

Contents

1	Introduction	1
1.1	A brief history of drones	1
1.2	Motivation	2
1.2.1	Ethical considerations	2
1.3	Context	2
1.4	Objectives	3
1.5	Ressources	3
1.5.1	Hardware	3
1.5.2	Software	4
1.6	Outline	5
2	State of the art	6
2.1	Drones	6
2.1.1	Numer of rotors	6
2.2	Computer Vision	6
2.2.1	Keypoint Detection, Description, and Matching	6
2.2.2	Bundle Adjustment	8
2.3	Simultaneous Localization And Mapping	8
2.3.1	Parallel Tracking and Mapping (PTAM)	8
2.3.2	ORB-SLAM	8
2.3.3	Large-Scale Direct Monocular SLAM (LSD-SLAM)	8
2.3.4	Direct Sparse Odometry (DSO)	9
2.3.5	SVO	9
2.3.6	CNN-SLAM	9
3	Computer Vision	10
3.1	Scale-state representation	10
3.2	SURF Keypoint Detector	11
3.3	SIFT Keypoint Descriptor	12
3.4	Keypoint Tracking	12
3.5	Detection and Tracking hybrid	13
3.6	Keypoint Matching	14
4	Mapping	16
4.1	Structure of the map	16
4.2	Evaluation procedure	17
4.3	Triangulation	18
4.3.1	Midpoint Method	18
4.3.2	Linear least squares	18
4.3.3	Optimal Correction	18
4.3.4	Comparison of triangulation methods	18

4.4	Bundle Adjustment	19
4.4.1	How it works	19
4.4.2	Performance	19
4.4.3	Tuning the Bundle Adjustment	19

Acronyms

UAV Unmanned Aerial Vehicle.

Chapter 1

Introduction

1.1 A brief history of drones

UAVs (Unmanned Aerial Vehicles), more commonly called drones, are defined as flying vehicles without human operators on board. They can be remote-controlled, or controlled by on-board computers. The earliest recorded use of UAVs dates back to 1849, when Austria launched about 200 unmanned balloons armed with bombs against the city of Venice [20]. Due to unfavorable wind conditions, this attack failed, and the experiment was not repeated. The first functional UAVs were made towards the end of World War 1 and their use was, like the Austrian balloons, military. One example is the Kettering Bug (Figure 1.1), which was a torpedo with wings and a propeller developed by the US Army in 1918 [25].

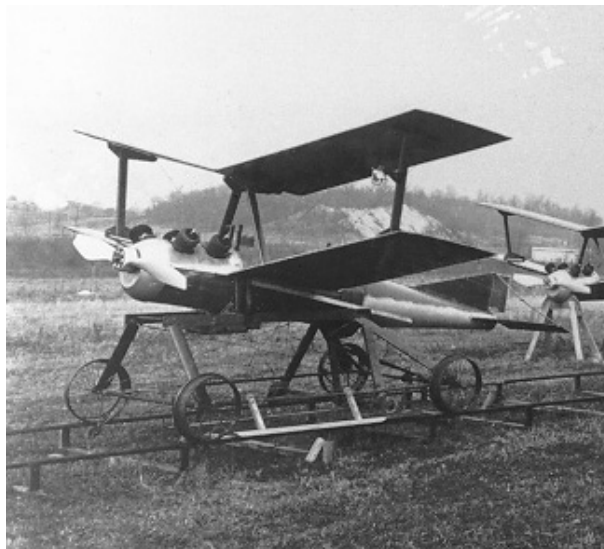


Figure 1.1: The Kettering Bug (1918)

Throughout the 20th century, UAVs became more and more sophisticated, and were used more and more, but always for military purposes. In the more recent years, civilian UAVs have started to appear on the markets and their number quickly exceeded that of military UAVs. In February 2017, the FAA (Federal Aviation Administration) of the United States estimated that around 1.1 million units were in use in the US alone, and expected that number to rise to 3.55 million by 2021 [9]. These civilian drones are very different from military drones, in both their form and their function: civilian drones are usually smaller, and use rotors to take off vertically. They are used in a wide variety of applications.

1.2 Motivation

The ability to remote-control small and agile flying objects over large distances through the air, and to bring them to previously inaccessible locations, makes many new things possible. With the increasingly lower prices and better performances of civilian UAVs, people keep finding more and more uses for these high-tech gadgets. Some examples of these applications are: crop monitoring in agriculture [3], delivery of mail or parcels, construction [5], cinematography, entertainment, or search and rescue operations. In all these applications, the more autonomous a drone is, the more efficient it will be at its task. One of the main challenges to achieve autonomy is for an UAV to be able to correctly identify its surroundings, and localize itself within them. In outdoor environments, GPS systems allow UAVs to know their position with great accuracy, but this is not possible in GPS-denied environments, like inside buildings. The main subject of this thesis will be fully autonomous navigation by a quadcopter in a GPS-denied environment.

1.2.1 Ethical considerations

The new possibilities brought by drones also pose ethical questions about security and privacy. Even though this technology can improve people's quality of life, it also has the potential to diminish it. If drones start to be widely used commercially, we could reach a point where the sound nuisances that they cause seriously impacts people who live in densely populated areas. Also, they can make us feel less at home, knowing that we could be observed from the sky at any moment. For this reason, it is important to adopt strict regulations regarding the use of drones in public spaces. Fortunately, many countries are already adopting legislation in this direction.

1.3 Context

This thesis is part of a project at the UCL that spans over several years and several masters theses. This project was launched by professor Julien Hendrickx in the 2012-2013 academic year, and had as long-term goal to develop a program that would enable low-cost UAVs to navigate autonomously in indoor environments. This means creating a map of their environment, localizing themselves in this map, and avoiding obstacles during exploration, using only on-board sensors. Another goal is to allow several drones to collaborate to speed up exploration. Five theses have already been written on this subject, each taking the work of the previous a little further.

2012-2015: First three theses In each of the three academic years (2012-2013, 2013-2014, 2014-2015), one masters thesis on the subject of indoor navigation for autonomous low-cost drones was written. These masters theses formed the base of the future work. They implemented visual SLAM methods to allow drones to build a two-dimensional map based on keypoints (first red pucks, then visual landmarks that the drone detected from a textured scene), and to localize itself within this map. During this time, inter-drone communication was also established, and was used to allow a drone to communicate the location of a target to another drone.

2015-2016: Recent work Last year, two groups of students simultaneously wrote theses on this subject. Before doing so, they joined forces to reimplement what had been done previously, but using ROS, an toolkit to develop software for robots that would make many things simpler, and allow more flexibility (see section 1.5.2). The work of the first group of students allowed a drone to search and follow a mobile target, and call a second drone to continue this task when its battery was low.

The second group of students extended the SLAM algorithm to allow to use a 3D map to localize the drone. Unfortunately, they did not implement triangulation to allow to project seen points

into 3D space, but rather made the assumption that all points were located on the ground when building the map. The end result was a drone capable of using a 3D map to localize itself, but not capable of building one from its observations.

