



**A Comparative Study on the Effect of Interest  
Point Methods and Descriptors on  
Feature-Feature Matching for Object Query in  
Point Clouds**

MICHAL STANIASZEK



## Abstract

Brief description of the aim of the project, the data used, some of the techniques used, some experimental results.

# Contents

<b>Contents</b>	<b>3</b>
<b>1 Introduction</b>	<b>5</b>
<b>2 Background</b>	<b>9</b>
2.1 Segmentation . . . . .	9
2.2 Interest Points and Saliency . . . . .	11
2.3 Descriptors . . . . .	12
2.3.1 2D . . . . .	12
2.3.2 3D . . . . .	13
2.3.2.1 Descriptors With Interest Point Extraction . . . . .	16
2.4 Storing and Querying Descriptors . . . . .	17
<b>3 System Development</b>	<b>19</b>
3.1 Preprocessing . . . . .	19
3.1.1 Downsampling . . . . .	19
3.1.2 Transformation and Trimming . . . . .	20
3.1.3 Normal Estimation . . . . .	22
3.1.4 Plane Extraction . . . . .	25
3.2 Interest Point Selection . . . . .	25
3.3 Feature Extraction . . . . .	25
3.4 Object Query . . . . .	25
<b>4 Experimental Results</b>	<b>27</b>
4.1 Segmentation . . . . .	27
4.2 Feature Extraction . . . . .	27
4.3 Object Query . . . . .	28
<b>5 Conclusion and Further Work</b>	<b>29</b>
<b>A Implementation</b>	<b>31</b>

Bibliography	<b>33</b>
--------------	-----------

# Chapter 1

## Introduction

Having a large amount of data is in most cases a good thing. Data, in an abstract sense, is the driving force behind the actions of every living thing, and as such holds great power. However, in order to make use of data, it is necessary to have some way of interacting with it in a useful way, and further processing it. While technologies for storing and interacting with data have been around for millennia, for the most part they were inconvenient and cumbersome. Writing allowed practically lossless transfer of information between generations, and is no doubt one of the most important inventions in the history of humanity. That being said, books are quite limited, especially when one wants to investigate a specific topic. Libraries are partial solutions to the problem, but most libraries don't possess all books in existence, and while indexes exist for a reason, it is still not easy to find what one is looking for.

With the internet and the immense amount of data available to its users, this problem of finding what one is looking for has been compounded, and good ways of getting around the problem have launched one of the most successful companies in history. At first, listing all of the early internet to create a database was a realistic proposition, and for some time this was a satisfactory solution. However, as the number of accessible data on the internet grew, the system became gradually more impractical. It could take minutes or even hours to get a result for a query, and the trawling of content caused network slowdowns [62, 67, 11]. Subsequent work in the area led to the development of search engines which were able to search for words in pages, and various innovations led this to become the very effective way of searching that we know today [12, 50].

While images have been on the internet since the early days, in recent years the advent of affordable digital cameras and the ubiquity of mobile phone cameras has led to hundreds of thousands of photographs being uploaded to the internet every minute [19, 29]. At its most basic, image search utilises the same techniques as text search, with information being extracted from metadata like tags, descriptions and keywords [31]. More recently, reverse image search has become more popular, allowing users to find similar images to an example by extracting information from

textures and trying to find other images which contain similar information [36]. There is still much information present in images that cannot currently be extracted and represented using image processing techniques, and this is an active research area.

An emerging method of data storage that will need to be searchable in the near future is 3D models and point clouds. 3D models have been used for a long time in computer games, medical imaging, and animation. More recently, developments in 3D printing have led to a growing number of websites which distribute models to use for printing [73]. For many years these sorts of models have been created using CAD programs, or in the case of object scanning, expensive time-of-flight cameras. The release of the Microsoft Kinect in late 2010 marked a turning point in the realm of 3D image processing, creating an affordable and effective method of gathering 3D data. Many research groups quickly purchased the hardware, and much work has been done in the area since. A 3D equivalent of the popular 2D image processing library OpenCV quickly came into existence for use with point clouds, as the data which comes from such 3D sensors is known [45, 51].

In this report, we will describe our approach the problem of retrieving from a data set objects that are similar to some object that we have provided, which we will call object query. In essence, we need to extract information from the data set and the objects that we are interested in which describes their properties in such a way that we can compare the descriptions to see if there are any similarities that imply the presence of an object in the data set. While the data set could be anything, in our particular case we have data from a project which studies long-term robot operation in office environments. The data set contains point clouds of a single office taken from the same position over a period of approximately a month. Some objects in the data set have also been labelled, so there is information about the positions of objects in the clouds.

While this project is not aiming to perform a specific task on an actual robotic system, within this context there are applications to which an object query system could be applied. Given a data set over long periods of time with clouds taken at various locations, the system could be used to track the motion of objects over time, and to provide information about where an object is likely to be at a certain time. One potential application is to help people find objects that have been misplaced.

The project will focus in particular on the implementation of a system which can perform object query. It will evaluate a number of standard methods for describing objects, and finding parts of objects that are particularly discriminative. While it is possible to describe objects as a whole, we will investigate the efficacy of using descriptions of small parts of the object to retrieve matches from the data set.

Something about why why the problem is interesting as opposed to model matching or somesuch

In chapter 2, we explain some concepts that are important to understand the work, provide background information on relevant parts of the image processing literature, and attempt to introduce the reader to previous work in similar areas. The specific techniques used in our system are described in chapter 3, and we also

discuss our reasoning behind certain steps of the process. Chapter A deals with the software, and how the system was implemented, along with some explanation of how the system works. Our experimental setup and the results of the experiments are described in chapter 4. We compare the quality of retrieval when different methods are used, and also investigate the time taken by the system under varying conditions. Finally, we summarise the system and our results in chapter 5.



# Chapter 2

## Background

In this chapter we will introduce some key ideas relating to the project, and papers which are related to what we are interested in doing. While some of the techniques mentioned here are not directly used in the implementation of our system, they can be useful for context, or to give examples of different ways of approaching problems in this area. We discuss methods of finding interesting regions in image and point cloud data, and how these regions can be represented using descriptors, along with some methods for storing descriptors in ways that make it easy and efficient to find similarities.

### 2.1 Segmentation

Segmentation encompasses techniques for splitting an image or a point cloud into different parts, or grouping similar parts — this is essentially two sides of the same coin. In terms of images, segmentation might be used to try to find background and foreground pixels, or for point clouds, to separate objects from the surfaces on which they are resting. There are many different types of methods in the area, which approach the problem from different starting points.

Superpixel clustering is the most common technique used for segmenting images. The intent is to create regions in which all pixels have some sort of meaningful relationship. Graph based algorithms treat pixels as nodes in a graph, where the weights on edges between nodes are related to the similarity between the connected pixels — intensity, proximity and so on [1]. The most simple method is to use a threshold on the edge weights to create superpixels. Fulkerson et al. use superpixel methods to identify object classes in images [24]. An algorithm which applies the idea of superpixels to point clouds to create supervoxels (3D pixels) has also been developed [48].

