



AGH

**AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W
KRAKOWIE**

**WYDZIAŁ ELEKTROTECHNIKI, AUTOMATYKI,
INFORMATYKI I INŻYNIERII BIOMEDYCZNEJ**

KATEDRA AUTOMATYKI I INŻYNIERII BIOMEDYCZNEJ

Praca dyplomowa magisterska

*Rozpoznawanie obiektów na obrazach RGB-D pod kątem
zastosowań w robotyce*

*Object recognition in RGB-D images for robotic
applications.*

Autor:

Jakub Olesiński

Kierunek studiów:

Automatyka i Robotyka

Opiekun pracy:

dr inż. Paweł Rotter

Kraków, 2016

Oświadczam, świadomy(-a) odpowiedzialności karnej za poświadczanie nieprawdy, że niniejszą pracę dyplomową wykonałem(-am) osobiście i samodzielnie i nie korzystałem(-am) ze źródeł innych niż wymienione w pracy.



AGH

**AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W
KRAKOWIE**

**FACULTY OF ELECTRICAL ENGINEERING, AUTOMATICS, COMPUTER SCIENCE AND
BIOMEDICAL ENGINEERING**

DEPARTMENT OF AUTOMATICS AND BIOENGINEERING

Master of Engineering Thesis

*Object recognition in RGB-D images for robotic
applications.*

Author:	<i>Jakub Olesiński</i>
Degree programme:	<i>Automatics and Robotics</i>
Supervisor:	<i>dr inż. Paweł Rotter</i>

Kraków, 2016

Contents

1. Introduction	7
1.1. Problem statement	7
1.2. Related work	8
1.3. Software tools	8
1.4. Test environment	8
2. Pre-processing	11
2.1. Representation	11
2.2. Neighbourhood	12
2.3. Noise filtering	13
2.4. Segmentation	13
3. Sparse feature matching	15
3.1. Shape description	15
3.2. Texture description	17
3.3. Clustering	17
3.4. Alignment	17
4. Dense template matching	19
4.1. Ransac	19
4.2. Linemod	19
4.3. Neural networks	19
5. Post-processing	21
5.1. Pose refinement	21
5.2. Hypothesis verification	21
6. Applications	23
6.1. Scenario 1	23
6.2. Scenario 2	23

6.3. Scenario 3	23
6.4. Scenario 4	23
Summary	25

1. Introduction

The aim of this work is to survey modern RGB-D image processing algorithms for model-based object recognition. Analysis of each method is focused on their usability in time and resource constrained robotic environment. Based on provided performance evaluation, selected methods are used to develop a complete, applicable in robotics, object recognition system.

The first chapter provides project background. A formal problem statement is included and followed by categorized solution proposals found in literature, outlining the scientific context of this work. The RGB-D imaging basics are subsequently introduced, together with software tools and test environment utilized throughout the project.

Second chapter introduces preprocessing techniques. Firstly, different data representation and conversion methods are discussed. Basic depth processing operations are introduced afterwards, including spatial transformations, neighbourhood selection and normal estimation. Different noise types encountered in RGB-D images and their corresponding filtering methods are further discussed.

Model fitting with sparse feature matching algorithms is presented in chapter three. Selected keypoint detectors and descriptors of both color and shape modalities are compared and an efficient matching technique is provided. Further, correspondence clustering and pose estimation methods are evaluated.

The following chapter is about

1.1. Problem statement

Given an input image and an object of interest, the task of object recognition is to provide answer to the following questions:

- Is the queried object *present* in the image?
- If it is, then what is the *pose* of the object?

If the object has the properties of a rigid body, its pose is given by the linear and angular position.

Object recognition is a process that comprises determination of both object instance presence and its pose in a given scene image.

Robotics: why object recognition Object recognition: why RGB-D

1.2. Related work

Survey of local methods [1].

1.3. Software tools

Throughout this work, most of the software was developed with an extensive use of the Point Cloud Library (PCL)[2], which is an open source C++ library for three dimensional image processing. Since 2011, it is developed by a large scientific community, maintained by the Open Perception foundation and delivered under BSD license. PCL provides implementations of a multitude of novel algorithms for 3D filtering, feature estimation, segmentation, registration and model fitting, together with tools for visualization and camera interfacing.