1.4 Objectives

For this thesis, the goal is to continue the work of last year's second group to allow true 3-D SLAM: to build a 3D map based on observations by the monocular camera. To achieve this goal we will first follow these intermediate objectives:

- Research the current state of the art for 3D Keyframe based monocular visual SLAM
- Implement a way to triangulate points based on multiple observations
- Implement a way to use new visual information to correct previous errors and build a consistent map

1.5 Ressources

1.5.1 Hardware

For the practical implementation of this project, a Parrot AR.Drone 2.0 was used. This quadrotor was commercialized in 2012 and is an updated version of the original AR.Drone that was launched in 2010. This drone was marketed as a high tech toy, and is designed to be controlled from a smartphone application (connected to the drone via Wi-Fi). A few augmented reality games are available for the AR.Drone, in which it can visually recognise some predefined tags using its cameras, and interact with objects or other drones with a tag. To encourage the creation of more games for their drones, Parrot has released an open SDK that allows to effectively reprogram the drones. This early release of an open SDK has made it quite popular in the scientific community to do research on autonomous flight. The drone consists of 4 rotors, each with their own electric motor and microcontroller, an internal computer with a 1GHz ARM Cortex A8 processor and 1GB DDR2 RAM at 200MHz, and various sensors.

Sensors

The AR.drone has the following sensors:

- 3 axis accelerometer with ± 50 mg accuracy
- 3 axis gyroscope with $\pm 2000^\circ/s$ accuracy
- Pressure sensor with ± 10 Pa accuracy
- 3 axis magnetometer with $\pm 6^\circ$ accuracy
- Ultrasound sensor (facing downwards)
- Front camera (HD 720p 30fps)
- Bottom camera (QVGA, 60 fps)

The accelerometer and gyroscope form the IMU, or Inertial Measurement Unit. The IMU is used to blindly estimate the displacement of the drone by integrating its speed, and is also used to estimate the speed and stabilize the drone. The pressure sensor has a precision of 10 Pa, which is equivalent to about 83 cm, so it is not precise enough to be a useful sensor indoors. Instead, the ultrasound sensor gives a very precise measure of the altitude of the drone. The

magnetometer is used along with the IMU to estimate the drone's orientations. The bottom camera is used to estimate the drone's speed, by using optical flow methods. Finally, the front camera is not used for navigation, but for the user to see through the drone, and make films or take pictures. In this work, however, we will use the front camera to create a map.

1.5.2 Software

Closed-source internal code

The code that runs on the drone during normal execution is unfortunately closed-source. This code receives a velocity command from the user via a smartphone app, and drives the drone according to this input. When no input is given, the drone stays stationary, in "hover mode". To achieve this, it fuses sensor information from the IMU, the ultrasound, pressure sensor, magnetometer, as well as odometry from the visual flow of the bottom camera to estimate its displacement, and uses this fused pose estimate to stabilize itself.

Parrot SDK

The software development kit released by Parrot allows to send commands and receive information from the drone. It is possible to send the drone commands to take off, land, emergency stop, hover, or move in a certain direction, but not to directly control the command sent to the motors, or to have access at the control law used for the hover mode. Similarly, it is possible to read the IMU, ultrasound, odometry and video data from the drone, but not to access the data fusion that it does internally.

ROS

For the practical programming of the drone, this project uses the ROS (robot operating system) toolbox since last year. This toolbox provides useful abstraction to program robots: the program is decomposed into "nodes", each node running simultaneously in a separate thread. Each node implements part of the drone's task, and the nodes can communicate between each other by sending structured messages. The objective is to make the code modular, to be able to use existing nodes uploaded by other developers as part of packages. One such package, called "ardrone autonomy" developed by the Autonomy Lab of Simon Fraser University handles communication with a physical AR.Drone, so by using this node, we can completely disregard Parrot's specific SDK, and send all commands in ROS format. This way, we hope to make it easy to reuse code in case of a change of hardware. ROS also provides a number of convenient tools for developing software to be used on robots: it provides ways to record and replay data, to easily set and change parameters, and to cherry pick which parts of the program to run when doing tests.

OpenCV

OpenCV is an open-source computer vision library. It provides functions and data structures to manipulate images, and implements many state-of-the-art computer vision algorithms. It can work with ROS messages. We will use this library for all of our computer vision implementation: keypoint detection, tracking, description, and matching.

PCL

PCL stands for Point Cloud Library and is also an open-source library for 3D point cloud processing. The map created by the drone is a point cloud, and is saved as a PCL data structure. Currently, PCL is only used to visualize this map, but in the future, it can be used for higher-level functions like combining maps created by multiple robots, or reconstructing a dense representation of the map from a point cloud.

Ceres Solver

Ceres solver [2] is an open-source library for modeling and solving non-linear least squares problems. It is developed and maintained by Google, where it is used among other applications to estimate the pose of the Google street cars. Google released it to the public as open-source software in 2012 [1]. It will be used in this thesis to solve Bundle Adjustment problems (see section).

1.6 Outline

First, the state of the art of miniature UAV's, computer vision, and Simultaneous Localisation and mapping will be explored in chapter 2. Chapters 3 through 5 will present the main subject of this work: computer vision (Chapter 3), localisation (Chapter 4), Mapping (Chapter 5). The results will then be presented in chapters 6: Simultaneous Localisation and Mapping. Finally, chapter 7 will expose the main conclusions of this thesis, and describe the main challenges for further work.

Chapter 2

State of the art

2.1 Drones

When talking about autonomous drone navigation, it is important to be aware of what the current hardware is capable of doing, and how we can expect it to evolve in the near future. There is quite a large choice of civilian drones available to choose from on the market. Here, we will look at the main differences between the available choices.

2.1.1 Numer of rotors

Multirotors, or multicopters, are defined as rotorcrafts with three or more rotors. Having more rotors enables them to maneuver in 3D space with fixed-pitch rotors, unlike helicopters, which have articulations at the bases of their rotors. The most common multirotors have 3, 4, 6, or 8 rotors, and are respectively called tricopters, quadcopters (or quadrotors), hexacopters, and octocopters. Having more rotors has the advantage of giving more agility, at the cost of more energy consumption, and therefore a shorter battery life. A free solid object in 3D space, such as a multicopter, has 6 degrees of freedom: 3 for translation and 3 for rotation. To be able to directly control each of these 6 degrees of freedom, it must be possible to give 6 independent controls to the drone. This means that tricopters and quadcopters are always under-actuated: they can't directly control all 6 degrees of freedom. For example, quadrotors whose rotors are all in the same plane (as is almost always the case), can only directly control their translational movement along the axis parallel to the rotation of their rotors, and their roll, pitch and yaw angles, so to control their position in the plane perpendicular to the direction of gravity, they have to first adapt their roll and pitch, so that the resulting force of gravity and the thrust of their motors points inside that plane. Most hexacopters also work this way, as their rotors are also often in the same plane. Some hexrotors, however, have tilted rotors, and are fully actuated [26].

