



Thick boundaries in binary space and their influence on nearest-neighbor search

Tomasz Trzcinski*, Vincent Lepetit, Pascal Fua

Ecole Polytechnique Fédérale de Lausanne (EPFL), Computer Vision Laboratory, CH-1015 Lausanne, Switzerland

ARTICLE INFO

Article history:

Received 4 November 2011

Available online 29 August 2012

Communicated by G. Borgefors

Keywords:

Approximate nearest neighbor search

Binary vectors

Locality sensitive hashing

Hierarchical *k*-means

ABSTRACT

Binary descriptors allow faster similarity computation than real-valued ones while requiring much less storage. As a result, many algorithms have recently been proposed to binarize floating-point descriptors so that they can be searched for quickly. Unfortunately, even if the similarity between vectors can be computed fast, exhaustive linear search remains impractical for truly large databases and approximate nearest neighbor (ANN) search is still required. It is therefore surprising that relatively little attention has been paid to the efficiency of ANN algorithms on binary vectors and this is the focus of this paper.

We first show that binary-space Voronoi diagrams have thick boundaries, meaning that there are many points that lie at the same distance from two random points. This violates the implicit assumption made by most ANN algorithms that points can be neatly assigned to clusters centered around a set of cluster centers. As a result, state-of-the-art algorithms that can operate on binary vectors exhibit much lower performance than those that work with floating point ones.

The above analysis is the first contribution of the paper. The second one is two effective ways to overcome this limitation, by appropriately randomizing either a tree-based algorithm or hashing-based one. In both cases, we show that we obtain precision/recall curves that are similar to those that can be obtained using floating point number calculation, but at much reduced computational cost.

© 2012 Elsevier B.V. All rights reserved.

1. Introduction

The problem of matching high-dimensional descriptors against large databases is pervasive in Computer Vision, e.g., in image-retrieval, or pose-estimation. When there are millions of such descriptors, linear search becomes prohibitively expensive, even after dimensionality reduction (Mikolajczyk et al., 2004; Brown et al., 2011) and no generic, exact, and more efficient algorithm is known.

Approximate Nearest Neighbor (ANN) search constitutes one effective approach to overcoming this limitation and there are many algorithms that can handle real-valued descriptors such as the Scale Invariant Feature Transform (SIFT) (Lowe, 2004) or Speeded Up Robust Feature (SURF) (Bay et al., 2006) descriptors. These algorithms rely on modified kd-trees (Beis and Lowe, 1997; Arya et al., 1998), multiple randomized kd-trees (Silpa-Anan and Hartley, 2008), hierarchical *k*-means (HKM) trees (Fukunaga and Narendra, 1975; Nister and Stewenius, 2006), spill trees (Liu et al., 2004), vantage-point trees (Yianilos, 1993), or hashing functions (Andoni and Indyk, 2008). A different approach to speeding up nearest-neighbor search is to binarize the real-valued descriptors using techniques such as Boosting (Shakhnarovich, 2005),

hashing (Andoni and Indyk, 2008; Kulis and Darrell, 2009), Principal Component Analysis (PCA) or Linear Discriminant Analysis (LDA) based methods (Raginsky and Lazebnik, 2009; Strecha et al., 2012), quantization (Gong and Lazebnik, 2011) and Semantic or Spectral Hashing (Salakhutdinov and Hinton, 2009; Weiss et al., 2009). Because the similarity between the resulting binary vectors can be evaluated using the Hamming distance, which can be computed much faster than the Euclidean one on modern CPUs, linear search is more efficient but remains too slow for large-scale applications. In favorable cases, the binary vectors can be used as indices to directly access their nearest neighbors (Weiss et al., 2009) which provides sub-linear complexity of the search. Unfortunately, this stops being possible when the typical Hamming distance between nearest neighbors is larger than a few units.

To get the best of both worlds under general conditions and to exploit the potential of binary descriptors, ANN search is necessary. Little attention has been paid to the performance of ANN algorithms on binary, as opposed to real-valued, vectors. Some of the algorithms discussed above such as Spectral Hashing are not adapted to binary vectors because they involve a PCA decomposition. Other methods can be used by treating binary descriptors as vectors of zeros and ones encoded as floating-point numbers. Even with the same search-accuracy, this encoding negates the advantages of binary vectors over real-valued ones: their compactness and the fact that the Hamming distance can be computed faster than the Euclidean one. Finally, there are algorithms such as vantage-point

* Corresponding author. Tel.: +41 216936862.

E-mail addresses: tomasz.trzcinski@epfl.ch (T. Trzcinski), vincent.lepetit@epfl.ch (V. Lepetit), pascal.fua@epfl.ch (P. Fua).

trees and HKM that can be modified to only deal with binary vectors and use the Hamming distance as a similarity measure. However, as we show, their accuracy is much lower.

The first contribution of this paper is to show that this performance loss can be traced to the fact that in Hamming spaces, unlike in Euclidean ones, the number of points that lie at the same distance from two random points, *i.e.* the points lying at the boundary of a Voronoi diagram, encompass a large proportion of the space. In other words, the Voronoi diagram has thick boundaries. This breaks the assumption made by many ANN algorithms that points can be unambiguously clustered with their closest neighbors. This phenomenon is different from the well-known *curse of dimensionality* which becomes apparent only for the high dimensional data (Indyk and Motwani, 1998). In the case of binary spaces, the thick boundaries of the Voronoi diagram influence the search regardless of the data dimensionality.

