

3D Binary Signatures

Siddharth Srivastava
Department of Electrical Engineering
Indian Institute of Technology, Delhi, India
eez127506@ee.iitd.ac.in

Brejesh Lall
Department of Electrical Engineering
Indian Institute of Technology, Delhi, India
brejesh@ee.iitd.ac.in

ABSTRACT

In this paper, we propose a novel binary descriptor for 3D point clouds. The proposed descriptor termed as 3D Binary Signature (3DBS) is motivated from the matching efficiency of the binary descriptors for 2D images. 3DBS describes keypoints from point clouds with a binary vector resulting in extremely fast matching. The method uses keypoints from standard keypoint detectors. The descriptor is built by constructing a Local Reference Frame and aligning a local surface patch accordingly. The local surface patch constitutes of identifying nearest neighbours based upon an angular constraint among them. The points are ordered with respect to the distance from the keypoints. The normals of the ordered pairs of these keypoints are projected on the axes and the relative magnitude is used to assign a binary digit. The vector thus constituted is used as a signature for representing the keypoints. The matching is done by using hamming distance. We show that 3DBS outperforms state of the art descriptors on various evaluation metrics.

CCS Concepts

•Computing methodologies → Interest point and salient region detections; Matching; Computer graphics; Object detection; Object recognition;

Keywords

3D Descriptors; 3D Matching; 3D Object Recognition

1. INTRODUCTION

Humans visualize a three dimensional world. Therefore, the ability to obtain and process information especially from visual senses in three dimensions has always been an exciting and potentially promising area of research. The success of various feature extraction and classification related tasks in 2D image analysis can be attributed to the availability of large scale annotated datasets such as ImageNet [5, 28]. With the growing availability of 3D scanning technologies such as Kinect [40], LIDAR etc., the availability of large and quality 3D datasets is also increasing. This has opened many research areas from 2D image analysis such as Object Recognition, Object Retrieval, Object Classification, Segmentation etc. for three dimensional data as well.

Local features and Deep Learning based approaches especially Convolutional Neural Networks (CNNs) have been successfully used for solving many research problems involving 2D images [32, 4, 17, 35, 15]. Similar to 2D domain, various local feature and deep learning based models have

been proposed for 3D point clouds. Although, CNN based techniques have been successfully proposed for 3D object classification [34, 14] and recognition [22], but 3D local features are state of the art in many tasks such as 3D object recognition and classification [9, 31], 3D scene reconstruction [12], 3D model retrieval [8] etc. Comparisons of descriptors on various benchmarks [10] show that high descriptiveness with low storage requirements and reasonable computational complexity depending upon the number of points in the point cloud constitute key requirements for a descriptor. Moreover, for practical applications involving 3D object matching, recognition, classification and reconstruction, the efficiency in feature matching and lower storage requirements become important [1]. In order to address these issues, we propose a novel 3D binary descriptor deriving motivation from binary descriptors for 2D images [3, 19, 27]. These descriptors have reported to match or outperform traditional SIFT like descriptors [21, 2] with significantly lower computational complexity and storage requirements. The proposed technique encodes the differences in the projection of normals into binary vectors. The projection is computed for nearest neighbours of a keypoint and aligned with a local reference frame. The nearest neighbours are chosen based on their angular orientation on a 2D projection plane. The binary descriptor thus generated is matched using Hamming Distance. We show that the proposed binary descriptor outperforms the state of the art on various benchmarks.

The closest work to ours is B-SHOT [26] which generates a binary vector from the popular Signature of Histograms of Orientations (SHOT) [31] descriptor. It quantizes the real valued SHOT descriptor to a binary vector. The major difference in the proposed technique is that the binary vector is generated directly from the point cloud data instead of first computing a real valued vector and quantizing it which results in an additional overhead. For highlighting the novelty of the technique, we refer to [10] where authors categorize 3d descriptors into two classes based on the methodology of generating histograms. First class of descriptors computes histograms either by computing a local reference frame and accordingly divide the support region spatially. Various spatial distribution measurements are then accumulated into histograms. The second class of descriptors, encodes geometric attributes such as normals, principal curvatures etc. of the points on the surface in the local neighbourhood the keypoint. The proposed approach is a hybrid of these two approaches. We define a local surface with the nearest neighbours of a keypoint with an angular constraint. Then similar to the first class we use a Local Reference Frame (LRF) for

aligning this local surface followed by computing normals and encoding the projective difference among them which is similar to the second class of descriptors. Since the proposed technique does not require computation of histograms, therefore to highlight this difference, we term the extracted descriptors as *signatures*. In view of the above discussion, the main contributions of this paper are:

- We propose a hybrid approach for constructing a highly distinctive yet compact 3D binary descriptor based on encoding differences among normal projections among nearest neighbours of a keypoint. To the best of our knowledge, this is the first work to directly generate 3D descriptors from point cloud data.
- We show that the proposed 3D Binary Signature (3DBS) outperforms state of the art descriptors on common evaluation benchmarks.

