordingly, the proposed sample-specific component selection leverages local statistics of each individual sample to discard the redundant Gaussian components<sup>3</sup>.

Specifically, we process the FV at the granularity of Gaussian components and select all the bits of those components with high importance. Here, the importance measured by the amplitude of their responses determines which Gaussian components are activated. In particular, we adopt the soft assignment coefficients  $\gamma(\mathbf{f}_t, i)$  of local features to indicate the importance of Gaussian  $i$ , which is employed to adaptively select parts of discriminative components for each sample. The importance, *i.e.*,  $I(\lambda_i)$ , is defined as the sum of posterior probabilities of local features  $\{\mathbf{f}_1, \dots, \mathbf{f}_T\}$  being assigned to the  $i$ -th Gaussian component in the continuous FV model, which is given by

$$I(\lambda_i) = \sum_{t=1}^T \gamma(\mathbf{f}_t, i). \quad (4)$$

Since each local feature can be soft quantized to multiple Gaussian components, we rank all the Gaussian components based on the accumulated soft assignment statistics of local features to Gaussian components. The Gaussian component ranked at the  $t$ -th position is called the  $t$ -th nearest neighbor Gaussian component. That is, sample-specific component selection can be implemented by sorting the set  $\{I(\lambda_i), i = 1, 2, \dots, K\}$ , and then the subvector of the binary Fisher vector (BFV)  $\mathbf{b}(i)$  with the largest  $I(\lambda_i)$  is first selected to generate Fisher codes, followed by the  $\mathbf{b}(j)$  with the second largest  $I(\lambda_j)$ . In this way, we compute the importance of individual components for each sample using Eq. (4) and select a subset of components with high importance. The rest of components are discarded, as shown in Figure 2. It is necessary to maintain a Gaussian selection mask for each sample. Accordingly, the sample-specific component selection is nearly memory free, and the computational cost of Gaussian selection is  $\mathcal{O}(KD)$ .

<sup>3</sup>Here, the redundant Gaussian components are those Gaussian components with low occurrence.

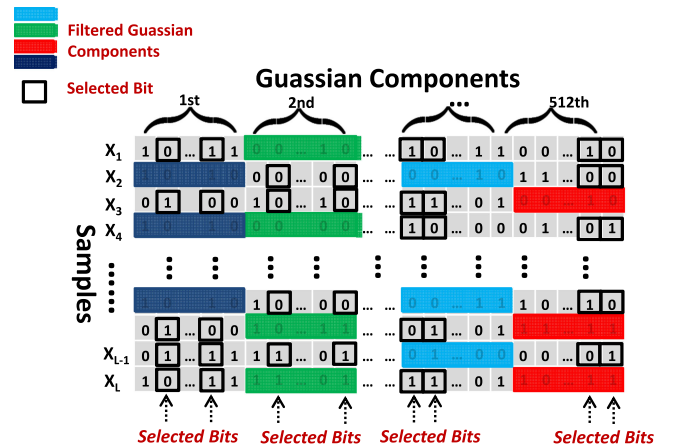


Fig. 3. Component-specific bit selection. The bits that carry as much information as possible are selected based on the global structure preserving sparse subspace learning. For each sample, the bits in black bounding boxes are the final compact Fisher codes.  $L$  denotes the number of training samples.

### B. Component-Specific Bit Selection

After sample-specific component selection, there probably exists bit-level redundancy within each selected component. Therefore, we may further improve the compactness of each selected component, while maintaining search performance. A naive solution is to apply random bit selection, however, it ignores the bit correlation within a component. In this work, we introduce a component-specific bit selection scheme by global structure preserving sparse subspace learning to select the most informative bits.

Without loss of generality, let  $\mathcal{B} \in \{0, 1\}^{\Theta \times D}$  denote a subset of the binary aggregated descriptor (after component selection) associated with the  $i$ -th Gaussian component over the whole dataset. Here  $\Theta$  denotes how many samples select the  $i$ -th component as an *important* one, as shown in Figure 3.  $\Theta \leq L$ , where  $L$  is the number of training samples. The goal of bits selection is to find a small set of bits that can capture most useful information of  $\mathcal{B}$ . One natural way is to measure how close the original data samples are over the learned subspace  $\mathbf{H}$  spanned

by the selected bits. Mathematically, the component-specific bit selection is formulated as

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \|\mathcal{B} - \mathcal{B}\mathbf{W}\mathbf{H}\|_F^2. \\ \text{s.t.} \quad & \mathbf{W} \in \{0, 1\}^{D \times D'} \\ & \mathbf{W}^\top \mathbf{1}_{D \times 1} = \mathbf{1}_{D' \times 1} \\ & \|\mathbf{W}\mathbf{1}_{D' \times 1}\|_0 = D' \end{aligned} \quad (5)$$

Here,  $\mathbf{W}$  is a selection matrix with entries of 0 or 1.  $\mathbf{W}^\top \mathbf{1}_{D \times 1} = \mathbf{1}_{D' \times 1}$  enforces that each column of  $\mathbf{W}$  has a single 1, *i.e.*, at most  $D'$  bits are selected for each Gaussian component.  $\|\mathbf{W}\mathbf{1}_{D' \times 1}\|_0 = D'$  guarantees that  $\mathbf{W}$  has the  $D'$  nonzero rows, and thus  $D'$  bits will be selected. Hence, Eq.(5) denotes the distance of  $\mathcal{B}$  to the learned subspace  $\mathbf{H}$ .

A major difficulty in solving Eq. (5) lies in handling the discrete constraints imposed on  $\mathbf{W}$ . In this work, we relax the 0-1 constraint of  $\mathbf{W}$  by introducing nonnegativity constraint and reduce the hard constraints of both  $\mathbf{W}^\top \mathbf{1}_{D \times 1} = \mathbf{1}_{D' \times 1}$  and  $\|\mathbf{W}\mathbf{1}_{D' \times 1}\|_0 = D'$  to a  $L_{2,1}$  norm constraint. Therefore, optimizing Eq. (5) is equivalent to solving

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \|\mathcal{B} - \mathcal{B}\mathbf{W}\mathbf{H}\|_F^2 + \beta \|\mathbf{W}\|_{2,1}, \\ \text{s.t.} \quad & \mathbf{W} \in \mathbb{R}_+^{D \times D'} \end{aligned} \quad (6)$$

where  $\mathbb{R}_+^{D \times D'}$  denotes a set of  $D \times D'$  nonnegative matrices and  $\beta$  is a parameter. Based on the resulting solution  $\mathbf{W}$ , we choose the bits corresponding to the  $D'$  rows of  $\mathbf{W}$  with the largest norms.

To solve this optimization, we employ the accelerated block coordinate update (ABCU) method [47] to alternately update  $\mathbf{W}$  and  $\mathbf{H}$  with

$$\begin{cases} f(\mathbf{W}, \mathbf{H}) = \|\mathcal{B} - \mathcal{B}\mathbf{W}\mathbf{H}\|_F^2 \\ g(\mathbf{W}) = \beta \|\mathbf{W}\|_{2,1}. \end{cases} \quad (7)$$

At the  $\tau$ -th iteration, we need to solve the following optimization problems

$$\begin{cases} \mathbf{W}^{\tau+1} = \arg \min_{\mathbf{W} \geq 0} \left\langle \nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau), \mathbf{W} - \widehat{\mathbf{W}}^\tau \right\rangle \\ \quad + \frac{\mathbf{L}_{\mathbf{W}}^\tau}{2} \left\| \mathbf{W} - \widehat{\mathbf{W}}^\tau \right\|_F^2 + \beta \|\mathbf{W}\|_{2,1} \\ \mathbf{H}^{\tau+1} = \arg \min_{\mathbf{H}} f(\mathbf{W}^{\tau+1}, \mathbf{H}), \end{cases} \quad (8)$$

where  $\mathbf{L}_{\mathbf{W}}^\tau = \|\mathbf{H}^\tau (\mathbf{H}^\tau)^\top\|_s \|\mathcal{B}^\top \mathcal{B}\|_s$  is the Lipschitz constant of  $\nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau)$  with respect to  $\mathbf{W}$ . Note that  $\|\Psi\|_s$  denotes the spectral norm and equals the largest singular value of  $\Psi$ .  $\widehat{\mathbf{W}}^\tau$  is given by

$$\widehat{\mathbf{W}}^\tau = \mathbf{W}^\tau + \omega_\tau (\mathbf{W}^\tau - \mathbf{W}^{\tau-1}), \quad (9)$$

