

# Efficient Conversion of Deep Features to Compact Binary Codes Using Fourier Decomposition for Multimedia Big Data

Jamil Ahmad, *Student Member, IEEE*, Khan Muhammad, *Student Member, IEEE*,  
Jaime Lloret, *Senior Member, IEEE*, and Sung Wook Baik, *Member, IEEE*

**Abstract**—Exponential growth of multimedia data has been witnessed in recent years from various industries, such as e-commerce, health, transportation, and social networks, etc. Access to desired data in such gigantic datasets require sophisticated and efficient retrieval methods. In the last few years, neuronal activations generated by a pre-trained convolutional neural network (CNN) have served as generic descriptors for various tasks including image classification, object detection and segmentation, and image retrieval. They perform incredibly well compared to hand-crafted features. However, these features are usually high dimensional, requiring a lot of memory and computations for indexing and retrieval. For very large datasets, utilization of these high dimensional features in raw form becomes infeasible. In this paper, a highly efficient method is proposed to transform high dimensional deep features into compact binary codes using bidirectional Fourier decomposition. This compact bit code saves memory and eases computations during retrieval. Further, these codes can also serve as hash codes, allowing very efficient access to images in large datasets using approximate nearest neighbor (ANN) search techniques. Our method does not require any training and achieves considerable retrieval accuracy with short length codes. It has been tested on features extracted from fully connected layers of a pretrained CNN. Experiments conducted with several large datasets reveal the effectiveness of our approach for a wide variety of datasets.

**Index Terms**—Deep learning, Fourier transform, hash codes, image retrieval, industrial informatics.

## I. INTRODUCTION

**B**IG data has recently emerged as a key concept, denoting the gigantic volume of data generated at a rapid pace due to the progress in sensing, communication, storage, cloud computing technologies, and algorithms. Recent statistics reveal

Manuscript received September 29, 2017; revised December 29, 2017; accepted January 22, 2018. This work was supported by the National Research Foundation of Korea grant funded by the Korea government (MSIP) (2016R1A2B4011712). Paper no. TII-17-2275. (*Corresponding author: Sung Wook Baik.*)

J. Ahmad, K. Muhammad, and S. W. Baik are with the Intelligent Media Laboratory, Digital Contents Research Institute, Sejong University, Seoul 143-747, South Korea (e-mail: jamilahmad@ieee.org; khan.muhammad@ieee.org; sbaik@sejong.ac.kr).

J. Lloret is with the Instituto de Investigación para la Gestión Integrada de Zonas Costeras, Universitat Politècnica de Valencia, València 46022, Spain (e-mail: jlloret@dcom.upv.es).

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/TII.2018.2800163

that 1200 Exabytes of data is generated annually and the rate is growing rapidly [1]. A huge fraction of this data is multimedia data (images and videos), generated by various industries, such as health, surveillance, agriculture, social web, online streaming services, movies, games, and internet protocol television (IPTV) industry [2]. For example, Facebook alone contains more than 40 billion photos [3]. Similarly, more than 500 h of videos are uploaded to YouTube every minute [4]. These massive amounts of data present enormous challenges for businesses and industries. At the same time, it provides opportunities for impressive future growth, based on effective utilization of the data for analysis. For instance, progress in medical imaging technologies allows visual analysis of patient through a variety of means including endoscopy, magnetic resonance imaging, radiography, ultrasonography, and many others. It causes huge amounts of data to be generated, which is stored for immediate or future use. Similarly, surveillance cameras deployed in wake of the recent security concerns throughout the globe, also generate huge amounts of multimedia data, required to be stored and properly indexed for possible future use. Major issues with these gigantic multimedia repositories include transmission, management, storage, and their efficient indexing and retrieval.

Providing reliable and efficient access to relevant data in large image repositories based on their contents is a highly challenging task which has been studied over the course of almost three decades. Content-based image retrieval (CBIR) methods allow retrieval of relevant images based on the content similarity between the query and target images [5], [6]. A core component of CBIR systems aims to represent images as feature vectors or feature histograms that correspond to the color or texture content of the image [7]. These systems can also be used to personalize and recommend contents for IPTV delivery services [8]. Traditionally, CBIR relied on hand-engineered features, such as scale invariant features transform [9], bag-of-visual-words histograms [10], [11], fisher vectors [12], vectors of locally aggregated descriptors [13], GIST [14], and CENsus TRansform hISTogram [15]. Each of these methods represented images in terms of low-level features; however, these features often fail to model high-level semantics in images. Therefore, their performance in large and challenging datasets was not very satisfactory [16]. In recent years, the hand-engineered feature extraction methods have been overshadowed by the feature learning based methods including

83 deep convolutional neural networks (CNN), and deep denoising  
 84 auto-encoders [17], [18]. They automatically extract features  
 85 from images, which have been used in a variety of tasks, such as  
 86 image classification, object localization, recognition, segmen-  
 87 tation, and image retrieval [19]. CNNs have been widely used  
 88 by the image retrieval community and have achieved state-of-  
 89 the-art performance [16], [20]–[22]. These architectures have  
 90 several convolutional, pooling, and fully connected (FC) lay-  
 91 ers, arranged in a hierarchy where successive layers learn com-  
 92 plex features of the input [23]. Deep features are usually ex-  
 93 tracted from the FC layers of CNN which correspond to acti-  
 94 vation values of the neurons in those layers. In a typical CNN,  
 95 these features often have thousands of dimensions. Though,  
 96 these features are capable of representing images effectively,  
 97 image indexing, and matching using these features become in-  
 98 feasible for large datasets [21].

99 Hash-based image retrieval methods aim at allowing efficient  
 100 access to relevant data in large datasets using approximate near-  
 101 est neighbor (ANN) search approaches. In wake of the growing  
 102 demands for efficient access to large image repositories, these  
 103 methods have appealed significant attention in recent years [24].  
 104 They work on the principle of locality sensitive hash functions  
 105 that transform high dimensional features to low-dimensional  
 106 hamming space (binary codes) and attempt to preserve origi-  
 107 nal neighbors in the hamming space [25]. These compact codes  
 108 are then used to directly retrieve nearest neighbors of the query  
 109 image from the hamming space without exhaustive search. A  
 110 large variety of hashing methods have been proposed in recent  
 111 years, which attempt to derive compact binary codes from  
 112 image features. A few notable methods include locality sensi-  
 113 tive hashing (LSH) [25], [26], principal component analysis  
 114 based hashing (PCAH) [27], spectral hashing (SH) [28], spheri-  
 115 cal hashing (SpH) [29], and density sensitive hashing (DSH)  
 116 [30], etc. Hash methods may be data-dependent or data-  
 117 independent. They may be trained in either supervised or un-  
 118 supervised manner. Typically, these methods are trained for a  
 119 particular dataset to generate hash codes of a certain length. If  
 120 the data changes or the length of the hash code needs to be modi-  
 121 fied, the training procedure has to be rerun. These characteristics  
 122 limit their utilization in real applications.

123 In this paper, we propose an efficient method to transform  
 124 selected deep features directly into compact binary codes. It  
 125 does not require any training and can be efficiently executed on  
 126 a graphics processing unit (GPU) to quickly convert deep fea-  
 127 tures to binary codes. We show that deep features from the FC  
 128 layers of CNNs are highly redundant, hence, we propose a fea-  
 129 ture selection algorithm to identify effective deep features based  
 130 on neuronal sensitivity and diversity. The proposed hash codes  
 131 yield considerable retrieval performance for 256 and 512 bit  
 132 codes. Major contributions in this work are summarized as  
 133 follows.

134 1) We show that the high dimensional deep features ex-  
 135 tracted from FC layers of a pretrained CNN are redun-  
 136 dant and a significant number of activation features can  
 137 be removed without any loss in retrieval performance,  
 138 particularly when dealing with images of a particular cat-  
 139 egory such medical or surveillance.

2) An effective feature selection algorithm is proposed for  
 deep feature based on neuronal sensitivity and diversity  
 measures.

3) A highly efficient method is proposed for transforming  
 deep features into compact binary codes, which can be  
 used as hash codes for efficient image search. Our method  
 uses bidirectional fast Fourier transform (BD-FFT) which  
 allows hash codes of desired length to be computed di-  
 rectly without requiring any training. The method can be  
 easily implemented on a GPU for significant speedup in  
 hash code computation at large scale.

4) We also show that the selected deep features yield better  
 hash codes with the proposed BD-FFT method, and offer  
 better locality sensitivity with 256 and 512 bit codes.