From this understanding comes the second contribution of the paper, an effective way to overcome the above mentioned problems inherent to Hamming spaces by creating multiple randomized data structures. Randomization produces structures that are independent from each other and therefore complementary. It solves the thick boundary problem and yields results similar to those obtained by converting the vectors to floating point values, but at a fraction of the computational cost. We instantiate this idea in two different ways, the first inspired by HKM trees and the second by the Locality-Sensitive Hashing scheme originally proposed for integer vectors (Gionis et al., 1999). In the first case, we replace the cluster centroids computed at each level of the tree by randomly chosen points and create multiple trees in this manner. In the second, we introduce an improved mechanism for selecting random subsets of coordinates used to index the vectors.

2. Related work

Local descriptors are high-dimensional vectors describing local regions extracted from images, and used in many applications of Computer Vision. In many of those applications, we need to match keypoints extracted from images against a large database as we do in our experiments: The distances between descriptors of keypoints corresponding to the same 3D point should be small. While the first descriptors were real-valued vectors, a recent trend focuses on *binary* descriptors, as they are more compact and distances between them can be evaluated efficiently. They can be computed directly from the images (Calonder et al., 2010), or from real-valued descriptors as in Locality Sensitive Hashing (Andoni and Indyk, 2008), Semantic or Spectral Hashing (SH) (Salakhutdinov and Hinton, 2009; Weiss et al., 2009), or LDAHash (Strecha et al., 2012). SH was designed to create binary vectors that can be used as table indices to directly access their nearest neighbors. However, as shown in Fig. 1, applying it to SIFT descriptors yields too large average Hamming distances between nearest neighbors to be practical. LDAHash (Strecha et al., 2012) produces average distances that are smaller but still too large.

In short, even when using sophisticated binary descriptors, quickly querying large databases still requires effective ANN methods. Even though Nearest Neighbor search has been widely discussed in the literature, no known generic algorithm is both exact and more efficient than brute force search. Many efficient approximate algorithms have been proposed for large-scale search. According to a recent comparative study (Muja and Lowe, 2009), the best ones for querying large databases are the randomized kd-trees (Silpa-Anan and Hartley, 2008) and hierarchical k-means tree algorithm (Fukunaga and Narendra, 1975; Nister and Stewenius, 2006).

The randomized kd-trees (Silpa-Anan and Hartley, 2008) are a recent modification of the original kd-trees (Friedman et al., 1977), which involved building a tree by recursively splitting in half along the dimension in which it exhibits the greatest variance. This performs well in low-dimensional spaces but loses its effectiveness as dimensionality increases (Amit et al., 1998). To prevent this, sets of randomized kd-trees can be built by recursively splitting along dimensions randomly chosen among the first D dimensions of greatest variance. Combining several trees with different splits mitigates the effects of quantization errors. Unfortunately, as we show in Section 3, this is a brittle technique when applied to binary vectors because a query vector can be moved to the wrong branch if only one of its bit is flipped, *e.g.* due to noise.

The hierarchical k-means tree (Fukunaga and Narendra, 1975; Nister and Stewenius, 2006) represents another successful alternative to brute force search. It recursively uses the k-means algorithm to split the data into k clusters. At run-time, a query vector follows the branch that corresponds to the closest centroid and back-tracking can be invoked to explore several leaves. Hierarchical k-means rely on means of vectors, which is problematic when dealing with binary vectors as we also see in Section 3.

Vantage-point trees (Yianilos, 1993) avoid the need to compute means by recursively picking a single vector among the data that reaches a node and splitting the others into those that are closer and those that are further. As we will show in Section 3, this method also performs poorly on binary vectors.

In short, state-of-the-art ANN techniques work well on real-valued vectors but not on binary ones. Furthermore, there is little work in connection to the latter. In (Cha and Srihari, 2000), an Additive Binary Tree (ABT) is associated to each binary vector. Each one of its nodes contains the frequency of 1's in a sub-part of the vector and this structure is used to stop the computation of the distance between two vectors early when the match is not promising. This approach, however, is still linear in the size of the database, and the speed gain is not clear compared to the full computation of the Hamming distance on modern hardware. In (Miller et al., 2005), the database is represented by a 256-ary tree in which each node corresponds to one byte of the vector. The parts of the tree that contain only one vector are pruned and replaced by a single leaf. This approach is sensitive to noise as changing a single bit may change how the branches are explored. In (Charikar, 2002), vectors are represented by a number of random permutations of bits. For each permutation, the vectors are sorted in a lexicographic order and when the query comes, the binary search is find the closest vectors. Although this method provides a sub-linear complexity, the memory required to store the sorted lists is a multiple of the dataset size. Moreover, it is reported to provide identical performance to the original LSH (Gionis et al., 1999) which is much more memory efficient.

3. Thick borders and performance loss

In this section, we first demonstrate that state-of-the-art ANN algorithms directly applied to binary vectors perform worse than when applied to floating-point ones. We then show that Voronoi diagrams have thick boundaries in binary spaces, which is what causes this performance drop.

3.1. ANN on binary vectors

To perform the comparison, we collected many images of Venice from the Flickr¹ database and created a first dataset containing 500k feature points and their SIFT descriptors. We then binarized

¹ <http://www.flickr.com>.