---

**Algorithm 1:** Proximal Operator for Solving  $\mathbf{W} = \text{Prox}(\mathbf{Y}, \eta)$

---

**Input:**  $\mathbf{Y} = \widehat{\mathbf{W}}^\tau - \frac{1}{\mathbf{L}_{\mathbf{W}}^\tau} \nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau)$  and  $\eta = 0.1$

- 1 **for**  $i = 1 \rightarrow D$  **do**
- 2     Let  $\mathbf{y}$  is the  $i$ -th row of  $\mathbf{Y}$  and  $\mathcal{I}$  the index set of positive components of  $\mathbf{y}$  ;
- 3     Set  $\mathbf{w} = \mathbf{0}$ ;
- 4     **if**  $\|\mathbf{y}_{\mathcal{I}}\|_2 > \eta$  **then**
- 5          $\mathbf{w}_{\mathcal{I}} = (\|\mathbf{y}_{\mathcal{I}}\|_2 - \eta) \frac{\mathbf{y}_{\mathcal{I}}}{\|\mathbf{y}_{\mathcal{I}}\|_2}$ ;
- 6         Set the  $i$ -th row of  $\mathbf{W}$  to  $\mathbf{w}$ .
- 7     **end**
- 8 **end**

---

where

$$\begin{cases} \omega_\tau = \min \left( \widehat{\omega}_\tau, \theta_\omega \sqrt{\frac{\mathbf{L}_{\mathbf{W}}^{\tau-1}}{\mathbf{L}_{\mathbf{W}}^\tau}} \right) \\ \widehat{\omega}_\tau = \frac{\delta_{\tau-1} - 1}{\delta_\tau} \\ \delta_\tau = \frac{1 + \sqrt{1 + 4\delta_{\tau-1}^2}}{2}. \end{cases} \quad (10)$$

$0 < \theta_\omega < 1$  and  $\delta_{\tau-1}$  are predetermined parameters. In our experiment, we set  $\delta_0 = 1$ ,  $\theta_\omega = 0.5$ .  $\widehat{\omega}_\tau$  has been widely applied to the accelerate proximal gradient method for the convex optimization problem [47].

In Eq. (8),  $\mathbf{W}^{\tau+1}$  is obtained by equivalently solving

$$\begin{aligned} \min_{\mathbf{W} \geq 0} \quad & \left\| \mathbf{W} - \left( \widehat{\mathbf{W}}^\tau - \frac{1}{\mathbf{L}_{\mathbf{W}}^\tau} \nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau) \right) \right\|_F^2 \\ & + \frac{\beta}{\mathbf{L}_{\mathbf{W}}^\tau} \|\mathbf{W}\|_{2,1}. \end{aligned} \quad (11)$$

Equation (11) can be decomposed into  $D$  smaller independent problems, each of which involves one row of  $\mathbf{W}$  and is then solved by a Proximal operator

$$\prod(\mathbf{y}) = \min_{\mathbf{w} \geq 0} \|\mathbf{w} - \mathbf{y}\|_2^2 + \eta \|\mathbf{w}\|_2. \quad (12)$$

where  $\mathbf{y}$  is the  $i$ -th row of  $\widehat{\mathbf{W}}^\tau - \frac{1}{\mathbf{L}_{\mathbf{W}}^\tau} \nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau)$ .  $\prod(\mathbf{y})$  is able to be implemented efficiently according to the Theorem 1 [48], [49], which is described in Algorithm 1.

*Theorem 1:* Given  $\mathbf{y}$  and  $\mathcal{I}$  the index set of positive components of  $\mathbf{y}$ , *i.e.*,  $\mathcal{I} = \{i : \mathbf{y}_i > 0\}$ , the solution of Eq.(12) is expressed as

$$\begin{cases} w_i = 0, & \forall i \notin \mathcal{I}; \\ \mathbf{w}_{\mathcal{I}} = \mathbf{0}, & \|\mathbf{y}_{\mathcal{I}}\|_2 \leq \eta; \\ \mathbf{w}_{\mathcal{I}} = (\|\mathbf{y}_{\mathcal{I}}\|_2 - \eta) \frac{\mathbf{y}_{\mathcal{I}}}{\|\mathbf{y}_{\mathcal{I}}\|_2}, & \text{otherwise.} \end{cases} \quad (13)$$

Readers are referred to [48], [49] for the detailed proof. In addition,  $\mathbf{H}^{\tau+1}$  can be obtained in a closed form, *i.e.*,

$$\mathbf{H}^{\tau+1} = [(\mathbf{W}^{\tau+1})^\top \mathcal{B}^\top \mathcal{B} \mathbf{W}^{\tau+1}]^{-1} (\mathbf{W}^{\tau+1})^\top \mathcal{B}^\top \mathcal{B}. \quad (14)$$

The accelerated block coordinate update method for solving Eq. (6) is described in Algorithm 2. The main cost lies in the update of  $\mathbf{W}$  and  $\mathbf{H}$ . For updating  $\mathbf{W}$ , we should compute the first order partial derivative of  $f(\mathbf{W}, \mathbf{H})$  with respect to  $\mathbf{W}$ , *i.e.*,  $\nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}) = \mathcal{B}^\top (\mathcal{B}\mathbf{W}\mathbf{H} - \mathcal{B})\mathbf{H}^\top$ . Computing

---

**Algorithm 2:** Accelerated Block Coordinate Update for Solving Eq.(6)

---

**Input:**  $\mathcal{B} \in \mathbb{R}^{L \times D}$ , the number of selected bits  $D'$ .  
**Output:** Index set of selected bits  $\mathcal{I} \subseteq \{1, 2, \dots, D\}$  with  $|\mathcal{I}| = D'$

- 1 **Initialize**  $\tau = 1, \delta_0 = 1, \theta_\omega = 0.5$
- 2 **for**  $\tau = 1, 2, \dots$  *until convergence do*
- 3     Compute  $\delta_\tau = \frac{1 + \sqrt{1 + 4\delta_{\tau-1}^2}}{2}$ ;
- 4     Compute  $\hat{\omega}_\tau = \frac{\delta_{\tau-1} - 1}{\delta_\tau}$ ;
- 5     Compute  $\omega_\tau = \min\left(\hat{\omega}_\tau, \theta_\omega \sqrt{\frac{L_{\mathbf{W}}^{\tau-1}}{L_{\mathbf{W}}^\tau}}\right)$ ;
- 6     Let  $\widehat{\mathbf{W}}^\tau = \mathbf{W}^\tau + \omega_\tau(\mathbf{W}^\tau - \mathbf{W}^{\tau-1})$ ;
- 7     Update  $\mathbf{W}^{\tau+1} \leftarrow \text{Prox}\left(\widehat{\mathbf{W}}^\tau - \frac{1}{L_{\mathbf{W}}} \nabla_{\mathbf{W}} f(\widehat{\mathbf{W}}^\tau, \mathbf{H}^\tau), \frac{\beta}{L_{\mathbf{W}}}\right)$   
     according to Algorithm 1;
- 8     Update  $\mathbf{H}^{\tau+1}$  according to Eq. (14);
- 9     **if**  $f(\mathbf{W}^{\tau+1}, \mathbf{H}^{\tau+1}) \geq f(\mathbf{W}^\tau, \mathbf{H}^\tau)$  **then**
- 10          $\widehat{\mathbf{W}}^\tau = \mathbf{W}^\tau$ ;
- 11         **else**  $\tau = \tau + 1$ ;
- 12     **end**
- 13 **end**
- 14 Normalize each column of  $\mathbf{W}$ ;
- 15 Sort  $\|\mathbf{W}_{i,\cdot}\|_2, i = 1, 2, \dots, D$  and select bits corresponding to the  $D'$  largest ones.