Gradient ascent based algorithms iteratively improve clusters until some criterion for convergence is reached [1]. Popularised by Comaniciu [17], mean shift was first introduced by Fukunaga [23] in 1975, and rediscovered by Cheng [13] in 1995. The technique finds stationary points in a density estimate of the feature space, for



(a) Superpixels size 64, 256 and 1024 computed using SLIC [1] (b) Supervoxel oversegmentation [48]

Figure 2.1: Examples of 2D and 3D superpixel segmentations

example pixel RGB values, and uses those points to define regions in the space by allocating pixels to them. One common way of computing a density estimate is to place Gaussians at the location of each pixel, and then to sum the values of all the Gaussians over the entire space. Pixels which follow the gradient of the density to the same stationary point are part of the same segment. An example can be seen in Figure 2.2.

Random Sample Consensus (RANSAC) is a technique which uses shape models to find ideal models in noisy data. Points in the data set are randomly sampled, and used to construct a shape. For example, in the case of a line, two points are sampled, and define the line. Distances from points in the data set to the model defined by the randomly sampled points are then computed to find points which are inliers to the model. This number is stored, and the process repeated a number of times. At the end of the process, the model with the largest number of inliers is returned [20]. RANSAC can be applied to segmentation tasks by using it to find planes, cylinders, spheres and so on in point clouds. In the case of planes this is particularly useful, as they are usually not part of objects of interest, mostly making up walls, floors or surfaces on which interesting objects rest. By removing the points corresponding to these uninteresting surfaces, it should be possible to work only with parts of clouds that contain objects of interest.

Several extensions to RANSAC have been proposed. Maximum Likelihood Estimation Sample Consensus (MLESAC) chooses a solution that maximises the likelihood of the model instead of just the number of inliers [72]. M-estimator Sample Consensus (MSAC) uses a different cost function to the original implementation, additionally scoring the inliers depending on how well they fit the data [72]. The Progressive Sample Consensus (PROSAC) uses prior information about the like-



Figure 2.2: Visualisation of mean shift [17]. a) First two components of image pixels in LUV space. b) Decomposition found by running mean shift. c) Trajectories of mean shift over the density estimate.

lihood of input data being an inlier or an outlier to limit the sampling pool and greatly reduce computation cost [15].

## 2.2 Interest Points and Saliency

Sipiran and Bustos extend the popular Harris detector [28] to 3D [60]. Knopp et al. extend the SURF detector to 3D [34].

Shilane and Funkhouser introduce a distinctiveness measure over classes of meshed objects [58].

A multi-scale signature defined by the heat diffusion properties of an object called the Heat Kernel Signature (HKS) [68] is used in [47] to retrieve shapes.

[76]

[63] extended to 3d by adding normal direction variation to intensity variation

## 2.3 Descriptors

The problem of describing regions of an image in a compact and useful manner has been studied for a long time in the computer vision community. For any given point in an image, we would like to create a description which can be used to represent the region around the point in some way. This descriptor, or feature, can then be compared to other descriptors to see if there is some similarity. If the similarity is within a given threshold, then we can assume that the points represented by the two descriptors come from the same object, or represent the same thing in both images. Thus, it is important to create features which are distinct for different regions. In addition, since objects move around and can be seen from different sides, or in different lights, an attractive property of descriptors is to give similar results for the same region which has been transformed in some way. In practice, this is quite difficult to achieve.

### 2.3.1 2D

While 2D descriptors are not directly usable on point clouds, the ideas that they use to give effective results can be transferred over to use for 3D description.

The Laplacian of Gaussians was introduced by Lindeberg, and uses derivatives combined with some other techniques to select interest points. [37]. This paper also introduces the concept of automatic scale selection for feature detection, which has played an important part in the field since then. The scale of features can be investigated by blurring an image using a Gaussian kernel — higher standard deviation blurs the image more, resulting in the removal of small scale features.

Even today the Scale Invariant Feature Transform (SIFT) is among the most popular descriptors for 2D images. It is invariant to scale and rotation, and is robust to some variation in affine distortion, viewpoint and illumination, and is distinctive, allowing for correct matching of single features in large databases. There are several stages of computation. Extrema are found in different scales to find points invariant to scale and orientation. Keypoints are selected at the extrema based on their stability. Image gradients at the keypoint are used to define its orientation for future computations. The image gradients are then transformed into a local descriptor vector with length 128 [38].

Mikolajczyk and Schmid [40] introduce the Harris-Laplace detector which is an improvement on SIFT [38] and the Laplacian of Gaussians [37] in the sense that it is able to deal with affine transformations. They do not, however, introduce a new type of descriptor to go with the point selection.

Speeded-Up Robust Features (SURF) is a more recent descriptor which can be computed and compared much faster than most other descriptors. It makes use of integral images, which replace pixels in an image or image patch with a

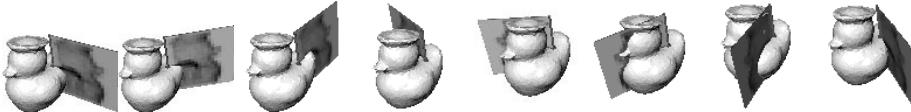


Figure 2.3: Frames from construction of a spin image [33]. The image plane spins around the oriented point normal and accumulates points.

cumulative sum of the pixel intensities over the rows and columns. This allows for fast computation of pixel intensities in an area of the image. SURF takes some ideas from SIFT, using the spatial distribution of gradients as a descriptor, but integrates over the gradients instead of using individual values, which makes it more robust to noise. The resulting descriptor is a 64 element vector, which means that it is also faster to compare than SIFT [4].

### 2.3.2 3D

One early descriptor which remains popular is the spin image. The descriptor is generated from a mesh model at oriented points with a surface normal. A plane intersecting the normal with a certain width and height is rotated around the normal, forming a cylinder. The plane is separated into bins. The bins accumulate the number of points which pass through a certain bin during the rotation. The resulting 2D image is the descriptor. By varying the width of the plane the region which defines the descriptor can be modified. A small width will give a local descriptor, while a large width will give a descriptor for the whole image [33, 32]. Figure 2.3 shows a visualisation of how the image is generated.

The Ensemble of Shape Functions (ESF) descriptor introduced in [74] by Wohlkinger and Vincze combines the Shape Distribution approach introduced by [46] along with some extensions proposed in [30]. It also makes use of their voxel-based distance measure from [75]. Pairs or triples of points are sampled from segmented partial clouds of objects, and histograms are created by extracting information such as distance, angle, ratios, and whether points are inside or outside (or a mix) of the model. See Figure 2.5.

The Point Feature Histogram (PFH) was introduced by Rusu et al. in [56]. It creates descriptors based on the angles between a point on a surface and  $k$  points close to it. The Fast Point Feature Histogram (FPFH) improved the speed of computation, and allowed the use of the descriptor in real time [54]. The Viewpoint Feature Histogram (VFH) extended the FPFH by adding viewpoint information to the histogram by computing statistics of surface normals relative to the viewpoint [55]. It also improved the speed of the FPFH. The clustered version (CVFH) further improved the viewpoint technique by mitigating the effect of missing parts and extending it to facilitate estimation of the rotation of objects [2].