The paper is organized as follows: In Section 2, we discuss the techniques that are similar to ours and have motivated the construction of the proposed 3D Binary Signature. In Section 3, we detail the generation of 3D Binary Signatures (3DBS) followed by experimental results in Section 4. Finally, Section 5 concludes this paper.

2. RELATED WORK

Numerous 3D descriptors have been proposed in the past two decades. Broadly these descriptors fall into two categories. First, the descriptors that encode the geometric attributes of a spatial region as histograms. Second, the descriptors that encode as histograms the statistical properties of a region. Another way of looking at these descriptors is whether a Local Reference Frame (LRF) has been used for forming the descriptor. The descriptors without LRF utilize local statistics such as normals etc. to form the descriptors while those using LRF encode information from spatial distribution or geometric properties of neighbouring points to form the descriptors. Authors in [10] provide an exhaustive evaluation of various 3D binary descriptors. In the rest of this section, we discuss a few techniques which are either (i) similar to the proposed technique or (ii) use surface normals for generating the descriptors. We first discuss a few techniques which particularly focus on speed and memory efficiency, followed by robust and general purpose techniques.

B-SHOT [26] is a 3D binary feature descriptor based on the popular SHOT descriptor. It is generated by quantizing the 352 length real valued SHOT descriptor to 352 length binary vector. The quantization begins by dividing the SHOT descriptor into sets of four values. The relative magnitudes of various combinations of values in each set are compared to assign a corresponding four-bit binary vector. There is a loss of information due to this quantization but it is compensated by significant gains in the descriptor matching efficiency and storage requirement. Signature of Histograms of Orientations (SHOT) [31], which is the basis behind B-SHOT, is based upon encoding into histograms, the surface normals in a spatial distribution. It works by constructing an LRF for a keypoint and the points in the support region are aligned with this LRF. The support region is then divided into several volumes. Each volume results in a local histogram by accumulation of point counts into bins as per the angles between normal of points in the support region

and the keypoint. The final descriptor is computed by concatenating all the local histograms. Spin Image (SI) [16] descriptor uses the surface normal at a keypoint as Local Reference Axis. Then in-plane and out-plane distances are computed for each point in the support region which is discretized into a 2D array. The final histogram is generated by binning the points from the support region to the 2D array constructed earlier. The dimension of the SI descriptor is equal to the number of bins across each dimension of the in-plane and out-plane space. Rotational Projection Statistics (RoPS) [11] constructs an LRF for each keypoint aligning it with the local surface to achieve invariance to transformations. The points on the local surface are rotated around x , y , and z axis and the corresponding support region are further projected onto the coordinate planes (xy , xz and yz). The planes are then divided into several bins and the number of points in each bin is counted. Various statistics are calculated on these bins and they are concatenated to form the final descriptor. 3D Shape Context (3DSC) [7] uses normal as Local Reference Axis. Further, it divides the support region into many bins along azimuth, elevation and radial dimensions. The descriptor is generated by weighted accumulation of the points lying in each bin. Unique Shape Context (USC) [37] improves the memory requirement and computational efficiency of 3DSC by avoiding computation of multiple descriptors at a keypoint. This is achieved by defining an LRF (same as that used in SHOT) at a keypoint and aligns the local surface accordingly. The support region is then divided into bins and the weighted accumulation is performed for the points lying in corresponding bins.