---

$\mathcal{B}\mathbf{W}, \mathbf{H}\mathbf{H}^\top, \mathcal{B}\mathbf{H}^\top, \mathcal{B}\mathbf{W}(\mathbf{H}\mathbf{H}^\top)$  and left multiplying  $\mathcal{B}^\top$  to  $\mathcal{B}\mathbf{W}(\mathbf{H}\mathbf{H}^\top) - \mathcal{B}\mathbf{H}^\top$  will take about  $3\Theta DD' + DD'^2 + \Theta D'^2$  flops. Similarly, updating  $\mathbf{H}$  will take  $2\Theta DD' + DD'^2 + \Theta D'^2 + D^3$  flops. Therefore, the computational complexity of Algorithm 2 is  $\mathcal{O}(\Theta DD')$ . Since  $D' < \min(\Theta, D)$ , the ABCU algorithm is scalable.

In summary, since the Gaussian components are independent of each other, the component-specific bit selection can be implemented in a parallel fashion. Figure 4 visualizes the bit correlations. Figure 4(a) shows the exemplar statistics of bit correlations between any two bits within a component, where 32 bits are randomly selected from a 64-dimensional binarized FV descriptor. By comparing Figure 4(a) and Figure 4(b), one can see that bits selected by the global structure preserving sparse subspace learning are much less correlated than random bit selection. The number of selected bits  $D'$  is fixed and applied to all samples. That is, the component-specific bit selection can be carried out offline such that the complexity does not affect online retrieval.

## V. HAMMING DISTANCE MATCHING

In online search, we apply the dual selection scheme to query samples as well and perform search using the selected bits. As the collection of selected components may vary in different samples, we have to solve the issue of matching descriptors across different Gaussian components. That is, given a query  $\mathbf{x}_q$  and a dataset sample  $\mathbf{x}_r$ , if the selected Gaussian components are different, the similarity cannot be computed directly using standard Hamming distance. So we adopt a normalized cosine similarity score  $S$  in [11] to calculate the distance, given by

$$S = \frac{\sum_{i=1}^K s_i^q s_i^r (D' - 2 * h(\text{sgn}(\mathbf{g}_i^{\mathbf{x}_q}), \text{sgn}(\mathbf{g}_i^{\mathbf{x}_r})))}{D' \sqrt{\|s^q\|_0 \|s^r\|_0}}, \quad (15)$$

where  $s_i^q$  and  $s_i^r$  denote whether the  $i$ -th Gaussian component is chosen for  $\mathbf{x}_q$  and  $\mathbf{x}_r$ , respectively, and  $h(\cdot, \cdot)$  is the Hamming

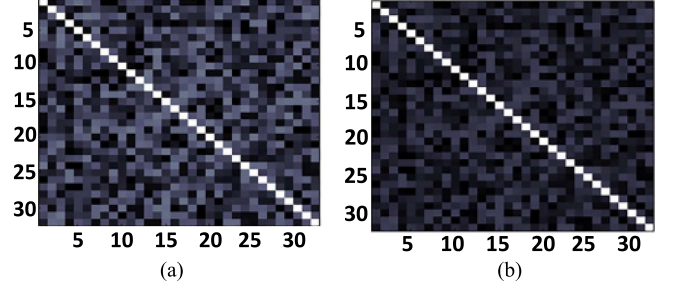


Fig. 4. Visualization of bit correlations of (a) randomly selected 32 bits and (b) optimized selected 32 bits from a selected component with of binarized FV descriptors over the INRIA Holidays dataset. Bright colors indicate strong bit correlation.

distance between binarized Fisher subvectors. In practice,  $S_c$  is computed based on the overlapping Gaussians  $s^q \cap s^r$  between  $\mathbf{x}_q$  and  $\mathbf{x}_r$ .

## VI. EXPERIMENTS

In this section, extensive experiments are conducted to evaluate the proposed method in both computational efficiency and search performance. Our approach is implemented in C++. The experiments are performed on an Dell Precision workstation 7400-E5440 with 2.83 GHz Intel XEON processor and 32 GB RAM in a mode of single core and single thread.

### A. Image Retrieval

To compare with baselines, we carry out retrieval experiments on MPEG CDVS datasets [50] and the publicly available INRIA Holidays dataset [51]. The CDVS dataset consists of five data sets used in the MPEG-CDVS standard: Graphics, Paintings, Frames, Landmarks and Common Objects. (1) The *Graphics* dataset depicts CD/DVD/book cover, text document and business card. There are 1,500 queries and 1,000 dataset images. (2) The *Painting* dataset contains 400 queries and 100 dataset images of paintings (say history, portraits, etc.) (3) The *Frame* dataset contains 400 queries and 100 dataset images of video frames captured from a range of video contents like movies and news. (4) The *Landmark* dataset contains 3,499 queries and 9,599 dataset images from building benchmarks, including the ZuBuD dataset, the Turin buildings, the PKUbench, etc. (5) The *Common Object* dataset contains 2,550 objects, each containing four images taken from different viewpoints. For each object, the first image is query and the rest are dataset images. (6) The *Holidays* dataset is a collection of 1,491 holiday photos, there are 500 image groups where the first image of each group is used as a query. To fairly evaluate the performance over a large-scale dataset, we use *FLICKRIM* as the distractor dataset, containing 1 million distractor images collected from Flickr. The retrieval performance is measured by mean Average Precision (mAP).

We evaluate the proposed method against several representative baselines including LSH [31], BP [30], and ITQ [20]. For the LSH, BP and ITQ methods, we randomly chose 20K images from Flickr1M dataset to train the projection matrices. The SIFT features from each image are extracted and then the



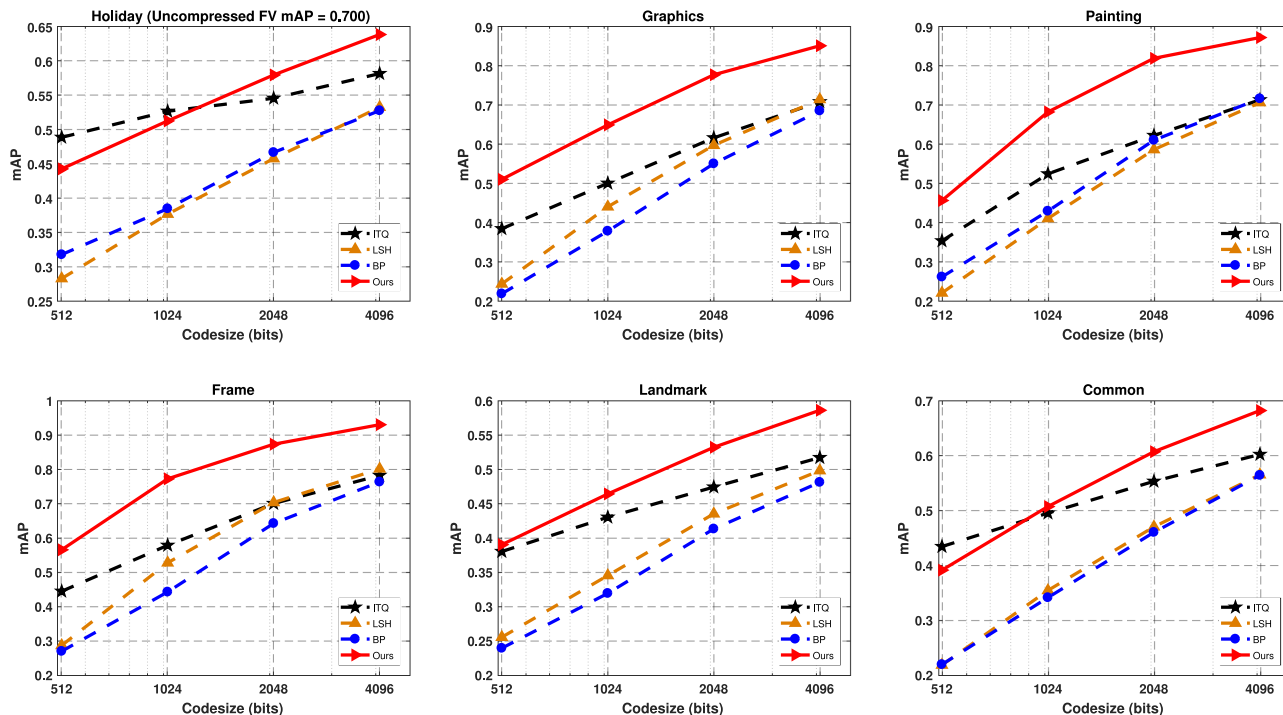


Fig. 5. Retrieval performance in terms of mAP vs. descriptor codesize over various benchmark datasets. (best viewed in color).