Bo et al. develop the kernel descriptor initially created for RGB images for use on depth images and point clouds. The kernels are used to describe size, shape

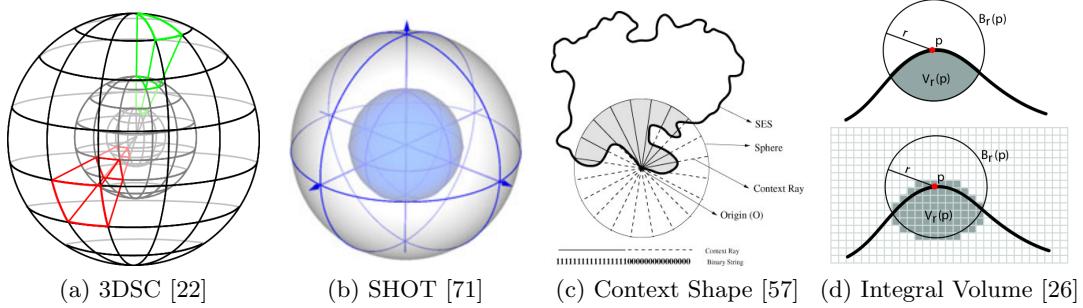


Figure 2.4: Visualisation of spherical descriptors.

and edge features. Local features are combined to object-level features . Kernel descriptors avoid the need to quantise attributes. Similarity is instead defined by a match kernel [9], which improves recognition accuracy [8].

The point pair feature describes the relation between two oriented points on a model. This means that it does not depend so much on the quality and resolution of the model data. The model is described by grouping the point pair features of the model, providing a global distribution of all the features on the model surface [18].

3D Shape Context (3DSC) is an extension of the original Shape Context descriptor for 2D images [6]. A sphere is placed at a point, and its “top” is oriented to match the direction of the normal at the point. Bins are created within the sphere by equally spaced boundaries in the azimuth and elevation, and logarithmically spaced boundaries in the radial dimension (Figure 2.4a). The logarithmic spacing means that shape distortions far from the basis point have less effect on the descriptor. Each bin accumulates a weighted count based on the volume of the bin and local point density [22]. 3DSC does not compute a local reference frame — the vector of the azimuth is chosen randomly, and subdivisions computed from that. This means that a number of descriptors equal to the number of azimuth divisions need to be computed and stored in order to compensate, and the matching process is complicated as a result. The Unique Shape Context (USC) solves this problem by defining a local reference frame and using the directions of that reference frame to subdivide the sphere [70].

The Signature of Histograms of Orientations (SHOT) descriptor improves on 3DSC by taking inspiration from SIFT and making extensive use of histograms. The sphere is split into 32 volumes: 8 azimuth regions, 2 elevation and 2 radial (Figure 2.4b). A local histogram is computed in each of the regions, using the angle between the normal of points and the feature point. The local histograms are then combined to form the final descriptor [71]. The authors also extend the descriptor to include colour (COLORSHOT) [69].

The Rotation Invariant Feature Transform (RIFT) is a generalisation of SIFT. Using intensity values computed at each point from the RGB values, a gradient is

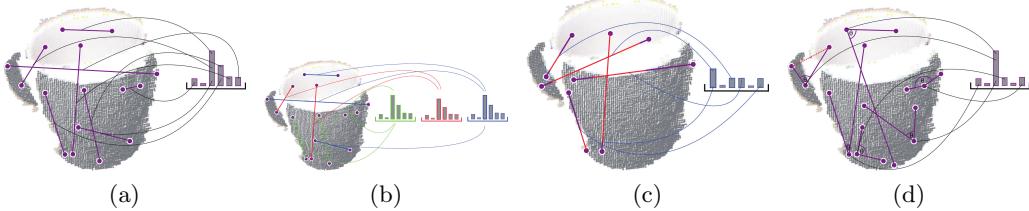


Figure 2.5: Examples of the measures used to construct the Ensemble of Shape Functions histograms of [74]. a) Distance between points. b) Whether the points are on or off the model, or mixed. c) Ratio of line segments on and off the surface of the model. d) Angle between pairs of lines.

computed. Concentric rings are placed around the initial point, and a histogram of the gradient orientations is created for points within each ring. The orientation of the gradient is computed relative to the line from the central point so that the descriptor is rotation invariant. The descriptor is 2D — one dimension is the distance, the other the gradient angle. The distance between two descriptors is measured using the Earth Mover’s Distance (EMD), which is a measure of the distance between two probability distributions [35].

Multi-scale descriptors are useful as they can be used to characterise regions of varying size. Cipriano et al. introduce such a descriptor for use on meshes [16]. It captures the statistics of the shape of the neighbourhood of a vertex by fitting a quadratic surface to it. Vertices in the region are weighted based on distance from an initial vertex, and a plane is constructed using a weighted average of the face normals. The parameters of the quadratic are then used to find its principle curvatures, which make up the descriptor.

Work in protein-protein docking also uses 3D descriptors to help with simulations of an otherwise lengthy and complex process. The Surface Histogram is introduced by Gu et al. [27], and uses the local geometry around two points with specific normals on the surface of a protein. A coordinate system is defined by the two points and the line between them, and a rectangular voxel grid is defined around the points. The grid is then marked in locations where the surface crosses the grid, and a 2D image is constructed by squashing the data onto one of the axes. The descriptor is designed to immediately give a potential pose for the docking.

Another example of a shape descriptor from biology is the Context Shape [57]. A sphere is centred on a point, and rays are projected from this point to points evenly distributed on the surface of the sphere (Figure 2.4c). Each of the rays is divided into segments, with a binary value associated with each segment depending on whether the segment is inside or outside the protein. To compare the descriptor, a rotation is applied to match the rays, and a volume of overlap is computed based on matching bits in the rays.

The splash descriptor was introduced by Stein et al. [66]. A point on the surface

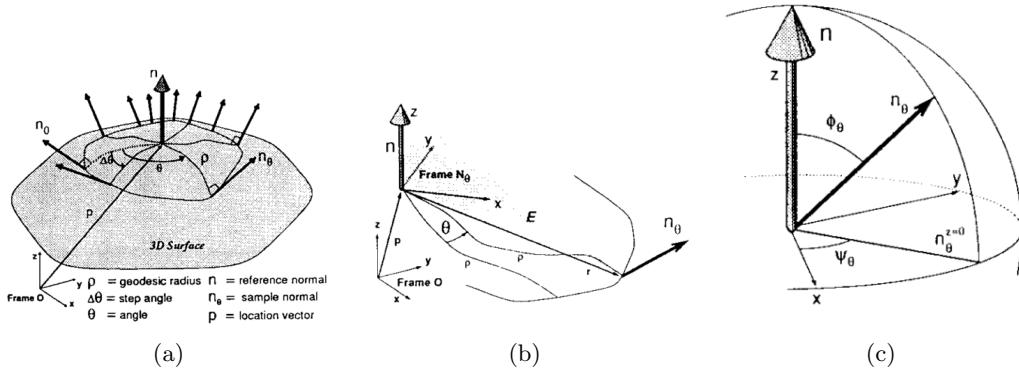


Figure 2.6: Splash descriptor [66]. a) shows the splash and normals around it. b) and c) show how the additional angles are defined.

with a given surface normal (the reference normal) is chosen, and a slice around that with some geodesic radius (distance along the surface) is computed. Points on the circle are selected using some angle step, and the normal at that point is determined. A super splash is when this process is repeated for several different radii. For each normal on the circle, additional angles between it and a coordinate system centred on the reference normal are computed. These angles and the angle around the circle are then mapped into a 3D space, where polygonal approximation is made, connecting each point with a straight line. Some additional computation is done to allow the encoded polygons to act as a hash. Figure 2.6 shows part of the formulation.