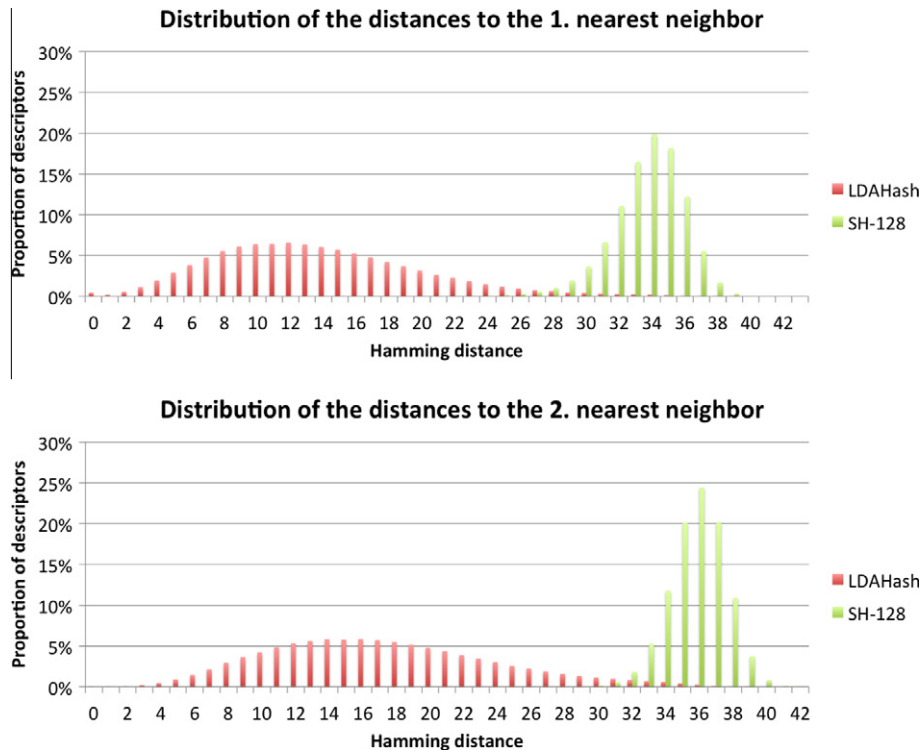


Fig. 1. Comparison of the distributions of distances from the descriptor to its first and second nearest neighbors in the Hamming spaces generated using LDAHash and SH on our 500k database. The average Hamming distance between descriptors is larger than 1 for both LDAHash and SH-generated 128-bit descriptors. However, because the distances are spread more widely for LDAHash vectors, all ANN algorithms tend to perform better on those, which is not all that surprising since SH was designed for a different purpose.

these using the publicly available implementation of LDAHash into 128-bits vectors (Strecha et al., 2012) whose length was shown to provide a good compromise. We used exhaustive linear-search to find the closest neighbors of each descriptor and we use this information as a ground-truth.

Table 1 summarizes our precision results for the first and second positions. The first simply is the the percentage of correct nearest neighbor that are retrieved. The second is computed by retrieving two nearest neighbors and checking whether both, only one, or none are the correct first two nearest neighbors of the query. The average proportion of correct matches divided by 2 is then taken to be the precision at the second position.

The results for the kd-trees and HKM algorithms were obtained using the publicly available code of the FLANN library (Muja and

Lowe, 2009), which automatically optimizes the algorithms parameters. We used our own implementation of the vantage-point trees.

The kd-trees and vantage-point trees can work on binary vectors without any modification since they do not involve averaging. By contrast, HKM involve computing centroids. We therefore tested two different versions of the algorithm, either rounding the coordinates of the centroids so that they remain binary vectors or using the floating-point coordinates.

As a baseline, we plot in the first column the results for matching the SIFT floating-point vectors. In the second column, we plot the systematically worse equivalent results using binary vectors. The degradation is noticeable for kd-trees and HKM, even if we treat binary vectors as floating-point ones. The performance drop is even more noticeable for the vantage-point trees.

Table 1

Precisions when looking for the first and second nearest neighbors for different methods and comparable query times, approximately 0.2 ms per query, for a dataset of 500k descriptors. The performances of state-of-the-art methods drop when they are applied on the binary 128-vectors obtained by running LDAHash (Strecha et al., 2012) on SIFT vectors. This is especially noticeable for the HKM algorithm when the centroids are forced to be binary vectors. The Parc-trees and Uniform LSH are two methods we introduce in Section 4 to avoid this loss of accuracy.

	Precision for			
	First position		Second position	
	SIFT descriptors	Binary descriptors	SIFT descriptors	Binary descriptors
Hierarchical k -means With real-valued centroids	0.94	0.82	0.93	0.79
kd-Trees	0.98	0.78	0.97	0.75
Hierarchical k -means With binary centroids	–	0.32	–	0.29
Vantage-point trees	0.35	0.17	0.31	0.15
Parc-trees	0.94	0.91	0.91	0.92
Original LSH for binary vectors (Gionis et al., 1999)	–	0.92	–	0.95
Uniform LSH	–	0.93	–	0.96

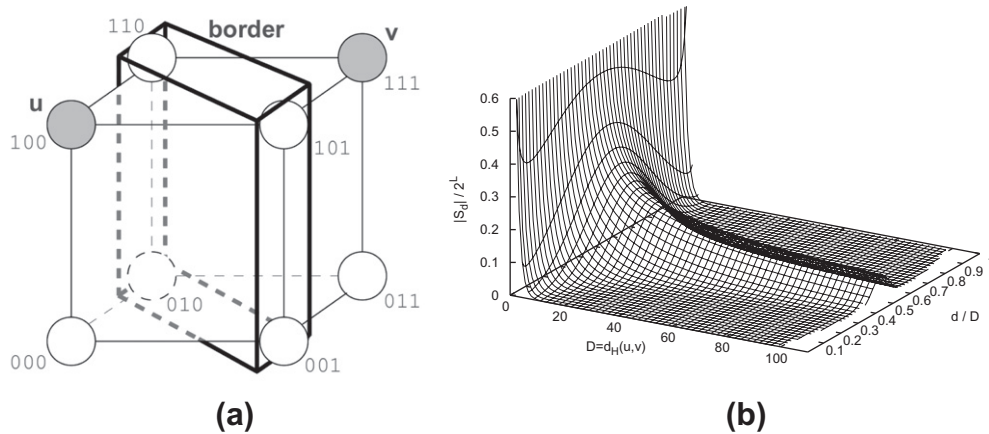


Fig. 2. Thick borders of Voronoi diagrams in binary space. (a) A significant proportion of the space is equidistant from two arbitrary points. In this example, four vectors are equidistant to the vectors u and v which accounts for half of the population. This makes the HKM algorithm fail on binary vectors. (b) Proportion of the binary vectors w that belong to the sets S_d defined in Eq. (1) as a function of the distance $D = d_H(u, v)$ and d . It is maximal for $d = D/2$, which corresponds to the set of vectors equidistant to u and v , and remains large even for large values of D . This phenomenon differs from the *curse of dimensionality* as it affects the data regardless of its dimensionality.