dimensionality of SIFT is reduced to 32 by PCA. We employ FV encoding to aggregate the dim-reduced SIFT features with 512 Gaussian components. Therefore, the total dimensionality of FV is  $512 \times 32 = 16384$ . Then, we apply power law [52], [53] ( $\alpha = 0.5$ ) followed by  $L2$  normalization to the raw FV feature.

Fig. 5 presents comparison results between the proposed method and hashing algorithms, in terms of mAP vs. different codesize over different datasets. As it is computationally expensive for  $L2$  distance between uncompressed FV, we present the retrieval accuracy of uncompressed FV on the INRIA Holidays dataset only. The proposed method achieves superior accuracy than other methods, especially for small codesizes. Note that ITQ yields better mAP than other baseline methods. This is reasonable since ITQ is fine tuned for projection learning to minimize the mean square error. However, ITQ suffers from huge memory footprint and computational cost because of the projection matrix computation.

### B. Video Retrieval

For the extension of video retrieval, we perform evaluation on the face retrieval task with two challenging TV-Series [55], *i.e.*, the Big Bang Theory (BBT), and Buffy the Vampire Slayer (BVS). The 3341 face videos are collected from the first six episodes from the first season of the BBT, and 4779 face videos of the first six episodes are acquired from BVS. As discussed in [7], we simply treat the video as a set of frames. Each face frame is represented with a 240-dimensional discrete cosine transformation feature. To reduce the dimensionality of original features, in this experiment, we employ FV encoding to aggregate the discrete cosine transformation features with 128

Gaussian components. Therefore, the  $240 \times 128 = 30,720$  dimensional aggregated FV is used to represent each face video. Following the training and testing dataset partitions introduced in [7], we randomly select we randomly selected 300 face videos for training on BBT and BVS datasets, and then selected 100 videos from the rest as the query for the retrieval task. We treat video retrieval as an ANN search problem. Given a query video, we find the videos that are the nearest neighbors of the query based on the Euclidean distance. The ground-truth is defined by the 50 nearest neighbors of each query video in the Euclidean space. We repeat the experiments 10 times and take the average mAP as the final result.

Although supervised binary coding methods achieve higher accuracy than those unsupervised and semi-supervised ones [7], [56], we compare our unsupervised method with several state-of-the-art unsupervised methods including LSH [31], BPBC [30], SGH [57], ITQ [20], SP [22], DBQ [58], CBE [59] and PCA-RR [20]. We utilize the published codes and suggested parameters of these methods from the corresponding authors. Performance comparisons are shown in Figure 6. Our method clearly outperforms other approaches on the long codes. Although our method yields comparable performance as ITQ especially when the codesize is larger than 2048, the proposed method runs much faster than ITQ.

### C. Effectiveness of Dual Selection

In this section, we further analyze the impact of dual selection on retrieval performance on the MPEG CDVS datasets and the INRIA Holidays dataset, *i.e.*, sample-specific component selection and component-specific bit selection. Note that the proposed method degenerates to two basic models: Ours\_S

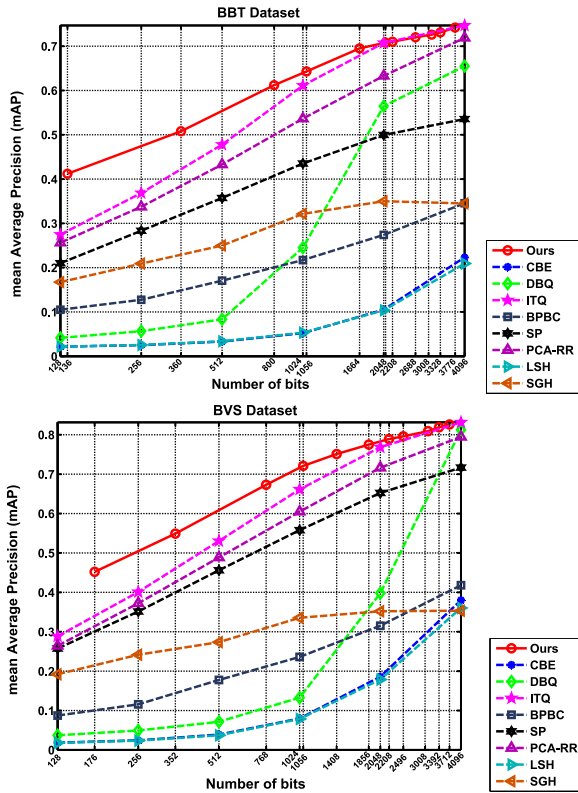


Fig. 6. Comparisons of the mean average precision on BBT and BVS datasets (best viewed in color).

denotes that only the sample-specific component selection scheme is employed. Ours\_C denotes that only the component-specific bit selection is applied. From Figure 7, the Ours\_S performs consistently better than Ours\_C at the same codesize over all datasets. In contrast, the proposed full-fledged method fusing Ours\_S and Ours\_C significantly outperforms both Ours\_S and Ours\_C (especially at small codesize). For example, at codesize of 2048 bits, Ours\_S and Ours\_C yield mAP 74.0% and 33.9% on the Graphics dataset, which is inferior to 77.8%. This gain shows the dual selection scheme is complementary to each other, which can select more informative bits. With more bits, Ours\_S, Ours\_C and the dual selection scheme can progressively improve the retrieval accuracy. In particular, Ours\_S and our full-fledged method approach to the retrieval mAP of uncompressed FV at large codesize. They obtain mAP 65.8% at codesize of 8192 bits on the Holidays dataset, while the original FV yields mAP 70.0%.

Noted that we only show the results of Ours\_C at the codesizes from 2048 to 8192, rather than presenting the results at the codesize of both 512 and 1024. This is because our FV encoding aggregates SIFT features with 512 Gaussian components. If the codesize is set to 512, just one bit is chosen from each Gaussian component, which is not reasonable. Moreover, Ours\_C is notably worse than Ours\_S and our full-fledged method. The main reason can be explained as follows. The original FV signal presents a sort of Gaussian basis sparsity. Indeed, if none of local feature was assigned to Gaussian  $i$ , then all the elements of the corresponding Fisher sub-vector in Eq. (2) and Eq. (3) are

supposed to be zeros. Although Ours\_C alone does not work well, we may yield satisfactory improvements by combining Ours\_C with Ours\_S.

#### D. Computational Complexity

Table I compares the compression ratio, memory and time complexity of the proposed method and other baselines. With comparable retrieval mAP, the compression ratio of our method is 2 to 6 times larger than the baseline schemes, resulting in much smaller compact codes. The memory footprint of our method is extremely low, *i.e.*, 0.015 MB (16 KB) for the globally bit selection mask. By contrast, RR+PQ and LSH cost over hundreds of megabytes to store the projection matrix. In addition, our method is super fast, as only binarization and selection operations are involved, while hashing methods and PQ often involve heavy floating point multiplications.