The rest of the paper is organized as follows: Section II  
 highlights strengths and weaknesses of recent hash-based re-  
 trieval methods. Section III explains the proposed method in  
 detail, highlighting the key features of the presented algorithms.  
 Section IV reports evaluation results of the proposed method on  
 several popular datasets. The paper is concluded in Section V  
 with some future research directions.

## II. RELATED WORK

Extraction of discriminative features is a primary factor in  
 the success of CBIR systems. The recent deep learning based  
 methods, especially CNNs yield highly discriminative features,  
 which achieve state-of-the-art performance in CBIR. Several  
 frameworks have been proposed for utilizing deep features  
 for image retrieval in challenging scenarios. For instance,  
 Krizhevsky *et al.* [23] showed that neuronal activations  
 extracted from FC layers can be used as feature descriptors and  
 image matching can be performed using standard Euclidean  
 distance. They also showed that these high dimensional features  
 can be easily compressed with dimensionality reduction  
 methods, such as principal component analysis (PCA), sacri-  
 ficing accuracy for some degree of efficiency. Razavian *et al.*  
 [17], [18] and Babenko *et al.* [21], [22] showed that features  
 from a pretrained CNN can be used as generic descriptors  
 for image retrieval and other related tasks. They showed  
 that features from a pretrained CNN, trained on a very large  
 dataset (ImageNet [31]) achieve state-of-the-art performance,  
 surpassing traditional hand-engineered features by a huge  
 margin. Deep features from FC layers are very powerful global  
 representations, however, they are high dimensional and directly  
 utilizing them becomes inefficient, particularly for large scale  
 datasets [32].

Large scale datasets demand efficient methods for storing  
 millions of images in memory and quickly finding relevant im-  
 ages to a query image. ANN search methods like LSH have  
 shown promising results in recent years. Typically images are  
 represented as features vectors in high dimensional Euclidean  
 space, such that the Euclidean distance corresponds to image  
 similarity. The main objective of hashing methods is to generate  
 a low-dimensional embedding in hamming space while preserv-  
 ing the neighborhood. Hence, when a query is issued, the hash  
 code of the query image is used to efficiently access nearest

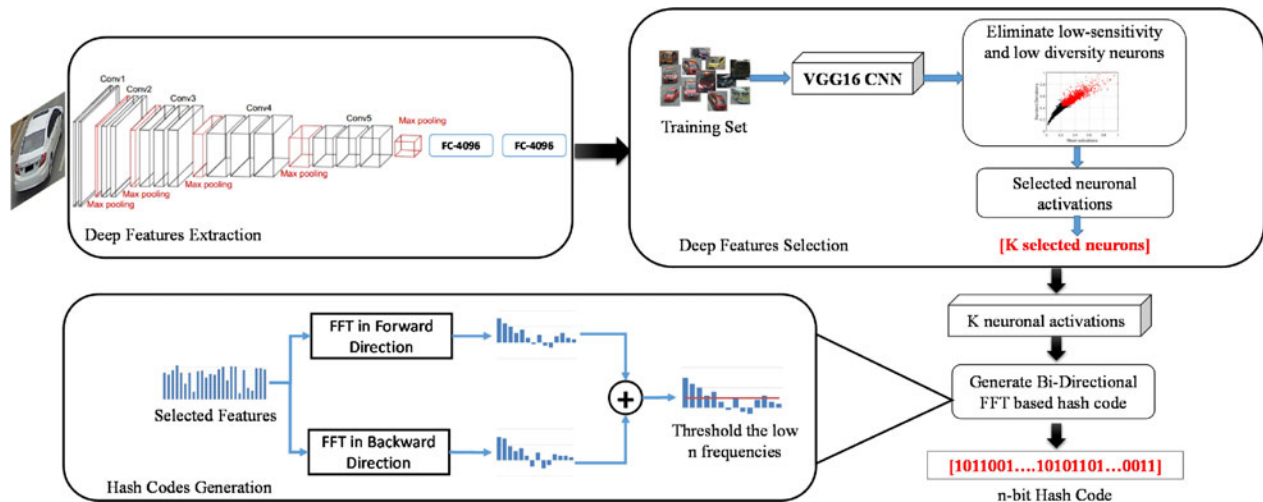


Fig. 1. Proposed framework.

195 neighbors of the query image using hamming distance. Based  
 196 on this idea, several approaches have been presented in the re-  
 197 cent years. For instance, PCAH [33] used principle directions  
 198 of data as the projection vectors to transform features to binary  
 199 codes. In LSH [25], [34], the binary code is computed through  
 200 random linear projection with a random threshold. In theory,  
 201 hamming distance between LSH codes and Euclidean distance  
 202 between image pairs are highly correlated, however, in practice  
 203 it can lead to very inefficient codes. SH [28] selects binary code-  
 204 words through minimum distance between similar points, where  
 205 similarity is defined by an approximate proximity matrix. Theo-  
 206 retically, it performs better than LSH, however, its optimization  
 207 is difficult to generalize for new data points. This problem is  
 208 solved with SpH [29] which uses Eigen functions of weighted  
 209 Laplace–Beltrami operations with the assumption of having a  
 210 multidimensional uniform distribution. It is highly efficient than  
 211 SH for hash code generation, however, its optimization is compu-  
 212 tationally expensive. DSH [30] is an extension of LSH which  
 213 utilizes random projections and also uses geometrical structure  
 214 of the data to guide the projections. It partitions the data points  
 215 into  $k$ -groups and splits each pair of adjacent groups with a  
 216 projection vector. From all such projections, DSH selects the  
 217 vectors based on the maximum entropy principle.

218 Hash-based image retrieval methods significantly improve re-  
 219 trieval efficiency in large scale datasets. However, these methods  
 220 are difficult to implement in real applications and some of them  
 221 require sufficient training data and time, while others are slow  
 222 at transforming feature vectors to hash codes. An ideal hash-  
 223 ing method is computationally efficient, simple to implement  
 224 and yield state-of-the-art performance for a variety of datasets.  
 225 In this paper, we present a simple and highly efficient way  
 226 of transforming deep features to compact binary codes using  
 227 BD-FFT.

### 228 III. PROPOSED METHOD

229 The proposed framework consists of two modules, feature  
 230 selection and hash code generation as shown in Fig. 1. First, we  
 231 studied deep features from FC layer of a pretrained VGG-16

232 CNN [35] in order to determine optimal set of features for a  
 233 particular type of data. Once the optimal features are selected,  
 234 they are converted to binary codes of different lengths using  
 235 bidirectional FFT. Details of both modules are provided in the  
 236 subsequent sections.

#### 237 A. Deep Features for Image Retrieval

238 Informatics and analytics systems make use of efficient ways  
 239 to access relevant information from large datasets. Visual data  
 240 constitute a large fraction of the data generated by different in-  
 241 dustries, where accurate and efficient access will allow analysts  
 242 and experts make better and timely decisions. Features extracted  
 243 from deep CNNs have shown state-of-the-art performance in  
 244 image retrieval from large datasets due to their impressive repre-  
 245 sentational capabilities. We used features from FC-4096 layer of  
 246 the VGG16 model [35] which was trained on ImageNet. These  
 247 features are regarded as generic descriptors for visual recogni-  
 248 tion tasks including image classification and retrieval [17], [18].  
 249 However, we argue that these features are highly powerful,  
 250 capable of representing a huge variety of visual data, and a  
 251 subset of these features will be sufficient to effectively represent  
 252 images of a particular type like medical radiographs or surveil-  
 253 lance images of vehicles, etc. In such specific datasets, subsets  
 254 of these generic features can prove to be more appropriate than  
 255 the full set of features. For this purpose, we propose an efficient  
 256 method to select deep features from a pretrained CNN for repre-  
 257 senting images of a particular type. Deep features from the FC  
 258 layer are constructed as global representations by combining the  
 259 local features extracted by various convolutional layers. VGG16  
 260 contains three FC layers having 4096, 4096, and 1000 neurons,  
 261 respectively. We used activation values of the second FC-4096  
 262 layer in our experiments because of their superior performance.  
 263 Each of these neurons are sensitive to particular objects or parts  
 264 of objects [36]. When a particular object appears in an image,  
 265 a subset of these neurons generate high activations indicating  
 266 its presence. Though these features are considered generic and  
 267 high level, their high dimensionality hinder their use in practical  
 268 applications.

