

Artificial Eye Orientation Estimation with Camera-Based Measurements of Natural Visual Features

Adding a camera to obtain an unbiased real-time estimate for the three-dimensional orientation of a humanoid eye

Mariana Ribeiro dos Reis do Vale Martins

Report for

Introduction to the Research in Electrical and Computer Engineering

Advisor(s)/Supervisor(s): Prof./Dr. Alexandre Bernardino
Prof./Dr. José Santos-Víctor
Prof./Dr. John van Opstal

January 2019

Contents

List of Tables	i
List of Figures	ii
1 Introduction	1
1.1 Motivation	1
1.2 Problem definition	2
1.3 Report Outline	3
2 Background and State of the art	4
2.1 Camera Model	4
2.2 Orthogonal Procrustes Problem	6
2.2.1 Translation	7
2.3 Epipolar Geometry	7
2.4 Feature Detection and Matching	11
2.4.1 SIFT - Scale Invariant Feature Transform	12
2.4.2 SURF - Speedup Robust Features	13
2.4.3 Matching step	15
2.5 Representation of eye orientations in 3 dimensions	15
2.5.1 Rotation axis and angle	15
2.5.2 Rotation vector in head-fixed coordinates	16
2.5.3 Quaternions	16
2.5.4 Coordinate system	16
2.5.5 Donders' and Listing's laws	17
3 Method and Proposal	19
4 Preliminary Work	22
4.1 Specifications	22
4.1.1 Hardware	22
4.1.2 Software	22
4.2 Comparison of rotations around each axis between the IMU sensor and the camera	22
4.2.1 Rotation around the torsional axis	23
4.2.2 Rotation around the horizontal axis	24
4.2.3 Rotation around the vertical axis	24
5 Thesis Development Plan	26
Bibliography	27

List of Tables

4.1	Rotation angle obtained when rotating the structure around the torsional component. . . .	23
4.2	Rotation angle obtained when rotating the structure around the horizontal component. . .	24
4.3	Rotation angle obtained when rotating the structure around the vertical component. . . .	24

List of Figures

1.1	Current mechanical eye model	2
1.2	IMU sensor's drift	3
2.1	Pinhole Camera Model	5
2.2	Epipolar geometry	8
2.3	Four possible solutions retrieved from E	10
2.4	Difference of Gaussians (DoG) by octave	12
2.5	Search for the maxima and minima of the DoG	12
2.6	SIFT feature descriptor	13
2.7	Comparison between Gaussian second order derivatives and box filter	14
2.8	SURF feature descriptor	14
2.9	Neuroscience vs Computer vision usual coordinate systems	17
2.10	Listing's law, illustrated for 3500 eye orientations made by a head-restrained monkey	18
3.1	Hypothetical coordinate reference frames of the world, camera from view 1 and camera from view 2	20
4.1	Structure used for making the experiments.	23
4.2	Matches obtained for rotation 90° around the torsional axis	24
4.3	Matches obtained for rotation 30° around the horizontal axis	25
4.4	Matches obtained for rotation 5° around the horizontal axis	25
4.5	Matches obtained for rotation 20° around the vertical axis	25
4.6	Matches obtained for rotation 5° around the vertical axis	25
5.1	Gantt chart of the thesis development	26

Chapter 1

Introduction

1.1 Motivation

Even after years of research, the brain remains to this day a mysterious information-processing biological system. For example, from many of its puzzling abilities, it is still not completely understood how the brain determines which muscles to control in order to fulfill a particular movement task. This task could be something as simple as grasping an object, or even before that, looking at (and identifying) that object.

The required number of degrees of freedom to perform a particular movement is typically much smaller than the ones made available by the muscular apparatus, thus yielding a redundant control problem with infinitely many possibilities. This aspect of motor control has become known as the "Degrees of Freedom Problem" (DOF)¹, and was first articulated by Nikolai Bernstein in its current form. For example, the eye has six extra-ocular muscles (6 DOF), but to point the eye in any given direction, only two coordinates are needed (the azimuth and the elevation angles). Bernstein theorized that the brain gradually tries to find the optimal control solution for a certain task, to constrain the DOF to the required motor space, which would result in consistent performance. Indeed, when different subjects perform the same task, their muscular activations are remarkably similar. As it will be seen below, the eye orientation thus appears to be constrained by the so-called Donders' law.