Point Signatures are similar to the splash descriptor in the sense that they both sample points on a circle [14]. This descriptor again selects a reference normal, and has a specific radius. This time, the radius defines a sphere around the point. The intersection of the surface with the sphere is a 3D space curve. The orientation of the curve is defined by fitting a plane to it. The distances between the space curve and the fitted plane at sampled points define the signature of the reference point. These signatures can be compared by lining them up and checking whether the query falls within the tolerance band of previous signatures. Figure 2.7 shows signatures from various surfaces.

### 2.3.2.1 Descriptors With Interest Point Extraction

While many descriptors designed for 2D applications also select interest points during an initial step in the process, the 3D descriptors that we have mentioned above do not automatically find locations in the cloud which are good points at which to compute descriptors.

The Normal Aligned Radial Feature (NARF) is an interest point extraction method with a feature descriptor. A score for the image points is determined based

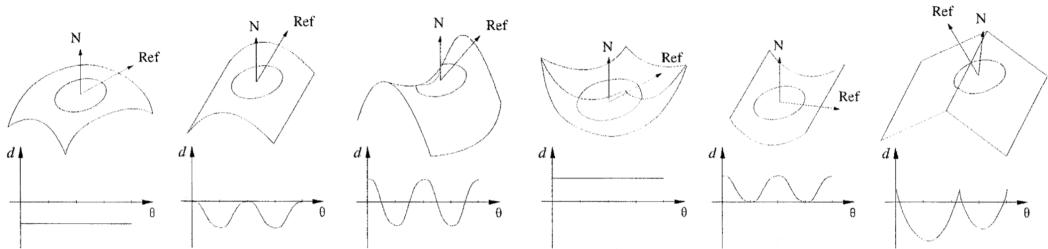


Figure 2.7: Examples of the point signature responses to different surfaces [14].  $d$  is the distance from the reference vector to the space curve defined by the intersection of the surface with a sphere centred at  $N$ . Ref rotates about  $N$ .

on the surface changes at the point, and information about borders. An interest value is computed from this based on the score of the surrounding points. Smoothing is applied, and non-maximum suppression is applied to find the final interest points. To compute the descriptor, rays are projected over the range image from the centre at certain intervals. The intensities of cells lying under the ray are weighted based on their distance from the centre, and a normalised weighted average of the pairwise difference of cells is used to define each element of the descriptor vector, which has a length equal to the number of rays [64]. The method is an improvement on a previous paper by the authors [65]. A problem with this method is that it uses range images directly. Point clouds can be used to generate range images by looking at them from different viewpoints, but this adds complexity to the method.

The integral volume descriptor is interesting as it combines interest point selection and description into one. The descriptor is defined as the volume of the intersection of a sphere centred at a point on the surface of an object with the inside of the object (Figure 2.4d). Interest points are selected by histogramming the descriptor values, identifying bins with a number of points less than a specified values, and selecting points from these bins. To ensure features are properly spaced, points in a certain radius of already selected points cannot be used. By modifying the radius of the sphere used to generate the descriptor, interest points at different scales can be selected [26].

## 2.4 Storing and Querying Descriptors

There are several techniques for storing and querying descriptors, mostly based on some form of tree. Recently, the k-d tree[7, 21] has been used for efficient approximate matching with either an error bound [3], where there is a bound placed on the error between the true nearest neighbour and the one found, or a time bound [5], where the search is stopped after examining a certain number of leaf nodes. Further improvements on the k-d tree are introduced in [59], where multiple randomised trees are used to optimise the search. A priority search tree algorithm is introduced in [42], which appears to be very effective. This may be the same

one as in [41]. The algorithm in the last two papers has been integrated into PCL, which is useful.

A different approach to nearest neighbour search is the balltree, which uses hyperspheres in a hierarchy to enclose points in the space [44]. Unlike the k-d tree, regions on the same level of the tree are allowed to intersect, and do not need to partition the whole space, which gives the balltree its representative power.

The vocabulary tree [43] makes use of techniques from document search to index images. Using  $k$ -means clustering, construction stage creates a hierarchical quantisation of the image patch descriptors. In the query phase, descriptors are compared to the cluster centres, and go down the tree until a leaf is reached. The path through the tree is used as a scoring measure to present retrieval results.

Philbin et al. [49] show that flat (single-level)  $k$ -means clustering can be scaled to large vocabulary sizes if approximate nearest neighbour methods are used. Early systems for image retrieval used a flat clustering scheme, which could not scale to large vocabularies [61]. The paper also introduces a re-ranking method which uses spatial correspondences, which improves the retrieval quality.

Boiman et al. [10] introduce the Naive Bayes Nearest Neighbour (NBNN) classifier. It uses nearest neighbour distances in the space of descriptors instead of images, computing “image-to-class” distances without quantising the descriptors. In general, quantisation allows for dimensionality reduction, at the expense of the discriminative power of descriptors. NBNN “can exploit the discriminative power of both (few) high and (many) low informative descriptors”. The problem here is that the classes must be known beforehand, and in our case we do not have that information. The local NBNN [39] does not do the search based on classes. Instead, all the descriptors are merged into a k-d tree on which approximate  $k$ -NN is run to find descriptors in the local region of a query descriptor. A distance to classes not present in the  $k$ -NN region is approximated by the distance to the  $k+1$ th neighbour.

Funkhouser and Shilane present a method for querying a database of 3D objects represented by local shape features [25]. Partial matches (correspondences) are stored in a priority queue sorted by geometric deformation and the feature similarity. This means that only objects in the database with a high probability of being a match need to be processed.

Some work has been done on optimising the retrieval of relevant images by learning from user input [52]. When retrieved images are presented, the user ranks them in terms of relevance, and this rank is then used to improve the relevance of future searches.

# Chapter 3

# System Development

In this chapter, we will describe the development of the object query system, and explain the reasoning behind our choices. The specific methods that are used will also be described in more detail, and examples of the output of different parts of the system will be shown.

## 3.1 Preprocessing

The first step in the object query system is to perform some preprocessing on the clouds in the data set — while not strictly necessary, there are some benefits to doing so, chief of which is a reduction in computation time. The data set that we have consists of around 80 clouds of a single room, taken at different times during different days of approximately a month of time. The clouds are made up of a number of intermediate frames, which are registered into a complete cloud. The robot used to collect the clouds takes several sweeps of the room, changing the angle of the camera after each sweep. The clouds are constructed using frames taken from a sweep where the camera is pointing slightly below the horizontal. Examples of the raw clouds can be seen in Figure 3.1.

### 3.1.1 Downsampling

In their merged forms, the clouds on average contain approximately 4,300,000 points for a room which is around 4m wide, 5.5m deep and 3m high. This number of points does not actually provide us with much additional information, since the intermediate frames all have the same resolution. As such, we can safely downsample the cloud to get a more reasonable number of points.

To downsample, we make use of a voxel grid, which splits the 3D space in which the cloud sits into smaller subspaces of equal size called voxels. The width, height and depth of voxels in the space can be specified, but we are interested in keeping all dimensions the same resolution, and so we specify the parameters so that each voxel is a cube. At this stage, we would like to perform a simple downsampling



Figure 3.1: Sample raw cloud viewed from several angles

to reduce the number of points, but we wish to keep small details in the cloud — something in the realm of a 1cm resolution is ideal in this case.