## 269 B. Optimal Deep Features Selection

270 Feature reduction offers improvements in efficiency and ac-  
 271 curacy as it helps in getting rid of the less useful and often mis-  
 272 leading features [37]. We propose an efficient method to select  
 273 optimal features from a pretrained CNN. An input image is usu-  
 274 ally feed-forwarded through a deep CNN (e.g., VGG16) and the  
 275 activation values from the FC-4096 layer are extracted, which  
 276 are then used to index or retrieve images. In hash-based retrieval  
 277 systems, these features are transformed to compact binary codes  
 278 and then images are retrieved using hamming distance. How-  
 279 ever, utilizing all these features for hash code generation and  
 280 retrieving images of specific type is ineffective.

281 Deep features from FC layers are global high level features  
 282 where particular neurons are sensitive to particular objects or  
 283 their parts. They respond actively when that particular part ap-  
 284 pears somewhere in the image. For a dataset consisting of a  
 285 particular type of images, e.g., medical, it is highly unlikely  
 286 that object parts belonging to other categories, such as sports,  
 287 surveillance, or animals, may be encountered. In such a case,  
 288 utilizing all the features to represent images become ineffec-  
 289 tive which may lead to decreased performance. In recent works,  
 290 we have seen that fine-tuning pretrained CNNs on particular  
 291 datasets yield better results, which is also a verification of the  
 292 fact that specific features perform better than generic ones [16],  
 293 [38]. Instead of fine-tuning, we propose to discard irrelevant  
 294 features before using them for image retrieval tasks in specific  
 295 datasets. For this purpose, we selected a representative set of  
 296 images from a target dataset and extracted deep features from  
 297 them. We eliminated those neurons which generated negligible  
 298 activations (low sensitivity to objects of interest) or similar acti-  
 299 vations (less discriminative) for the training set. Mean activation  
 300 values  $\mu$  and standard deviations  $\sigma$  were computed for all 4096  
 301 neurons over the entire training set, where the training set  $T_s$   
 302 consisted of randomly chosen images from all the datasets we  
 303 used in the experiments and were represented by  $R^{4096}$  vectors  
 304 of deep features. Neurons having  $\mu_i$  greater than the threshold  
 305  $t_\mu$ , and  $\sigma_i$  greater than the threshold  $t_\sigma$  were selected as the data  
 306 specific discriminative features in a set  $F_s$ . This process can be  
 307 performed for selecting specific features for representing im-  
 308 ages of a particular category. The feature selection mechanism  
 309 is presented in Algorithm 1.

## 310 C. Conversion to Compact Binary Codes

311 In this paper, we consider the selected feature vector as a  
 312 one-dimensional signal, and construct its frequency domain  
 313 representation using FFT. During this transformation, the time-  
 314 domain signal is represented as a combination of different fre-  
 315 quencies. These frequencies correspond to the activation pat-  
 316 terns of neurons in the selected feature set. The Fourier spectrum  
 317 effectively captures those patterns and represents them as fre-  
 318 quencies with different amplitudes. The original signal can be  
 319 reconstructed using a certain representative frequencies of this  
 320 spectrum as shown in Fig. 2. Each frequency component will  
 321 indicate the presence or absence of a certain frequency content  
 322 (i.e., neuronal activation pattern) in the features. Based on this  
 323 idea, we select low n frequency components of the spectrum  
 324 (excluding the dc component) and transform them into binary

---

### Algorithm 1: Selection of optimal deep features.

---

**Input:** Training feature vectors  $Tf_i$  having size  $T \times R^{4096}$   
 extracted from FC-4096 (VGG16)

**Output:** Indices of selected deep features  $F_s$

**Steps:**

1. Calculate mean activation values  $\mu_i$  and standard deviation  $\sigma_i$  for all 4096 neurons across the entire training set  $T$

For  $i = 1$  to 4096

$$\mu_i = \sum_{t=1}^T Tf_i$$

$$\sigma_i = \sqrt{\frac{\sum_{t=1}^T (Tf_i - \mu_i)^2}{T}}$$

End for

2. Keep the neurons whose  $\mu_i$  are greater than  $t_\mu$  and  $\sigma_i$  is greater than  $t_\sigma$ .

$$F_{s_i} = \begin{cases} \text{Select neuron,} & \mu_i > t_\mu \text{ and } \sigma_i > t_\sigma \\ \text{Discard neuron,} & \text{otherwise} \end{cases}$$

where  $t_\mu$  and  $t_\sigma$  are selected empirically.

3. Return the indices of selected neurons in  $F_s$ .
- 

---

### Algorithm 2: Conversion of deep features to binary codes.

---

**Input:** Deep feature vector  $f_i$  having  $R^d$

**Output:** n-bit binary code

**Steps:**

1. Compute FFT of  $f_i$  in forward direction to obtain a Fourier spectrum  $F_f$

$$F_f = \sum_{j=0}^{d-1} f_i e^{-i2\pi kj/n}, \quad k = 0, \dots, d-1$$

2. Compute FFT of  $f_i$  in backward direction to obtain  $F_b$

$$F_b = \sum_{j=d-1}^0 f_i e^{-i2\pi kj/n}, \quad k = 0, \dots, d-1$$

3. Compute the sum of  $F_f$  and  $F_b$  to obtain  $F$ .

$$F = F_f + F_b$$

4. Calculate the real part of  $F$

$$F' = \text{real}(F)$$

5. Calculate the mean frequency component  $f_m$  from  $F'$  without considering the DC component ( $F'_0$ )

$$f_m = \frac{1}{d} \sum_{i=1}^{d-1} F'_i$$

6. Convert the low-n frequencies in  $F'$  to binary codes  $H$  using the  $f_m$  as a threshold

$$H = \begin{cases} 1, & F'_i > f_m \\ 0, & \text{Otherwise} \end{cases}$$

7. Return the n-bit binary code  $H$ .
- 

codes as illustrated in Algorithm 2. Frequencies that are less 325  
 than certain threshold are converted to zero bits and the rest are 326  
 converted to ones. Though some information is lost during this 327  
 conversion, the main gist of the spectrum is somehow retained 328  
 which leads to high performing binary codes. Since each neuron 329  
 represent a semantic concept (such as object part), a sufficiently 330

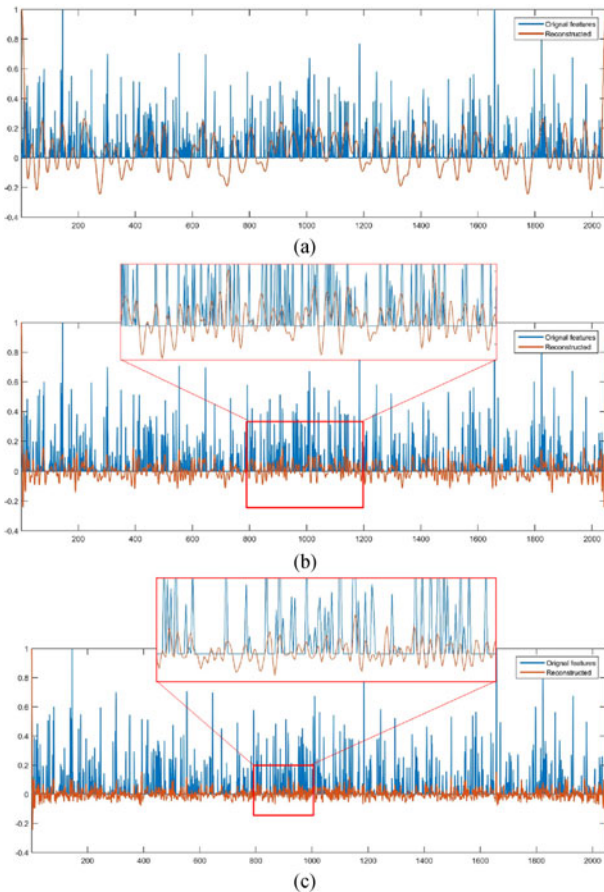


Fig. 2. Reconstruction of features from (a) 64 bits, (b) 256 bits, and (c) 512 bit hash codes generated using BD-FFT.

331 strong activation usually refer to the presence of that object part.  
 332 With such high level representation, if the reconstructed signal  
 333 adequately identify the high activation neurons, the code will be  
 334 an effective representation of the original features. The procedure  
 335 for conversion of deep features to binary codes is provided in  
 336 Algorithm 2.

#### 337 D. Bidirectional Fourier Decomposition

338 Though the simple FFT based binary conversion yield strong  
 339 representative codes [39], their quality can be further improved  
 340 with bidirectional FFT. In this case, we compute FFT of the  
 341 features in both forward and backward directions and then add  
 342 the corresponding frequency spectra. The dc component is ig-  
 343 nored and the subsequent  $n$  frequency components are binarized  
 344 to obtain the  $n$ -bit binary codes. Since the deep features are not  
 345 time-dependent, the bidirectional FFT actually helps capture the  
 346 patterns in neuronal activations more effectively, thereby yield-  
 347 ing better codes. Experimental results revealed that the BD-FFT  
 348 based codes perform much better than the regular FFT based  
 349 codes as reported in the experiments section.