Moreover, navigating through an environment in which multiple stimuli compete for attention, and being equipped with multiple sensory and motor systems to do so, severely complicates the challenges for the brain to rapidly select the optimal plan for a response to achieve a particular goal [1]. This leads to several questions: How does the normal brain decide which signals are "goals", and which are "distractors" in those complex environments? How is the plan to reach the goal truly defined? And how does this differ for sensory-impaired brains?

When focusing on audiovisual stimuli, it is important to apply the scientifically acquired concepts from neuroscience and psychophysics to an approximate model of the human sensory-motor system (eyes, head, and ears), which is the main topic of this research project. The idea is to create an autonomous humanoid eye-head robot with foveal vision, realistic auditory inputs, three-dimensional nested eye and head motor systems, and rapid sensory-motor feedback controls and learning algorithms. So far, a working mechanical model of a biologically inspired eye with six muscles has been built in a previous project [2], and is shown in Figure 1.1.

This model was constructed with Donders' and Listing's Laws in mind. Donders' law states that the 3D orientation of the eye, when looking in a specific direction, is always the same. Listing propounded

¹N. Bernstein, "The Coordination and Regulation of Movements. Oxford : Pergamon Press," 1967.

that the law could be further constrained by the discovery that the rotation axes corresponding to all possible eye orientations all lie in a common plane, called Listing's plane. The component of the rotation axis normal to this plane quantifies the eye's torsional orientation component, and results to be close to zero in these conditions. The normal to Listing's Plane is directed into the so-called *primary direction*, and points about 15 to 30 degrees upwards relative to the straight-ahead viewing direction (head fixed). It is still not well understood which aspects of the ocular motor system determine the primary orientation (and hence, Listing's law).

Whether the rotation axis of eye orientation is programmed by the brain, fully determined by the mechanical properties of the eye muscles, or by both factors, is still not known for sure, and is highly debated in the literature. Yet, implementing this system in a humanoid robot might help at understanding this problem better. Note that for active eye-head movements, for vergence eye movements (near viewing), for vestibular eye movements (head rotations), and for eye-movements under tilted head orientations, Listing's law no longer holds true, which provides a strong counter argument against the eye muscle (pulley) hypothesis. [3]

The current prototype uses an Inertial Measurement Unit (IMU) sensor to estimate the eye's pose. Although IMU units are good on acceleration and velocity measurements, they tend to have a significant position drift when determining the orientation (see Figure 1.2), which is of course a considerable problem when the eye's orientation needs to be established with better than 0.5 deg resolution. Therefore, this project's focus is to study the addition of a camera to the system, in order to determine the orientation of the eye in a more reliable manner for this model. This will hopefully contribute to a better understanding of the eye's control system, and consequently help at restoring impaired vision, which was, after all, Franciscus Donders' ultimate desire.



Figure 1.1: The current mechanical eye model. It is composed of an IMU sensor, seen on the left, and connected to a supportive green eye mounted on a global joint, in turn connected to six elastics, representing the eye muscles, and finally controlled by three motors that pull at those elastics (paired by the aluminum strips). The three motors are controlled by the computer, which is seen lying on the table.

1.2 Problem definition

As the eye model will eventually rely on accurate camera output, and to foster solid neuroscientifically inspired study with this model, it's indispensable to have the best possible estimates of the camera's

translation and orientation (and relative image transformations). As well as that, the computational speed versus accuracy trade off will also have to be dealt with. The objective of this study is to approximate the current model to a biological eye as well as possible. Thus, determining the instantaneous orientation of the eye in 3D when the gaze direction changes suddenly, which happens, for example, during rapid saccadic eye movements, should be calculated at the appropriate speed.

The aim of this work is to find the accurate camera poses within the desired time constraints. The accuracy and computational time mainly depend on the complexity of feature detection, on feature matching, and on the algorithm responsible for calculating the pose, making each of these factors an important sub-topic for this study.

This project intends to go through a detailed analysis of the advantages and disadvantages of using the camera over the IMU, and what may be gained by the combination of both.