#### E. Discussion

*Impact of Key Parameters:* We evaluate the impact of key parameters of the proposed algorithm, including the number of selected Gaussian components  $M$  and the number of selected bits  $D'$ , in terms of retrieval mAP on BBT and BVS datasets, as shown in Table III. In our experiments, different numbers of selected bits  $D' = \{8, 16, 32, 64\}$  are applied to different configuration settings. For example, to obtain a binary code with about 1024 bits, we employ the cross-validation algorithm to obtain the combination  $\{M, D'\}$  that produces the best performance on the evaluation dataset. From Table III,  $M = 66$  and  $D' = 16$  provides the best results. In practice, we employ cross-validation to determine the best parameters.

*Performance Impact of Combining Deep Features and Hand-crafted Features:* To further improve the performance of our model, we incorporate CNNs and our hand-crafted compact aggregated descriptors into the well-established CDVS evaluation framework to study the effectiveness of our method. In particular, we use a Nested Invariance Pooling (NIP) method [19], [60] to derive the compact and robust CNNs descriptor. Given a query  $\mathbf{x}_q$  and a dataset sample  $\mathbf{x}_r$ , instead of simply concatenating the NIP derived deep descriptors and hand-crafted descriptors, we apply the weighted similarity scores as follows:

$$S(\mathbf{x}_r, \mathbf{x}_q) = \alpha * S_c(\mathbf{x}_r, \mathbf{x}_q) + (1 - \alpha)S_t(\mathbf{x}_r, \mathbf{x}_q),$$

where  $\alpha$  is the weighting factor.  $S_c$  and  $S_t$  represent the matching score of NIP and hand-crafted descriptors, respectively. In this work,  $\alpha$  is empirically set to 0.75.

It is shown that the combination descriptors, referred to as Ours+NIP, can greatly improve performance, as depicted in Figure 7. For example, we have achieved more than 20% mAP improvement over Ours on the Holidays dataset. However, the performance growth trend of Ours+NIP is less apparent on the Holidays and Common datasets. This is because the NIP can get remarkable results with the codesize of 512, *e.g.*, 0.885 mAP on the Holidays dataset, 0.964 mAP on the Common dataset. In this case, the incremental performance improvements of our scalable descriptor is inhibited. However, the performance growth trend

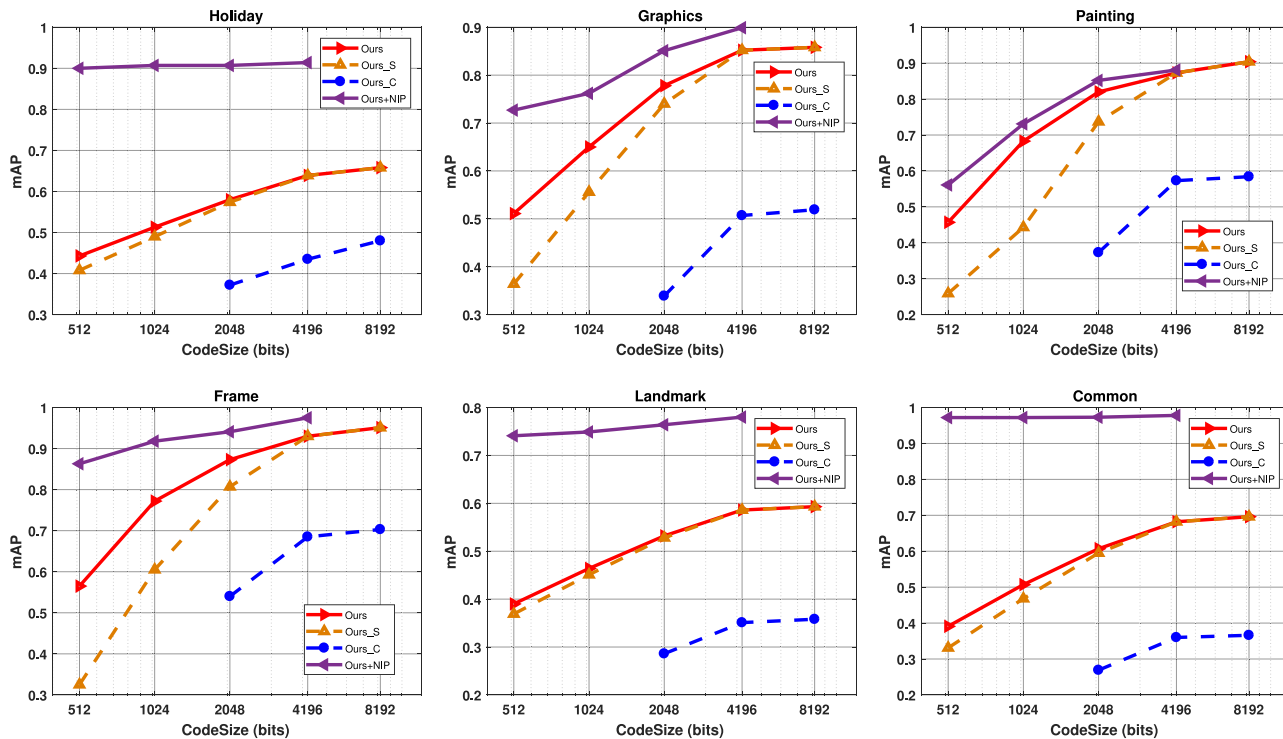


Fig. 7. Effectiveness of Dual Selection. Retrieval performance in terms of mAP vs. descriptor codesize over various benchmark datasets (best viewed on high-resolution display). It is noted that the resulting dimensions of the NIP descriptor of VGG-16 is 512 [19], [60].

TABLE I  
COMPARISONS OF DESCRIPTOR COMPRESSION RATIO, COMPRESSION TIME AND MEMORY FOOTPRINT FOR FV COMPRESSION WITH COMPARABLE RETRIEVAL MAP, *e.g.*, ABOUT 51% ON THE INRIA HOLIDAYS DATASET

Method	Compression Ratio	Memory Cost		Compression Time	
		Theoretical	Practice (MB)	Theoretical	Practice (ms)
LSH [31]	93.6	$\mathcal{O}(NP)$	91.8	$\mathcal{O}(NP)$	151
BPBC [30]	154.2	$\mathcal{O}(\sqrt{N}\sqrt{P})$	0.031	$\mathcal{O}(N\sqrt{P} + P\sqrt{N})$	17
PQ [24]	65.5	$\mathcal{O}(NQ)$	8.4	$\mathcal{O}(NQ)$	25
RR+PQ [24], [38], [54]	65.5	$\mathcal{O}(NQ + N^2)$	277	$\mathcal{O}(NQ + N^2)$	278
ITQ [20]	256	$\mathcal{O}(N^2)$	254	$\mathcal{O}(N^2)$	257
Ours	512	$\mathcal{O}(N)$	0.015	$\mathcal{O}(N)$	< 1

Note that  $N, P, Q$  denote the dimensionality of FV, the target codesize and the size of vector quantization codebooks for the PQ method, respectively.

TABLE II  
PERFORMANCE OF VIDEO RETRIEVAL ON THE MPEG CDVA DATASET

	mAP	Precision@ $R$	Descriptor Size	Search time(s)
Pool5 [61]	0.670	0.638	2KB	2.3
R-MAC [62]	0.771	0.738	2KB	2.3
NIP [60]	0.801	0.767	2KB	2.3
Ours+NIP	0.849	0.824	4KB	4.9

of Ours+NIP is comparable with Ours as the number of codesize grows on the Graphics, Painting and Frame datasets.