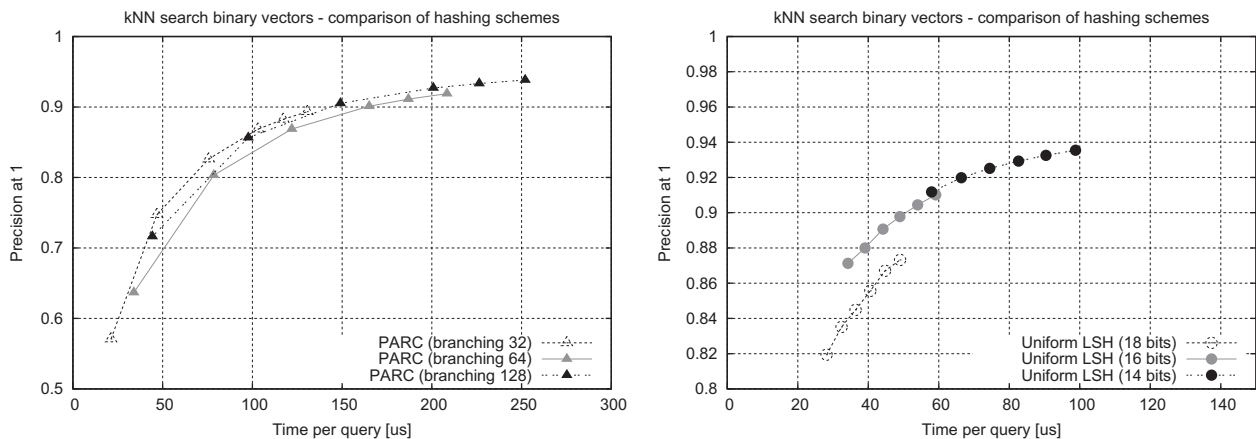


Fig. 3. Comparison of precision for first position with different parameters for the parc-trees and Uniform LSH on the 500k binary vectors dataset. For PARC trees (left) we plot the results for $T = 8, 16, 24, 32, 36, 40$ trees and different branching factors. For Uniform LSH (right) we plot the results for 30, 35, 40, 45, 50, 55, 60 hashing tables with various key lengths.

As a sanity check, even though Spectral Hashing (Weiss et al., 2009) is not truly designed to produce vectors that can be searched by ANN but rather used as indices to directly access their nearest neighbors, we ran the same series of tests on the vectors it produces. The ANN precision rates are globally lower but we observed the same behavior.

3.2. Interpretation

That kd-trees perform poorly on binary vectors is not that surprising since the splits are performed one dimension at a time and binary vectors can take only two values per dimension. Hence this method is sensitive to flip noise.

To understand the performance drops for the HKM and the vantage-point tree, one must consider that the topology of the Hamming space is different than the Euclidean one. This is because of the discrete nature of the binary spaces where many vectors are equidistant to two random points.

This affects the vantage-point trees because many vectors may lie on the splitting sphere: If the dimensionality of the binary space is L and the sphere radius is D , the proportion of uniformly distributed vectors that lie on the sphere boundary is $\frac{1}{2^L} \binom{L}{D}$. For exam-

ple, for $L = 16$ and $D = 8$, this represents 20% of the binary space, an enormous fraction. This is problematic because the algorithm depends on the assumption that the splits separate the data well.

The same thing happens with the HKM trees, especially when one binarizes the centroids. As explained below, the boundaries of the Voronoi diagram defined by such binarized centroids contain a significant proportion of the binary space. This is detrimental to the algorithm because points in those thick boundaries can be arbitrarily assigned to one or the other cluster and can fall down the wrong branch of the tree at run-time.

Let us consider two L -dimensional binary centroids u and v . We would like to evaluate the number of vectors around the boundary defined by u and v , that is, around the hyperplane made of w vectors equidistant from u and v . To this end, let us consider the cardinality of the sets S_d defined by:

$$S_d = \{w \text{ such that } d_H(v, w) = d_H(u, w) + D - 2d\}, \quad (1)$$

where $d_H(\cdot, \cdot)$ is the Hamming distance, and D the Hamming distance between u and v . The S_d family spans the Hamming space, with $u \in S_0, v \in S_D$, and $S_{D/2}$ the set of vectors vectors equidistant from u and v .

u and v have $L - D$ bits in common and D bits that are different. Let us first consider the case when D is even. For a vector w to

