

Quadra-Embedding: Binary Code Embedding with Low Quantization Error

Youngwoon Lee, Jae-Pil Heo, and Sung-Eui Yoon

KAIST

Abstract. Thanks to compact data representations and fast similarity computation, many binary code embedding techniques have been recently proposed for large-scale similarity search used in many computer vision applications including image retrieval. Most of prior techniques have centered around optimizing a set of projections for accurate embedding. In spite of active research efforts, existing solutions suffer both from diminishing marginal efficiency as more code bits are used, and high quantization errors naturally coming from the binarization.

In order to reduce both quantization error and diminishing efficiency we propose a novel binary code embedding scheme, *Quadra-Embedding*, that assigns two bits for each projection to define four quantization regions, and a novel binary code distance function tailored specifically to our encoding scheme. Our method is directly applicable to a wide variety of binary code embedding methods. Our scheme combined with four state-of-the-art embedding methods has been evaluated with three public image benchmarks. We have observed that our scheme achieves meaningful accuracy improvement in most experimental configurations under k - and ϵ -NN search.

1 Introduction

Scalable and efficient similarity search plays a key role in many large-scale computer vision applications dealing with high-dimensional data space. One example is the web-scale image retrieval. Common image descriptors, e.g., Bag-of-visual-Words or GIST, used for image retrieval have hundreds or thousands of dimensionality, and there are billions of images available on the web.

Traditional solutions adopting hierarchical structures (e.g., kd-trees) [1, 2] do not provide sufficient scalability in terms of both computational time and storage costs for high-dimensional, large-scale data sets. Recently, embedding high-dimensional data to short binary codes has been recognized as one of the most promising approaches to address both high-dimensionality and large-scaleness of data, since it can provide a compact representation of data and efficient similarity search. Consequently, a lot of binary code embedding algorithms have been studied lately [3–12]. The core goal of binary code embedding methods is to preserve similarities among original high-dimensional data in the corresponding binary codes, *i.e.* neighbor data points in the original high-dimensional space should be encoded to similar binary codes.

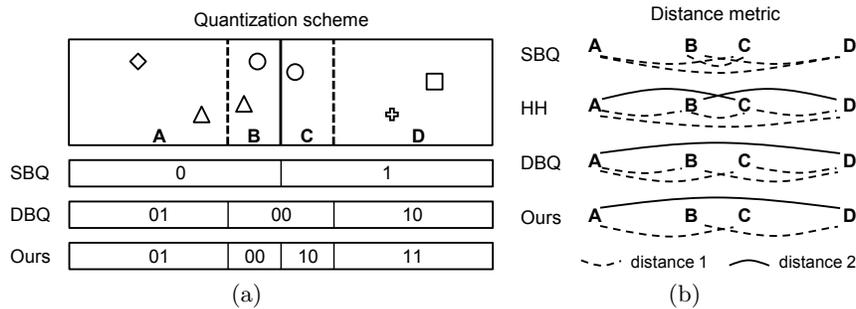


Fig. 1. (a) Different quantization schemes for a data set shown in the top row given a projection. The solid vertical line is a hyperplane associated with the projection, while two dashed lines are additional partitioning planes, i.e. offset surfaces used for defining the buffer area. *SBQ* and *DBQ* indicate the Single-Bit Quantization and Double-Bit Quantization [13], respectively. (b) Different distance metrics used in different encoding schemes. Any pair that does not have distance of 1 or 2 has zero distance. *HH* is Hierarchical Hashing [11], which uses the same encoding scheme with *Ours*.

Most binary embedding methods compute a binary code of each data element using multiple projections, which can be categorized into two groups: data-independent and data-dependent methods. Data-independent methods construct the projections based on vectors randomly drawn from some specific distributions. The well-known work in this category is locality-sensitive hashing (LSH) [3]. LSH is extended to various similarity metrics [4, 6, 7]. However, recent researches give more attentions to data-dependent methods for designing more data-friendly projections in order to achieve the higher accuracy. Notable examples include spectral hashing [5], sequential projection [8], joint optimization [9], and iterative quantization [10].

Despite the intensive research efforts to obtain effective projections, there are still remaining issues; 1) diminishing returns as the number of projections increases, and 2) quantization errors in high-density regions. The main cause of the diminishing return is the growing difficulty of defining both independent and informative set of projections as the number of projections increases. In regard to the quantization error, quantization boundaries are usually within dense regions, which is causing neighbor data points can be assigned to different binary codes [13]. In this paper we aim to resolve these two issues. More precisely, our contributions are:

- For each projection, we assign two bits to define four quantization regions as shown in Fig. 1(a) instead of the conventional binary regions based on a single bit (Sec. 3.2). In addition we propose a novel distance measure between two binary codes tailored to our binary code embedding scheme (Sec. 3.3).
- We formulate an optimization problem to decide partitioning boundaries of four regions suitable for our distance metric by minimizing the quantization error (Sec. 3.4).