Furthermore, we extend experiments in the evaluation framework of emerging MPEG CDVA to validate the effectiveness of Ours+NIP. The MPEG CDVA dataset<sup>4</sup> includes 9974 query and 5127 reference videos. For video retrieval, the performance is evaluated by the mean Average Precision (mAP), the precision at a given cut-off rank  $R$  for a single query (Precision@ $R$ ) and

search time. For each query, we retrieve its  $R$  nearest items and compute the ratio of the number of retrieved relevant points to  $R$ . In our experiment, we set  $R = 50$ . Table II shows that NIP achieves the promising retrieval performance over both pool5 (CNNs features) and R-MAC. Furthermore, Ours+NIP has significantly improved the retrieval performance against Pool5 (CNNs features) by 17.9% in mAP and 18.6% in Precision@ $R$ . Compared with the state-of-the-art R-MAC, 7.8% in mAP and 8.6% in Precision@ $R$  are achieved. Overall, it is shown that a combination of our global CDVS descriptor and NIP CNNs descriptors can greatly improve performance, in which the positive complementary effects of hand-crafted descriptors and deep invariant descriptors have been well demonstrated.

*Performance impact of Codebook Size:* We evaluate our method for different numbers of Gaussian components (*i.e.*, the codebook size) and present performance in Figure 8. For raw FV, the longer codesize is, the better performance is. Moreover, the compressed descriptors of different lengths yield the best results

<sup>4</sup><http://www.cldatlas.com/cdva/dataset.html>

TABLE III  
IMPACT OF KEY PARAMETERS OF THE PROPOSED ALGORITHM, INCLUDING THE NUMBER OF SELECTED COMPONENTS  $M$ , THE NUMBER OF SELECTED BITS  $D'$

Feature	$M \times D'$						
	$128 \times 8$	$64 \times 16$	$65 \times 16$	$66 \times 16$	$32 \times 32$	$33 \times 32$	$16 \times 64$
BBT	0.580	0.617	0.629	<b>0.643</b>	0.616	0.632	0.591
BVS	0.638	0.674	0.683	<b>0.721</b>	0.697	0.705	0.644

Assume that the code length is about 1024.

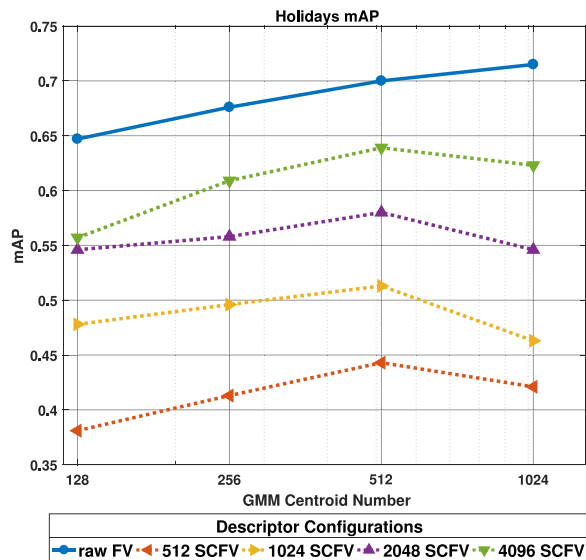


Fig. 8. Impact of the number of Gaussian components (*i.e.*, codebook size) measured on the Holidays dataset.

over benchmarks when the number of Gaussian components is 512 and the performance is further boosted as the code length grows. However, performance drops when the codebook size is set to 1024, which is mainly attributed to the fewer or even none overlapping Gaussians between the query and database images. In addition, we simply use  $sgn(g)$  to binarize each element  $g_j$  in Eq. (2) and Eq. (3), which may introduce more noise as the dimensionality of the original FV or VLAD grows. So ours is worse than the original uncompressed FV.

*Comparisons with Binary VLAD Code:* VLAD encoding is regarded as a simplified version of FV encoding [16]. We can obtain the compact binary VLAD code using the proposed scheme as well. Table IV compares the retrieval accuracy in terms of mAP, STM<sup>5</sup> and search time on the “Holidays + 1M Flickr” dataset [51] mentioned above. The proposed method is compared against several state-of-the-art methods including LSH [31], PQ [24], HKM [63], EWH [64] and MIH [65]. The compact descriptors derived from both VLAD and FV are evaluated. Both are significantly faster than linear search with comparable retrieval accuracy. For instance, our method obtains around 30 (*resp.* 16) times speedup than linear search at the cost of a minor mAP drop, *i.e.*, 0.18% (*resp.* 0.09%) for the VLAD feature (*resp.* FV feature). The proposed method performs significantly better than HKM (*i.e.*, +18% mAP) with comparable search time. Although the mAP of ours is slightly worse than

<sup>5</sup>STM denotes the Success rate for Top Match to measure the precision at rank 1, which is defined as (the number of times the top retrieved image is relevant)/(the number of queries).

TABLE IV  
COMPARISONS OF THE COMPACT BINARY CODES DERIVED FROM BOTH VLAD AND FV ON THE “HOLIDAYS + 1M FLICKR” DATASET [51]. THE CODE LENGTH IS 4096

Feature	Method	INRIA Holidays		
		mAP	STM	time (s)
Binary codes obtained from FV	Linear search	63.59	69.85	2.23
	HKM [63]	45.20	47.80	0.12
	PQ [24]	62.53	68.15	2.98
	LSH [31]	62.69	68.88	1.06
	EWH [64]	61.23	67.99	0.43
	MIH [65]	63.50	69.80	0.53
	<b>Ours</b>	<b>63.90</b>	<b>69.50</b>	<b>0.14</b>
Binary codes obtained from VLAD	Linear search	62.23	68.2	2.41
	HKM [63]	43.70	45.60	0.11
	PQ [24]	60.15	65.23	2.95
	LSH [31]	61.16	67.50	1.25
	EWH [64]	60.46	67.10	0.47
	MIH [65]	62.20	68.20	0.50
	<b>Ours</b>	<b>62.05</b>	<b>67.90</b>	<b>0.08</b>

MIH, search is faster. Overall, the FV codes outperform VLAD codes.

## VII. CONCLUSION

We have proposed an effective solution for learning compact binary codes to address the fast ANN problem with high-dimensional aggregated descriptor. Both the sample-specific Gaussian component selection and the component-specific bit selection are proposed to produce a codebook-free compact descriptor, which exhibits extremely low compression memory and time complexity, and supports fast Hamming distance matching. Moreover, our approach has been validated in the evaluation framework of the MPEG CDVS standard, and provided a groundwork of compact hand-crafted features for the emerging MPEG CDVA standard. Extensive experimental results demonstrate superior retrieval performance against the state-of-the-art methods such as hashing and PQ.

## REFERENCES

- [1] H. Müller and D. Unay, “Retrieval from and understanding of large-scale multi-modal medical datasets: A review,” *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 2093–2104, Sep. 2017.
- [2] F. Shen *et al.*, “Asymmetric binary coding for image search,” *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 2022–2032, Sep. 2017.
- [3] Z. Chen, J. Lu, J. Feng, and J. Zhou, “Nonlinear sparse hashing,” *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 1996–2009, Sep. 2017.
- [4] L. Liu, M. Yu, and L. Shao, “Learning short binary codes for large-scale image retrieval,” *IEEE Trans. Image Process.*, vol. 26, no. 3, pp. 1289–1299, Mar. 2017.
- [5] R. Ji *et al.*, “When location meets social multimedia: A survey on vision-based recognition and mining for geo-social multimedia analytics,” *ACM Trans. Intell. Syst. Technol.*, vol. 6, no. 1, Art. no. 1, 2015.