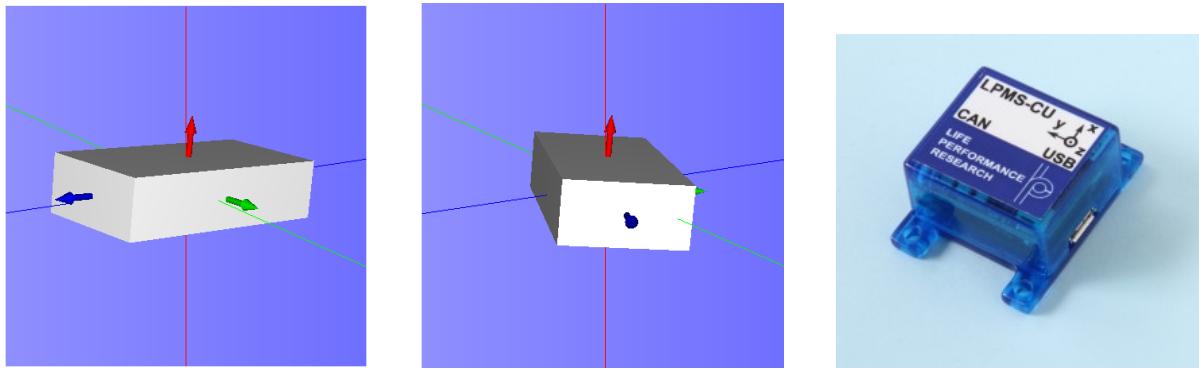


Figure 1.2: IMU sensor's drift. It can be observed that, in 80 seconds, the horizontal angle drifts about 60 degrees. The coordinate system of the IMU sensor is defined on the image on the left, where the axis in red is the torsional component, z, green is y, and blue is x.

1.3 Report Outline

The next chapters of this report are structured as follows. Chapter 2 gives insight on the initial but essential concepts necessary to launch this project, along with an overview of the state of the art on some of the matters described. Chapter 3 explains the methods that will be experimented and later on compared, in order to accomplish the project's goal. Chapter 4 describes the preliminary work on some of the project's subjects that have taken place during the last months. And finally, chapter 5 defines the master thesis development plan for the months ahead, with the respective Gantt chart.

Chapter 2

Background and State of the art

2.1 Camera Model

The first step in this project is to understand and apply the pinhole camera model. Digital cameras are equipped with a sensor (mostly charge-coupled devices (CCD) and CMOS') that transforms light into discrete electrical signals. When taking a picture, some of the light reflecting on the object is directed towards the camera pinhole (tiny aperture on the camera), forming a 180° inverted image of the object on the sensor. All the light rays entering the camera converge into the pinhole and diverge on the other side (with a unique one-to-one correspondence). Therefore the pinhole is also named the **center of projection**, principal point, or focal point.

In cameras the aperture size and light-exposure time can be controlled. A smaller aperture produces a sharper image but needs a longer exposure time to collect enough light, whereas a larger pinhole allows more light to be captured, but at the expense of producing blurred images.

The image of the object is projected onto the so-called **image plane**. When moving this plane closer to the pinhole (zoom out), the projected object gets smaller, and vice-versa. The distance between the plane and the aperture is the **focal length**, or focal distance. The view-angle of a camera depends on the focal distance, and on the sensor size. The image resolution depends solely on the number of pixels that fit on the sensor.

A useful coordinate system, called the **camera reference system**, is centered at the pinhole, and has its z-axis perpendicular to the image plane. The line defined between the center of projection and the center of the image plane is called the **optical axis**.

Figure 2.1 summarizes the concepts explained above, and will help to understand the upcoming details. As may be observed, a correspondence between the real world (3D) points $M = [X \ Y \ Z]^T$ and the points projected onto the image plane $m = [x \ y]^T$ can be established by the following triangular similarity as

$$x = \frac{f \cdot X}{Z}, \quad y = \frac{f \cdot Y}{Z}, \quad (2.1)$$

where f is the focal length (in meters), (x, y) are the projected points (in meters) on the image plane, and (X, Y, Z) (in meters) are the 3D points of the object. In an ideal case, where the camera is normalized,