#### 350 E. Locality Sensitivity of the Binary Codes

351 In LSH, the distance between the original features must corre-  
 352 late with the distance between the computed binary codes.  
 353 To evaluate locality sensitivity of the proposed binary codes,

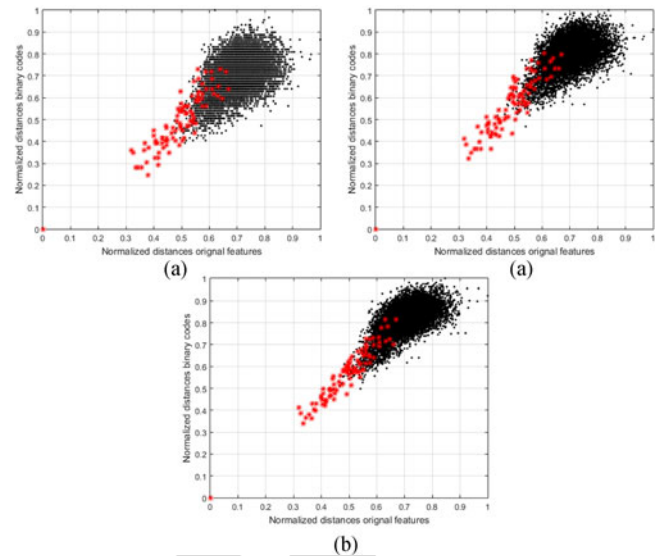


Fig. 3. Locality sensitivity of the proposed binary codes (a) 128 bits, (b) 256 bits, and (c) 512 bits.

we compared the normalized distances between deep features and their corresponding binary codes. Fig. 3 (a)– (c) reports the correlation among the distances between deep features and their corresponding binary codes. The distances of the query image with the rest of the images are shown on the  $x$ - and  $y$ -axis using deep features and binary codes, respectively. The red dots correspond to the relevant images and the black dots represent the irrelevant images in the dataset. Visualization of the distances reveal that the binary codes strongly correlate with the original deep features, especially for 256 and 512 bit codes, achieving correlation scores of 0.8975 and 0.9447, respectively. Increase in the distance between the original features is appropriately reflected by the distance between the binary codes. The relevant images have relatively smaller distances than the irrelevant ones which shows that those images will be retrieved at higher ranks. This characteristic of the proposed binary codes will help it achieve almost similar performance as the deep features.

## IV. EXPERIMENTS AND RESULTS

In this section, we present a detailed evaluation of the proposed method on a number of datasets used for benchmarking image retrieval methods. Different experiments were designed to measure performance of the proposed scheme and the effects of deep feature selection. All the experiments were carried out in MATLAB [40] environment on a Windows 7 PC equipped with 16 GB RAM. All the hashing methods were implemented and evaluated in MATLAB.

### A. Datasets

A number of datasets have been used to evaluate retrieval performance of the proposed method, including Corel-10 K, Holiday [41], IRMA-2009 [42], vehicle reidentification (VeRI) dataset [43], and stanford online products (SOP) dataset [44]. Each of these datasets contain thousands of images and are widely used to benchmark CBIR systems. Corel-10 K and Holiday datasets consist a variety of natural images whereas

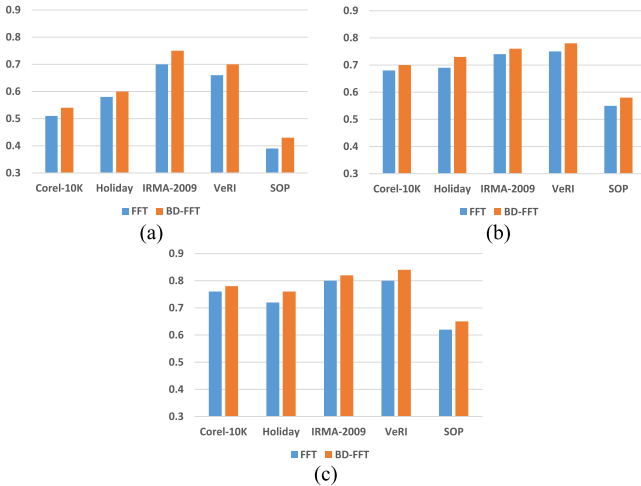


Fig. 4. Retrieval performance comparison FFT and BD-FFT based hash codes for (a) 128-bit, (b) 256-bit, and (c) 512-bit hash codes.

388 IRMA-2009, VeRI, and SOP contain images of particular cat-  
 389 egories including medical radiographs, vehicles, and products,  
 390 respectively.

### 391 B. Retrieval Performance of FFT Versus BD-FFT

392 A bidirectional Fourier decomposition of the feature vector  
 393 allowed us to capture patterns in the neuronal activations in a  
 394 much better way. Each bit in the hash code indicate either the  
 395 presence (1-valued bits) or absence (0 valued bits) of activa-  
 396 tion pattern in the original features. With BD-FFT, certain high  
 397 frequency patterns are captured in a much better manner than  
 398 the regular FFT based codes which leads to its superior perfor-  
 399 mance as reported in Fig. 4. The precision scores for various  
 400 datasets have been computed at recall = 0.2. The results reveal  
 401 that BD-FFT yield 3% to 10% better performance in terms of  
 402 precision scores as compared to FFT for all datasets at different  
 403 code lengths.

### 404 C. Retrieval Performance With Hash Codes Using 405 Different Subsets of Deep Features

406 In these experiments, we evaluated retrieval performance us-  
 407 ing hash codes of different lengths, computed from different  
 408 subsets of deep features. Hash codes of 128, 256, and 512 bits  
 409 were generated for five different sets of features, which con-  
 410 tained 4096, 1816, 1366, 820, and 585 neuronal activations.  
 411 These subsets were obtained by varying the threshold values in  
 412 Algorithm 1. Several images were selected at random from each  
 413 dataset and top ranked images were retrieved using hamming  
 414 distance between the query code and codes in the database. The  
 415 commonly used metrics including precision and recall were used  
 416 to report retrieval performance for each dataset. Fig. 5 shows  
 417 retrieval results in Corel-10 K dataset with 128, 256, and 512  
 418 bit codes for five different subsets of features. For each subset  
 419 of features, the precision-recall curves are presented for hash  
 420 codes of different lengths. In all of these results, the subset with  
 421 1816 activations yield better performance than the other subsets,  
 422 even the full-feature set. The margin is clearly visible in 128-bit  
 423 codes and gradually reduces for 256 and 512 bit hash codes, yet

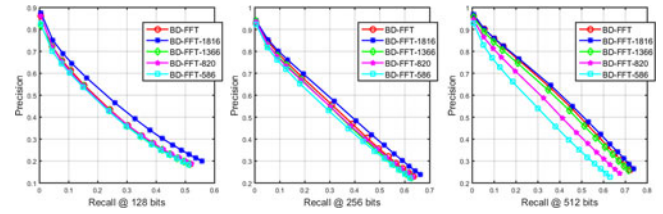


Fig. 5. Retrieval performance with hash codes generated from varying subsets of deep features for Corel-10 K dataset.

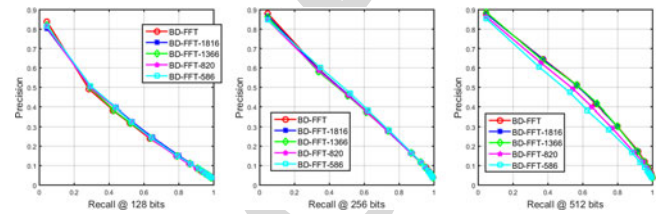


Fig. 6. Retrieval performance with hash codes generated from varying subsets of deep features for holiday dataset.