Octocopters on the other hand, are always over-actuated. One example, the Omnicopter, developed at ETH Zurich [7] can perform a 360° rotation along any axis, and move in a straight line in any direction, which enables it to perform complex and precise manoeuvres. It is able to catch thrown ping-pong balls with a little net.

2.2 Computer Vision

2.2.1 Keypoint Detection, Description, and Matching

Keypoint detection will be at the basis of our work, as we will build a map to localize the drone, and this map will be made up of landmarks, visual features observed with the camera of the drone. Good keypoints need to be recognizable from a wide variety of viewpoints, and ideally be invariant to changes of illumination, scale, rotation, or occlusion. Here is a quick overview of

some of the most important keypoint detectors and descriptors that exist in the literature.

Scale-Invariant Feature Transform (SIFT) Published in 1999 by David Lowe from Columbia University, SIFT [17] has long been the reference in keypoint description, because it is quite robust. It uses the maxima and minima of a difference of Gaussian function of the image, rescaled at different levels. The image is blurred by gaussian filters at different scales, and the difference between blurred images is taken. The result is an image without its highest spacial frequencies (noise) and lowest spacial frequencies (untextured parts of the image), leaving only a certain range of frequencies, that correspond to specific detail levels. The maxima and minima of the resulting image are considered to be corners, and become keypoints. Each keypoint is then described by a vector of size 128, representing histograms of orientation in the pixels around the keypoint. SIFT keypoints are very robust to changes of viewpoint, occlusion and illumination. Their main drawback is the computational cost to find them and to extract their descriptor.

Speeded Up Robust Features (SURF) Inspired by SIFT and first published in 2006, Speeded-Up Robust Features [bay_surf] was a new kind of keypoint detector and descriptor. SURF uses square shaped filters to approximate gaussian smoothing, which can be computed much faster, and then selects points where the determinant of the hessian matrix is maximal. Its performance in terms of robustness to changes in viewpoint, occlusion, and illumination is similar to that of SIFT but it is computationally much faster.

Features from Accelerated Segment Test (FAST) Published in 2005, the FAST detector [22] uses a different approach than SIFT and SURF, and is even faster than SURF. FAST compares the intensity of a central pixel with the intensities of the 16 pixels forming a Bresenham circle of radius 3 around the central pixel. If the central pixel is brighter or darker by a certain threshold than N contiguous pixels in the circle, the pixel is considered a corner. This test is very fast, because it is possible to reject points that do not match the criteria without comparing the intensities of all points of the circle. The threshold and the number N are parameters that have to be fixed by the user. In the original version of FAST, $N = 12$ and the four pixels at the cardinal points of the circle (up, down, left, right) are first tested, and the point is rejected if the values of these four points make it impossible for N consecutive brighter or darker pixels to exist. Expanding on this idea, the creators of FAST proposed an upgrade to their algorithm in [23] using machine learning. Their upgrade uses the ID3 algorithm to learn a decision tree to decide whether a point is a keypoint or not using the intensities of the 16 pixels. This algorithm automatically selects the best tests to perform on these 16 pixels to quickly decide whether they are corners or not. With this detector, the natural descriptor to use are the intensities of the 16 surrounding pixels, as well as whether the keypoint is a minimum or a maximum. We can already see that the dimension of the descriptor vector of FAST is 16, where SIFT's descriptor is of dimension 128, so we can expect this descriptor to be much less robust than SIFT.

Binary Robust Independent Elementary Features (BRIEF) BRIEF [8] is a descriptor introduced in 2010 that represent keypoints with binary strings. Each bit in the string is the result of a test testing which one of two specific pixels in the neighborhood of the keypoint is brightest in the smoothed image. Different lengths can be used for the string, but the creators of this descriptor found that a size of only 256 or even 128 bits is often enough for accurate matching. Matching between BRIEF features is done using the Hamming distance, which can be computed very quickly on modern computers.

Oriented FAST and Rotated BRIEF (ORB) In 2011 Rublee et al. [24] proposed a combination of a modified FAST detector, and a modified BRIEF descriptor, to obtain the

ORB detector and descriptor. One of the main drawbacks of FAST is that it has no orientation component, so it is unable to recognise keypoints when these are rotated in the camera plane. ORB computes an orientation component to each keypoint detected by FAST by using the intensity centroid of the patch, and then computes a rotated BRIEF descriptor, so that the resulting keypoint is invariant to rotations. Because the FAST detector and BRIEF descriptor are both very fast methods, and result in very small descriptors, ORB keypoints are very efficient computationally.

2.2.2 Bundle Adjustment

2.3 Simultaneous Localization And Mapping

Simultaneous Localization and Mapping (SLAM) refers to the joint task of creating a map of a robot's surroundings, while also keeping track of the robot's location in this map. The word "robot" should be understood very broadly in this context, for example it could be a simple handheld camera. Because there are countless different types of robots that do SLAM, SLAM is also a very diverse field, with different algorithms for different kinds of sensors.

2.3.1 Parallel Tracking and Mapping (PTAM)

PTAM [14] was one of the first applications of Bundle Adjustment in real time. It is able to track the movement of a handheld camera by constructing a map of visual features (it uses FAST). It is called parallel tracking and mapping, because the tracking and mapping tasks are completely decoupled, and happen simultaneously in different threads. However, it is quite expensive computationally, so it only works in real time with small workspaces, and lacks loop closure (the ability to correct accumulated errors when revisiting a previously visited location). PTAM has been adapted to work with Parrot AR.Drones at the Technische Universität München [10] with impressive results, however this solution suffers from PTAM's problems: the size of the map is quite limited, and there is no loop closure. As a result, the implementation in [10] builds an initial map, and then stops mapping, only using this initial map for tracking, which makes it unusable for large environments.

2.3.2 ORB-SLAM

ORB-SLAM [18], created in 2015 is a keyframe-based monocular visual SLAM that builds a map of ORB features. The creators of ORB-SLAM extended it to ORB-SLAM 2 [19], which works with stereo and RGB-D cameras for better accuracy. It uses bundle adjustment to build a consistent map and for loop closure, and works in real time in large environments, by using a covisibility graph to focus the mapping and localization tasks on a small portion of the map at a time. It also uses a survival of the fittest approach to remove redundant keyframes and points.

2.3.3 Large-Scale Direct Monocular SLAM (LSD-SLAM)

Developed at the Technische Universität München in 2014, LSD-SLAM [12] uses a semi-dense approach to SLAM. Unlike the previously mentioned methods, it does not use keypoints, but rather, the entire image. Each new image is used to estimate a similarity transform from the previous image to estimate the position of the camera, and is then used to either refine the last keyframe, or create a new one. Each keyframe consists of an image and a depth map created with per-pixel stereo comparisons between consecutive images. The result works in real time on a CPU, and is able to obtain accurate maps of large environments, and to correct accumulated errors after loop closure.