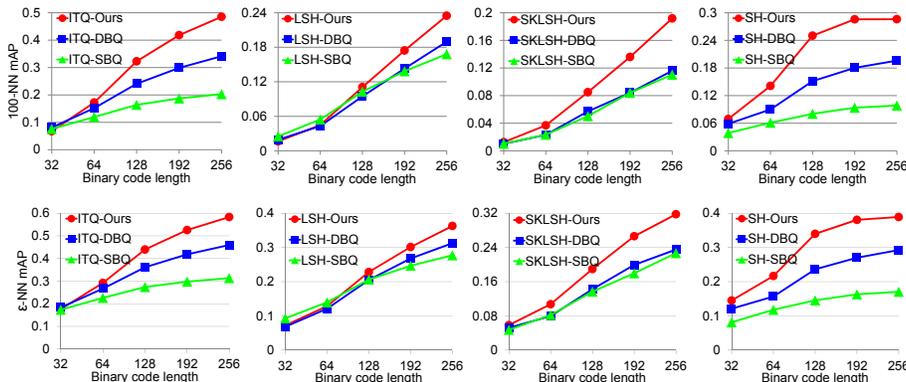


Fig. 2. Results on *CIFAR-60K-512D* with k -NN (top) and ϵ -NN search (bottom). Our method improves accuracy over *SBQ* and *DBQ* with different hashing methods. Our method shows the best results when combined with *ITQ* [10].

According to our best knowledge, only two existing approaches, Hierarchical Hashing (HH) [11] and Double-Bit Quantization (DBQ) [13] are aimed for similar goals that our method strives for. HH allows each projection to have four quantization states, but used the common Hamming distance that does not exploit all the benefits of having four states. DBQ quantized projection values into three different states with two bits and used the Hamming distance (see Fig. 1). In contrast, our method fully utilizes four states that two bits can encode, and adopts a specialized distance metric that further lowers down the quantization error caused by having four regions. To demonstrate benefits of our method, we have tested our method in three well-known image retrieval benchmarks in the context of two different types of nearest neighbor search, k -NN and ϵ -NN, in Sec. 4. Our method achieves significant improvements in accuracy over other prior encoding schemes across different hashing methods including LSH [4], spectral hashing [5], shift-invariant kernel-based LSH [7], and iterative quantization [10].

2 Related Work

2.1 Image Descriptors

To compactly and robustly represent images, a lot of image descriptors have been developed. Some well-known examples include Bag-of-visual-Words (BoW) [14] and GIST [15]. Finding similar images is then reduced to nearest neighbor search among image descriptors. Since these image descriptors are defined as high-dimensional vectors (e.g., larger than a few hundreds), performing nearest neighbor search for such high-dimensional image descriptors in an efficient and effective manner is very challenging [3]. Conventional approaches such as using

hierarchical data structures (e.g., kd-trees) [16, 1] based on recursive space partitioning have been known to be inefficient for such high-dimensional data in terms of search time and storage costs.

2.2 Binary Code Embedding Methods

To overcome the difficulty of handling large-scale and high-dimensional data sets, a significant amount of researches has been put on embedding high-dimensional data into compact binary codes preserving similarity among original data.

One of the most popular embedding methods is Locality-Sensitive Hashing (LSH) [3], which uses random projections drawn from a specific distribution. Many extensions of LSH have been proposed such as LSH with p -stable distributions [4], kernelized LSH [6], shift-invariant kernel-based LSH [7], etc.

Data-independent methods [4, 7] based on the random projections do not exploit data distributions. Departing from data-independent techniques, data-dependent approaches have been more widely studied recently because they can achieve higher accuracy given a fixed code length over data-independent techniques. Spectral hashing [5] derived projection directions based on spectral graph partitioning. He et al. [9] introduced a hashing technique that jointly optimizes both search accuracy and time. Heo et al. [12] proposed spherical hashing that partitions data using hyperspheres instead of commonly adopted hyperplanes. These data-dependent techniques share similar optimization goals such as balanced partitioning for each hashing function and independence between hashing functions.

In addition to considering aforementioned optimization goals there are some efforts to reduce the quantization error of a binary hashing function. Wang et al. [8] sequentially constructed each hyperplane to reduce the quantization error. Gong et al. [10] computed hyperplanes based on principal component analysis and then rotated them to minimize the quantization error iteratively. Joly et al. [17] used hyperplanes computed by support vector machine that produces maximum margins between data points and hyperplanes.