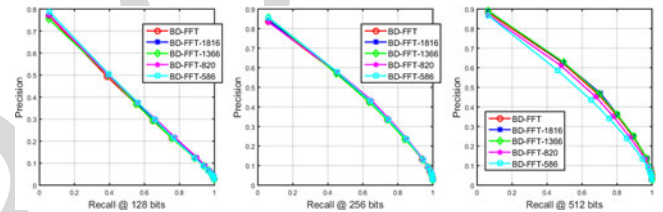


Fig. 7. Retrieval performance with hash codes generated from varying subsets of deep features for IRMA-2009 dataset.

it performs better than the other sets of features. Interestingly, 424  
 the performance of other reduced feature sets remains almost 425  
 the same as the full feature set, especially at 128 and 256 bit 426  
 codes. However, the 820 and 586 dimensional features failed to 427  
 catchup to the performance with other subsets in 512 bit codes. 428  
 It is important to observe here that performance remains almost 429  
 unchanged even if significant number of neuronal activations are 430  
 dropped. In 512 bit code, the scores for 4096, 1816, and 1366 431  
 features are almost the same. These results reveal the redundant 432  
 nature of deep features extracted from the FC layer. 433

The same experiments were carried out for Holiday image 434  
 datasets and the results presented in Fig. 6 reveal similar results 435  
 as compared to Corel-10 K. Features with 820 and 586 scores 436  
 slightly lower at 128 bits than the other subsets. However, the 437  
 performance with 4096, 1816 and 1366 features remains the 438  
 same for all hash codes. Though we did notice slightly better 439  
 performance at low recall for 1816 and 1366 subsets, the reduced 440  
 feature set performed almost the same as the full feature set. The 441  
 same results were observed with IRMA-2009 dataset as shown 442  
 in Fig. 7, where the reduced feature sets perform slightly better 443  
 at low recall and yield similar performance to the full feature 444  
 set for the rest of recall values with 128 and 256 bit codes. 445  
 However, with 512 bits, 1816-d, and 1366-d features achieve 446  
 better precision than the full feature set at all recall settings. 447

The VeRI dataset is quite challenging due to its large 448  
 volume and diversity. Carefully chosen subsets of features 449  
 either perform better than the full feature set or yield identical 450

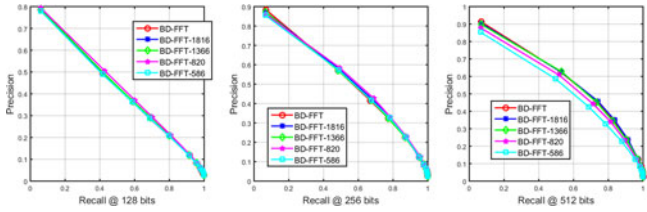


Fig. 8. Retrieval performance with hash codes generated from varying subsets of deep features for VeRI dataset.

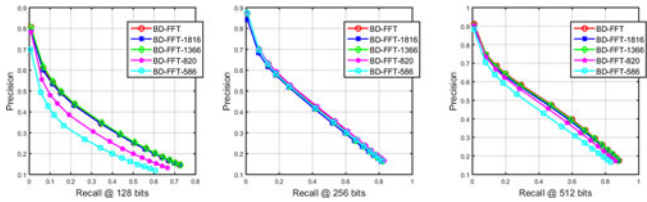


Fig. 9. Retrieval performance with hash codes generated from varying subsets of deep features for SOP dataset.

451 performance. In this dataset, we observed similar performance  
 452 for all subsets with 128 and 256 bit codes. With 512 bit codes,  
 453 820-d, and 586-d features scored slightly lower precision at all  
 454 recall settings as shown in Fig. 8. Finally, same experiments  
 455 were run for the SOP dataset which is the most challenging  
 456 dataset with huge volume and large number of product  
 457 categories. Precision scores dropped significantly when recall  
 458 rates are increased, particularly for 128 bit codes. At this length,  
 459 the hash codes generated for 4096, 1816, and 1366 features  
 460 yield similar retrieval performance, whereas the other subsets  
 461 achieve very low precision scores. With 256 bit codes, all the  
 462 subsets achieve similar precision scores at all recall rates. At  
 463 512 bits, 1816-d, and 1366-d features score almost the same as  
 464 the 4096-d features as presented in Fig. 9.

465 With these results, we can conclude that the FC layer features  
 466 are highly redundant and can be substantially reduced without  
 467 any loss in performance. Even in some cases, may get improved  
 468 retrieval results. Through these experiments, we decided to uti-  
 469 lize the selected 1816 neuronal activations from the FC-7 layer  
 470 instead of the 4096 features to generate hash codes for efficient  
 471 image retrieval in large datasets.

#### 472 D. Retrieval Performance With State-of-the-Art Hashing 473 Schemes

474 In this section, we compare the retrieval performance of the  
 475 proposed hash codes with five other schemes including LSH  
 476 [25], [34], SH [28], PCAH [33], DSH [30], and SpH [29]. In  
 477 these experiments, query images were randomly chosen from  
 478 each dataset and top ranked images were retrieved using hash  
 479 codes of 128, 256, and 512 bits. Precision-recall scores are  
 480 reported for each experiment. Fig. 10 presents the retrieval per-  
 481 formance of various hashing methods for Corel-10 K dataset.  
 482 The proposed method performed better than LSH at 128 bits,  
 483 however, it achieved low precision scores compared to other  
 484 methods. At 256 bits, BD-FFT outperformed LSH and PCAH  
 485 at low recalls, and LSH, PCAH, and SH at high recall rates. At

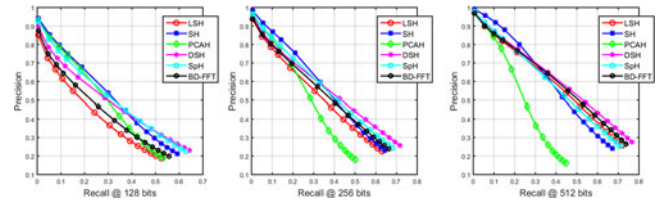


Fig. 10. Retrieval performance with hash codes compared with state-of-the-art methods for Corel-10 K dataset.

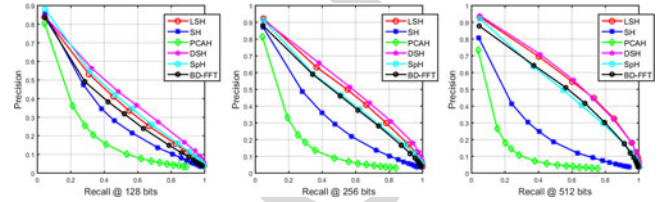


Fig. 11. Retrieval performance with hash codes compared with state-of-the-art methods for holiday dataset.

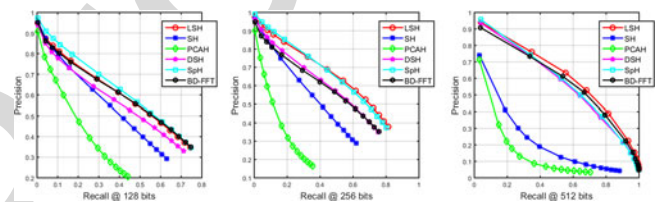


Fig. 12. Retrieval performance with hash codes compared with state-of-the-art methods for IRMA-2009 dataset.

low recall rates, BD-FFT performed similar to DSH. The per-  
 487 formance of BD-FFT improved significantly with 512 bit codes  
 488 where it outperformed LSH and PCAH at low recalls and LSH,  
 489 PCAH, SH, and SpH at all recall rates above 0.35.

490 In Holiday dataset, BD-FFT performed better than PCAH  
 491 and SH at 128 and 256 bit codes (see Fig. 11). At 512 bits,  
 492 it significantly outperformed PCAH, SH, and yielded slightly  
 493 better precision scores than SpH at most recall settings. How-  
 494 ever, the performance of LSH and DSH was relatively better  
 495 for this dataset. In IRMA-2009 dataset, BD-FFT yielded better  
 496 results than PCAH, SH, DSH, and LSH at 128 bit codes. Only  
 497 SpH performed slightly better than our method. With 256 bit  
 498 codes, BD-FFT scored better than PCAH and SH, however it  
 499 performed slightly poor than the rest of the methods. Increasing  
 500 the hash code length to 512 bits resulted in much better per-  
 501 formance of our method, surpassing SpH, SH, PCAH, and DSH  
 502 for recall rates above 0.4 as shown in Fig. 12.

503 In the VeRI dataset, BD-FFT significantly outperformed  
 504 PCAH, SH, and DSH in all experiments. With 512 bits, it per-  
 505 formed better than LSH at high recalls and reached the perfor-  
 506 mance of SpH (see Fig. 13). Similarly in SOP dataset, BD-FFT  
 507 outperformed PCAH and SH at 128, 256, and 512 bit codes.  
 508 However the other methods LSH, SpH, and DSH performed  
 509 much better at low recall rates as shown in Fig. 14. This is the  
 510 most challenging dataset and that is why its precision scores are  
 511 much lower than the other datasets.