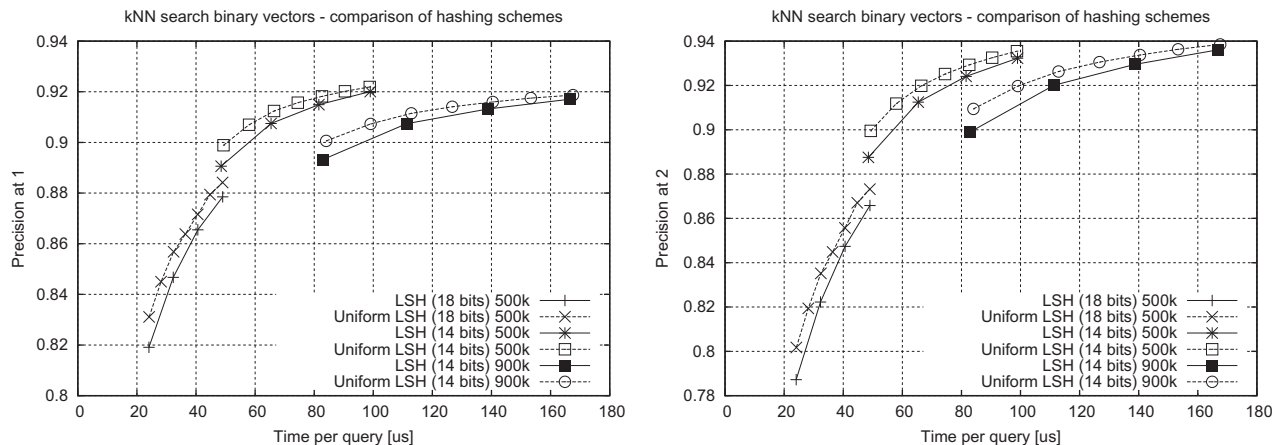


Fig. 4. Comparison of precision for first (left) and second (right) positions for the original LSH and Uniform LSH with hashing keys optimized as explained in Section 4.2. We varied the number of keys, their lengths and the sizes of the datasets. While the improvement is limited, this demonstrates that the hashing keys can be optimized in Hamming spaces.

belong to S_d , d bits among the D different bits must be changed between u and w . In addition any number n of bits among the $L - D$ common bits can also be flipped between u and w . We then have $d_H(u, w) = n + d$ and $d_H(v, w) = n + D - d$, and w indeed belongs to S_d .

The number of possible such w vectors is therefore $2^{L-D} \binom{D}{d}$ and their proportion of the full space is $\frac{2^{L-D}}{2^L} \binom{D}{d} = \frac{1}{2^D} \binom{D}{d}$. Remarkably this expression does not depend on the dimension L of the space but only on D and d . We plot its values in Fig. 2(b).

This expression reaches its maximum for $d = D/2$, that is, for the set of vectors that lie at equal distance from vectors u and v . Using Stirling's approximation (Stirling, 1764), this expression for $d = D/2$ can be approximated by $\sqrt{\frac{2}{\pi D}}$ when D increases, and therefore slowly decreases towards 0 (see Fig. 2(b)). For example, when $D = 2$, 50% of the space lies at equal distance from the 2 centroids! For $D = 64$, 10% of the space is still equidistant from the centroids. As a result, the borders of the Voronoi diagram defined by u and v contain a significant proportion of the binary space which leads to a severe performance drop of the HKM algorithm.

When D is odd, no vector is equidistant from u and v . However, we can derive a similar expression for the number of points for which the distances to u and v differ by 1. The number of such points remains large. The borders of the Voronoi diagram defined by the centroids in the HKM algorithm therefore contain a significant proportion of the binary space.

4. Randomized data partitioning

In this section, we address the above-mentioned shortcoming of state-of-the-art ANN search algorithms in Hamming spaces, and describe two simple yet effective ANN search methods that work by randomizing data partitioning.

4.1. Parc-trees

This first algorithm relies on multiple trees. Like the nodes of a hierarchical k -means (HKM) (Fukunaga and Narendra, 1975; Nister and Stewenius, 2006) tree, the parc-tree nodes split the data into k parts by storing k vectors we call *centroids*, and associating the data with the closest centroid. Each non-terminal node has k children, corresponding to the different parts. By contrast with HKM, we do not optimize on the centroids but randomly select them among

the data vectors that reach the node, except those which have previously been used. The recursion stops when the number of data vectors is less than k . Because of the randomization, the trees are independent from each other.

At run-time, a query vector recursively follows the branch associated to the closest node vector until it reaches a leaf, as in HKM. In HKM however, the leaves have to store all the data that reach them, and a linear search over this data is required to find the vector closest to the query vector as the candidate nearest neighbor. In parc-trees, the centroids belong to the dataset and when the query vector reaches a leaf, we already computed its distances to the centroids of the nodes it visited. Hence, a candidate nearest neighbor is chosen to be the closest vector among those in the leaf and the centroids of the visited nodes.

The query operation is repeated over all T trees and the best match is retained. This allows the parc-trees to mitigate the quantization error introduced by the thick Voronoi boundary: We can find the correct nearest neighbor even if it is present in only one visited node among all the trees.

The influence of parameters T and k on the obtained precision and computational time can be seen in Fig. 3. The performance increases with the number of trees T , until it saturates, linearly with T . Increasing the branching factor k also improves the performance. The average tree depth then decreases, but the computation time still increases: A tree of depth d contains $k + k^2 + k^3 + \dots + k^d \approx k^d$ vectors, so a tree of S data vectors is approximately of depth $\log S / \log k$. The number of distance computations required when dropping a query vector into the tree is therefore $k \log S / \log k$, which increases with k sublinearly.

Most of the operations involved by this approach are Hamming distance computations. They amount to an XOR operation followed by a popcount instruction present on modern CPUs, and are much faster to evaluate than the Euclidean distance between floating-point vectors. Moreover, the T trees can be simultaneously queried on a multi-core machine, which means we incur only a limited penalty for using several trees.