2.3.4 Direct Sparse Odometry (DSO)

Also produced at the TUM, DSO [11] was created in 2016. Like LSD-SLAM, DSO is a direct method: it works directly on the image pixels by computing an associated depth field, instead of extracting keypoints from the image. Working directly on the image allows to bypass the costly point detection and description steps, and also allows to use points that are not recognisable by themselves, making more data useable (edges, weak intensity regions) and to work in less textured environments. Unlike LSD-SLAM, DSO is a sparse method, it samples a limited number of points on each keyframe, and does not use a smoothness geometry prior. Because DSO permanently marginalizes old points, it cannot detect a loop closure and correct accumulated errors, so it is more a pure visual odometry than a complete SLAM system, and cannot be used to obtain a consistent map (if it visits the same location twice, all points will be doubled). However, it is very accurate, and even after very large loops, the accumulated error remains small.

2.3.5 SVO

2.3.6 CNN-SLAM

Chapter 3

Computer Vision

Our approach to the SLAM problem will expand on the work done in previous years, and therefore we will continue using keyframe-based SLAM. This method builds a map made out of points observed from a small number of previous poses of the drone (keyframes). We will use the same terminology as in [10]: detected points in images will be referred to as keypoints, and keypoints whose 3D position we have an estimation of will be referred to as landmarks. As we will see later, at least two observations of a same keypoint are necessary to estimate its 3D position and use it as a landmark. Being able to recognize points in different images that correspond to the same world location lies at the basis of both the mapping and localization tasks. To do this, we will proceed in three steps. First, we need to determine which points of an image are particular enough to be recognizable from many different viewpoints. Then we need to characterize and save those points in the form of a descriptor. Finally, when we detect more points later on, we need to be able to compare these descriptors to determine which ones describe the same keypoint.

We already outlined the principal keypoint detectors and descriptors in the state of the art (see section 2.2.1). For this work, we will continue building on what was started last year, and we will use a SURF detector and a SIFT descriptor for the reasons they explained in their thesis (see [4]).

3.1 Scale-state representation

One of the problems when trying to recognize points from a variety of viewpoints, is that depending on the distance between the point and the camera, the point can have many different sizes in the image. The scale-state representation is an attempt to solve this problem. An image can be represented as a two variable function $f(x, y)$, where the value of f is the intensity of the pixel at location (x, y) . The scale-space representation of f is a family of signals obtained by convoluting the original signal with a gaussian filter g : $L(x, y; \sigma) = g(x, y; \sigma) \star f(x, y)$. For different values of the parameter σ (the scale), the result is a blurred image where finer and finer details are indistinguishable. When $\sigma = 0$, the gaussian becomes an impulse function and the result of the convolution is the original image. The scale-state representation of an image, $L(x, y, \sigma)$ can be seen as a three dimensional image, made by stacking images obtained by blurring the source image more and more. By working on the scale-state representation of images, we ensure that our methods are scale invariant. [15]

To illustrate this with an example, we look at figure 3.1. Images 3.1a and 3.1b have been observed and we would like to match them. Zooming in on image 3.1a gives 3.1c, but this image is quite different from 3.1b. If we look at the scale-space representation of 3.1b we might find a match at higher scales, when the image becomes blurred to look like image 3.1d. Note that we are matching individual keypoints, which are usually smaller than this image (the corners of this image could each be keypoints).



(a) View from far away



(b) View from close by



(c) From far away, zoomed in



(d) From close by, filtered by a gaussian

Figure 3.1: Usefulness of the scale-space representation

3.2 SURF Keypoint Detector

The SURF detector uses the determinant of the Hessian matrix to measure local change around points. The Hessian matrix is defined as follows:

$$\mathcal{H}(\mathbf{x}, \sigma) = \begin{bmatrix} L_{xx}(\mathbf{x}, \sigma) & L_{xy}(\mathbf{x}, \sigma) \\ L_{xy}(\mathbf{x}, \sigma) & L_{yy}(\mathbf{x}, \sigma) \end{bmatrix} \quad (3.1)$$

where L_{xx} , L_{xy} and L_{yy} are the convolution of the image with the second order derivatives of the gaussian $g(\mathbf{x}; \sigma)$. The creators of SURF decided to approximate the second order Gaussian derivatives with box filters. This decision stemmed from the fact that even though Gaussian filters are optimal, they are already modified quite heavily, as they have to be discretized and cropped, and the result of the convolution is then resampled. The big advantage of using square filters is that they can be evaluated very quickly by using integral images. To obtain a scale-space representation, the SURF method uses box filters of different sizes: for example, a 9x9 filter, approximates gaussian with parameter $\sigma = 1.2$, with larger filters corresponding to bigger values of σ . Local maxima on 3x3x3 neighborhoods of the scale-space representation are then taken as points of interest.

Finally, so that the detected keypoints are rotationally invariant, we need to assign an orientation to each keypoint that can be recognized from subsequent viewpoints. To do this, the SURF method computes the response to vertical and horizontal Haar wavelets in a circular neighborhood of radius $6s$ around the detected points, s being the scale at which that point was detected. The Haar wavelets are rectangle shaped, so they can also be computed efficiently using integral images. The responses are plotted in a plane, where the abscissa is the horizontal response, and the ordinate is the vertical response. The responses are summed within a sliding orientation

window, giving an orientation vector for each position of the sliding window. The orientation vector with the largest norm is taken to be the orientation of the keypoint. [bay_surf]

3.3 SIFT Keypoint Descriptor

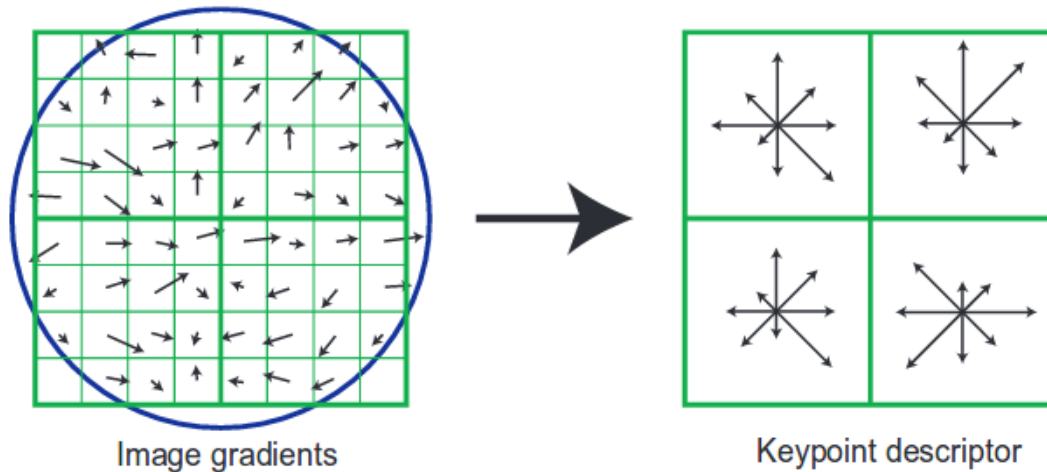


Figure 3.2: Illustration of a SIFT descriptor. Each subregion containing 4x4 points defines a histogram with 8 bins. Note that SIFT uses 4x4 subregions, not 2x2 like in this illustration. [16]

To create a SIFT descriptor, a 16x16 square region is taken around the keypoint. This region is rotated and scaled according to the orientation and scale of the detected keypoint. The 16x16 region is subdivided into 16 4x4 subregions. In each subregion, a histogram of local gradients is computed and quantized into 8 bins. The magnitudes of these histograms form the descriptor. Figure 3.2 shows an illustration of the histograms. As there are 16 histograms of 8 bins each, the total length of the descriptor vector is 128. This is quite a large dimension, which has the advantage of being highly discriminative, even with many different keypoints, but the disadvantage of requiring more space to be saved, and to slow down the matching algorithms that will be used later on. [17]

3.4 Keypoint Tracking

Because the above algorithm to extract a descriptor of a keypoint is quite slow (it takes approximatively to extract all descriptors from a new image, see), it is not practical to run these algorithms on each new image, as this would greatly reduce the framerate. For this reason, another, faster alternative has been found : keypoint tracking. This method finds keypoints by looking for keypoints that were observed in the previous image, and assuming that their displacement inside the image was small.