3. PROPOSED METHODOLOGY

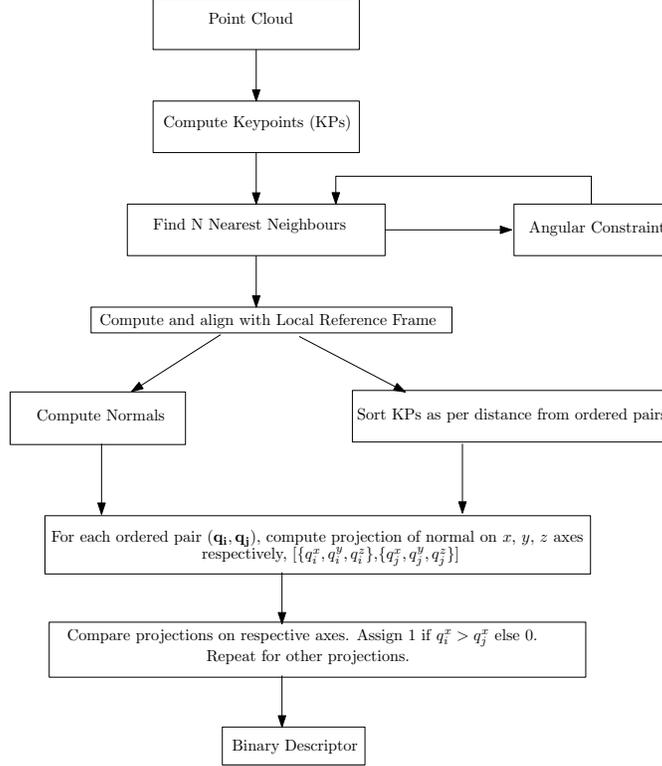
In this section, we discuss the methodology for generating the 3D Binary Signatures. Figure 3 shows the construction pipeline of the proposed binary descriptor.

The proposed methodology is abstractly inspired from the descriptor construction pipeline of 2D binary descriptors. Generation of 2D binary descriptors in general involves three steps: a) Choosing a sampling pattern b) Orientation Compensation and c) Selection of Sampling pairs. The sampling pattern in the current methodology is obtained with the help of nearest neighbours while Local Reference Frame compensates for invariance to orientation and other transformation. The sampling pairs are the ordered nearest neighbour pairs which are compared for encoding the difference in projections of the normals to a binary vector. We describe these steps in detail in the following subsections.

3.1 Keypoint Detection

The keypoints can be detected using any of the standard keypoint detection techniques. From our experiments on various keypoint detectors, namely, Intrinsic Shape Signature (ISS) [40], MeshDoG [39], Keypoint Quality (KPQ) [23], Harris3D [33], Heat Kernel Signatures [36] and 3D SIFT [6], it was observed that ISS and Harris3D performed best in our experiments. Authors in [38] present a comprehensive survey on various keypoint detectors. They show that Intrinsic Shape Signatures(ISS) [41] demonstrate the highest repeatability while being the most efficient keypoint detector. Moreover, the work in [1] showed that selection of an appropriate keypoint detector has a reasonable impact on the performance of a descriptor. Therefore, for further experiments in this work, ISS has been chosen as the keypoint

Figure 1: Construction of 3D Binary Signature



detector.

3.2 Nearest Neighbours with Angular Constraint

The objective of this criteria is to uniformly distribute the nearest neighbours along the surface of the point cloud. The process is pictorially shown in Figure 2(a). Instead of a spherical support region, we define the local surface for each keypoint with its N nearest neighbours. Nearest neighbours are traditionally computed based on a distance criterion. Authors in [20] observe that computing nearest neighbour in such a manner may not be an optimal choice for local surface representation and hence introduce an angle criterion. We therefore use the angle criterion to distribute neighbours around the keypoint \mathbf{p} , by projecting the neighbours of the points on a plane ϕ which is best fitting plane of neighbours of \mathbf{p} . The plane ϕ is obtained as a least squares best fitting plane which turns the approximation into a quadratic form and therefore can be solved efficiently. The points thus obtained are sorted as per distance and for each neighbour of \mathbf{p} , if the angle between \mathbf{p} , its neighbour \mathbf{q}_m and its successor \mathbf{q}_{m+1} , does not exceed a threshold θ i.e. $\angle \mathbf{q}_m \mathbf{p} \mathbf{q}_{m+1} \leq \theta$, it is considered added to the set of nearest neighbours. This process is repeated till N neighbours have been found. Angle criterion also helps in avoiding ambiguity in the order of neighbours having the same distance from a keypoint. For instance, if $\mathbf{q}_1, \mathbf{q}_2$ and \mathbf{q}_3 are neighbours of keypoint \mathbf{p} with the following distance relation, $\|\mathbf{p} - \mathbf{q}_2\|_2 = \|\mathbf{p} - \mathbf{q}_3\|_2 < \|\mathbf{p} - \mathbf{q}_1\|_2$. If $\angle \mathbf{q}_1 \mathbf{p} \mathbf{q}_2 < \angle \mathbf{q}_1 \mathbf{p} \mathbf{q}_3$, \mathbf{q}_2 is listed before \mathbf{q}_3 avoiding any ambiguity.