4.2. Uniform LSH

We also developed an ANN method inspired by the Hashing-based method of Gionis et al. (1999), which involves converting integer vectors into binary ones and randomly selecting and concatenating bits from them to generate multiple hashing keys. A query vector is then matched against the vectors in the buckets corresponding to its keys values by linear search. The smaller the

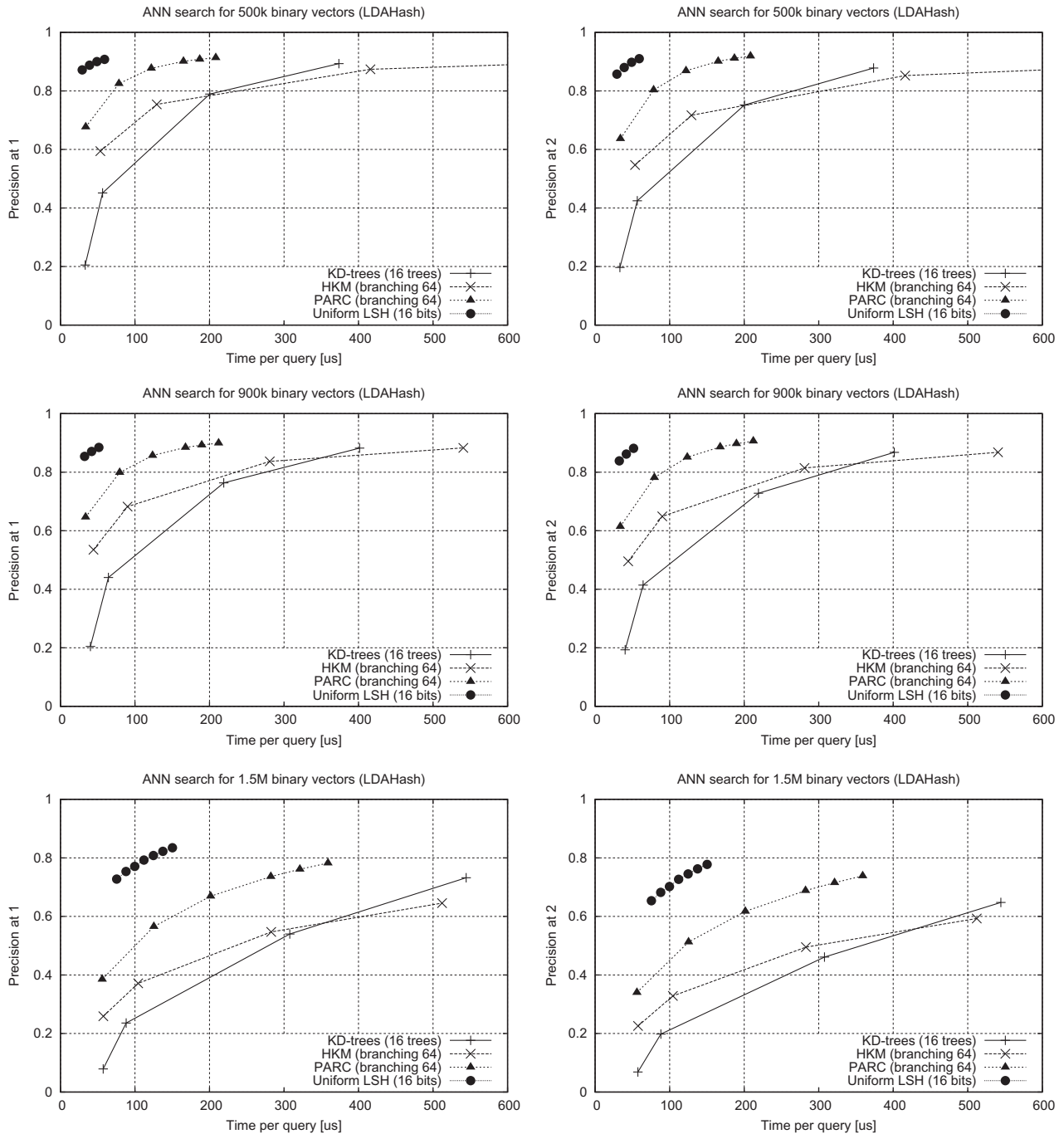


Fig. 5. Comparison of precision for first and second positions for different ANN search algorithms on 500k, 900k and 1.5 M binary vector datasets. Uniform LSH outperforms the parc-trees, which themselves outperform all the other state-of-the-art methods for all configurations.

lengths of the keys, the greater the size of the buckets becomes, which yields higher precision at the cost of increased computational time. This simple scheme performs well, as our experiments show. As for parc-trees, most of the operations are Hamming distance evaluations, which can be performed efficiently, and searches, which can be parallelized.

However, the random selection of coordinates may lead to unnecessary overhead, as some coordinates may be selected more frequently than others, whereas some of them may not be picked at all. This problem can be solved by increasing the number of keys, but to run fast, it is desirable to use as few keys as possible.

To resolve this dilemma, we optimize the keys so that the bits selected to generate the keys are distributed more uniformly. This way, the keys generate more various partitioning of the database. More formally, we can define the keys K_i as sets containing the selected bits coordinates: $K_i = \{c_{ij} \in [1; L] \mid j \in [1; n]\}$ where L is the dimensionality of the binary vectors, n is the number of bits in a key. We also define N_k as the number of times a given coordinate is used in a key: $N_k = |\{c_{ij} = k \mid i \in [1; m] \text{ and } j \in [1; n]\}|$, where m is the number of keys. Then we optimize the keys to minimize

$$\min_{\{K_i\}} \sum_k (N_k - N)^2 \quad (2)$$



Fig. 6. Top: Two out of the 25 13824×7680 pixels images from the Marseilles dataset. Bottom: Two out of the 68 11500×7500 pixels images from the Zwolle dataset.