2.3 Encoding Schemes with Low Quantization Error

All the prior binary code embedding techniques mentioned in the last section aim to optimize projections for achieving the similarity-preserving property. On the other hand, our approach targets for reducing the quantization error and thus maximizing the discriminative power of given a set of projections.

There have been only a few research efforts on designing encoding schemes that can reduce the quantization error. Liu et al. [11] interpreted an approximate nearest neighbor structure with an anchor graph and then applied the graph Laplacian technique to the graph, in order to construct discriminative projections. In order to attain higher accuracy, Liu et al. [11] proposed Hierarchical Hashing (HH) that uses an additional hash bit for each projection that minimizes the cut value of the graph Laplacian. Although HH achieved significant improvement for Anchor Graph Hashing (AGH) [11], it is not straightforward

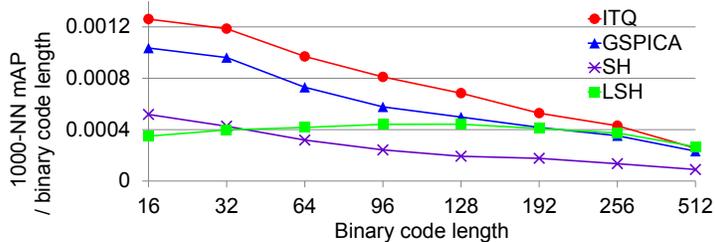


Fig. 3. This figure shows the diminishing return of having more projections for hashing. Y-axis is the mAP (mean Average Precision) of nearest neighbor search divided by the number of projections. In the case of data-dependent methods (ITQ [10], GSPICA [9], and spectral hashing (SH) [5]), mAP values per each hash bit consistently decrease as the total number of projections increases. However a data-independent technique, LSH, shows the diminishing return after 128 bit code lengths.

to apply its encoding scheme to other hashing methods. On the other hand, our technique is orthogonal and can be used together with all the prior embedding methods mentioned in the last section.

Kong and Li [13] will present a double-bit quantization, a similar strategy to our method for reducing the quantization error. This technique quantizes projection values into three states with two bits to inform whether a point is near a partitioning hyperplane or not. This technique uses the common Hamming distance for their encoding scheme. It was applied to prior binary code embedding methods such as iterative quantization [10] and achieved meaningful improvement. Unlike this method, our approach utilizes four states, the maximum number of states that two bits can represent, and adopts a novel distance metric that can further reduce the quantization error caused by additional hashing boundaries. We will demonstrate that our method shows noticeable improvements over this technique in a wide variety of settings.

3 Our Approach

In this section we elaborate the motivation of our approach, followed by explaining components of our method.

3.1 Motivation

Most binary code embedding techniques do not differentiate whether data points are located closely or far away, when they have the same binary code. This phenomenon is prominent when we encode data points with short binary code sizes. As increasing their code sizes one can achieve higher accuracy. Nonetheless, the benefit realized by having more projections, resulting in allocating more bits, diminishes quickly in most prior binary code embedding techniques, as demonstrated in Fig. 3. This is mainly because it is difficult to decorrelate hashing functions and to find useful projections for hashing [11, 13].

Before presenting our method, we first illustrate the quantization error of projection-based binary code embedding methods in Fig 1. At a high level, the quantization error is caused by two opposing cases: nearby points with a high Hamming distance, and faraway points with a low Hamming distance. We name these two as *inter-quantization* and *intra-quantization* errors, respectively. As an example of the inter-quantization error, two nearby circular points in the top row of Fig. 1(a) are partitioned by a hyperplane shown in the solid line and thus have two different binary codes. Therefore, the Hamming distance between them is not zero, even though these two circular points are closely located. On the other hand, the right circular point is far from the rectangular point, but they are partitioned together in the right side of the hyperplane. This is an example of the intra-quantization error, since they have the same binary code and thus zero Hamming distance even though they are located faraway.

Although both cases generating the quantization error are important, we focus on minimizing the inter-quantization error caused by nearby points with high Hamming distances. This decision is based on two reasons: first, most hashing techniques compute a short list of candidate nearest neighbor points by identifying data points that have zero or low Hamming distance. Once two nearby points have a high Hamming distance, these nearby points may not be included in the short list, unless computing an excessively long candidate list. Second, we can effectively control intra-quantization error by culling out faraway points that have low Hamming distances from the query point by adopting more expensive distance computation or re-ranking methods on the short list.

3.2 Binary Code Embedding Function

In order to reduce the inter-quantization error, we introduce a buffer area along a hyperplane. The buffer area is defined by two offset surfaces that are constructed by offsetting the hyperplane (or hypersphere) with offset distance r , called *margin*. Two dashed lines shown in the top row of Fig. 1(a) are two offset surfaces defining the buffer area given a hyperplane denoted by the solid line.