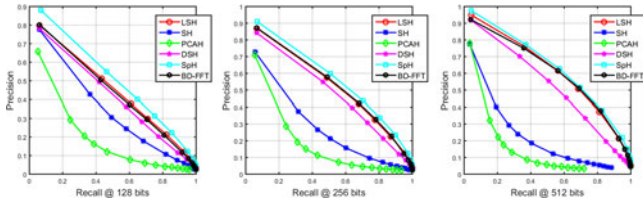


Fig. 13. Retrieval performance with hash codes compared with state-of-the-art methods for VeRI dataset.

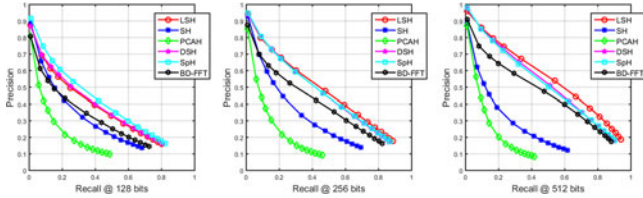


Fig. 14. Retrieval performance with hash codes compared with state-of-the-art methods for SOP dataset.

512 In most of the datasets, BD-FFT outperformed majority of the  
 513 methods, and achieved impressive performance especially with  
 514 256 and 512 bit hash codes. Moreover, the proposed method  
 515 yields more significant performance gains than the other com-  
 516 peting methods when size of the hash code increases. Keeping  
 517 in view the simplicity of our method, these results are very  
 518 promising. From these results, we can conclude that the pro-  
 519 posed method is capable of transforming high dimensional deep  
 520 features to compact binary codes of any length. We recommend  
 521 hash codes of length 256 or 512 bits to be used for image index-  
 522 ing and retrieval in large datasets. Though higher length codes  
 523 can also be generated in the same efficient manner, which may  
 524 yield performance improvements in most cases.

### 525 E. Qualitative Retrieval Performance Using the 526 Proposed Hash Codes

527 In this experiment, randomly chosen query images were used  
 528 to retrieve top-ranked images from each of the five datasets us-  
 529 ing hash codes generated with the proposed BD-FFT method  
 530 having 512-bit length. Results of two queries have been shown  
 531 for each dataset in terms of top 20 retrieved images in Fig. 15.  
 532 Results reveal that the proposed hash codes is capable of retriev-  
 533 ing relevant images at top ranks despite the huge volume and  
 534 diversity within these datasets, particularly IRMA-2009, Stan-  
 535 ford Online Products, and VeRI. The proposed hash codes can  
 536 effectively represent deep features, allowing almost the same  
 537 retrieval results as the raw features. The top two queries were  
 538 taken from Corel-10 K dataset where all relevant images have  
 539 been retrieved at top ranks. The next two rows contain results  
 540 from Holiday dataset where the first query image had three  
 541 other relevant images in the dataset, which have been success-  
 542 fully retrieved at top ranks. It is important to note here, that the  
 543 rest of the images, though irrelevant, resemble the query image  
 544 in visual appearance. Similar is the case with the other query  
 545 where the images at ranks 1, 2, 3, and 5, have been correctly  
 546 retrieved. The other images are also visually similar to the query  
 547 image. In the third pair of queries, visually similar images have

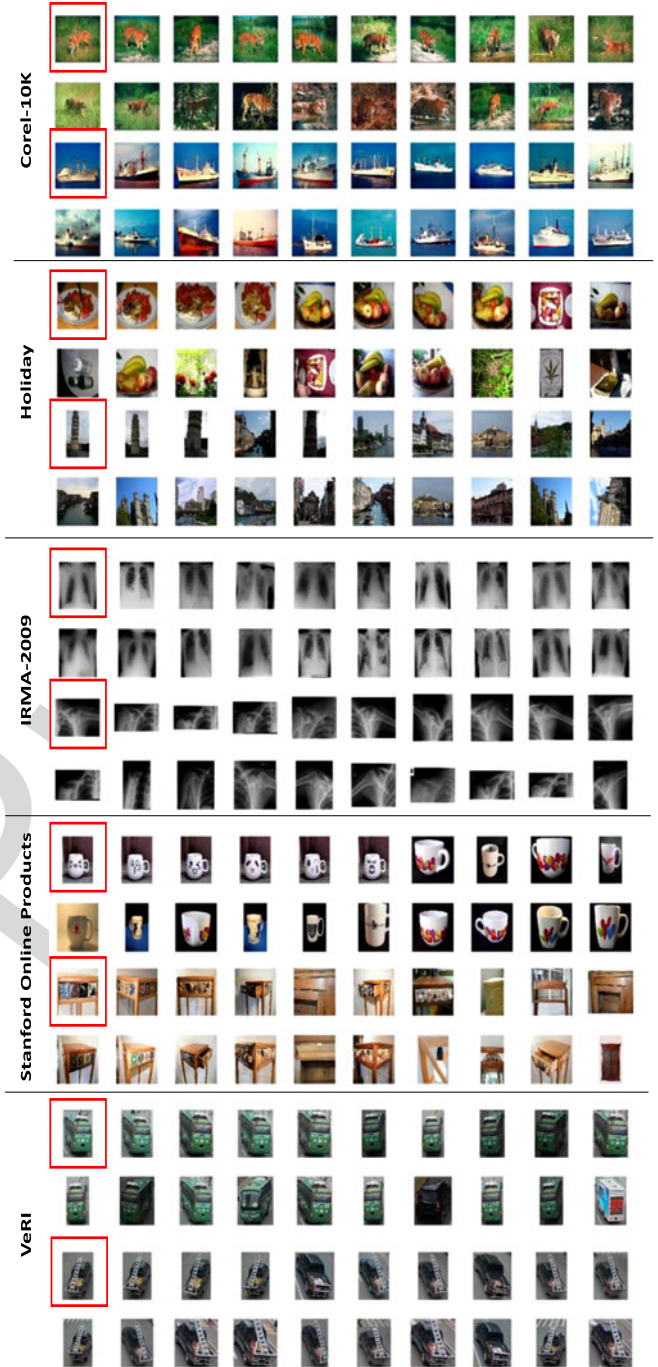


Fig. 15. Retrieval results using BD-FFT based 512-bit hash codes.

TABLE I  
 TRAINING TIME REQUIRED (IN SECONDS) FOR VARIOUS HASHING METHODS

Method	Training Time (20000 × 4096 features) 512-bits
LSH	0.03
SH	20.6
PCAH	19.7
DSH	30.2
SpH	252.1
BD-FFT	0.00



TABLE II  
TIME REQUIRED (IN SECONDS) FOR TRANSFORMING FEATURES TO HASH CODES USING VARIOUS METHODS

Method	Feature Size								
	10000 × 4096			20000 × 4096			200000 × 4096		
	128-bit	256-bit	512-bit	128-bit	256-bit	512-bit	128-bit	256-bit	512-bit
LSH	0.30	0.31	0.40	0.60	0.61	0.62	1.68	2.92	5.50
SH	1.10	4.47	16.18	2.20	8.32	33.3	24.76	84.86	341.9
PCAH	0.07	0.13	0.27	0.16	0.33	0.55	1.53	2.78	5.49
DSH	0.08	0.14	0.29	0.16	0.34	0.59	2.18	3.09	6.2
SpH	0.22	0.33	0.63	0.47	0.71	1.44	0.51	0.98	1.82
BD-FFT (CPU)		0.55			1.2			13.9	
BD-FFT (GPU)		<b>0.02</b>			<b>0.041</b>			<b>0.43</b>	

TABLE III  
STORAGE SPACE REQUIREMENTS FOR 1 MILLION IMAGES WITH DEEP FEATURES AND PROPOSED HASH CODES

Features	Storage required (MB)	Storage required (GB)	Retrieval performance % of original features
Raw (4096 deep features)	31250	30.51758	100
<b>512-bit</b>	<b>61.03516</b>	<b>0.059605</b>	<b>97.02</b>
<b>256-bit</b>	<b>30.51758</b>	<b>0.029802</b>	<b>92.09</b>
128-bit	15.25879	0.014901	86.10
64-bit	7.629395	0.007451	64.25
32-bit	3.814697	0.003725	40.31

548 been successfully retrieved at top ranks for both queries. The  
549 last two pairs of queries are from the most challenging datasets  
550 SOP, and VeRI. Despite the challenging nature and large size  
551 of these datasets, the proposed codes were able to retrieve the  
552 relevant images at top ranks. These results show the promis-  
553 ing performance of the proposed codes. With sufficiently sized  
554 codes, almost the same retrieval results can be achieved with the  
555 proposed method.