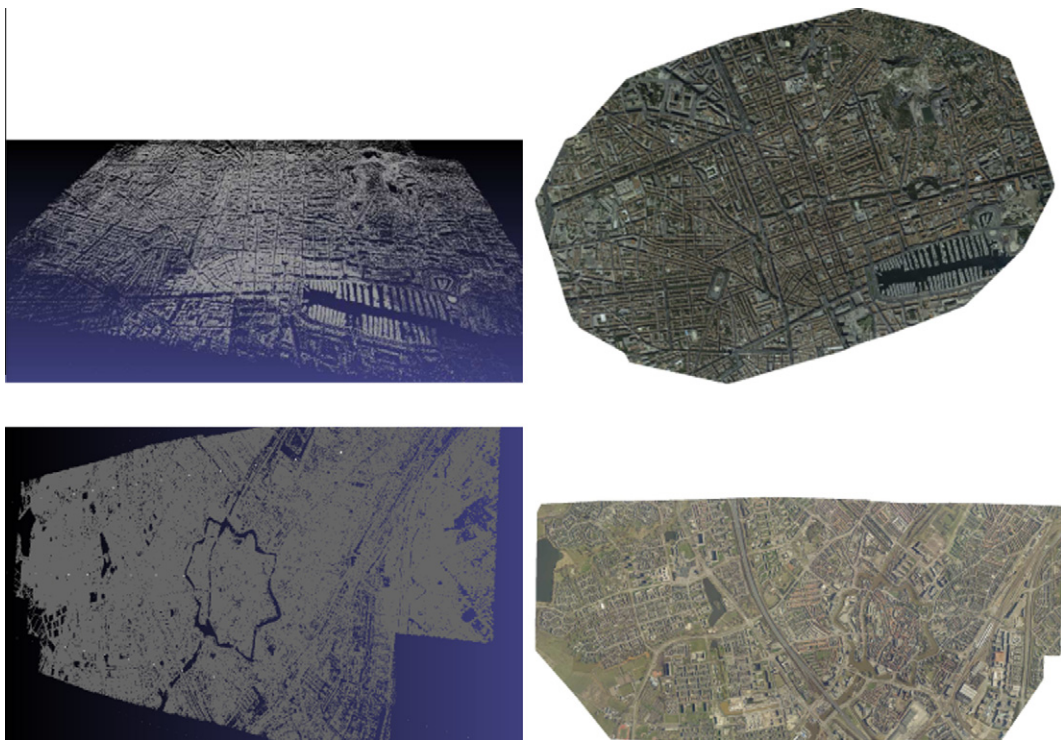


Fig. 7. Top left: The aerial triangulation of 1.1M 3D points and top right: the generated ortho-image for the Marseilles dataset. Bottom left: The aerial triangulation of 2.1M 3D points and bottom right: the ortho-image made of 68 individual images (right) for the Zwolle dataset.

with $N = n \cdot m/L$, the ideal number of times a coordinate should be picked. To do this, we use a simple greedy algorithm that generates the keys one by one, by randomly selecting the bits among those which were used less often for the previous keys. We call this modification Uniform LSH as the distribution of the bits used in the keys is optimized to be more uniform and hence partition the dataset better.

We experimented with different numbers of keys and numbers of bits per key used. The results of those experiments are shown in Fig. 3. Computation times increase linearly with the number of keys, but the precision of the search also increases. Similarly, the lower the number of bits per key used, the bigger the data partitions and hence the longer the search time. However, since we perform a linear search within the selected data partition, shorter

keys (and bigger data partitions) lead to performance improvement.

Overall, as it can be seen in Fig. 4, the resulting Uniform LSH algorithm improves performances over those of the original LSH.

5. Results

In this section, we compare parc-trees and uniform LSH against kd-trees and HKM trees. The results were obtained for the 500k Venice dataset of 128-bit binary descriptors from Section 3.1. To see if the results hold also for bigger datasets, we created two more datasets containing 900k and 1.5 M binary descriptors which were generated by binarizing the SIFT descriptors extracted from more Flickr images of Venice. We draw the plots by setting the parameters of all algorithms so that the query time is approximately the same. To produce the different points in the plots, the number of hashing tables of LSH varied from 30 to 60. The test datasets are larger than the ones used to evaluate the recent FLANN library (Muja and Lowe, 2009) which mostly contained only 100k vectors. Furthermore, datasets of comparable sizes are frequently used for real-life applications, such as image-based 3D reconstruction. The results presented here are the average over three runs. The computation times were evaluated on a computer with two Intel Xeon E5620 2.4 GHz CPUs and 48 Gb RAM.

Fig. 5 presents the comparison of different ANN search algorithms applied to binary descriptors. LSH outperforms all tree-based methods. Out of those, parc-trees remain the best. The speed-up of LSH and parc-trees over the other algorithms is especially visible for higher precision levels. For instance, for 500k dataset KD-trees needs approximately 300 μ s to reach the precision at second position equal to 0.85, whereas it takes parc-trees and Uniform LSH less than 100 μ s and 50 μ s, respectively. As the dataset size grows, it takes more time to find the nearest neighbors, but the relative ordering of the performances remains the same for all the methods: LSH performs the best, followed by the parc-trees. For the 1.5 M dataset, LSH achieves a precision of 0.7 at first position about an order of magnitude faster than KD-trees and HKM.

To verify our approach, we applied it to Aerial Triangulation. We extracted feature points from aerial images, and match them using the same binary descriptors as before. Matched points correspond to the same 3D points, and we use these matches to jointly optimize the 3D points and the camera parameters by bundle block adjustment (Hartley and Zisserman, 2000).

We tested our binary search strategy on two datasets of large aerial images. The first dataset contains 25 high resolution (13824×7680) images of Marseilles,² two of which are shown in Fig. 6. The second dataset consists of 68 11500×7500 aerial images of the Dutch city of Zwolle.