3.3 Alignment with Local Reference Frame

To infuse invariance to various transformations such as

translation and rotation along with robustness to noise and clutter, the local surface formed previously is aligned with a Local Reference Frame. Authors in [25] show that repeatability of an LRF has a direct impact on the robustness and descriptive ability of a descriptor. Therefore, we construct an LRF using the technique of [31] which is the basis for the popular SHOT descriptor. It computes a weighted covariance matrix \mathbf{M} given by Eq. 1 around the keypoint \mathbf{p} where the distant points are assigned smaller weights.

$$\mathbf{M} = \frac{1}{\sum_{i:d_i \leq R} (R - d_i)} \sum_{i:d_i \leq R} (R - d_i) (\mathbf{p}_i - \mathbf{p})(\mathbf{p}_i - \mathbf{p})^T \quad (1)$$

where $d_i = \|\mathbf{p}_i - \mathbf{p}\|_2$ and R is the spherical support region. As the computation of the covariance matrix assumes a spherical support region, we set the radius of the support region, R , as the farthest neighbour of the point \mathbf{p} , as given by Eq. 2.

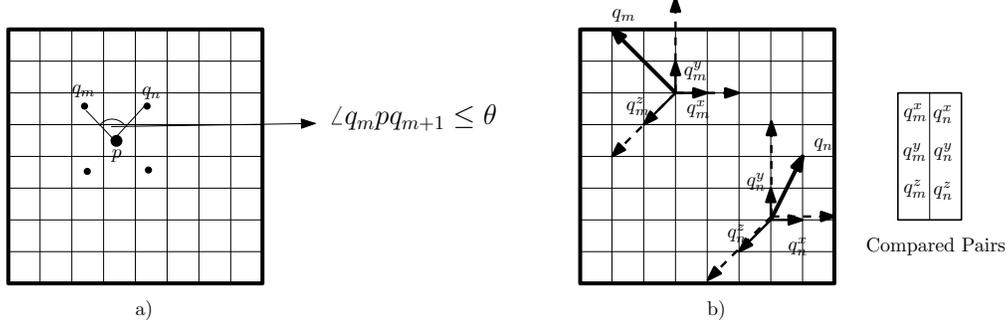
$$R = \max_{1 \leq i \leq N} \|\mathbf{p}_i - \mathbf{p}\|_2 \quad (2)$$

The region from the previous step is then aligned with this LRF for further processing. We denote this collection of points excluding the keypoint as C .

3.4 Descriptor Generation

For generating the descriptor, the projection of surface normals on each of the three axes, x, y, z are computed. Let us denote the projection of a point \mathbf{q} on axis a where $a \in \{x, y, z\}$ as q^a . These are computed for the points in C and are mathematically given by Eq. 3.

Figure 2: Visualization of a) Angular Constraint b) Comparison of projection of Normals



$$q^x = \langle \mathbf{q}, \hat{\mathbf{i}} \rangle, \quad q^y = \langle \mathbf{q}, \hat{\mathbf{j}} \rangle, \quad q^z = \langle \mathbf{q}, \hat{\mathbf{k}} \rangle \quad \forall \mathbf{q} \in C \quad (3)$$

where $\langle \cdot \rangle$ represents inner product and $\hat{\mathbf{i}}, \hat{\mathbf{j}}, \hat{\mathbf{k}}$ are the unit vector in the direction of x,y and z axes i.e. $(1, 0, 0)$, $(0, 1, 0)$ and $(0, 0, 1)$ respectively.

Then the respective projections on each axis for pairs of ordered points $(\mathbf{q}_m, \mathbf{q}_n)$ s.t. $1 \leq m, n \leq N$ are compared. The ordered pairs are constructed by first sorting the points in C based upon the distance from the keypoint \mathbf{p} . Let us denote the sorted set of points by $\mathbf{q}_1, \mathbf{q}_2 \dots \mathbf{q}_N$ such that

$$\|\mathbf{q}_m - \mathbf{p}\|_2 \leq \|\mathbf{q}_n - \mathbf{p}\|_2 \quad \forall m < n \text{ and } m, n \in (1, N) \quad (4)$$

Then the ordered set of points consist of the pairs

$$(\mathbf{q}_m, \mathbf{q}_n) \text{ s.t. } m < n \quad (5)$$

Each ordered pair thus obtained, results in three comparisons corresponding to projections on three axes. Each comparison is assigned a binary bit $b_a \in \{0, 1\}$ for an axis a based upon the relative magnitude of the projections, as given in Eq. 6.

$$b_a = \begin{cases} 1, & \text{if } q_m^a \geq q_n^a \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

The comparison for an ordered pair, i , results in a bit vector of size 3 given by Eq. 7.

$$\mathbf{b}(i) = b_x b_y b_z \quad (7)$$

Finally, the descriptor is constructed by concatenating the binary strings $\mathbf{b}(i)$ from all the comparisons resulting in a binary vector of size $3 * \frac{N*(N-1)}{2}$. The process is pictorially visualized as in Figure 2(b).