Two offset surfaces with the given hyperplane define four disjoint regions. In order to encode where a data point x is located among these four regions, our Quadra-Embedding method allocates two binary bits, $h_1(x)$ and $h_2(x)$, per one projection $f(x)$. In the case of using a simple hyperplane for the projection $f(x)$ is defined by $w^T x$, where w is a normal of the hyperplane; projections based on kernel techniques [6, 7, 9] and hyperspheres [12] are defined similarly. Specifically, the first and second hash bits are defined as follows:

$$h_1(x) = \text{sign}(f(x)) \quad \text{and} \quad h_2(x) = \text{sign}(|f(x)| - r) \quad (1)$$

Instead of defining the buffer area with the margin r we can define four regions in a more general manner with three thresholds $\{t_1, t_2, t_3\}$ indicating positions of the left offset surface, the hyperplane, and the right offset surface respectively, as follows:

$$h_1(x) = \text{sign}(f(x) - t_2) \quad \text{and} \quad h_2(x) = \begin{cases} 0 & \text{if } t_1 \leq f(x) \leq t_3 \\ 1 & \text{otherwise} \end{cases} \quad (2)$$

Intuitively, the first hash bit $h_1(x)$ indicates which side of the hyperplane contains the data point in the same manner of hash bits used in most prior hashing techniques. On the other hand, the second hash bit $h_2(x)$ indicates whether the data point is inside the buffer area. For a given data point $x \in \mathbb{R}^d$, the binary code $X \in \{0, 1\}^m$ of a length m is defined by concatenating two vectors of the first and second hashing functions, $H_1(x)$ and $H_2(x)$, each of which is based on $m/2$ projections as follows:

$$X = (H_1(x), H_2(x)) = (h_1^1(x), \dots, h_1^{m/2}(x), h_2^1(x), \dots, h_2^{m/2}(x)) \quad (3)$$

3.3 Hamming Distance for Quadra-Embedding

In order to maximize the benefit of our quantization scheme using two hash bits, we propose a Quadra-Embedding Distance (QED) function tailored to our method. The QED between two binary codes, $X = (H_1(x), H_2(x)) = (X_1, X_2)$ and $Y = (H_1(y), H_2(y)) = (Y_1, Y_2)$, is defined as follows:

$$d_{QED}(X, Y) = 2|(X_1 \oplus Y_1) \wedge (X_2 \wedge Y_2)| + |(X_1 \oplus Y_1) \wedge (X_2 \oplus Y_2)| \quad (4)$$

Intuitively QED between X and Y measures how many regions we should cross to reach a region of the data point X (or Y) from the region of the query point Y (or X). Therefore QED imposes zero distance on neighboring regions. It results in low inter-quantization errors, since data points are not discriminated by a single boundary unlike most prior methods.

QED among four disjoint regions is shown in Fig. 1(b). For example, given the right circular point (having the binary code of 10) in the top row of Fig. 1(a), the QED against itself (10), the left circular point (00), the rectangular point (11), and the diamond point (01) are 0, 0, 0, and 1 respectively. Computing our distance function has only minor overhead compared to the Hamming distance as discussed in Sec 4.5.

3.4 Threshold Optimization Process

Given our encoding scheme and distance function, it is critical to optimize the margin of the buffer area or more generally three positional values $\{t_1, t_2, t_3\}$ of two offset surfaces and the hyperplane. As we have larger buffer area, the inter-quantization error given its projection boundary decreases. On the other hand, a large buffer area makes a hash bit underutilized in terms of discriminative power. Our goal of optimizing thresholds is minimizing quantization errors along boundaries while keeping sufficient discriminative power.

For t_2 we can simply use the original quantization boundary, t_2^o , of a projection computed by any hashing methods. However, we have found that we can achieve higher accuracy by mildly optimizing the value of t_2 that are not too far away from t_2^o . This is mainly because the original hashing method may not consider well to lower down the quantization error while constructing projections. Nonetheless, adjusting the original quantization boundary has a potential

risk to destruct benefits of the original hashing method. As a result, we adjust the value of t_2 such that it does not deviate much from t_2^o , while aggressively attempting to optimize the locations (t_1 and t_3) of two offset surfaces.