The Open Computer Vision (OpenCV)[3] library was utilized as another source of implementation for the analysed algorithms. OpenCV, which is a more mature project, developed since 1999 is likewise delivered under open source BSD license. While PCL is mainly focused on volumetric data analysis, OpenCV provides a wide array of functions for 2D intensity and color image processing.

1.4. Test environment

Dataset Willow

Recognition system was validated against sample dataset from XXX, which provides full object models together with real world scene scans and ground truth position annotations. To

Dataset own Xtion sensor

Timing perf hardware

Testing of implemented algorithms was done on two hardware platforms: an embedded computing board NVidia Jetson TK1 and personal laptop Lenovo Y50-70. Parameters of both computers are summarized in table [x].

NVidia Jetson TK1		Lenovo Y50-70	
CPU:	ARM Cortex-A15 2.32GHz x4	Intel Core i7-4720HQ 2.60GHz x4	
GPU:	NVIDIA Kepler GK20a with 192 SM3.2 CUDA cores (upto 326 GFLOPS)	NVidia GeForce GTX 960M with 640 SM5 CUDA Maxwell cores (upto xx GFLOPS)	
RAM:	2GB DDR3L 933MHz 64 bit	16GB DDR3L 1600MHz (4GB GPU 2.5Mhz 128bit)	

Table 1.1: Test hardware specification

2. Pre-processing

This chapter describes basic RGB-D image processing tool set, that prepares an input image for extracting relevant information during object recognition.

2.1. Representation

Depth and colour information of the environment can be represented in computer memory by the means of various data structures, which differ in their complexity and applicability. In case of typical, low-cost RGB-D cameras, the raw data that is retrieved from the device firmware and operating system drivers come in a form of two images, one from a conventional colour camera and the other from the depth sensor, as depicted on Figure x.

Figure XXX - drawing of camera with 2 raycasts, 2 images

Data retrieved from a depth sensor have a form of one channel, two dimensional image. Each pixel contains measurement of the distance between depth camera reference frame origin and a corresponding point on the reflected surface. In case of Asus Xtion sensor, this data is stored as 32 bit floating point values. It is important to note, that depending on the acquisition method, not all of the surface point distances are measured. For example, if the surface is translucent, most sensors will fail to retrieve proper distance. Such situation is reflected in data by the use of *Not a Number (NaN)* pixel values.

Due to the displacement between color and depth sensor [reference], RGB-D cameras require a calibration step to align the measurements. For this purpose, open source packages such as [XXX] can be used. Finally, successful RGB-D data alignment yields a four channel digital image, which can be defined more formally as:

$$I : \{1, \dots, M\} \times \{1, \dots, N\} \rightarrow [0, 1]^4, \quad I(u, v) = \begin{bmatrix} I_R(u, v) \\ I_G(u, v) \\ I_B(u, v) \\ I_D(u, v) \end{bmatrix}, \quad (2.1)$$

where:

- $M \in \mathbb{N}, N \in \mathbb{N}$ - row and column image size,
- $I_R(u, v), I_G(u, v), I_B(u, v)$ - red, green and blue colour channel intensities,
- $I_D(u, v)$ - depth channel intensity

RGB-D image formulation 2.1 is mostly suitable for per channel, classic two-dimensional image processing, such as applying morphological operators or linear filters [XXXImProc]. It does, however, limit the three dimensional surface information to a single viewpoint, which is undesirable for some applications in robotics, especially for object modelling and environment mapping. A generalized data representation that overcomes this limitation is called a *point cloud*. It can be defined as a point set $C \subset \mathbb{R}^N$ of N -dimensional vectors. Depending on the usage, the points can represent different image modalities and features, such as point location in some reference frame, colour, surface normal vector or more advanced descriptors (see chap. 3). Direct conversion from a depth map to point cloud, requires projective transformation [XXX] of the depth channel to render Cartesian x, y, z point coordinates in camera reference frame. In this manner, multiple RGB-D image frames, taken from different viewpoints can be stored within a single point cloud, by registering their points to a common coordinate system.