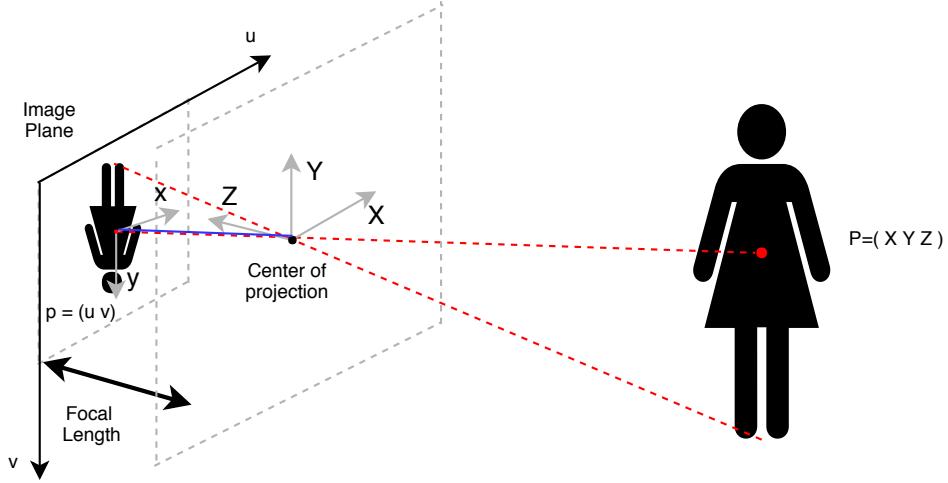


Figure 2.1: Pinhole Camera Model. The girl defined by the generic set of points M is projected onto the image plane (defined by the set of generic points m_{pixels}) through the light rays passing through the camera's pinhole, called the center of projection. The focal length is the distance from the center of the image plane to the center of projection. The digital image's coordinate system defined by the camera is not centered on the image plane. The optical axis is the blue line.

i.e., $f = 1$ m, the perspective projection can be defined, using homogeneous coordinates, as

$$\lambda \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \Leftrightarrow \lambda \tilde{m} = [I \ 0] \tilde{M} = M. \quad (2.2)$$

(with $\lambda = Z/f = Z$). However, in reality the image plane reference system is usually different from the digital image one, as shown in Figure 2.1, thus necessitating to translate it by $(c_x \ c_y)$ pixels. Also, the focal length is not normalized to one meter. Furthermore, the digital image obtained is actually expressed in pixel units, whereas the camera model is described in metric units, requiring the points to be converted by scalings (s_x, s_y) . And finally, due to sensor errors, there might be a deformation skew, s , associated to the camera coordinate system, which means that it may not have exactly perpendicular axes. All these aspects are taken into consideration through the so-called **intrinsic parameters** representation, which is described as,

$$K = \begin{bmatrix} fs_x & s & c_x \\ 0 & fs_y & c_y \\ 0 & 0 & 1 \end{bmatrix}, \quad (2.3)$$

resulting in,

$$\begin{aligned} u &= fs_x x + sy + c_x \\ v &= fs_y y + c_y \end{aligned}, \quad (2.4)$$

where (u, v) are the shifted image coordinates in pixels, denoted by $m' = [u \ v]^T$ and (x, y) are in metric units, denoted earlier by m .

On top of this, the world frame, where object M is located, might at an arbitrary position and orientation, so it is necessary to convert it to the camera coordinate frame. This conversion consists on expressing the origin of the world coordinate frame in camera coordinates, through a translation t , and the rotation between the two frames, through a rotation matrix, R . These are called the **extrinsic parameters**.

Combining this, the full perspective projection, or camera model, is defined as

$$\begin{aligned}\lambda \tilde{m}' &\sim K[R|t]\tilde{M} \\ \lambda \tilde{m}' &\sim P\tilde{M}\end{aligned}, \quad (2.5)$$

where P is a 3×4 matrix, named **camera matrix**, and λ is a scale factor. From now on m' will be noted just as m for simplicity, and its coordinates may be expressed as,

$$u = \frac{p_1^T \tilde{M}}{p_3^T \tilde{M}}, \quad v = \frac{p_2^T \tilde{M}}{p_3^T \tilde{M}}, \quad (2.6)$$

where p_i^T denotes the row i of P .

In practice, the camera model contains several nonlinear effects, generically denoted as "distortion", since real cameras have lenses instead of a pinhole. The lens allows all the emitted rays of light to be refracted into one converging point, the focal point. Although the focal length will now increase to be the sum of the distance from the lens to the converging point, and from that point to the image plane, the camera model explained above still applies, but the distortions must be accounted for. These distortions are incorporated in (2.7). Radial distortion is modeled by the k_i coefficients on the left-hand side of the polynomial transformation in (2.7) (taken from Zernike's model of aberrations), and the tangential distortion (which occurs when the image plane is not parallel to the camera lens) is defined by the p_j coefficients in the right-hand term.¹

$$\begin{aligned}u_d &= u \frac{1+k_1r^2+k_2r^4+k_3r^6}{1+k_4r^2+k_5r^4+k_6r^6} + (2p_1uv + p_2(r^2 + 2u^2)) \\ v_d &= v \frac{1+k_1r^2+k_2r^4+k_3r^6}{1+k_4r^2+k_5r^4+k_6r^6} + (p_1(r^2 + 2v^2) + 2p_2uv)\end{aligned}, \quad (2.7)$$

where $r^2 = u^2 + v^2$ and u_d and v_d , are the distorted coordinates of the image point.