3.5 Descriptor Matching

The descriptors are matched using Hamming Distance. There are efficient algorithms to compute hamming distances especially on CPUs having hardware support to count bits in a word (for ex: POPCNT instruction). In practical applications, there are usually millions of such descriptors, say, D , which are needed to be matched. Performing a linear search over D descriptors would be of the order of $O(D^2)$, which would be computationally very expensive. For binary descriptors from 2D images, the popular way of matching these descriptors is to hash the binary bit vectors using techniques

such as Locality Sensitive Hashing (LSH)[27] and if needed, perform a linear search in each bucket. But, it is important to note here that the matching complexity of such techniques is affected by two important factors. Firstly, by the number of descriptors in the search space i.e. D and secondly by the length of the binary bit vectors. Typically if the length of the binary vector is greater than 32, various matching techniques use linear scan [18]. These limitations are aggravated in the context of 3D point clouds as the number of points could be of the order of millions resulting in a huge number of keypoints. Moreover, the size of the proposed 3D Binary Signature is $O(N^2)$ for each keypoint, where N is the number of nearest neighbour of a keypoint. Therefore, we adopt a fast binary feature matching technique proposed in [24]. The technique has been chosen for its low memory footprint and the ability to scale to large datasets. The technique is briefly described below.

3.5.1 Fast Feature Matching

- *Building search tree of features*: The input data (descriptors) is divided into K clusters by randomly selecting K data points. The remaining data points are assigned to the closest cluster center (similar to k-medoids clustering). If the number of data points in a cluster (DP_C) is above a certain threshold i.e. maximum number of leaf nodes (S_L), then the algorithm is recursively repeated until each cluster has $DP_C < S_L$.
- *Search for nearest neighbours*: The search is performed in parallel on multiple hierarchical clustering trees. The search begins by recursively exploring the node nearest to the query descriptor while the unexplored nodes are added to a priority queue. Once a node has been completely traversed, the next nearest node from the priority queue is extracted and is again explored recursively. This stopping criteria for this recursive search is based upon a search precision i.e. the fraction of exact neighbours discovered in the total number of returned neighbours.

4. EXPERIMENTS AND RESULTS

In this section, we describe the experimental details and the results thus produced.

4.1 Experimental Setup

4.1.1 System Specification

Table 1: Datasets used in the study

Dataset Name	#Model	#Scene
Random Views	6	36
Laser Scanner	5	10
LIDAR	5	10
Retrieval	6	18

The experiments were performed on a system having an Intel i7 processor with 128GB RAM and Ubuntu 14.04 Operating System. The implementation has been done in C++ (g++ 4.8) using Point Cloud Library (PCL) [30] with OpenMP enabled.

4.1.2 Datasets

The experiments were performed on four publicly available datasets [13] with the details as shown in Table 1. The chosen datasets allow us to evaluate the proposed descriptor on various application contexts. The *Random Views*, *Laser Scanner* and *LIDAR* datasets are suited for object recognition. In these datasets the model are full 3D meshes while the scenes are 2.5D views from specific viewpoints. On the other hand, the *Retrieval* dataset is for 3D shape retrieval where the scenes are built by introducing rigid transformations and noise. Moreover, these datasets also allow for evaluating the descriptor on varying quality of point clouds. The *Random Views* and *Retrieval* have been derived from high resolution models. Comparatively, the *Laser Scanner* and *LIDAR* datasets can be categorized as medium and low quality respectively.

4.2 Performance Evaluation

We evaluate the proposed descriptor for descriptiveness, compactness and efficiency against the best performing descriptors in the comparative analysis of [10]. These descriptors are Fast Point Feature Histogram (FPFH) [29], Signature of Histogram of Orientations (SHOT) [31, 37], Unique Shape Context (USC) [37] and Rotational Projection Statistics (RoPS) [11]. As discussed in Section 3.1, Intrinsic Shape Signature (ISS) [40] keypoint detector is used. We use the implementations available in Point Cloud Library (PCL) for the keypoint detector ISS and descriptors FPFH, SHOT, USC and RoPS. Unless specified, the default parameters from the corresponding implementation are used for various techniques. The computation of best-fitting plane was performed in parallel on a GPU. For performing the fast feature matching (Section 3.5.1), the number of parallel search trees has been fixed to 3, branching factor to 16 and maximum leaf nodes to 150.