Downsampling with a 1cm resolution gives a reduction in size of the clouds of on average 78%, to approximately 950,000 points. Figure 3.2 shows the effect of the downsampling. While there is slight degradation of the textures, this is to some extent a visual effect which is viewpoint dependent. Most of the structure in the cloud is retained, which is key. This step is important, as it greatly affects the speed of computation of subsequent steps in the system, but it is a trade off. If the downsampling resolution is too low, then we lose a lot of information about the surface structure of parts of the cloud, and this is likely to lead to worse performance when trying to find matches. How tolerant we are to low resolution also depends on the kinds of objects that we are interested in finding. If we do not care about smaller objects, then even with a lower resolution the results should still be satisfactory. However, a lower resolution likely means that it will be necessary to look at larger regions of space in order to describe points. We will investigate the effects of this in chapter 4.

Show some of the objects which are made up of slices due to viewing angle - downsampling can reduce this slicing effect.

### 3.1.2 Transformation and Trimming

Once the cloud has been downsampled, there is a little more that needs to be done in order to get the cloud into a convenient form. The raw data that we have has clouds which have their origin at the position of the camera while the room was being scanned. Our data is a subset of a larger dataset which contains clouds of

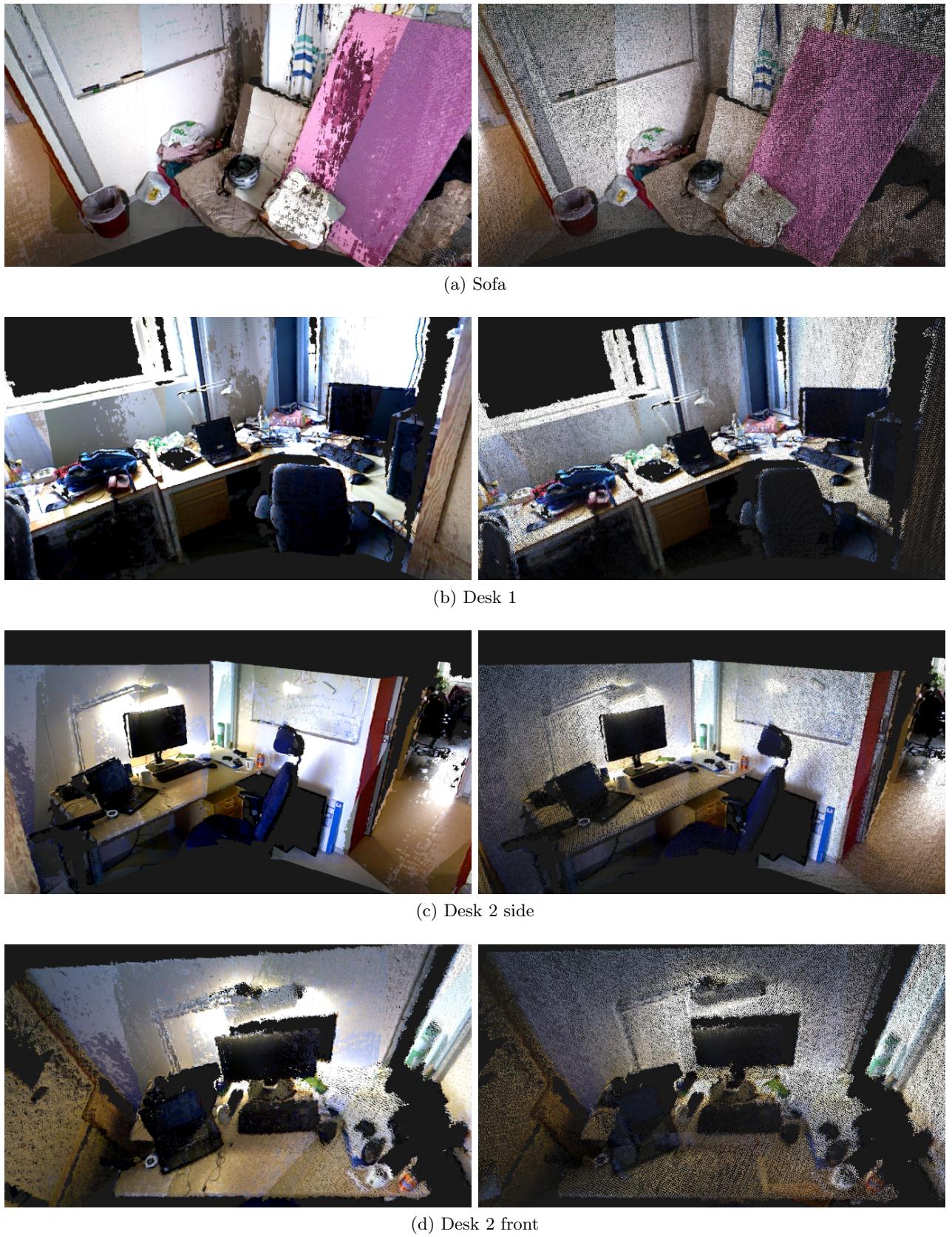


Figure 3.2: The effect of downsampling. The left column shows the original clouds, the right column clouds downsampled with voxel size of  $1\text{cm}^3$ .

more than one room — if we were to use the data without applying any additional transformations, all the clouds would sit on top of each other at the origin, whereas we would ideally like to have them in their true position relative to the origin. The robot collecting data knows its position, so this information is stored.

As mentioned before, each cloud is a combination of a number of intermediate frames, each of which has corresponding information about the pose of the camera when the frame was taken, which we can use to transform the complete cloud into its actual position in space.

An added benefit of this transformation is that it allows us to remove the floor and ceiling by using a simple thresholding filter on the  $z$  axis, as the floor of the cloud is now aligned with the  $x$ - $y$  plane of the global reference frame, as opposed to being aligned with the cloud's rotated reference frame (Figure 3.3) The threshold for the ceiling can be determined by measuring the ceiling height, and the floor is assumed to be a  $z = 0$ . We add a small offset to each of the values to ensure that the parts are correctly removed even if there is some noise.

Although we would like the system to be as generic as possible, the particular subset of clouds that we are using have a large number of points outside the room which do not give any useful information. To this end, we also include additional filters on the  $x$  axis to remove these points. Figure 3.4 shows the end result of this step.

### 3.1.3 Normal Estimation

In this step, normals are estimated for each point in the cloud. The normal at a point is the vector which is perpendicular to the curvature of the surface at that point. By estimating normals for clouds, we can get some more information about the surface structure of the cloud. Normals are used in several parts of the system, including by feature selection methods and features. They are also used in the plane extraction step to increase accuracy.

There are many ways of estimating normals, but the method we use is formulated as a least squares plane fitting problem, which is used to estimate the normal of the plane tangent to the surface at the point at which the normal is to be computed [53]. The computation gives an ambiguous result in terms of the sign of the normal. To correct for this, a viewpoint is needed, which serves to define what sign is used. Perhaps the most important thing to note is that the normal must be computed using points in a neighbourhood; either within a certain radius, or the nearest  $k$  points. The neighbourhood determines the scale factor that results. A small neighbourhood gives a small scale factor, and a large neighbourhood a large scale factor. A large scale factor can be bad if the objects that one is trying to examine have regions where the rate of change of surface curvature is high, such as at the corners of tables. It results in the smearing of edges and the suppression of fine detail [53]. Figure 3.5 shows an example of the effect of different neighbourhood sizes on the results.