2.2 Orthogonal Procrustes Problem

After having established the camera model, it is now time to discuss some potential methods of determining image transformations, as needed for this project. The first of these methods is called the Orthogonal Procrustes Problem. Procrustes, a son of Poseidon in Greek mythology, made his victims fit in a wonderful all-fitting bed, either by stretching their limbs, or by cutting them off. Ultimately, he was fitted into his own bed by Theseus. Similar to the analysis held here, there are 3 elements to Procrustes's problem: M_1 , the bed, M_2 , the unlucky adventurer, and T , the fitting treatment. The solution to the problem is a matrix, T , that minimizes

$$\|M_1 - TM_2\|^2, \quad (2.8)$$

which transforms one cloud of points, describing M_1 , into another cloud, M_2 . Considering pure rotations only, T is an orthogonal matrix, defined as R . Through the Frobenius norm,

$$\|M_1 - RM_2\|^2 = \text{trace}(M_1^T M_1 + M_2^T M_2) - 2 \text{trace}(M_2^T M_1 R), \quad (2.9)$$

¹Born, Max; Wolf, Emil. Principles of Optics: Electromagnetic Theory of Propagation, Interference and Diffraction of Light. Chapter 5.

so minimizing (2.8) with respect to R is equivalent to maximizing the second term of (2.9). By applying singular value decomposition (SVD) to $M_2^T M_1$, the latter term can be further simplified into

$$\text{trace}(M_2^T M_1 R) = \text{trace}(U \Sigma V^T R) = \text{trace}(\Sigma V^T R U) = \text{trace}(\Sigma H) = \sum_{i=1}^N \sigma_i h_{ii}. \quad (2.10)$$

The singular values of σ_i are all non-negative, and so the expression becomes maximum when $h_{ii} = 1$ for $i = 1, 2, \dots, N$, since H , a product of orthogonal matrices, is an orthogonal matrix itself, thus having its maximal value when $H = I$. This results in $I = V^T R U$, and so $R = V U^T$.

2.2.1 Translation

Often, M_1 and M_2 may be associated with different origins. Hence, it is customary to consider better fits of the configurations by allowing shape-preserving translations of the origin. When translating M_1 and M_2 by t_1 and t_2 , respectively, (2.8) becomes

$$\|(M_1 - \mathbf{1}t_1^T) - R(M_2 - \mathbf{1}t_2^T)\|, \quad (2.11)$$

which can also be written as

$$\|P_1 - RP_2 - \mathbf{1}t^T\|, \quad (2.12)$$

where $t^T = t_1^T R - t_2^T$. The expression is minimized when t^T corresponds to the column-means of $M_1 R - M_2$, therefore a simple way of effecting this translation is to remove the column-means of M_1 and M_2 , separately, before minimizing. This is called centering. The data then becomes expressed in terms of the mean deviations. [4]

2.3 Epipolar Geometry

Epipolar geometry provides an alternative yet powerful tool to obtain image transformations. It describes the relation between two views of a single scene through a 3x3 singular, non-invertible, matrix called the **essential matrix**, E , if the camera matrix is known, or the **fundamental matrix**, F , otherwise. The essential matrix is the product of a rotation matrix and a skew-symmetric matrix, whose determinant is zero, as it will be shown below. If a point in the (3D) world, M , is projected as m_1 in the first view, and m_2 in the second, then those points satisfy the relation

$$\tilde{m}_1^T F \tilde{m}_2 = 0. \quad (2.13)$$

In Figure 2.2, c_1 and c_2 are the centers of projection of the first and second views of the scene from one camera that has been displaced by (R, t) (or by two cameras staring at the same scene). Given the point on the first view, m_1 , its correspondence to the second view is constrained to a line, which is called the epipolar line, l_{e2} . This line is defined as the intersection of the plane L (the epipolar plane, containing M , m_1 , and m_2) and the second view's image plane. Furthermore, plane L is delineated by the centers of projection of both views and m_1 . The epipoles, e_1 and e_2 , are born from the intersection of the line $c_1 - c_2$ (baseline) with the epipolar lines, l_{e1} and l_{e2} .