Each image contains approximately 400k binary keypoints which makes exhaustive feature matching, even on binary vectors, excessively slow. Our approach reduces the matching time by a factor 20 over linear search with a 95% accuracy, which is consistent with the results reported in Section 5. The final aerial triangulations and the combined ortho-images for the Marseilles and Zwolle datasets are shown in Fig. 7.

6. Conclusion

We showed that Voronoi diagrams in Hamming spaces have thick borders, which reduces the precision of state-of-the-art

ANN algorithms. We then proposed two techniques that rely on randomized data partitioning to overcome this problem and yield precisions that are comparable to those obtained using floating-point vectors at a fraction of the computational cost.

References

- Amit, Y., Geman, D., Jedynak, B., 1998. Efficient focusing and face detection. *Face Recognition: From Theory to Applications*, 143–158.
- Andoni, A., Indyk, P., 2008. Near-Optimal Hashing Algorithms for Approximate Nearest Neighbor in High Dimensions. *Commun. ACM* 51 (1), 117–122.
- Arya, S., Mount, D., Netanyahu, N., Silverman, R., Wu, A., 1998. An optimal algorithm for approximate nearest neighbor searching fixed dimensions. *J. ACM* 45, 891–923.
- Bay, H., Tuytelaars, T., Van Gool, L., 2006. SURF: speeded up robust features. *Eur. Conf. Computer Vision.*, 404–417.
- Beis, J., Lowe, D., 1997. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. *Conf. on Computer Vision and Pattern Recognition*, 1000–1006.
- Brown, M., Hua, G., Winder, S., 2011. Discriminative Learning of Local Image Descriptors. *IEEE Trans. Pattern Anal. Machine Intell.* 33 (1), 43–57.
- Calonder, M., Lepetit, V., Strecha, C., Fua, P., 2010. BRIEF: Binary Robust Independent Elementary Features. In: *Eur. Conf. Computer Vision*. Heraklion, Greece.
- Cha, S.-H., Srihari, S., 2000. Nearest neighbor search using additive binary tree. In: *Conference on Computer Vision and Pattern Recognition*.
- Charikar, M., 2002. Similarity estimation techniques from rounding algorithms. *Proc. 34th Ann. ACM Sympos. Theory Comput.*, 380–388.
- Friedman, J., Bentley, J., Finkel, R., 1977. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Software* 3 (3), 209–226.
- Fukunaga, K., Narendra, P., 1975. A branch and bound algorithm for computing k-nearest neighbors. *IEEE Trans. Comput.* 24, 750–753.
- Gionis, A., Indyk, P., Motwani, R., 1999. Similarity search in high dimensions via Hashing. In: *International Conference on Very Large Databases*.
- Gong, Y., Lazebnik, S., 2011. Iterative quantization: a procrustean approach to learning binary codes. In: *Conference on Computer Vision and Pattern Recognition*.
- Hartley, R., Zisserman, A., 2000. *Multiple View Geometry in Computer Vision*. Cambridge University Press.
- Indyk, P., Motwani, R., 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In: *Proceedings of the 30th Symposium on Theory of Computing*, pp. 604–613.
- Kulis, B., Darrell, T., 2009. Learning to hash with binary reconstructive embeddings. *Adv. Neural Informat. Proc. Systems*, 1042–1050.
- Liu, T., Moore, A., Gray, A., Yang, K., 2004. An investigation of practical approximate nearest neighbor algorithm. *Advance in Neural Information Processing Systems*.
- Lowe, D., 2004. Distinctive image features from scale-invariant keypoints. *Internat. J. Computer Vision* 20 (2), 91–110.
- Mikolajczyk, K., Schmid, C., Zisserman, A., 2004. Human Detection Based on a Probabilistic Assembly of Robust Part Detectors. *Eur. Conf. Computer Vision*, 69–81.
- Miller, M., Rodriguez, M., Cox, I., 2005. Audio fingerprinting: nearest neighbor search in high dimensional binary spaces. *J. VLSI Signal Process. Systems* 41 (3), 285–291.
- Muja, M., Lowe, D., 2009. Fast approximate nearest neighbors with automatic algorithm configuration. In: *International Conference on Computer Vision*.
- Nister, D., Stewenius, H., 2006. Scalable recognition with a vocabulary tree. In: *Conference on Computer Vision and Pattern Recognition*.
- Raginsky, M., Lazebnik, S., 2009. Locality-sensitive binary codes from shift-invariant kernels. *Adv. Neural Inform. Proc. Syst.*, 1509–1517.
- Salakhutdinov, R., Hinton, G., 2009. Semantic hashing. *Internat. J. Approx. Reason.* 50 (7), 969–978.
- Shakhnarovich, G., 2005. *Learning Task-Specific Similarity*, Ph.D. thesis, Massachusetts Institute of Technology.
- Silpa-Anan, C., Hartley, R., 2008. Optimised kd-trees for fast image descriptor matching. In: *Conference on Computer Vision and Pattern Recognition*.
- Stirling, J., 1764. *Methodus Differentialis* [electronic Resource]: Sive Tractatus De Summatione Et Interpolatione Serierum Infinitarum, J. Whiston and B. White.
- Strecha, C., Bronstein, A., Bronstein, M., Fua, P., 2012. LDHash: improved matching with smaller descriptors. *IEEE Trans. Pattern Anal. Mach. Intell.* 34.
- Weiss, Y., Torralba, A., Fergus, R., 2009. Spectral hashing. *Adv. Neural Inform. Process. Systems* 21, 1753–1760.
- Yianilos, P., 1993. Data structures and algorithms for nearest neighbor search in general metric spaces. In: *Fourth ACM-SIAM Symposium on Discrete Algorithms*, 1993.

² From “Benchmarking of Image Matching Approaches for DSM computation” <http://eurosdrrbenchmarkofimagematching.ign.fr/>.