- [6] R. Tao, A. W. Smeulders, and S.-F. Chang, "Attributes and categories for generic instance search from one example," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 177–186.
- [7] Y. Li, R. Wang, Z. Huang, S. Shan, and X. Chen, "Face video retrieval with image query via hashing across Euclidean space and Riemannian manifold," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 4758–4767.
- [8] L.-Y. Duan *et al.*, "Overview of the MPEG-CDVS standard," *IEEE Trans. Image Process.*, vol. 25, no. 1, pp. 179–194, Jan. 2016.
- [9] J. Lin *et al.*, "Rate-adaptive compact fisher codes for mobile visual search," *Signal Process. Letters*, vol. 21, no. 2, pp. 195–198, 2014.
- [10] L.-Y. Duan, J. Lin, J. Chen, T. Huang, and W. Gao, "Compact descriptors for visual search," *IEEE MultiMedia*, vol. 21, no. 3, pp. 30–40, Jul.–Sep. 2014.
- [11] L.-Y. Duan, J. Lin, Z. Wang, T. Huang, and W. Gao, "Weighted component hashing of binary aggregated descriptors for fast visual search," *IEEE Trans. Multimedia*, vol. 17, no. 6, pp. 828–842, Jun. 2015.
- [12] S. Paschalakis *et al.*, "Information Technology - Multimedia content descriptor interface - Part 13: Compact Descriptors for Visual Search", International Standard ISO/IEC15938-13, First Edition, 2015-09-01.
- [13] L.-Y. Duan *et al.*, "Compact descriptors for video analysis: The emerging MPEG standard," 2017, arXiv:1704.08141.
- [14] A. Redondi, L. Baroffio, M. Cesana, and M. Tagliasacchi, "Compress-then-analyze vs. analyze-then-compress: Two paradigms for image analysis in visual sensor networks," in *Proc. IEEE Int. Workshop Multimedia Signal Process.*, 2013, pp. 278–282.
- [15] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," *Int. J. Comput. Vision*, vol. 105, no. 3, pp. 222–245, 2013.
- [16] H. Jégou *et al.*, "Aggregating local image descriptors into compact codes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 9, pp. 1704–1716, Sep. 2012.
- [17] H. F. Yang, K. Lin, and C. S. Chen, "Supervised learning of semantics-preserving hash via deep convolutional neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 2, pp. 437–451, Feb. 2017.
- [18] V. E. Liong, J. Lu, Y. P. Tan, and J. Zhou, "Deep video hashing," *IEEE Trans. Multimedia*, vol. 19, no. 6, pp. 1209–1219, Jun. 2017.
- [19] Y. Lou *et al.*, "Compact deep invariant descriptors for video retrieval," in *Proc. Data Compression Conf.*, 2017, pp. 73–82.
- [20] 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.
- [21] X. Zhu *et al.*, "Graph PCA hashing for similarity search," *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 2033–2044, Sep. 2017.
- [22] Y. Xia, K. He, P. Kohli, and J. Sun, "Sparse projections for high-dimensional binary codes," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 3332–3339.
- [23] L.-Y. Duan *et al.*, "Minimizing reconstruction bias hashing via joint projection learning and quantization," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 3127–3141, Jun. 2018.
- [24] H. Jegou, M. Douze, and C. Schmid, "Product quantization for nearest neighbor search," *IEEE Trans. Pattern Analysis Mach. Intell.*, vol. 33, no. 1, pp. 117–128, Jan. 2011.
- [25] L. Jin *et al.*, "Online variable coding length product quantization for fast nearest neighbor search in mobile retrieval," *IEEE Trans. Multimedia*, vol. 19, no. 3, pp. 559–570, Mar. 2017.
- [26] S. Liu, J. Shao, and H. Lu, "Generalized residual vector quantization and aggregating tree for large scale search," *IEEE Trans. Multimedia*, vol. 19, no. 8, pp. 1785–1797, Aug. 2017.
- [27] T. Zhang, G.-J. Qi, J. Tang, and J. Wang, "Sparse composite quantization," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 4548–4556.
- [28] Y. Cao, M. Long, J. Wang, and S. Liu, "Deep visual-semantic quantization for efficient image retrieval," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2017, pp. 1328–1337.
- [29] J. Wang *et al.*, "Optimized cartesian k-means," *IEEE Trans. Knowledge and Data Engineering*, vol. 27, no. 1, pp. 180–192, Jan. 2015.
- [30] Y. Gong, S. Kumar, H. Rowley, and S. Lazebnik, "Learning binary codes for high-dimensional data using bilinear projections," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2013, pp. 484–491.
- [31] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, "Locality-sensitive hashing scheme based on p-stable distributions," in *Proc. 12th Annu. Symp. Comput. Geom.*, 2004, pp. 253–262.
- [32] M. Raginsky and S. Lazebnik, "Locality-sensitive binary codes from shift-invariant kernels," in *Proc. Advances Neural Inf. Process. Syst.*, 2009, pp. 1509–1517.
- [33] Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in *Proc. Advances Neural Inf. Process. Syst.*, 2009, pp. 1753–1760.
- [34] Z. Wang, L.-Y. Duan, J. Yuan, T. Huang, and W. Gao, "To project more or to quantize more: Minimizing reconstruction bias for learning compact binary codes," in *Proc. 25th Int. Joint Conf. Artif. Intell.*, 2016, pp. 2181–2188.
- [35] L. Liu, M. Yu, and L. Shao, "Learning short binary codes for large-scale image retrieval," *IEEE Trans. Image Process.*, vol. 26, no. 3, pp. 1289–1299, Mar. 2017.
- [36] K. He, F. Wen, and J. Sun, "K-means hashing: An affinity-preserving quantization method for learning binary compact codes," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2013, pp. 2938–2945.
- [37] Y. Guo, G. Ding, L. Liu, J. Han, and L. Shao, "Learning to hash with optimized anchor embedding for scalable retrieval," *IEEE Trans. Image Process.*, vol. 26, no. 3, pp. 1344–1354, Mar. 2017.
- [38] T. Ge, K. He, Q. Ke, and J. Sun, "Optimized product quantization for approximate nearest neighbor search," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2013, pp. 2946–2953.
- [39] T. Zhang, C. Du, and J. Wang, "Composite quantization for approximate nearest neighbor search," in *Proc. Int. Conf. Mach. Learn.*, 2014, no. 2, pp. 838–846.
- [40] K. Lin, J. Lu, C.-S. Chen, and J. Zhou, "Learning compact binary descriptors with unsupervised deep neural networks," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2016, pp. 1183–1192.
- [41] V. Erin Liang, J. Lu, G. Wang, P. Moulin, and J. Zhou, "Deep hashing for compact binary codes learning," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 2475–2483.
- [42] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *Proc. Eur. Conf. Comput. Vision*, 2014, pp. 392–407.
- [43] F. Perronnin, Y. Liu, J. Sánchez, and H. Poirier, "Large-scale image retrieval with compressed fisher vectors," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2010, pp. 3384–3391.
- [44] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *Proc. IEEE Int. Conf. Comput. Vision*, 2003, pp. 1470–1477.
- [45] D. Chen *et al.*, "Residual enhanced visual vector as a compact signature for mobile visual search," *Signal Process.*, vol. 93, no. 8, pp. 2316–2327, 2013.
- [46] O. M. Parkhi, K. Simonyan, A. Vedaldi, and A. Zisserman, "A compact and discriminative face track descriptor," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2014, pp. 1693–1700.
- [47] Y. Xu and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," *SIAM J. Imag. Sci.*, vol. 6, no. 3, pp. 1758–1789, 2013.
- [48] N. Parikh and S. Boyd, "Proximal algorithms," *Foundations Trends Optim.*, vol. 1, no. 3, pp. 123–231, 2013.
- [49] N. Zhou, Y. Xu, H. Cheng, J. Fang, and W. Pedrycz, "Global and local structure preserving sparse subspace learning: An iterative approach to unsupervised feature selection," *Pattern Recognit.*, vol. 53, pp. 87–101, 2016.
- [50] V. Chandrasekhar *et al.*, "Feature matching performance of compact descriptors for visual search," in *Data Compression Conf.*, 2014, pp. 3–12.
- [51] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *Proc. Eur. Conf. Comput. Vision*, 2008, pp. 304–317.
- [52] R. Arandjelović, "Advancing large scale object retrieval," Ph.D. dissertation, Dept. Eng. Sci., Univ. of Oxford, Oxford, U.K., 2013.
- [53] H. Jégou, and A. Zisserman, "Triangulation embedding and democratic aggregation for image search," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, Jun. 2014, pp. 3310–3317.
- [54] M. Norouzi and D. J. Fleet, "Cartesian k-means," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2013, pp. 3017–3024.
- [55] M. Bauml, M. Tapaswi, and R. Stiefelhagen, "Semi-supervised learning with constraints for person identification in multimedia data," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2013, pp. 3602–3609.
- [56] Y. Li, R. Wang, S. Shan, and X. Chen, "Hierarchical hybrid statistic based video binary code and its application to face retrieval in TV-series," in *Proc. IEEE Int. Conf. Workshops Autom. Face Gesture Recognit.*, 2015, pp. 1–8.