The algorithm used for this keypoint tracking is called the pyramidal Lucas-Kanade method [6]. This method assumes that the movement between successive frames is small and can be approximated by an affine transformation in local regions. It uses small dense pixel regions to estimate the movement of a keypoint between successive frames. The method was explained more in detail in last year's master thesis [4]. Keypoint tracking has the advantage of being much faster than detecting new keypoints, especially because the descriptor of each tracked keypoint is already known and does not have to be extracted again at each frame. The main drawback of this method is that it does not allow to detect any new keypoints, as all tracked keypoints have to be present in the previous frame.

3.5 Detection and Tracking hybrid

In the previous years, detection and tracking were used alternatively: each new frame was populated with keypoints using one of the two methods. As long as there were enough keypoints in the preceding frame, tracking was used on all these keypoints with no detection of new keypoints. While tracking is used, the number of keypoints in each successive frame is nonincreasing, as each keypoint can either be lost or tracked, and when the number of keypoints in the frame fell below some threshold, keypoint detection was used to populate the next frame with new keypoints. When the detection method is used, all previous keypoints are discarded, and keypoints are detected on the entire image. The reason for this choice was that detection was too slow to be used on each frame, but necessary to find new keypoints, and that if detection and tracking were used simultaneously, keypoints that are already being tracked would be re-detected, which would result in duplicate keypoints in the same frame. However, this method is suboptimal, because every time detection is used, all currently tracked keypoints are lost and have to be re-detected. It also causes quite an inconsistent framerate.

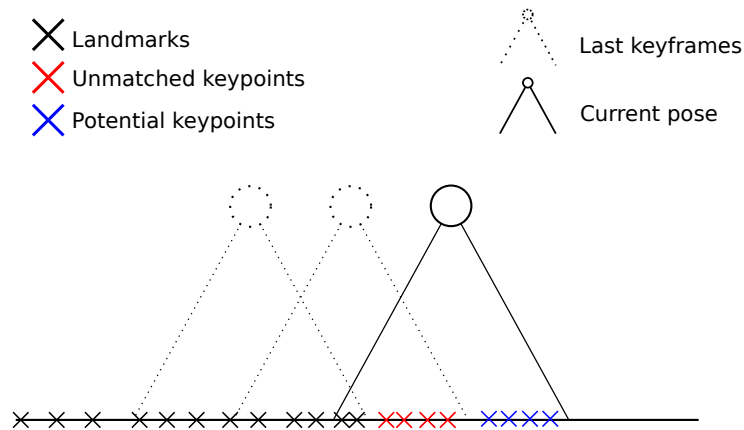


Figure 3.3: Delay between first observation of a keypoint, and mapping of this keypoint

In addition to being suboptimal, this method has a few drawbacks that make it hard to use for this project. As we will see later, at least two views of a single keypoint from different keyframes are required to estimate the 3D position of the corresponding landmark. Figure 3.3 is an illustration of this problem, in a one dimensional world for simplicity. When a keypoint is seen for the first time and saved in a keyframe, it becomes an unmatched keypoint (red cross), and its 3D location is unknown. It has to be observed from a second keyframe in order to become a landmark (black cross), with an estimated 3D position. Only then can it be used as a reference to localize the drone. On figure 3.3, if the drone creates a new keyframe at its current position, the unmatched keypoints can become landmarks, and the potential keypoints become new unmatched keypoints. If we had not detected and saved those new unmatched keypoints, when the drone continues advancing it will not be able to add any landmarks to the map. The difficulty lies in the fact that each time we create a new keyframe, we ideally need to see:

1. Enough already mapped landmarks to accurately estimate the position of the keyframe
2. Unmatched keypoints from previous keyframes so that we can put them into the map to allow further exploration
3. New, previously unobserved keypoints to ensure long-term viability (these points will be the unmatched keypoints of the next keyframe)

The need for these three ingredients means that keyframes have to be created more frequently than when we are mapping keypoints from single observations. In many cases, there are enough common keypoints between successive keyframes so that the old method did not require detection at any moment between the two keyframes. As a result, all keypoints of the most recent keyframe would already have been observed before. This means that they can be added to the map, but it seriously harms possible future exploration, as the third element of the above list is missing, so it won't be able to create any new landmarks when the next keyframe is created. Simply decreasing the threshold at which we used detection instead of tracking was not a good solution as it resulted in many more slow detection steps.

To solve this problem, a new technique had to be found. It is clear that we need new, unmatched keypoints in each keyframe, and these points cannot be tracked from previous images, so they have to be detected. However, keypoint detection takes so much time it's prohibitive to do it on each image. We noticed that the vast majority of keypoints that are lost during tracking are lost because they are moved out of the field of vision, so they leave the image from the sides. Similarly, new potential keypoints enter the field of view from the sides, so after some displacements, there will often be sections of the image with many keypoints (parts of the scene that were observed continuously since the last full keypoint detection), and sections of the image with no keypoints (parts of the scene that were not in the field of view when the last full detection took place). Our solution is to treat the part of the image with no keypoints as if it was a separate image, and do keypoint detection on that part, while keeping the part of the image with keypoints, and tracking its keypoints.



Figure 3.4: Illustration of the hybrid between tracking and detection

3.6 Keypoint Matching

The final computer vision task is keypoint matching. After keypoints from multiple viewpoints are saved, we need to know which observations correspond to the same point. This is needed when a keyframe is added to the map, because we need keypoint matches to place landmarks in the map, and also when trying to match an image with the map to localize the drone. We will compare the descriptors of the keypoints to determine which ones are similar enough to represent the same point. The Euclidean distance is a good metric of how different the descriptors are, because they are vectors of floating point numbers (if we were working with binary descriptors, like in the case of BRIEF, we should use the Hamming distance instead).