4.2.1 Descriptiveness

Descriptiveness is measured using Area under the *Precision-Recall* curve (PRC). The PRC is generated using the following steps. Firstly, the keypoints are detected from the considered scenes and models. The keypoints are then described using various descriptors. For a fair comparison between the feature matching capability of the floating point and binary descriptors, we index the features using the technique described in Section 3.5.1 which can be applied consistently across floating point and binary descriptors. Due to this, the results of our experiments are slightly different from those reported in [10]. However, the relative performance results

are still valid even though the absolute numbers change by a small margin (7.2% on an average). The number of matches returned from the search tree depends upon the search precision τ . A linear scan is performed on the matches obtained from the search tree and following [38], a match is considered correct if the matched features belong to the corresponding objects in the scene and model point clouds, and the matched keypoint lies within a small neighbourhood of the ground-truth. This neighbourhood is defined by a sphere of radius 2 mesh resolution (mr) [16] and centered at the ground-truth keypoint. The *precision* and *recall* are then computed as

$$Precision = \frac{\#CorrectMatches}{\#TotalMatches} \quad (8)$$

$$Recall = \frac{\#CorrectMatches}{\#CorrespondingMatches} \quad (9)$$

We then vary τ from 0 to 1 and compute the Area under the PR curve (AUC_{PR}). The results are shown in Table 2. We denote the proposed binary descriptor with N neighbours as 3DBS- N . The ranking for FPFH, USC, RoPS and SHOT are consistent with those reported in [10]. It can be observed that 3DBS-32 outperforms other descriptors on *Retrieval* dataset while 3DBS-64 outperforms on all datasets except *LIDAR* and *Random Views*. Moreover, the magnitude of difference between AUC_{PR} of RoPS and 3DBS with other descriptors is approximately 80%. This observation is important since it shows that the proposed technique has good performance on low resolution point cloud while also performing consistently for various transformations in other datasets. Another observation that can be made is that 3DBS-64 consistently performs better than 3DBS-32. This is expected since 64 nearest neighbours span a larger local surface around a keypoint making the descriptor more robust to occlusion and clutter. This observation is in line with the performance of the other descriptors when an increase in the radius of the support region increases the performance of the descriptor [10, 38]

4.2.2 Compactness

Compactness of a descriptor is a measure to compare descriptors when memory footprint and storage requirements become important. Compactness is given as the

$$Compactness = \frac{AverageAUC_{PR}}{\#Floats_{descriptor}} \quad (10)$$

The number of floats (length) in each of the considered descriptors is shown in Table 3 with 32 bits per float. The results are graphically shown in Figure 3. It can be seen that 3DBS-32 is highly compact being close to FPFH. Moreover, 3DBS-64 has lower compactness than FPFH and is close to SHOT and RoPS. It would be important to note here that 3DBS can be made more compact by using bit vector compression schemes. Although the compactness measure is provided for completeness of comparative analysis with other popular 3D descriptors, it must be noted that the binary descriptors are not stored in memory as floats for computation. Therefore, compactness does not impact the matching efficiency of the proposed binary descriptor.

4.2.3 Efficiency

Table 2: Area under the Precision-recall curve (AUC_{PR})

Dataset/Descriptor	FPFH	SHOT	USC	RoPS	3DBS-32	3DBS-64
Random Views	0.24334	0.24799	0.05982	0.20001	0.22341	0.24010
Laser Scanner	0.07341	0.05018	0.01103	0.16310	0.155953	0.16389
LIDAR	0.00198	0.00136	0.00164	0.00521	0.004682	0.004897
Retrieval	0.49319	0.56114	0.59521	0.52457	0.61109	0.64120