There are other possible representations of depth data. In contrast to point clouds, *Patch volumes* [XX] and *Signed distance functions* are dense representations of the three-dimensional environment, where each point represents distance to the nearest surface. Patch volumes use BLABLABLA. OctoMaps use BLABLABLA. Those representations, however, are mainly used for environment mapping and SLAM applications. This work will further focus on point clouds, as they are the most covered representation for object recognition. ??

2.2. Neighbourhood

Similarly to classical image processing, a wide range of point cloud processing algorithms rely on the notion of point neighbourhood. Given a metric function d and a query point p_q in the point cloud P , two commonly used types of neighbourhoods can be defined. A *r-neighbourhood* is a set composed of all the points $p_i \in P$ that lie within a sphere of a radius r and a center in p_q , thus satisfying the condition

$$d(p_q, p_i) \leq r. \quad (2.2)$$

The other type is *k-neighbourhood*, a set of k nearest points in the sense of given metric. Proper determination of neighbourhood size parameters, either r or k correspondingly, is crucial in further analysis. Too small values do not provide enough information about the surface. Too large,

on the other hand, will average the surface and skip small details. Complexity and efficiency of neighbourhood search depends on different internal data representation in a point cloud.

From computational point of view, a point cloud that is directly converted from RGB-D image can be internally stored in a two dimensional array of points, with elements corresponding to image pixels. Point clouds with such data arrangement are commonly referred to as *organized*. Image-like ordering is beneficial, because it reflects the spatial distribution of points directly into memory.

Organized neighbourhoods. KD-tree. Flann.

Boundary estimation

2.3. Noise filtering

Types of noise, frequency response, low pass filtering, nonlinear filtering

2.4. Segmentation

segmentation

3. Sparse feature matching

Solutions to problem stated in section 1.1 are commonly divided [?] into feature-based and template-based. The former approach, also referred to as local matching, is focused on comparing point neighbourhoods between the scene and model datasets.

The solution in such methods is composed of several steps. Firstly, an input subset with elements of rich and distinguishable neighbourhood information is selected. Points that belong to this subset are referred to as *keypoints*. For each keypoint, its neighbourhood information is expressed in the form of a *descriptor* vector. By the means of selected metric, commonly the L_2 norm, descriptors are further compared between the scene and model datasets, to find *correspondences* with minimal distance. Such point pairs are then grouped to share similar geometric constraints and finally, the largest clusters are used to calculate affine transformation, which renders the initial problem solution.

The inherent locality of feature matching methods has direct implications on the solution performance. Such methods are by design robust to occlusions [?]. Each point is processed independently, which enables data parallelization to boost time performance. Conversely, key-point identification requires costly analysis of the whole input dataset, thus a balance between the recognition performance and time effectiveness is required [?] for real-time applications.

3.1. Shape description

There is a multitude of proposals for shape key-point detectors existing in literature. An overview and performance evaluation of the most popular methods can be found in [4] and [5]. From both evaluations, the *Intrinsic Shape Signatures* (ISS) [6] detector is worth particular attention. [4] states that ISS, as a fixed-scale detector, copes well with full three-dimensional models and provides a proper balance between the repeatability rate and time efficiency. In [5], the ISS is evaluated with the best repeatability rate among detector implementations available in PCL. The ISS introduces a saliency measure, defined by the smallest eigenvalue of the neighbourhood scatter matrix. For a given point p and its neighbourhood N , the scatter matrix Σ is

given by

$$\Sigma(p, N) = \frac{1}{|N|} \sum_{q \in N} (q - \mu_p)(q - \mu_p)^T, \mu_p = \frac{1}{|N|} \sum_{q \in N} q \quad (3.1)$$

By denoting the eigenvalues of $\Sigma(p, N)$ as $\lambda_1 > \lambda_2 > \lambda_3$, the ISS detector classifies p as a keypoint, if the following condition is satisfied

$$\frac{\lambda_2}{\lambda_1} < \epsilon_1 \wedge \frac{\lambda_3}{\lambda_2} < \epsilon_2 \wedge \lambda_3 > \epsilon_3. \quad (3.2)$$