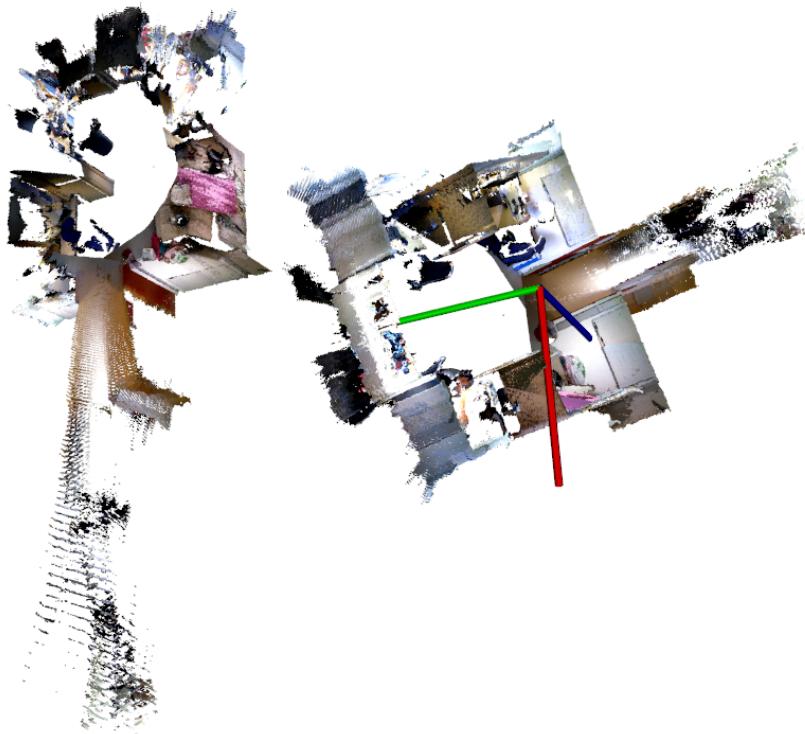


Figure 3.3: Original cloud and the transformed cloud. The original cloud is on the right. The coordinate axis shows the global reference frame — none of the axes are aligned for the original cloud, but the transformed cloud is well aligned with the  $x$ - $y$  plane.

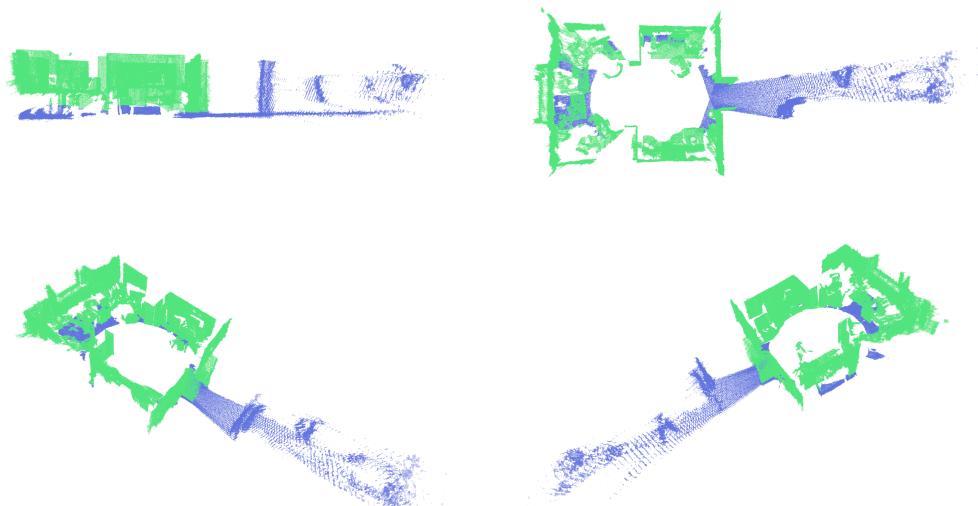


Figure 3.4: Result of trimming step. Transformed cloud is blue, trimmed is green.

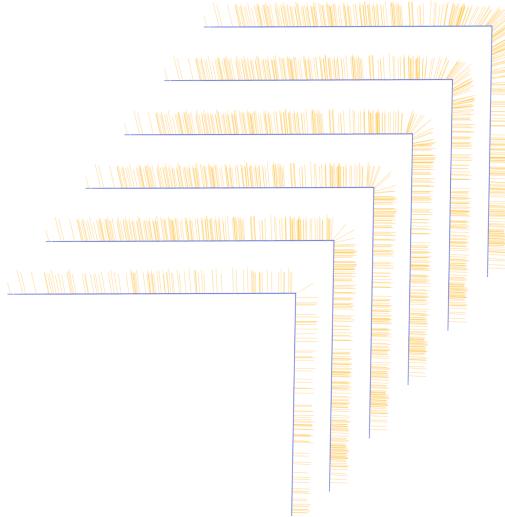


Figure 3.5: Example of the smoothing effect of normal computation radius. From bottom to top, 0.01, 0.025, 0.5, 0.2, 0.25, 0.5cm radius. Normals are indicated by orange lines. Note the tendency of normals with higher radius to tilt as they approach the corner. Normals on the top section are slightly skewed due to perspective.

Normal and plane extraction section order depends partly on the order in which it is done. Talk about how normal extraction can help with plane extraction, if I do that.

### 3.1.4 Plane Extraction

Basic description of RANSAC or whichever of its extensions I end up using.

## 3.2 Interest Point Selection

### 3.3 Feature Extraction

Talk specifically about each of the features used and their good/bad points, and how the features extraction is done (not particularly complex, probably)

### 3.4 Object Query

Query methods, how the database and query object are compared.

## Chapter 4

# Experimental Results

Description of experiments and experimental results with different configurations of the system.

### 4.1 Segmentation

Apply different segmentation methods/parts of the segmentation method, and see how they affect the overall system performance, as well as the time taken to complete the segmentation.

### 4.2 Feature Extraction

Same as above, but for feature extraction. Try different features and see which give best performance. Consider normal extraction cost for those features which require it.

	Original	Downsampled	Trimmed	Num planes	Points on planes	After preprocessing
Default	$4344879 \pm 125893$	$964004 \pm 114379$	$810469 \pm 97833$	$7.10 \pm 1.48$	$456273 \pm 56508$	$354195 \pm 77925$
dt 0.02	$4348571 \pm 124781$	$718345 \pm 246373$	$599503 \pm 208953$	$7.99 \pm 3.07$	$492154 \pm 152830$	$107348 \pm 64374$

Table 4.1: Average number of points in various stages of preprocessing.

	Load	Downsample	Trim	Normals	Planes	per plane time
Default	$1.91 \pm 0.58$	$1.08 \pm 0.31$	$20.12 \pm 13.56$	$30.95 \pm 10.88$	$221.62 \pm 130.12$	$31.36 \pm 16.09$
dt 0.02	$2.45 \pm 1.25$	$1.18 \pm 0.17$	$0.16 \pm 0.06$	$21.65 \pm 13.36$	$235.50 \pm 118.06$	$29.08 \pm 7.50$

Table 4.2: Average time taken for various stages of preprocessing.