Let $T \subset \mathbb{R}^d$ be a training set containing n points and P be a set of projected data points, i.e. $P = \{p | p = f(x), x \in T\}$. Our optimization adjusts thresholds such that projected data points in each region get away from thresholds, in order to reduce the quantization error. Let four regions divided by three thresholds be P_1, P_2, P_3 , and P_4 , i.e. $P_1 = \{p \in P | p \leq t_1\}$, $P_2 = \{p \in P | t_1 < p \leq t_2\}$, $P_3 = \{p \in P | t_2 < p < t_3\}$, $P_4 = \{p \in P | t_3 \leq p\}$. We also define μ_1, μ_2, μ_3 , and μ_4 to be the mean values of four subsets P_1, P_2, P_3 and P_4 respectively. Our objective is finding thresholds that minimize the following penalty function:

$$J = \sum_{p \in P_1} \ell(p - \mu_1)^2 + \sum_{p \in P_2} \ell(\mu_2 - p)^2 + \sum_{p \in P_3} \ell(p - \mu_3)^2 + \sum_{p \in P_4} \ell(\mu_4 - p)^2, \quad (5)$$

where $\ell(\cdot)$ is a filtering function and set to be $\max(\cdot, 0)$. $\ell(\cdot)$ ignores negative inputs and thus serves as the following two purposes: 1) the filtering function used in the second and third terms ignores projected data points near the given hyperplane, in order to less aggressively adjust t_2 , and 2) the filtering function used in the first and fourth terms ignores faraway points that do not contribute much to the change of quantization error, but affect heavily the penalty function.

Values of each term of Eq. (5) can be precomputed for all the possible pairs of two neighboring thresholds in $O(n^2)$ with dynamic programming. We then find the optimal solution by exhaustively searching all the possible pairs of t_1 and t_3 . Note that as we incrementally change the value of t_1 or t_3 , the optimal t_2 can be computed in a constant time by looking into only a fixed number of potential candidates. This incremental computation makes our algorithm to run in $O(n^2)$. In practice our algorithm takes approximately 1.44 s on average for each projection to compute three thresholds with 20 k training data points. For example, it took less than one minute for 64 bit code lengths.

4 Results and Discussions

In this section we explain our experiment protocols and compare our method against the state-of-the-art techniques.

4.1 Datasets

We evaluated our method on the following three different image datasets:

- *CIFAR-60K-512D*: A set of 512 dimensional GIST descriptors [15] from the CIFAR-10 dataset [18] sampled from Tiny Images [19].
- *GIST-1M-960D*: A set of 960 dimensional, one million GIST descriptors that were used by Jégou et al. [20].
- *GIST-75M-384D*: A set of 384 dimensional GIST descriptors [15] of 75 million images from Tiny Images [19].

4.2 Tested Hashing Methods

In order to test our method, we compared our method against the following state-of-the-art hashing methods:

- *LSH*: Locality-Sensitive Hashing [4] with the normalized data to have the zero mean, as discussed in [10].
- *SH*: Spectral Hashing [5].
- *SKLSH*: Shift-invariant Kernel-based Locality-Sensitive Hashing [7] with the bandwidth parameter, which is the inverse of the mean distance [21].
- *ITQ*: ITERative Quantization [10] with the normalized data to have the zero mean.

Since *ITQ* and *SH* are based on principal component analysis, the number of hashing functions cannot exceed the dimensionality of the original data. We therefore measured the performance up to 256 bits.

4.3 Quantization Strategies

Given the hashing methods mentioned above, we tested different bit allocation methods as follows:

- *SBQ*: Single-Bit Quantization which is commonly used in most prior hashing methods.
- *DBQ*: Double-Bit Quantization proposed by Kong and Li [13].
- *Ours*: Our binary code embedding method, Quadra-Embedding.

We tested different quantization schemes with hashing functions computed by hashing methods listed in Sec. 4.2. We combined one of quantization strategies with one of hashing methods. To concisely represent them, we use a naming scheme like *LSH-SBQ*, which uses *SBQ* with *LSH*. As another example, *LSH-Ours* denotes our embedding scheme combined with *LSH*.

We did not compare our method against hierarchical hashing [11], mainly because hierarchical hashing is not naturally applicable to different hashing methods listed in Sec. 4.2, and its simplified version (using four regions with *DBQ*'s thresholds) has been shown to be consistently inferior than *DBQ* [13].

4.4 Protocols

We followed the evaluation protocols that have been widely used in recent papers [12, 17]. We trained hashing functions with randomly chosen 20 k training points on *CIFAR-60K-512D* and *GIST-1M-960D*, 100 k points on *GIST-75M-384D*. Then we randomly selected 10 k, 1 k, and 500 queries on *CIFAR-60K-512D*, *GIST-1M-960D*, and *GIST-75M-384D* respectively.

All the experimental results on *CIFAR-60K-512D* and *GIST-1M-960D* are averaged over five independent training times, and results on *GIST-75M-384D* are averaged over three independent training times because of its long computational time.