- [57] Q.-Y. Jiang and W.-J. Li, "Scalable graph hashing with feature transformation," in *Proc. 24th Int. Joint Conf. Artif. Intell.*, 2015, pp. 2248–2254.
- [58] W. Kong and W.-J. Li, "Double-bit quantization for hashing," in *Proc. AAAI Conf. Artificial Intell.*, 2012, pp. 634–640.
- [59] F. X. Yu, S. Kumar, Y. Gong, and S.-F. Chang, "Circulant binary embedding," in *Proc. 31th Int. Conf. Mach. Learn.*, 2014, pp. 946–954.
- [60] J. Lin *et al.*, "HNIP: Compact deep invariant representations for video matching, localization, and retrieval," *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 1968–1983, Sep. 2017.
- [61] H. Azizpour, A. S. Razavian, J. Sullivan, A. Maki, and S. Carlsson, "From generic to specific deep representations for visual recognition," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit. Workshops*, 2015, pp. 36–45.
- [62] G. Toliás, R. Sivic, and H. Jégou, "Particular object retrieval with integral max-pooling of CNN activations," 2016, arXiv:1511.05879v2.
- [63] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2006, vol. 2, pp. 2161–2168.
- [64] M. M. Esmaili, R. K. Ward, and M. Fatouchechi, "A fast approximate nearest neighbor search algorithm in the hamming space," *IEEE Trans. Pattern Analysis Mach. Intell.*, vol. 34, no. 12, pp. 2481–2488, Dec. 2012.
- [65] M. Norouzi, A. Punjani, and D. J. Fleet, "Fast search in hamming space with multi-index hashing," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2012, pp. 3108–3115.



**Yuwei Wu** received the Ph.D. degree in computer science from the Beijing Institute of Technology (BIT), Beijing, China, in 2014. From August 2014 to August 2016, he was a Postdoctoral Research Fellow with the School of Electrical & Electronic Engineering, Nanyang Technological University, Singapore. He is now an Assistant Professor with the School of Computer Science, BIT. His research interests include computer vision and information retrieval. Dr. Wu received the Outstanding Ph.D. Thesis award from BIT, and the Distinguished Dissertation Award Nominee from China Association for Artificial Intelligence.



**Feng Gao** received the B.S. degree in computer science from University College London, London, U.K., in 2007, and the Ph.D. degree in computer science from Peking University, China, in 2018. His research interest is to work on the intersection of computer science and art, including but not limited to artificial intelligence and painting art, deep learning and painting robot, etc.



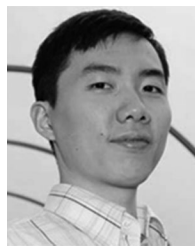
**Yicheng Huang** received the bachelor's degree in computer science and technology from Peking University, Beijing, China, in 2015, and is currently working toward the M.S. degree at the School of Electrical Engineering and Computer Science, Peking University, Beijing, China. His current research interests include large-scale image retrieval and fast nearest neighbor search.



**Jie Lin** received the B.S. and Ph.D. degrees from the School of Computer Science and Technology, Beijing Jiaotong University, Beijing, China, in 2006 and 2014, respectively. From 2011 to 2014, he was a visiting student with the Rapid-Rich Object Search Laboratory, Nanyang Technological University, Singapore, and the Institute of Digital Media, Peking University, Beijing, China. He is currently a Research Scientist with the Institute of Infocomm Research, A\*STAR, Singapore. His work on image feature coding has been recognized as core contribution to the MPEG-7 Compact Descriptors for Visual Search standard. His research interests include deep learning, feature coding, and large-scale image/video retrieval.



**Vijay Chandrasekhar** received the B.S and M.S. degrees from Carnegie Mellon University, Pittsburgh, PA, USA, in 2002 and 2005, respectively, and the Ph.D. degree in electrical engineering from Stanford University, Stanford, CA, USA, in 2013. His research interests include signal processing, mobile visual search, visual feature coding, etc. His Ph.D. work on feature compression led to the MPEG-CDVS (Compact Descriptors for Visual Search) standard, which he actively contributed from 2010 to 2013. He has authored or coauthored more than 80 papers/MPEG contributions in a wide range of top-tier journals/conferences such as the *International Journal of Computer Vision*, *ICCV*, *CVPR*, the *IEEE SIGNAL PROCESSING MAGAZINE*, *ACM Multimedia*, the *IEEE TRANSACTIONS ON IMAGE PROCESSING*, *Designs, Codes and Cryptography*, the International Society of Music Information Retrieval, and MPEG-CDVS, and has filed 7 U.S. patents (one granted, six pending). He is currently a research scientist with the Institute for Infocomm Research, Singapore. His research interests include mobile audio and visual search, large-scale image and video retrieval, machine learning, and data compression. Dr. Chandrasekhar was the recipient of the A\*STAR National Science Scholarship in 2002.



**Junsong Yuan** (M'08–SM'14) graduated from the Special Class for the Gifted Young of Huazhong University of Science and Technology, Wuhan, China, in 2002. He received the M.Eng. degree from the National University of Singapore, Singapore, in 2005, and the Ph.D. degree from Northwestern University, Evanston, IL, USA, in 2009. He was an Associate Professor with the School of Electrical and Electronics Engineering, Nanyang Technological University, Singapore. He is currently an Associate Professor with the Department of Computer Science and Engineering, State University of New York, Buffalo, NY, USA. Dr. Yuan received the 2016 Best Paper Award of the *IEEE TRANSACTIONS ON MULTIMEDIA*, the Doctoral Spotlight Award of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'09), the Nanyang Assistant Professorship from NTU, and the Outstanding EECS Ph.D. Thesis award from Northwestern University. He was a Guest Editor for the *International Journal of Computer Vision*. He is currently a Senior Area Editor for the *Journal of Visual Communication and Image Representation*, and an Associate Editor for the *IEEE TRANSACTIONS ON IMAGE PROCESSING* and the *IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY*. He is a Program Co-Chair of ICME'18 and VCIP'15, and an Area Chair of ACM MM'18, ACCV'18'14, ICPR'18'16, CVPR'17, ICIP'18'17, etc.



**Ling-Yu Duan** (M'06) received the M.Sc. degree in automation from the University of Science and Technology of China, Hefei, China, in 1999, the M.Sc. degree in computer science from the National University of Singapore (NUS), Singapore, in 2002, and the Ph.D. degree in information technology from The University of Newcastle, Callaghan, Australia, in 2008. He was the Associate Director of the Rapid-Rich Object Search Laboratory, a joint lab between Nanyang Technological University, Singapore, and Peking University (PKU), Beijing, China since 2012. He is currently a Full Professor with the National Engineering Laboratory of Video Technology, PKU, China. His research interests include multimedia indexing, search, and retrieval, mobile visual search, visual feature coding, video analytics, etc. Dr. Duan received the EURASIP Journal on Image and Video Processing Best Paper Award in 2015, the Ministry of Education Technology Invention Award (First Prize) in 2016, the National Technology Invention Award (Second Prize) in 2017, China Patent Award for Excellence (2017), and the National Information Technology Standardization Technical Committee "Standardization Work Outstanding Person" Award in 2015. He was a Co-Editor for the MPEG Compact Descriptor for Visual Search Standard (ISO/IEC 15938-13). He is currently an Associate Editor for the *ACM Transactions on Intelligent Systems and Technology* and the *ACM Transactions on Multimedia Computing, Communications, and Applications*. He is a Co-Chair of the MPEG Compact Descriptor for Video Analytics.