Figure 3: Compactness of Descriptors

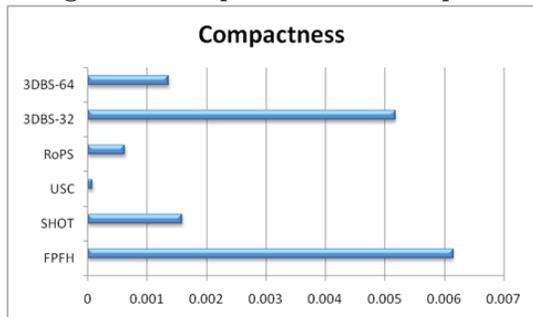


Figure 4: Matching time comparison

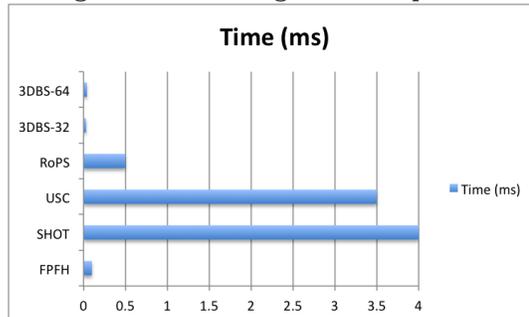


Table 3: Descriptor Length

Descriptor	# of floats
FPFH	33
SHOT	135
USC	1980
RoPS	352
3DBS-32	48
3DBS-64	192

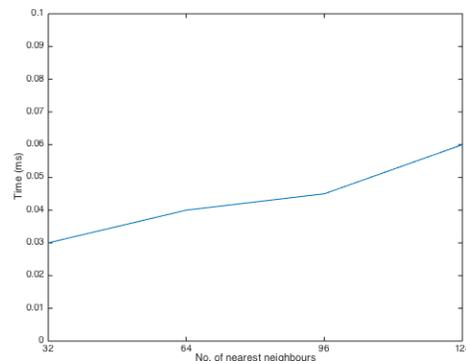
The major advantage of binary descriptors is that they can be matched extremely fast. To evaluate the descriptor matching time, the average matching time of keypoints on LIDAR dataset is reported in Figure 4. It can be seen that 3DBS is nearly 10 times as efficient than FPFH while almost three order of magnitude faster than other descriptors. This speed-up can be attributed to two factors. Firstly, the matching of binary vectors is by design faster than matching floating point vectors. Secondly, as discussed in Section 3.5, we leverage the built-in POPCNT instruction in GNU C Compiler providing tremendous computational efficiency in matching binary vectors.

It can also be observed that the size of the proposed descriptor quadruples when the number of nearest neighbours doubles. As discussed previously, the number of nearest neighbours impacts the descriptiveness. Therefore, in Figure 5 we show the descriptor retrieval and matching time by gradually increasing the number of nearest neighbours. As can be seen, the increase in matching time is nearly sublinear when the number of nearest neighbours are doubled.

5. CONCLUSION

A novel 3D binary descriptor termed as 3D Binary Signature was proposed. The descriptor was based upon aligning the local surface as per a Local Reference Frame. The local surface has been identified with a nearest neighbour approach with angular constraint. This is in contrast with previous descriptors where a spherical region was used. The neighbour-

Figure 5: Performance on increasing NN



hood was characterized with projections of surface normals and encoding them as binary vector. We showed that the proposed descriptor outperforms the state of the art methods on various standard datasets. It is highly compact and nearly 3 – 10 times faster than traditional descriptors, while demonstrating comparable or better descriptiveness.

6. REFERENCES

- [1] L. A. Alexandre. 3d descriptors for object and category recognition: a comparative evaluation. In *Workshop on Color-Depth Camera Fusion in Robotics at the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Vilamoura, Portugal, volume 1, page 7. Citeseer, 2012.
- [2] H. Bay, T. Tuytelaars, and L. Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006.
- [3] M. Calonder, V. Lepetit, M. Ozuysal, T. Trzcinski, C. Strecha, and P. Fua. Brief: Computing a local binary descriptor very fast. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1281–1298, 2012.