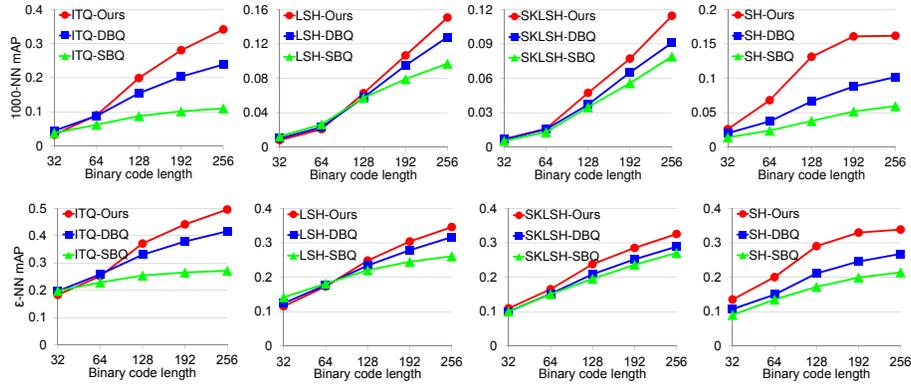


Fig. 4. Results on *GIST-1M-960D* with k -NN (top) and ε -NN search (bottom). Our method improves accuracy over *SBQ* and *DBQ* with different hashing methods. Our method shows the best results when combined with ITQ [10].

The performance is measured by the mean average precision (mAP), which is the average area under the recall-precision curves. The ground truths are computed based on the Euclidean distance by exhaustively testing all the data sets against query points. We adopted two major protocols for approximate nearest neighbor search, k -NN and ε -NN. The ground truths of a point x in k -NN and ε -NN are defined as follows:

$$g_{k\text{-NN}}(x, y) = \begin{cases} 1 & \text{if } y \in k\text{-NN}(x) \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad g_{\varepsilon\text{-NN}}(x, y) = \begin{cases} 1 & \text{if } d(x, y) < \varepsilon \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

In ε -NN, ε is determined by averaging distances of 50th nearest neighbors from query points, as used in [13]. In this setting, approximately 299 queries in *GIST-1M-960D* and 1731 queries in *CIFAR-60K-512D* have no nearest neighbors within ε and thus we ignore these queries for evaluation. For *GIST-75M-384D* we measure the performance only with k -NN since computing the ground truths of ε -NN takes long computational time. We compared different methods including ours given the same number of bits. Specifically *DBQ* and our method use only half of the number of hashing functions that *SBQ* uses, but allocate two bits for each hashing function.

4.5 Results

Figure 2 shows results of different methods for k -NN and ε -NN on the benchmark of *CIFAR-60K-512D*. Especially for the benchmark we used $k = 100$ and set ε to be the average distance of the 50th nearest neighbors, respectively; we have also tried different k and ε for k -NN and ε -NN, and observed similar results. Also, Table 1 shows mAP values at 64, 128, and 256 bits for two different types of nearest neighbor search.

Table 1. Results on *CIFAR-60K-512D* dataset.

100-NN mAP									
# bits	64			128			256		
	SBQ	DBQ	OURS	SBQ	DBQ	OURS	SBQ	DBQ	OURS
ITQ	0.1194	0.1516	0.1826	0.1639	0.2408	0.3228	0.2028	0.3411	0.4677
LSH	0.0545	0.0431	0.0452	0.1029	0.0945	0.1156	0.1680	0.1897	0.2317
SKLSH	0.0233	0.0231	0.0336	0.0503	0.0570	0.0875	0.1104	0.1159	0.1833
SH	0.0612	0.0898	0.1380	0.0804	0.1514	0.2194	0.0984	0.1961	0.2387

ϵ -NN mAP									
# bits	64			128			256		
	SBQ	DBQ	OURS	SBQ	DBQ	OURS	SBQ	DBQ	OURS
ITQ	0.2260	0.2668	0.2921	0.2738	0.3610	0.4395	0.3134	0.4597	0.5831
LSH	0.1396	0.1210	0.1285	0.2062	0.2050	0.2281	0.2771	0.3125	0.3627
SKLSH	0.0821	0.0797	0.1073	0.1366	0.1424	0.1895	0.2267	0.2351	0.3176
SH	0.1178	0.1566	0.2165	0.1454	0.2362	0.3397	0.1704	0.2925	0.3890

In both of k -NN and ϵ -NN our method shows better results, more than 100% improvement in some cases, over *SBQ* and *DBQ* under different hashing methods in most cases, especially for 64 or more bits code lengths. For example, our method achieved 42%, 24%, 65%, and 46% improvement over *DBQ* when using *ITQ*, *LSH*, *SKLSH*, and *SH* with 256 bit length for k -NN, respectively. Also, compared to *SBQ*, our method achieved 139%, 40%, 74%, and 191% improvement when using *ITQ*, *LSH*, *SKLSH*, and *SH* for the 256 bit length, respectively.