Thresholds ϵ_1 and ϵ_2 are meant to provide sufficient difference between variations along each principal direction, which aids in estimation of a repeatable reference frames for further description stages. The third threshold ϵ_3 ensures that the variations are large enough to consider the point as interesting. A further improvement to the ISS detection criteria is to apply non-maximum suppression over the saliency measure computed at each point. With this edge thinning technique, a point is classified as a keypoint, if it has the largest l_3 over its neighbourhood. An example of keypoints detected with ISS is presented on figure 3.1. The object (cleaning spray in the center of the figure) has marked points in areas of rich shape information, like the bottleneck or the nozzle. It is also visible, that ISS classifies wrong points in two cases. Firstly, when a point lies on a visibility boundary (i.e. the edges of flat surfaces in the figure), thus it should be extended with boundary detection as described in xx. Secondly, when the sensor noise is high (distant wall flat surface in upper left edge of the figure). This issue can be handled by ignoring points above certain z-distance from the sensor or by adaptatively selecting neighbourhood radius for saliency measure.

The overall time complexity of ISS classification is at the order of $O(nk)$, where n is the image size and k is the max neighbourhood size. This translates into average execution times on target platforms, provided in table 3.1.

	NVidia Jetson TK1	Lenovo Y50-70
PCL ISS	4s	4s
PCL ISSBD	4s	4s
CUDA ISS	4s	4s
CUDA ISSBD	4s	4s

Table 3.1: ISS detector time performance statistics.

It is important to note, that in practical applications the keypoint detection step is often replaced by volumetric down sampling. About voxel grid. Model is densely grained. This avoids computational costs of keypoint detection, but increases the complexity of further matching and clustering stages. How to choose this tradeoff.

Descriptor comparison. Take SHOT [7]. Complexity, time perf. Extension with color. [8]



Figure 3.1: ISS keypoints (marked with red) computed with $\epsilon_1 = \epsilon_2 = 0.75$ and $\epsilon_3 = 0.0005$

3.2. Texture description

matching [9]

3.3. Clustering

geometric consistency, hough

3.4. Alignment

umeyama, ransac

4. Dense template matching

4.1. Ransac

4.2. Linemod

4.3. Neural networks

Convolutional neural networks

5. Post-processing

5.1. Pose refinement

icp

5.2. Hypothesis verification

ghv, ghv from v4r

multi pipeline

6. Applications

6.1. Scenario 1

something simple with recognition
simulation of something simple

6.2. Scenario 2

something simple with classification
simulation of something simple

6.3. Scenario 3

Robocup@Work
simulation of Robocup@Work?

6.4. Scenario 4

APC?
Synthetic dataset from APC gazebo. simulation of APC?

Summary

Bibliography

- [1] 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*, vol. 36, no. 11, pp. 2270–2287, 2014.
- [2] 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.
- [3] G. Bradski, “The OpenCV Library,” *Dr. Dobb’s Journal of Software Tools*, 2000.
- [4] F. Tombari, S. Salti, and L. Di Stefano, “Performance evaluation of 3d keypoint detectors,” *International Journal of Computer Vision*, vol. 102, no. 1-3, pp. 198–220, 2013.
- [5] S. Filipe and L. A. Alexandre, “A comparative evaluation of 3d keypoint detectors in a rgb-d object dataset,” in *Computer Vision Theory and Applications (VISAPP), 2014 International Conference on*, vol. 1, pp. 476–483, IEEE, 2014.
- [6] Y. Zhong, “Intrinsic shape signatures: A shape descriptor for 3d object recognition,” in *Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on*, pp. 689–696, IEEE, 2009.
- [7] F. Tombari, S. Salti, and L. Di Stefano, “Unique signatures of histograms for local surface description,” in *European conference on computer vision*, pp. 356–369, Springer, 2010.
- [8] F. Tombari, S. Salti, and L. Di Stefano, “A combined texture-shape descriptor for enhanced 3d feature matching,” in *2011 18th IEEE International Conference on Image Processing*, pp. 809–812, IEEE, 2011.
- [9] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “Orb: An efficient alternative to sift or surf,” in *International Conference on Computer Vision*, (Barcelona), 11/2011 2011.