#### 556 F. Efficiency Analysis

557 In this section, we evaluate efficiency of the proposed scheme  
558 in terms of training time, hash code computation time, and  
559 storage requirements for the varying length hash codes. We aim  
560 to provide an insight into how efficient the proposed method  
561 is, compared to other similar approaches. In Table I, we listed  
562 the training times for various competing methods when 20 000  
563 features having 4096-dimensions were used for training the  
564 hashing functions. The training time mentioned in seconds, re-  
565 veal that the LSH method is the quickest to train and takes only  
566 0.03 s. The SH and PCAH methods take around 20 s, whereas,  
567 DSH require 30.2 s. The most computationally expensive  
568 method was found to be SpH which took 252.1 s to train for  
569 generating 512-bit hash codes. Though some of these methods  
570 are quite fast to train, they would require retraining when  
571 the hash code size gets changed. Further, the data-dependent  
572 methods like SH and SpH require to be trained each time  
573 when utilized for a different kind of dataset. Contrary to these  
574 methods, the proposed method do not require any training and  
575 can be used to directly transform deep features into binary hash  
576 codes of any length. Further, using specialized hardware (GPU),  
577 the proposed method can be executed in parallel, yielding very

high speeds for transforming features to hash codes. These  
characteristics make its implementation in real applications  
very easy. The proposed method can be easily implemented to  
transform the indexed features to binary codes which would  
allow efficiently locating similar images using ANN schemes.

Table II lists the hash code computation times for varying  
length codes using deep features. We used three test sets, having  
10 K, 20 K, and 200 K vectors of 4096-d to evaluate the con-  
version efficiency. Hash codes of 128, 256, and 512-bits were  
obtained using different hashing methods and the conversion  
times were recorded. The average conversion times reported in  
Table II reveal that majority of the methods including LSH,  
PCAH, DSH, and SpH are very efficient when shorter length  
hash codes are generated. The slowest method SH required  
1.10 s to convert 10 K features to 128-bit hash codes, however  
it took 341.9 s to convert 200 K features to 512-bit codes. In  
comparison, most of the hashing methods are more efficient  
than the proposed method on a CPU, which require 0.55, 1.2,  
and 13.9 s to convert 10 K, 20 K, and 200 K features into  
128, 256, and 512 bit hash codes, respectively. However, the  
advantage of the proposed method over other methods is that  
it can be easily computed on a GPU which yield significant  
gains in efficiency, reducing the computation times to 0.0002,  
0.041, and 0.43 s for 128, 256, and 512-bits, respectively.  
If the proposed method is implemented on a GPU, it can  
compute hash codes significantly faster than all the other  
competing methods. This characteristic also favors our method  
for implementation in practical applications.

In Table III, we show the amount of storage required for  
1 Million images when the raw features are stored to index  
images. We also show the amount of storage required to in-  
dex 1M images with 32, 64, 128, 256, and 512-bit codes. In

610 addition, we also report the relative image retrieval performance  
 611 to the original deep features for each code. With 32-bit codes, we  
 612 would require only 3.8 MB storage to index the images, however  
 613 we would only get 40.3% retrieval performance. Hash codes  
 614 greater than or equal to 128-bits, yield considerable retrieval  
 615 performance as well as saves storage space. The recommended  
 616 setting is to generate 256 or 512 bit codes for representing im-  
 617 ages because they would respectively yield 92% and 97% rela-  
 618 tive retrieval performance as compared to the original features.  
 619 Further, these hash codes reduce the storage requirements of the  
 620 index file from 30.5 GB to only 30 or 61 MB, which allow them  
 621 to be easily fit into memory. This would significantly improve  
 622 retrieval efficiency for large scale datasets.

## 623 V. CONCLUSION AND FUTURE WORK

624 In this paper, we presented an efficient method to directly  
 625 transform deep features into compact hash codes with locality  
 626 sensitivity property. These hash codes allow efficient retrieval  
 627 from large scale datasets utilizing ANN search procedures. The  
 628 proposed hash code conversion method require two steps. First,  
 629 salient deep features are selected using the proposed feature se-  
 630 lection algorithm, which analyzes the deep features and selects  
 631 features with higher diversity than a certain threshold. We an-  
 632 alyzed deep features and found that these features are highly  
 633 redundant and a significant number of these features can be  
 634 ignored without any loss in retrieval performance. Through ex-  
 635 periments, we determined 1816 features out of 4096 to represent  
 636 images. In the second step, we computed the FFT of these se-  
 637 lected features and binarized the top-n frequencies using mean  
 638 frequency as the threshold. The parameter n determined the  
 639 desired length of the hash code. The main idea behind the pro-  
 640 posed method is to represent the selected deep feature as a  
 641 signal and the FFT is used to approximate the feature vector  
 642 in the frequency domain. The computed hash codes have sig-  
 643 nificant representational capability with 128, 256, and 512 bit  
 644 codes, where the 512 bit codes yield almost the same retrieval  
 645 accuracy as the original deep features.

646 An essential characteristic of the proposed hashing method is  
 647 that it is completely data-independent and does not require any  
 648 training. Hash codes of any length can be directly computed very  
 649 efficiently. The implementation and operational simplicity of the  
 650 proposed scheme makes it very convenient to be implemented in  
 651 real-world applications. Further, GPU based acceleration of the  
 652 proposed method can substantially improve overall efficiency  
 653 of the retrieval system of large scale datasets. In this work, we  
 654 showed that the proposed method yield comparable performance  
 655 to the state-of-the-art for codes above 256 bits, however its  
 656 performance with smaller codes is relatively weak. Further, the  
 657 proposed method performs well for deep features, however, it  
 658 may not perform well for sparse features and further study is  
 659 needed to improve its performance for any type of features.

660 In future, we plan to study the effects of deep features on  
 661 its frequency spectrum and devise more effective ways of  
 662 capturing information in deep features into the compact binary  
 663 representations. Further, we will also evaluate wavelet based  
 664 methods to construct high performance short codes so that the  
 665 retrieval efficiency could be further enhanced.

## REFERENCES

- [1] P. Louridas and C. Ebert, "Embedded analytics and statistics for big data," *IEEE Softw.*, vol. 30, no. 6, pp. 33–39, Nov/Dec. 2013. 667
- [2] J. Lloret, M. Garcia, M. Atenas, and A. Canovas, "A QoE management system to improve the IPTV network," *Int. J. Commun. Syst.*, vol. 24, pp. 118–138, 2011. 668
- [3] B. Thomee *et al.*, "YFCC100M: The new data in multimedia research," *Commun. ACM*, vol. 59, pp. 64–73, 2016. 669
- [4] M. R. Robertson, "300+ hours of video uploaded to youtube every minute," *ReelSEO*, vol. 21, Nov. 2014. 670
- [5] J. Ahmad, M. Sajjad, S. Rho, and S. W. Baik, "Multi-scale local structure patterns histogram for describing visual contents in social image retrieval systems," *Multimedia Tools Appl.*, vol. 75, pp. 12669–12692, 2016. 671
- [6] J. Ahmad, M. Sajjad, I. Mehmood, S. Rho, and S. W. Baik, "Saliency-weighted graphs for efficient visual content description and their applications in real-time image retrieval systems," *J. Real-Time Image Process.*, vol. 13, pp. 431–447, Sep. 2017. 672
- [7] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 12, pp. 1349–1380, Dec. 2000. 673
- [8] M. Garcia, A. Canovas, M. Edo, and J. Lloret, "A QoE management system for ubiquitous IPTV devices," in *Proc. 3rd Int. Conf. Mobile Ubiquitous Comput., Syst., Service Technol.*, 2009, pp. 147–152. 674
- [9] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, pp. 91–110, 2004. 675
- [10] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2006, pp. 2169–2178. 676
- [11] H. Jégou, M. Douze, and C. Schmid, "Improving bag-of-features for large scale image search," *Int. J. Comput. Vis.*, vol. 87, pp. 316–336, 2010. 677
- [12] M. Douze, A. Ramisa, and C. Schmid, "Combining attributes and fisher vectors for efficient image retrieval," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2011, pp. 745–752. 678
- [13] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2010, pp. 3304–3311. 679
- [14] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, pp. 145–175, 2001. 680
- [15] J. Wu and J. M. Rehg, "CENTRIST: A visual descriptor for scene categorization," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 8, pp. 1489–1501, Aug. 2011. 681
- [16] J. Ahmad, M. Sajjad, I. Mehmood, and S. W. Baik, "SiNC: Saliency-injected neural codes for representation and efficient retrieval of medical radiographs," *PLoS One*, vol. 12, 2017, Art. no. e0181707. 682
- [17] A. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features off-the-shelf: An astounding baseline for recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2014, pp. 806–813. 683
- [18] H. Azizpour, A. Razavian, J. Sullivan, A. Maki, and S. Carlsson, "From generic to specific deep representations for visual recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 36–45. 684
- [19] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 580–587. 685
- [20] J. Ahmad, K. Muhammad, and S. W. Baik, "Data augmentation-assisted deep learning of hand-drawn partially colored sketches for visual search," *PLoS One*, vol. 12, 2017, Art. no. e0183838. 686
- [21] A. Babenko, A. Slesarev, A. Chigorin, and V. Lempitsky, "Neural codes for image retrieval," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 584–599. 687
- [22] A. Babenko and V. Lempitsky, "Aggregating local deep features for image retrieval," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1269–1277. 688
- [23] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. 25th Int. Conf. Inf. Process. Syst.*, 2012, pp. 1097–1105. 689
- [24] J. Wang, W. Liu, S. Kumar, and S.-F. Chang, "Learning to hash for indexing big data—A survey," *Proc. IEEE*, vol. 104, no. 1, pp. 34–57, Jan. 2016. 690
- [25] A. Gionis, P. Indyk, and R. Motwani, "Similarity search in high dimensions via hashing," in *Proc. 25th Int. Conf. Very Large Data Bases*, 1999, pp. 518–529. 691
- [26] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, "Locality-sensitive hashing scheme based on p-stable distributions," in *Proc. 20th Annu. Symp. Comput. Geom.*, 2004, pp. 253–262. 692
- [27] X. Yu, S. Zhang, B. Liu, L. Zhong, and D. Metaxas, "Large scale medical image search via unsupervised PCA hashing," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2013, pp. 393–398. 693