Even at a low bit codes (e.g., 32 bits) our method showed a bit lower or modest improvement in most cases. For example, at 64 bits our method showed improvement performance over *DBQ* in all the cases. However, our method showed a lower accuracy over *SBQ* only in a single case when combined with *LSH*. These results indicate that at low code lengths, more projections give higher discriminative power over allocating more bits to reduce the quantization error as we did in our method. On the contrary, as we use more bits to represent high dimensional data points, our method reduces the quantization errors and thus improve the overall discriminative power over using all the bits for having different projections.

While our method showed improved results across different hashing methods, it achieved the best results with *ITQ* across all the tested code lengths and different types of NN search. Since *ITQ* constructs hashing functions to minimize the quantization error, the goal of our method aligns best with *ITQ* and thus achieved the best results.

We also tested the methods with *GIST-1M-960D*, which is a much bigger dataset and has higher dimensionality than *CIFAR-60K-512D*. Nonetheless in this benchmark we observed similar results achieved with *CIFAR-60K-512D* (see Table 2).

Finally we tested *LSH* and *ITQ* in *GIST-75M-384D*, which is one of the largest image benchmarks available to computer vision community. Overall re-

Table 2. Results on *GIST-1M-960D* dataset.

1000-NN mAP									
# bits	64			128			256		
	SBQ	DBQ	OURS	SBQ	DBQ	OURS	SBQ	DBQ	OURS
ITQ	0.0620	0.0881	0.0899	0.0875	0.1545	0.1989	0.1101	0.2388	0.3415
LSH	0.0264	0.0230	0.0208	0.0573	0.0576	0.0626	0.0970	0.1279	0.1506
SKLSH	0.0129	0.0155	0.0160	0.0347	0.0370	0.0473	0.0790	0.0910	0.1145
SH	0.0239	0.0373	0.0681	0.0379	0.0665	0.1313	0.0595	0.1017	0.1618

ϵ -NN mAP									
# bits	64			128			256		
	SBQ	DBQ	OURS	SBQ	DBQ	OURS	SBQ	DBQ	OURS
ITQ	0.2289	0.2585	0.2520	0.2544	0.3304	0.3713	0.2715	0.4164	0.4958
LSH	0.1812	0.1774	0.1734	0.2206	0.2343	0.2484	0.2610	0.3167	0.3458
SKLSH	0.1508	0.1517	0.1651	0.1956	0.2088	0.2380	0.2706	0.2888	0.3257
SH	0.1353	0.1497	0.2003	0.1721	0.2112	0.2903	0.2143	0.2678	0.3385

sults (Fig. 6) on this benchmark are similar to those achieved in other benchmarks. In data-independent method *LSH*, our method showed comparable performance with *SBQ* and *DBQ* up to 128 bits, while outperforming them afterwards. With data-dependent method *ITQ*, our method showed better performance from 128 bits.

4.6 Discussions

Figure 5 shows additional gains brought by different components of our method. *Ours w/ DBQ-Op.* uses our encoding scheme and distance function, but with the threshold learning of *DBQ* [13], while *Ours w/ Our-Op.* uses our optimization process mentioned in Sec. 3.4. By using our encoding and distance function over *DBQ* we achieved 29% improvement at 128 bit length. In addition we achieve 5% additional gain by using our optimization method.

We also measured distance computation time of our proposed distance function QED. One million distance computations of QED for 256 bit code lengths took 8.3 ms on average, while the Hamming distance took 7.4 ms. Our distance function significantly improved the overall accuracy of nearest neighbor search with comparable computational time to the Hamming distance.

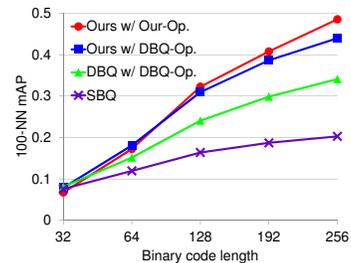


Fig. 5. This figure shows benefits of different components of our method on *CIFAR-60K-512D* with *k*-NN protocol and projections computed by *ITQ*.

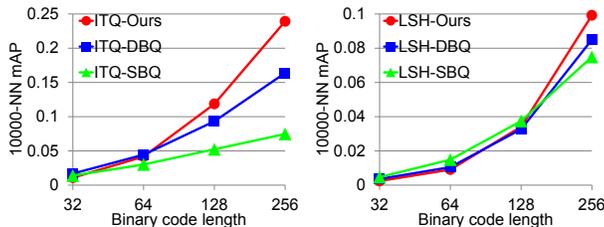


Fig. 6. Results on *GIST-75M-384D* with k -NN.