When matching a query image with a reference image, we find the nearest descriptor of the reference image for each descriptor of the query image. For each pair, we then determine if they are close enough to be considered a match. Finding the exact nearest neighbor would take $O(n * m)$ time for an exhaustive search, where n and m are respectively the number of keypoints in the query and in the reference images. A search with a k-d tree runs in $O(n \log(m))$ time, but suffers from the curse of dimensionality, so in high dimensional spaces it is not always faster than an exhaustive search in practice. To overcome these problems, we only look for

an approximate nearest neighbor. There exist many different approximate nearest neighbor methods, one of the most famous for high dimensional search is the best-bin first method, which was recommended by the creators of SIFT. It is based on a k-d tree and finds the nearest neighbor in most cases, and if not, it finds a close neighbor. Other, better approximate nearest neighbor methods have since been invented, so to marginalize the choice of an approximate nearest neighbor algorithm, we use the implementation of a FLANN matcher of OpenCV for the nearest neighbor computation. FLANN stands for Fast Library for Approximate Nearest Neighbors. It contains many state-of-the-art nearest neighbors algorithms and automatically chooses the most appropriate one depending on the data.

Chapter 4

Mapping

The task of mapping was the main challenge of this work. This task consists in determining the 3D coordinates of recognizable features, that can later be used as landmarks by the drone to estimate its own position. The main challenge when building a 3D map using a single monocular camera, is that a point needs to be observed from at least 2 different positions to be mapped. The simplest approach to map a point, is to simply triangulate its position from two different views. We will begin by exploring this approach. We will find that although this does work reasonably well when we are certain of the position of the cameras, it does not when this position is uncertain. In addition, this method does not allow to take more than two views into account. To remedy these problems we will implement a bundle adjustment step, that allows to build a map that is globally consistent. Throughout this section, we will have to make design choices to try to obtain a method that is both fast enough to work in real time, and accurate enough for the drone to control its position. To evaluate these performances, we will perform a standardized test.

4.1 Structure of the map

We will keep the keyframe-based structure for the map from the previous years, but this structure will need to be adapted so that each landmark does not belong to just one keyframe. The map will contain two main data structures: a list of keyframes, and a list of landmarks. Each keyframe will contain the following information: the drone's pose estimation at the moment the keyframe was created, the estimated corrected pose of the keyframe, a list of observed keypoints, meaning their descriptor and their 2D position in the image plane of the keyframe, and for each keypoint, whether it corresponds to a landmark of the map, and which one. Each landmark will contain a descriptor and estimated 3D coordinates.

The map grows every time we decide to add a keyframe. When we do, the keypoints of the current image seen by the camera, along with their 2D coordinates and descriptors, are saved in a new keyframe. The descriptors are then matched with descriptors of the other keyframes, and when there is a match between two keypoints of two different keyframes, we can triangulate a new landmark into the map.

The main challenges are :

- How to triangulate new points?
- How to take into account when more than two keyframes see a same keypoint ?
- How to prevent errors from accumulating over time?

4.2 Evaluation procedure

All along the creation of the mapping algorithm, we will need to make design choices and to set a number of parameters. To make these decisions, we will evaluate the performance of the map using a benchmark test. The evaluation will happen in two phases: in the map creation phase, the drone will initialize its map using four keyframes (a somewhat arbitrary number to ensure a large enough map), and in the map testing phase, it will be moved to various known locations, and we will measure the accuracy of its position estimation. The setup used for this evaluation is illustrated on figure 4.1. During the map creation phase, we will try to emulate the way the drone would initialize its map just after taking off during a real flight mission (see more in section). First the drone is placed in a know position on a table (position A in figure 4.1). There, the drone is turned on, and it begins initializing its map by creating a first keyframe. The drone is then successively placed in 3 other known locations (B, then C, then D), and at each of these locations, the drone manually receives its position. Every time this happens, the drone creates a keyframe, adds landmarks into the map, and optionally updates existing landmarks' positions, or even removes some landmarks. In reality, the drone would not know its exact position when making keyframes at points B, C, and D (position A is defined as the origin), especially at the beginning, as it would have to rely on its IMU and other internal sensors to estimate its position. To emulate this, we implement a second type of test: the robustness test, where the position that is communicated to the drone is slightly different from its real position, at each of the uncertain points (B, C, and D).

In the map testing phase, the drone is again placed at different known locations (positions 1 through 5 on figure 4.1). This time, no information is communicated to the drone from the outside, and the drone does not modify its map. The drone estimates its position using only its camera and the map that it built during the initialization phase. We compare the drone's estimated position with its real position to evaluate the quality of the map. The two quantities we will seek to optimize are the accuracy of the drone's visual pose estimation, and the time taken to create the map.

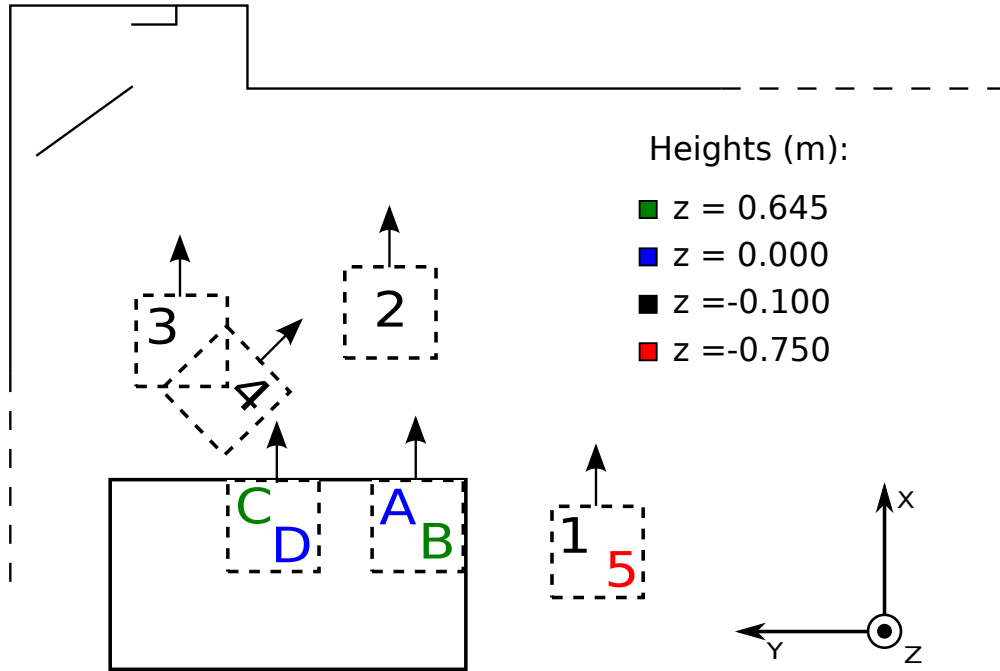


Figure 4.1: Different positions of the drone during the validation

4.3 Triangulation

Our first problem is where to put a landmark in the 3D world from two observations at two different keyframes. If all measurements were perfect, we could simply draw a ray at each keyframe that goes from the camera center and passes through the keypoint in the image plane, and those two rays would intersect at the position of the landmark.

Unfortunately, those lines never intersect in practice, due to various errors (measurement errors, errors in the model of the cameras, errors on the position estimation of the cameras), so we need to find a method to find the best possible location of a 3D point from the pair of images.