### **4.3 Object Query**

This is evaluated in the previous two sections, but maybe there is more than one method of finding the query object. Multi-index tree vs single index tree, for example.

## **Chapter 5**

# **Conclusion and Further Work**

Describe what the project was about, what was achieved, summarise the experiments, describe what can be improved.



# **Appendix A**

# **Implementation**

Description of any parts of the implementation which are non standard or worth a mention. Perhaps talk about PCL and version control, as well as documentation and some other general software approaches to the system?



# Bibliography

- [1] Radhakrishna Achanta et al. “SLIC superpixels compared to state-of-the-art superpixel methods”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 34.11 (2012), pp. 2274–2282.
- [2] Aitor Aldoma et al. “CAD-model recognition and 6DOF pose estimation using 3D cues”. In: *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*. IEEE. 2011, pp. 585–592.
- [3] Sunil Arya et al. “An optimal algorithm for approximate nearest neighbor searching fixed dimensions”. In: *Journal of the ACM (JACM)* 45.6 (1998), pp. 891–923.
- [4] Herbert Bay et al. “Speeded-up robust features (SURF)”. In: *Computer vision and image understanding* 110.3 (2008), pp. 346–359.
- [5] Jeffrey S Beis and David G Lowe. “Shape indexing using approximate nearest-neighbour search in high-dimensional spaces”. In: *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*. IEEE. 1997, pp. 1000–1006.
- [6] Serge Belongie, Jitendra Malik, and Jan Puzicha. “Shape matching and object recognition using shape contexts”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 24.4 (2002), pp. 509–522.
- [7] Jon Louis Bentley. “Multidimensional binary search trees used for associative searching”. In: *Communications of the ACM* 18.9 (1975), pp. 509–517.
- [8] Liefeng Bo, Xiaofeng Ren, and Dieter Fox. “Depth kernel descriptors for object recognition”. In: *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*. IEEE. 2011, pp. 821–826.
- [9] Liefeng Bo, Xiaofeng Ren, and Dieter Fox. “Kernel descriptors for visual recognition”. In: *Advances in Neural Information Processing Systems*. 2010, pp. 244–252.
- [10] Oren Boiman, Eli Shechtman, and Michal Irani. “In defense of nearest-neighbor based image classification”. In: *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE. 2008, pp. 1–8.
- [11] C Mic Bowman, Peter B Danzig, and Michael F Schwartz. *Research problems for scalable internet resource discovery*. Tech. rep. DTIC Document, 1993.

- [12] Sergey Brin and Lawrence Page. “The anatomy of a large-scale hypertextual Web search engine”. In: *Computer networks and ISDN systems* 30.1 (1998), pp. 107–117.
- [13] Yizong Cheng. “Mean shift, mode seeking, and clustering”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 17.8 (1995), pp. 790–799.
- [14] Chin Seng Chua and Ray Jarvis. “Point signatures: A new representation for 3d object recognition”. In: *International Journal of Computer Vision* 25.1 (1997), pp. 63–85.
- [15] Ondrej Chum and Jiri Matas. “Matching with PROSAC-progressive sample consensus”. In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 1. IEEE. 2005, pp. 220–226.
- [16] Gregory Cipriano, George N Phillips, and Michael Gleicher. “Multi-scale surface descriptors”. In: *Visualization and Computer Graphics, IEEE Transactions on* 15.6 (2009), pp. 1201–1208.
- [17] Dorin Comaniciu and Peter Meer. “Mean shift: A robust approach toward feature space analysis”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 24.5 (2002), pp. 603–619.
- [18] Bertram Drost et al. “Model globally, match locally: Efficient and robust 3D object recognition”. In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE. 2010, pp. 998–1005.
- [19] Facebook, Inc. *Form S-1 Registration Statement Under The Securities Act of 1933: Facebook, Inc.* 2012. URL: <http://www.sec.gov/Archives/edgar/data/1326801/000119312512034517/d287954ds1.htm> (visited on 05/12/2015).
- [20] Martin A Fischler and Robert C Bolles. “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography”. In: *Communications of the ACM* 24.6 (1981), pp. 381–395.
- [21] Jerome H Friedman, Jon Louis Bentley, and Raphael Ari Finkel. “An algorithm for finding best matches in logarithmic expected time”. In: *ACM Transactions on Mathematical Software (TOMS)* 3.3 (1977), pp. 209–226.
- [22] Andrea Frome et al. “Recognizing objects in range data using regional point descriptors”. In: *Computer Vision-ECCV 2004*. Springer, 2004, pp. 224–237.
- [23] Keinosuke Fukunaga and Larry Hostetler. “The estimation of the gradient of a density function, with applications in pattern recognition”. In: *Information Theory, IEEE Transactions on* 21.1 (1975), pp. 32–40.
- [24] Brian Fulkerson, Andrea Vedaldi, and Stefano Soatto. “Class segmentation and object localization with superpixel neighborhoods”. In: *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE. 2009, pp. 670–677.

- [25] Thomas Funkhouser and Philip Shilane. “Partial matching of 3 D shapes with priority-driven search”. In: *ACM International Conference Proceeding Series*. Vol. 256. Citeseer. 2006, pp. 131–142.
- [26] Natasha Gelfand et al. “Robust global registration”. In: *Symposium on geometry processing*. Vol. 2. 3. 2005, p. 5.
- [27] Shengyin Gu et al. “Surface-histogram: A new shape descriptor for protein-protein docking”. In: *Proteins: Structure, Function, and Bioinformatics* 80.1 (2012), pp. 221–238.
- [28] Chris Harris and Mike Stephens. “A combined corner and edge detector.” In: *Alvey vision conference*. Vol. 15. Manchester, UK. 1988, p. 50.
- [29] Stan Horaczek. *How Many Photos Are Uploaded to The Internet Every Minute?* 2013. URL: <http://www.popphoto.com/news/2013/05/how-many-photos-are-uploaded-to-internet-every-minute> (visited on 05/12/2015).
- [30] Cheuk Yiu Ip et al. “Using shape distributions to compare solid models”. In: *Proceedings of the seventh ACM symposium on Solid modeling and applications*. ACM. 2002, pp. 273–280.
- [31] Yushi Jing and Shumeet Baluja. “Pagerank for product image search”. In: *Proceedings of the 17th international conference on World Wide Web*. ACM. 2008, pp. 307–316.
- [32] Andrew E. Johnson and Martial Hebert. “Using spin images for efficient object recognition in cluttered 3D scenes”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 21.5 (1999), pp. 433–449.
- [33] Andrew Edie Johnson. “Spin-images: a representation for 3-D surface matching”. PhD thesis. Citeseer, 1997.
- [34] Jan Knopp et al. “Hough transform and 3D SURF for robust three dimensional classification”. In: *Computer Vision–ECCV 2010*. Springer, 2010, pp. 589–602.
- [35] Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. “A sparse texture representation using local affine regions”. In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 27.8 (2005), pp. 1265–1278.
- [36] Michael S Lew et al. “Content-based multimedia information retrieval: State of the art and challenges”. In: *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 2.1 (2006), pp. 1–19.
- [37] Tony Lindeberg. “Feature detection with automatic scale selection”. In: *International journal of computer vision* 30.2 (1998), pp. 79–116.
- [38] David G Lowe. “Distinctive image features from scale-invariant keypoints”. In: *International journal of computer vision* 60.2 (2004), pp. 91–110.
- [39] Sancho McCann and David G Lowe. “Local naive bayes nearest neighbor for image classification”. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE. 2012, pp. 3650–3656.