5 Conclusion and Future Work

We have introduced a novel binary code embedding method, Quadra-Embedding technique for efficient and effective nearest neighbor search for image retrieval. In order to reduce the quantization error, our method uses two different code bits for each projection and a specialized distance function tailored to our encoding scheme. We have explained an optimization method that adjusts the margin of the buffer area for our method. We have tested our method with two different types of nearest neighbor search and a diverse set of hashing methods under three different image benchmarks. We have observed that our method achieves significant improvement over prior quantization strategies such as single-bit and double-bit quantizations [13] in most experimental configurations. These results have demonstrated the usefulness and robustness of our approach.

Many interesting future research directions lie ahead. Our method showed superior results over *SBQ* in long code lengths, especially with data-dependent methods. In case of data-dependent methods, our method showed comparable accuracy even at 32 bit code lengths and outperformed *SBQ* afterwards. However, our method showed a bit lower or comparable performance over *SBQ* with *LSH* at 32 or 64 bit code lengths, while outperforming *SBQ* afterwards. These results indicate that it is better to use a single bit for each projection for short code lengths, i.e. 32 bit code lengths. Also, our approach can be easily extended to allocate more than two bits per one hashing function. The most important question that needs to be addressed is how to determine an appropriate number of bits for each projection. Furthermore, designing a distance metric for this encoding scheme would be very challenging and thus we may have to learn useful metric functions for maximizing the performance of this approach. Currently, we optimized our encoding scheme given a hashing method. As a next step, it would be very interesting to jointly optimize both encoding scheme and hashing functions. As demonstrated by the fact that our method achieved the best results with iterative quantization [10] that also aims to minimize the quantization error, this joint optimization can further improve the accuracy of nearest neighbor search.

Acknowledgement. This work was supported in part by MCST/KOCCA/CT/R&D 2011, MKE/KEIT [KI001810035261], MKE/MCST/IITA [2008-F-033-02], BK,

DAPA/ADD (UD110006MD), MEST/NRF/WCU (R31-2010-000-30007-0), KMCC, MSRA, and MEST/NRF (2011-0030822).

References

1. Nistér, D., Stewénus, H.: Scalable recognition with a vocabulary tree. In: CVPR. (2006)
2. Muja, M., Lowe, D.G.: Fast approximate nearest neighbors with automatic algorithm configuration. In: VISAPP. (2009) 331–340
3. Indyk, P., Motwani, R.: Approximate nearest neighbors: toward removing the curse of dimensionality. In: STOC. (1998)
4. Datar, M., Immorlica, N., Indyk, P., Mirrokni, V.S.: Locality-sensitive hashing scheme based on p-stable distributions. In: SoCG. (2004)
5. Weiss, Y., Torralba, A., Fergus, R.: Spectral hashing. In: NIPS. (2008)
6. Kulis, B., Grauman, K.: Kernelized locality-sensitive hashing for scalable image search. In: ICCV. (2009)
7. Raginsky, M., Lazebnik, S.: Locality sensitive binary codes from shift-invariant kernels. In: NIPS. (2009)
8. Wang, J., Kumar, S., Chang, S.F.: Sequential projection learning for hashing with compact codes. In: ICML. (2010)
9. He, J., Radhakrishnan, R., Chang, S.F., Bauer, C.: Compact hashing with joint optimization of search accuracy and time. In: CVPR. (2011)
10. Gong, Y., Lazebnik, S.: Iterative quantization: a procrustean approach to learning binary codes. In: CVPR. (2011)
11. Liu, W., Wang, J., Kumar, S., Chang, S.F.: Hashing with graphs. In: ICML. (2011)
12. Heo, J.P., Lee, Y., He, J., Chang, S.F., Yoon, S.E.: Spherical hashing. In: CVPR. (2012)
13. Kong, W., Li, W.J.: Double-bit quantization for hashing. In: AAAI. (2012)
14. Sivic, J., Zisserman, A.: Video google: A text retrieval approach to object matching in videos. In: ICCV. (2003)
15. Oliva, A., Torralba, A.: Modeling the shape of the scene: a holistic representation of the spatial envelope. IJCV (2001)
16. Friedman, J.H., Bentley, J.L., Finkel, R.A.: An algorithm for finding best matches in logarithmic expected time. ACM TOMS **3** (1977) 209–226
17. Joly, A., Buisson, O.: Random maximum margin hashing. In: CVPR. (2011)
18. Krizhevsky, A.: Learning multiple layers of features from tiny images. Technical report, University of Toronto (2009)
19. Torralba, A., Fergus, R., Weiss, Y.: Small codes and large image databases for recognition. In: CVPR. (2008)
20. Jégou, H., Douze, M., Schmid, C.: Product quantization for nearest neighbor search. IEEE TPAMI (2011)
21. Gordo, A., Perronnin, F.: Asymmetric distances for binary embeddings. In: CVPR. (2011)