4.3.1 Midpoint Method

The simplest and most obvious solution is to take the midpoint of the common perpendicular of the two rays. This method is intuitive to understand geometrically, and is quite easy to compute. In practice, however its results are not very good, as there is no theoretical reason for this point to be the best.

4.3.2 Linear least squares

A more advanced method for finding the coordinates of a point from a correspondence is called Linear Least squares, and is described in Hartley & Zisserman's *Multiple view geometry in computer vision* [21]. This method also has the advantage that it is already implemented in the OpenCV library.

4.3.3 Optimal Correction

If we assume that the error on observed points is random and follows a Gaussian distribution with zero mean, then the optimal solution would be to displace the pixels on both images until the resulting rays meet, keeping the displacement of the pixels as small as possible in the least squared sense. Such a solution would give the maximum likelihood estimator of the position of the 3D points, under those assumptions. There are several algorithms in the literature that triangulate the position of a point using optimal correction. The most popular one, proposed by Hartley and Sturm [13], computes the solution directly but requires finding the root of a 6th degree polynomial. Kanantani et. al.'s method [**kanatani**] finds a solution iteratively, but requires very few iterations to have an accurate solution, and in practice, is faster than the Hartley Sturm method. It also has better numerical properties, as unlike the Hartley-Sturm method, it does not have singularities at the epipoles.

4.3.4 Comparison of triangulation methods

Using the evaluation procedure described in section 4.2, we can compare these 3 triangulation methods. A comparison of their performances can be seen on table ??.

We can see that the optimal correction method is more performant, at the cost of being computationally heavier. We will see, however, that this computation time is negligible in comparison with the time taken for bundle adjustment. Because the time difference is not significant, and because optimal correction is better motivated theoretically while also giving better results in practice, we will use the optimal correction method to triangulate points.

The performance increase from optimal correction is not very big however, because the biggest source of error is not gaussian noise on the measurements of the points, but errors in the position estimate of the cameras at the moment of creating the keyframes (especially in the robustness test).

4.4 Bundle Adjustment

Having a bad estimation of the keyframe's pose is problematic, as it will result in badly located landmarks, which in turn will cause a bad estimation of the camera's position when future keyframes are created. In the long term, errors will accumulate, and the map will be completely distorted. Luckily, if we have enough point correspondences between two images, it is possible to deduce the relative displacement between the two images. This means that from a set of images, we can reconstruct a scene, without even needing a prior estimation of the position of the cameras that took the images. This is good news as it means that the images can give us some absolute information about the scene, that can be used to correct the errors on the camera poses of the previous keyframes. The problem of adjusting camera poses and 3D point locations in order to minimize the reprojection errors of the 3D points onto the image planes is known as bundle adjustment. As stated above, bundle adjustment has the advantage of being absolute with respect to the world, and so not having errors accumulate. Another advantage of bundle adjustment, is that it can easily take into account points that are seen by more than two cameras, which is not trivial for the triangulation techniques described above. The main disadvantage of bundle adjustment is that it is computationally heavy, so it is important to adapt it to be useable in real time.

To show the advantages of bundle adjustment, we compare it with triangulation using the optimal correction method. When using bundle adjustment, optimal correction triangulation is also used to obtain an initial solution, from which we optimize.

4.4.1 How it works

4.4.2 Performance

We see that in the accuracy test, bundle adjustment gives a significant improvement to the performance of the pose estimation. In the robustness test, however, bundle adjustment makes the results quite worse. This is surprising, as it is just in the case where the pose of the camera is uncertain that bundle adjustment should improve the result, by adjusting the pose of the camera. But because the problem is quite non linear, it is possible to fall in local minima. We will later see that by tuning the bundle adjustment phase, we can make it perform much better, and actually improve the results in the robustness test. The other big problem of bundle adjustment is the time it takes. In both tests, this time was over 13s, which is a prohibitively large amount of time for real time applications. In addition to improving the performance on the robustness test, we will try to reduce the time taken by the bundle adjustment.

4.4.3 Tuning the Bundle Adjustment

The main drawback of Bundle Adjustment is that it requires an iterative method to be solved and can take a lot of time, which of course is a limiting factor for a robot that builds a map in real time. Therefore it is important to optimize both the speed of the computations, and their accuracy. As is often the case, there will have to be a tradeoff between these two. To measure both the speed of the Bundle Adjustment and the accuracy of the map obtained, we will conduct some experiments using different parameters to initialize the map.

Convergence of the solver

The first element we can tune is the convergence criterion of the solver that solves the bundle adjustment problem. We will stop when the ratio of the change in the objective function to the value of this function arrives below some threshold. This ensures that the criterion scales with the problem, which is important as the size of the problem can vary during operation (as the map grows, for example). The default value of the Ceres solver is 10^{-6} , but experimentally,

we find that we can use a softer threshold to significantly speed up the computations, without impactig the quality of the results too much (and even improving them in the robustness test).

Table 4.1: Effect of convergence criterion on Bundle Adjustment speed and performance

Convergence Threshold	Accuracy			Robustness		
	Distance (m)	Angle (rad)	Time (s)	Distance (m)	Angle (rad)	Time (s)
1.000×10^1	0.099	0.041	0.529	0.776	0.219	0.527
5	0.099	0.041	0.533	0.776	0.219	0.560
1	0.099	0.041	0.537	0.776	0.219	0.529
5×10^{-1}	0.099	0.041	0.532	0.284	0.086	0.615
1×10^{-1}	0.116	0.037	0.671	0.244	0.098	1.176
5×10^{-2}	0.068	0.016	0.743	0.291	0.113	2.117
1×10^{-2}	0.090	0.032	1.421	0.823	0.316	2.770
5×10^{-3}	0.091	0.031	1.533	0.142	0.050	4.186
1×10^{-3}	0.085	0.023	3.620	1.035	0.397	4.907
1×10^{-4}	0.084	0.026	5.058	1.362	0.480	10.484
1×10^{-5}	0.090	0.030	12.453	1.483	0.528	13.206
1×10^{-6}	0.090	0.030	15.703	1.374	0.492	13.802
1×10^{-7}	0.088	0.028	15.795	1.559	0.567	12.935

As we can observe on table 4.1, the more severe the convergence threshold, the better the performance on the accuracy test. On the robustness test, however, this is true for high convergence thresholds, but when the convergence threshold is lower than 0.01, the performance starts to degrade as the convergence threshold becomes lower. In both the accuracy and robustness tests, the lower the threshold, the longer it takes. We will chose a threshold of 0.05, as it is a nice compromise of performance on the accuracy test, performance on the robustness test, and speed.