- [40] Krystian Mikolajczyk and Cordelia Schmid. “Scale & affine invariant interest point detectors”. In: *International journal of computer vision* 60.1 (2004), pp. 63–86.
- [41] Marius Muja and David G Lowe. “Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration.” In: *VISAPP (1) 2* (2009).
- [42] Marius Muja and David G. Lowe. “Scalable Nearest Neighbor Algorithms for High Dimensional Data”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36.11 (2014), pp. 2227–40.
- [43] David Nister and Henrik Stewenius. “Scalable recognition with a vocabulary tree”. In: *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*. Vol. 2. IEEE. 2006, pp. 2161–2168.
- [44] Stephen Malvern Omohundro. *Five balltree construction algorithms*. International Computer Science Institute Berkeley, 1989.
- [45] opencv.org. *OpenCV*. URL: <http://www.opencv.org> (visited on 05/13/2015).
- [46] Robert Osada et al. “Shape distributions”. In: *ACM Transactions on Graphics (TOG)* 21.4 (2002), pp. 807–832.
- [47] Maks Ovsjanikov et al. “Shape Google: a computer vision approach to invariant shape retrieval”. In: *Proc. NORDIA 1.2* (2009).
- [48] Jeremie Papon et al. “Voxel cloud connectivity segmentation-supervoxels for point clouds”. In: *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE. 2013, pp. 2027–2034.
- [49] James Philbin et al. “Object retrieval with large vocabularies and fast spatial matching”. In: *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*. IEEE. 2007, pp. 1–8.
- [50] Brian Pinkerton. “Finding what people want: Experiences with the WebCrawler”. In: *Proceedings of the Second International World Wide Web Conference*. Vol. 94. Chicago. 1994, pp. 17–20.
- [51] pointclouds.org. *Point Cloud Library*. URL: <http://www.pointclouds.org> (visited on 05/13/2015).
- [52] Yong Rui and Thomas Huang. “Optimizing learning in image retrieval”. In: *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on*. Vol. 1. IEEE. 2000, pp. 236–243.
- [53] Radu Bogdan Rusu. “Semantic 3D Object Maps for Everyday Manipulation in Human Living Environments”. PhD thesis. Computer Science department, Technische Universitaet Muenchen, Germany, 2009.
- [54] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. “Fast point feature histograms (FPFH) for 3D registration”. In: *Robotics and Automation, 2009. ICRA’09. IEEE International Conference on*. IEEE. 2009, pp. 3212–3217.

- [55] Radu Bogdan Rusu et al. “Fast 3d recognition and pose using the view-point feature histogram”. In: *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*. IEEE. 2010, pp. 2155–2162.
- [56] Radu Bogdan Rusu et al. “Persistent point feature histograms for 3D point clouds”. In: *Proc 10th Int Conf Intel Autonomous Syst (IAS-10), Baden-Baden, Germany*. 2008, pp. 119–128.
- [57] Zujun Shentu et al. “Context shapes: Efficient complementary shape matching for protein–protein docking”. In: *Proteins: Structure, Function, and Bioinformatics* 70.3 (2008), pp. 1056–1073.
- [58] Philip Shilane and Thomas Funkhouser. “Distinctive regions of 3D surfaces”. In: *ACM Transactions on Graphics (TOG)* 26.2 (2007), p. 7.
- [59] Chanop Silpa-Anan and Richard Hartley. “Optimised KD-trees for fast image descriptor matching”. In: *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE. 2008, pp. 1–8.
- [60] Ivan Sipiran and Benjamin Bustos. “Harris 3D: a robust extension of the Harris operator for interest point detection on 3D meshes”. In: *The Visual Computer* 27.11 (2011), pp. 963–976.
- [61] Josef Sivic and Andrew Zisserman. “Video Google: A text retrieval approach to object matching in videos”. In: *Computer Vision, 2003. Proceedings. Ninth IEEE International Conference on*. IEEE. 2003, pp. 1470–1477.
- [62] Bill Slawski. *Just What Was The First Search Engine?* May 2006. URL: <http://www.seobythesea.com/2006/02/just-what-was-the-first-search-engine/> (visited on 05/11/2015).
- [63] Stephen M Smith and J Michael Brady. “SUSANa new approach to low level image processing”. In: *International journal of computer vision* 23.1 (1997), pp. 45–78.
- [64] Bastian Steder et al. “Point feature extraction on 3D range scans taking into account object boundaries”. In: *Robotics and automation (icra), 2011 ieee international conference on*. IEEE. 2011, pp. 2601–2608.
- [65] Bastian Steder et al. “Robust on-line model-based object detection from range images”. In: *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*. IEEE. 2009, pp. 4739–4744.
- [66] Fridtjof Stein and Gérard Medioni. “Structural indexing: Efficient 3-D object recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 14.2 (1992), pp. 125–145.
- [67] Danny Stieben. *The Archie Search Engine — The World’s First Search!* May 2013. URL: <http://www.makeuseof.com/tag/the-archie-search-engine-the-worlds-first-search/> (visited on 05/11/2015).

- [68] Jian Sun, Maks Ovsjanikov, and Leonidas Guibas. “A Concise and Provably Informative Multi-Scale Signature Based on Heat Diffusion”. In: *Computer graphics forum*. Vol. 28. 5. Wiley Online Library. 2009, pp. 1383–1392.
- [69] Federico Tombari, Samuele Salti, and Luigi Di Stefano. “A combined texture-shape descriptor for enhanced 3D feature matching”. In: *Image Processing (ICIP), 2011 18th IEEE International Conference on*. IEEE. 2011, pp. 809–812.
- [70] Federico Tombari, Samuele Salti, and Luigi Di Stefano. “Unique shape context for 3D data description”. In: *Proceedings of the ACM workshop on 3D object retrieval*. ACM. 2010, pp. 57–62.
- [71] Federico Tombari, Samuele Salti, and Luigi Di Stefano. “Unique signatures of histograms for local surface description”. In: *Computer Vision–ECCV 2010*. Springer, 2010, pp. 356–369.
- [72] Philip HS Torr and Andrew Zisserman. “MLESAC: A new robust estimator with application to estimating image geometry”. In: *Computer Vision and Image Understanding* 78.1 (2000), pp. 138–156.
- [73] RepRap Wiki. *Printable Part Sources*. URL: [http://reprap.org/wiki/Printable\\_part\\_sources](http://reprap.org/wiki/Printable_part_sources) (visited on 05/13/2015).
- [74] Walter Wohlkinger and Markus Vincze. “Ensemble of shape functions for 3d object classification”. In: *Robotics and Biomimetics (ROBIO), 2011 IEEE International Conference on*. 2011, pp. 2987–2992.
- [75] Walter Wohlkinger and Markus Vincze. “Shape distributions on voxel surfaces for 3D object classification from depth images”. In: *Signal and Image Processing Applications (ICSIPA), 2011 IEEE International Conference on*. IEEE. 2011, pp. 115–120.
- [76] Yu Zhong. “Intrinsic shape signatures: A shape descriptor for 3D object recognition”. In: *Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on*. IEEE. 2009, pp. 689–696.