666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741

- [28] Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, pp. 1753–1760.
- [29] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, and S.-E. Yoon, "Spherical hashing," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2012, pp. 2957–2964.
- [30] Z. Jin, C. Li, Y. Lin, and D. Cai, "Density sensitive hashing," *IEEE Trans. Cybern.*, vol. 44, no. 8, pp. 1362–1371, Aug. 2014.
- [31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [32] K. Lin, H.-F. Yang, J.-H. Hsiao, and C.-S. Chen, "Deep learning of binary hash codes for fast image retrieval," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 27–35.
- [33] X.-J. Wang, L. Zhang, F. Jing, and W.-Y. Ma, "Annosearch: Image auto-annotation by search," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2006, pp. 1483–1490.
- [34] M. Slaney and M. Casey, "Locality-sensitive hashing for finding nearest neighbors [lecture notes]," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 128–131, Mar. 2008.
- [35] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Representations*, 2015.
- [36] A. B. Valdez, M. H. Papesh, D. M. Treiman, K. A. Smith, S. D. Goldinger, and P. N. Steinmetz, "Distributed representation of visual objects by single neurons in the human brain," *J. Neurosci.*, vol. 35, pp. 5180–5186, 2015.
- [37] G. Chandrashekar and F. Sahin, "A survey on feature selection methods," *Comput. Elect. Eng.*, vol. 40, pp. 16–28, 2014.
- [38] A. Babenko, A. Slesarev, A. Chigorin, and V. Lempitsky, "Neural codes for image retrieval," in *Proc. 13th Eur. Conf. Comput. Vis.*, 2014, pp. 584–599.
- [39] J. Ahmad, K. Muhammad, and S. W. Baik, "Medical image retrieval with compact binary codes generated in frequency domain using highly reactive convolutional features," *J. Med. Syst.*, vol. 42, p. 24, Dec. 19, 2017.
- [40] "MATLAB," 2016a ed: MathWorks, 2016.
- [41] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *Proc. 10th Eur. Conf. Comput. Vis.*, 2008, pp. 304–317.
- [42] H. Müller *et al.*, "Overview of the CLEF 2009 medical image retrieval track," in *Proc. Workshop Cross-Lang. Eval. Forum Eur. Lang.*, 2009, pp. 72–84.
- [43] X. Liu, W. Liu, H. Ma, and H. Fu, "Large-scale vehicle re-identification in urban surveillance videos," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2016, pp. 1–6.
- [44] H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 4004–4012.



**Jamil Ahmad** (S'16) received the BCS degree in computer science from the University of Peshawar, Peshawar, Pakistan, in 2008 with distinction. He received the Master's degree in computer science with specialization in image processing from Islamia College, Peshawar, Pakistan, in 2014. He is currently working toward the Ph.D. degree in digital contents at Sejong University, Seoul, South Korea.

He is also a Regular Faculty Member with the Department of Computer Science, Islamia College Peshawar, Peshawar, Pakistan. He has published several articles in these areas in reputed journals, including *Journal of Real-Time Image Processing*, *Multimedia Tools and Applications*, *Journal of Visual Communication and Image Representation*, *PLOS One*, *Journal of Oral and Maxillofacial Surgery*, *Computers and Electrical Engineering*, *SpringerPlus*, *Journal of Sensors*, and *KSII Transactions on Internet and Information Systems*. He is also an Active Reviewer for *IET Image Processing*, *Engineering Applications of Artificial Intelligence*, *KSII Transactions on Internet and Information Systems*, *Multimedia Tools and Applications*, *IEEE TRANSACTIONS ON IMAGE PROCESSING*, and *IEEE TRANSACTIONS ON CYBERNETICS*. His research interests include deep learning, medical image analysis, content-based multimedia retrieval, and computer vision.

787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811



**Khan Muhammad** (S'16) received the Bachelor's degree in computer science from the Islamia College Peshawar, Peshawar, Pakistan, in 2014, with a focus on information security. He is currently working toward the M.S. degree leading to Ph.D. degree in digital contents at Sejong University, Seoul, South Korea.

Since 2015, he has been a Research Associate with the Intelligent Media Laboratory. He has authored more than 40 papers in peer-reviewed international journals and conferences,

such as *IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS*, *Future Generation Computer Systems*, *Neurocomputing*, the *IEEE ACCESS*, the *Journal of Medical Systems*, *Biomedical Signal Processing and Control*, *Multimedia Tools and Applications*, *Pervasive and Mobile Computing*, *SpringerPlus*, *KSII Transactions on Internet and Information Systems*, *MITA 2015*, *PlatCon 2016*, *FIT 2016*, and *ICNGC 2017*. His research interests include image and video processing, information security, image and video steganography, video summarization, diagnostic hysteroscopy, wireless capsule endoscopy, computer vision, deep learning, and video surveillance.

812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833



**Jaime Lloret** (M'07–SM'10) received the B.Sc. and M.Sc. degrees in physics from University of Valencia, Valencia, Spain, in 1997, the B.Sc. and M.Sc. degrees in electronic engineering from University of Valencia, Valencia, Spain, in 2003 and the Ph.D. (Dr. Ing.) degree in telecommunication engineering from the Polytechnic University of Valencia, Valencia, Spain, in 2006.

He is currently an Associate Professor with the Polytechnic University of Valencia, Valencia, Spain. He is the Chair of the Integrated Management Coastal Research Institute and the Head of the "Active and collaborative techniques and use of technologic resources in the education" Innovation Group. He has authored 22 book chapters and has more than 400 research papers published in national and international conferences and journals. He is an Editor-in-Chief of the "Ad Hoc and Sensor Wireless Networks" (with ISI Thomson Impact Factor) and "Networks Protocols and Algorithms". He led many national and international projects. He is currently the Chair of the working group of the Standard IEEE 1907.1. He has been General Chair (or Co-Chair) of 38 International workshops and conferences. He is an IARIA Fellow.

834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855



**Sung Wook Baik** (M'16) received the B.S. degree in computer science from Seoul National University, Seoul, South Korea, in 1987, the M.S. degree in computer science from Northern Illinois University, DeKalb, DeKalb, IL, USA, in 1992, and the Ph.D. degree in information technology engineering from George Mason University, Fairfax, VA, USA, in 1999.

From 1997 to 2002, he was a Senior Scientist of the Intelligent Systems Group with Datamat Systems Research Inc. In 2002, he joined the

faculty of the College of Electronics and Information Engineering, Sejong University, Seoul, South Korea, where he is currently a Full Professor and Dean of Digital Contents. He is also the Head of Intelligent Media Laboratory, Sejong University, Seoul, South Korea. His research interests include computer vision, multimedia, pattern recognition, machine learning, data mining, virtual reality, and computer games.

856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873