- [4] T.-H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma. Pcanet: A simple deep learning baseline for image classification? *IEEE Transactions on Image Processing*, 24(12):5017–5032, 2015.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009.
- [6] A. Flint, A. R. Dick, and A. Van Den Hengel. Thrift: Local 3d structure recognition. In *dicta*, volume 7, pages 182–188, 2007.
- [7] A. Frome, D. Huber, R. Kolluri, T. Bülow, and J. Malik. Recognizing objects in range data using regional point descriptors. In *European conference on computer vision*, pages 224–237. Springer, 2004.
- [8] Y. Gao and Q. Dai. *View-based 3-D object retrieval*. Morgan Kaufmann, 2014.
- [9] Y. Guo, M. Bennamoun, F. Sohel, M. Lu, and J. Wan. 3d object recognition in cluttered scenes with local surface features: a survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(11):2270–2287, 2014.
- [10] Y. Guo, M. Bennamoun, F. Sohel, M. Lu, J. Wan, and N. M. Kwok. A comprehensive performance evaluation of 3d local feature descriptors. *International Journal of Computer Vision*, 116(1):66–89, 2016.
- [11] Y. Guo, F. Sohel, M. Bennamoun, M. Lu, and J. Wan. Rotational projection statistics for 3d local surface description and object recognition. *International journal of computer vision*, 105(1):63–86, 2013.
- [12] Y. Guo, F. Sohel, M. Bennamoun, J. Wan, and M. Lu. An accurate and robust range image registration algorithm for 3d object modeling. *IEEE Transactions on Multimedia*, 16(5):1377–1390, 2014.
- [13] Y. Guo, J. Zhang, M. Lu, J. Wan, and Y. Ma. Benchmark datasets for 3d computer vision. In *2014 9th IEEE Conference on Industrial Electronics and Applications*, pages 1846–1851. IEEE, 2014.
- [14] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik. Learning rich features from rgb-d images for object detection and segmentation. In *European Conference on Computer Vision*, pages 345–360. Springer, 2014.
- [15] K. Jiang, Q. Que, and B. Kulis. Revisiting kernelized locality-sensitive hashing for improved large-scale image retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4933–4941, 2015.
- [16] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. *IEEE Transactions on pattern analysis and machine intelligence*, 21(5):433–449, 1999.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [18] B. Kulis and T. Darrell. Learning to hash with binary reconstructive embeddings. In *Advances in neural information processing systems*, pages 1042–1050, 2009.
- [19] S. Leutenegger, M. Chli, and R. Y. Siegwart. Brisk: Binary robust invariant scalable keypoints. In *2011 International conference on computer vision*, pages 2548–2555. IEEE, 2011.
- [20] L. Linsen. *Point cloud representation*. Univ., Fak. für Informatik, Bibliothek, 2001.
- [21] D. G. Lowe. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. Ieee, 1999.
- [22] D. Maturana and S. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*, pages 922–928. IEEE, 2015.
- [23] A. Mian, M. Bennamoun, and R. Owens. On the repeatability and quality of keypoints for local feature-based 3d object retrieval from cluttered scenes. *International Journal of Computer Vision*, 89(2-3):348–361, 2010.
- [24] M. Muja and D. G. Lowe. Fast matching of binary features. In *Computer and Robot Vision (CRV), 2012 Ninth Conference on*, pages 404–410. IEEE, 2012.
- [25] A. Petrelli and L. Di Stefano. On the repeatability of the local reference frame for partial shape matching. In *2011 International Conference on Computer Vision*, pages 2244–2251. IEEE, 2011.
- [26] S. M. Prakhya, B. Liu, and W. Lin. B-shot: A binary feature descriptor for fast and efficient keypoint matching on 3d point clouds. In *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*, pages 1929–1934. IEEE, 2015.
- [27] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski. Orb: An efficient alternative to sift or surf. In *2011 International conference on computer vision*, pages 2564–2571. IEEE, 2011.
- [28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [29] R. B. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (fpfh) for 3d registration. In *Robotics and Automation, 2009. ICRA '09. IEEE International Conference on*, pages 3212–3217. IEEE, 2009.
- [30] R. B. Rusu and S. Cousins. 3D is here: Point Cloud Library (PCL). In *IEEE International Conference on Robotics and Automation (ICRA)*, Shanghai, China, May 9-13 2011.
- [31] S. Salti, F. Tombari, and L. Di Stefano. Shot: unique signatures of histograms for surface and texture description. *Computer Vision and Image Understanding*, 125:251–264, 2014.
- [32] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek. Image classification with the fisher vector: Theory and practice. *International journal of computer vision*, 105(3):222–245, 2013.
- [33] I. Sipiran and B. Bustos. Harris 3d: a robust extension of the harris operator for interest point detection on 3d meshes. *The Visual Computer*, 27(11):963–976, 2011.
- [34] R. Socher, B. Huval, B. Bath, C. D. Manning, and A. Y. Ng. Convolutional-recursive deep learning for 3d object classification. In *Advances in Neural*

- Information Processing Systems*, pages 665–673, 2012.
- [35] S. Srivastava, P. Mukherjee, and B. Lall. Characterizing objects with sika features for multiclass classification. *Applied Soft Computing*, 2015.
- [36] J. Sun, M. Ovsjanikov, and L. Guibas. A concise and provably informative multi-scale signature based on heat diffusion. In *Computer graphics forum*, volume 28, pages 1383–1392. Wiley Online Library, 2009.
- [37] F. Tombari, S. Salti, and L. Di Stefano. Unique shape context for 3d data description. In *Proceedings of the ACM workshop on 3D object retrieval*, pages 57–62. ACM, 2010.
- [38] F. Tombari, S. Salti, and L. Di Stefano. Performance evaluation of 3d keypoint detectors. *International Journal of Computer Vision*, 102(1-3):198–220, 2013.
- [39] A. Zaharescu, E. Boyer, K. Varanasi, and R. Horaud. Surface feature detection and description with applications to mesh matching. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 373–380. IEEE, 2009.
- [40] Z. Zhang. Microsoft kinect sensor and its effect. *IEEE multimedia*, 19(2):4–10, 2012.
- [41] Y. Zhong. Intrinsic shape signatures: A shape descriptor for 3d object recognition. In *Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on*, pages 689–696. IEEE, 2009.