Rejecting outliers

When looking closer at the results of bundle adjustment, we find that a large part of the errors come from a small number of points. Figure 4.2 shows the repartition of contribution to the total error between the landmarks as a logarithmic scale. We can see that 1 % of the landmarks account for more than 40 % of the total error, and that 10 % of the landmarks account for more than 80 % of the total error. One possible explanation is that these points do not correspond to real world points, so that even if we have the correct position of all cameras exatly, the rays corresponding to these points will not intersect, or even be close to intersection. For this reason, after every bundle adjustment pass, we could look at the contribution of every mapped point to the objective function, and eliminate points whose error is too high. Before doing this, there are some things that we have to take into account, becuase as a point is observed by an increasingly large number of cameras:

- Its total contribution to the cost function will increase, as each observtion of a point adds a term to the objective funtion
- Its contribution per observation to the cost function will also increase, because every time an observation is added, the location of the point will change a bit, and its position will become less optimal if we only consider the cameras that already saw the point before.

For these reasons, we will put a different threshold on the points depending on the number of cameras that see the point. If a point is seen by a large number of cameras, we will allow it to have larger errors in the objective function before removing it.

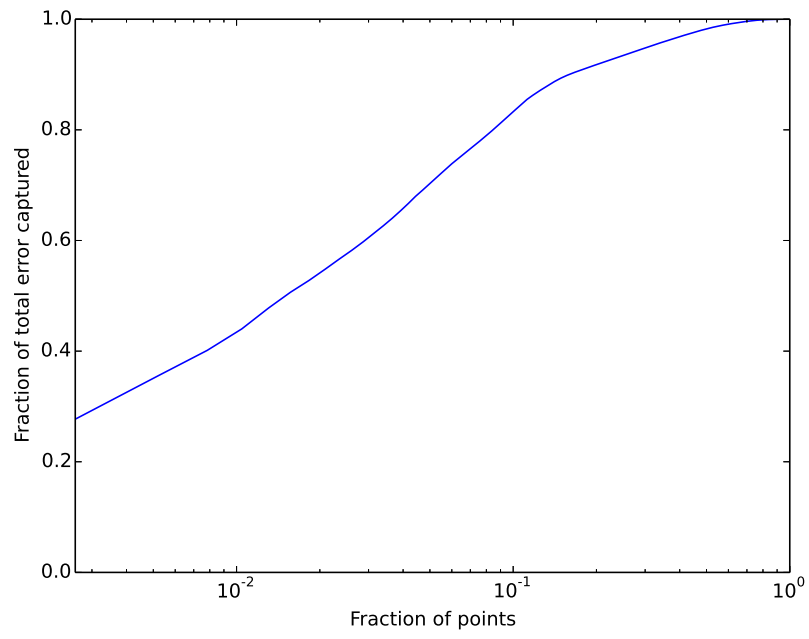


Figure 4.2: Repartition of error among points

Bibliography

- [1] Sameer Agarwal and Keir Mierle. “Introducing Ceres Solver - A Nonlinear Least Squares Solver”. In: *Google Open Source Blog* (2012).
- [2] Sameer Agarwal, Keir Mierle, et al. *Ceres Solver*. <http://ceres-solver.org>.
- [3] Chris Anderson. “Agricultural Drones”. In: *MIT Technology Review* vol. 117 | no. 3 (2014).
- [4] Alexandre Leclerc Arnaud Jacques. “Towards 3D Visual SLAM for an Autonomous Quadcopter Running ROS”. MA thesis. Université Catholique de Louvain, 2016.
- [5] “Bâtir avec des drones ?” In: <https://uclouvain.be/fr/sciencetoday/actualites/batir-avec-des-drones.html> (2016).
- [6] Jean-Yves Bouguet. “Pyramidal Implementation of the Lucas Kanade Feature Tracker”. In: *Intel Corporation* ().
- [7] Dario Brescianini and Raffaello D’Andrea. “Design, Modeling and Control of an Omni-Directional Aerial Vehicle”. In: *IEEE International Conference on Robotics and Automation* (2016).
- [8] Michael Calonder, Vincent Lepetit, and Pascal Fua. “BRIEF: Binary robust independent elementary features”. In: *European Conference on Computer Vision*. 2010.
- [9] Marcus Chavers. “Consumer Drones By the Numbers in 2017 and Beyond”. In: *newsledge.com* (2017).
- [10] J. Engel. “Autonomous Camera-Based Navigation of a Quadcopter”. MA thesis. Germany: Technical University Munich, Dec. 2011.
- [11] J. Engel, V. Koltun, and D. Cremers. “Direct Sparse Odometry”. In: *arXiv:1607.02565*. July 2016.
- [12] J. Engel, T. Schöps, and D. Cremers. “LSD-SLAM: Large-Scale Direct Monocular SLAM”. In: *European Conference on Computer Vision (ECCV)*. Sept. 2014.
- [13] Richard I. Hartley and Peter Sturm. “Triangulation”. In: *CVIU - Computer Vision and Image Understanding* Vol. 68, No. 2 (1997).
- [14] Georg Klein and David Murray. “Parallel Tracking and Mapping for Small AR Workspaces”. In: *Proc. Sixth IEEE and ACM International Symposium on Mixed and Augmented Reality (ISMAR’07)*. Nara, Japan, Nov. 2007.
- [15] Tony Lindeberg. “Scale-space theory: A basic tool for analysing structures at different scales”. In: *Journal of Applied Statistics* 21 (1994), pp. 224–270.
- [16] David G. Lowe. *Distinctive Image Features from Scale-Invariant Keypoints*. 2004.
- [17] David G. Lowe. “Object recognition from local scale-invariant features”. In: *Proceedings of the International Conference on Computer Vision*. 1999.
- [18] Raul Mur-Artal, J. M. M. Montiel, and Juan D. Tardós. “ORB-SLAM: a Versatile and Accurate Monocular SLAM System”. In: *CoRR* abs/1502.00956 (2015).

- [19] Raul Mur-Artal and Juan D. Tardós. “ORB-SLAM2: an Open-Source SLAM System for Monocular, Stereo and RGB-D Cameras”. In: *CoRR* abs/1610.06475 (2016).
- [20] Russell Naughton. *Remote Piloted Aerial Vehicles : An Anthology*. Tech. rep. Monash University, 2003.
- [21] Andrew Zisserman Richard Hartley. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2003.
- [22] Edward Rosten and Tom Drummond. “Fusing points and lines for high performance tracking.” In: *IEEE International Conference on Computer Vision*. Vol. 2. Oct. 2005, pp. 1508–1511. DOI: 10.1109/ICCV.2005.104. URL: http://www.edwardrosten.com/work/rosten_2005_tracking.pdf.
- [23] Edward Rosten and Tom Drummond. “Machine learning for high-speed corner detection”. In: *European Conference on Computer Vision*. Vol. 1. May 2006, pp. 430–443. DOI: 10.1007/11744023_34. URL: http://www.edwardrosten.com/work/rosten_2006_machine.pdf.
- [24] Ethan Rublee et al. “ORB: an efficient alternative to SIFT or SURF”. In: *In ICCV*. 2011.
- [25] Jimmy Stamp. “Unmanned Drones Have Been Around Since World War I”. In: *smithsonian.com* (2013).
- [26] Richard Voyles. *Dexterous Hexrotor Research*. <http://web.ics.purdue.edu/~rvoyles/research/Hexrotor.html>.