There is a unique correspondence between a point in one image plane and an epipolar line on another that can be defined by the cross product $l_{e2} = e_2 \times \tilde{m}_2$, as well as, an homographic mapping via L from the points of the first to the second view, characterized by $\tilde{m}_2 = H_M \tilde{m}_1$. The fundamental matrix

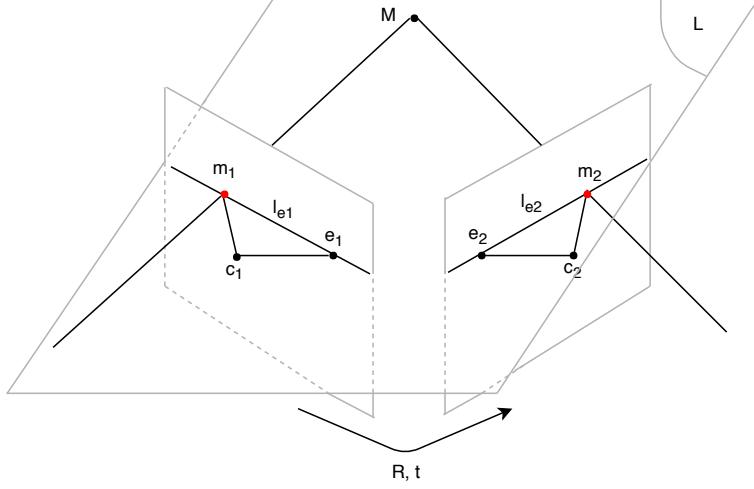


Figure 2.2: Epipolar geometry. The two small planes correspond to the image planes of two different views by the camera, separated by a rotation R and a translation t , of the same scene, with centers of projection given by c_1 and c_2 , respectively. M is the fixed point gazed at in the real world. m_1 and m_2 are the projections of point M into the respective image planes, whose rays are lying on the plane L . l_{e1} and l_{e2} are the epipolar lines, that also lie on the plane and e_1 and e_2 are the epipoles. The line defined by the centers is called the baseline.

may be constructed as $F = e_2 \times H_M$ through

$$l_{e2} = e_2 \times H_M \tilde{m}_1 = F \tilde{m}_1, \quad (2.14)$$

and if point m_2 lies on its epipolar line then $\tilde{m}_2^T l_{e2} = 0$, leading to the (2.13). F is a homogeneous matrix of rank-2 (since l_{e2} is rank-2 and H_M is rank-3) with 7 degrees of freedom, and $\det(F) = 0$.

Algebraically, having M expressed in the first camera's reference frame, under the camera model studied above, (2.15) can be deduced.

$$\lambda_1 \tilde{m}_1 = K_1 [I \ 0] \tilde{M}, \quad \lambda_2 \tilde{m}_2 = K_2 [R \ t] \tilde{M} \quad (2.15)$$

In the case of only one camera in two different views: $K_1 = K_2$. Now by eliminating \tilde{M} , the previous equation becomes

$$K_2^{-1} \lambda_2 \tilde{m}_2 = R K_1^{-1} \lambda_1 \tilde{m}_1 + t, \quad (2.16)$$

and because the cross product of two vectors is orthogonal to them both,

$$(K_2^{-1} \lambda_2 \tilde{m}_2)^T \cdot (K_2^{-1} \lambda_2 \tilde{m}_2 \times t) = 0 \quad (2.17)$$

$$(K_2^{-1} \lambda_2 \tilde{m}_2)^T \cdot ((R K_1^{-1} \lambda_1 \tilde{m}_1 + t) \times t) = 0, \quad (2.18)$$

the essential matrix $E = [t]_\times R$ and the fundamental matrix $F = K_2^{-T} E K_1^{-1}$ can be determined by

$$\begin{aligned} \tilde{m}_2^T K_2^{-T} [t]_\times R K_1^{-1} \tilde{m}_1 &= 0 \\ \tilde{m}_2^T K_2^{-T} E K_1^{-1} \tilde{m}_1 &= 0 \\ \tilde{m}_2^T F \tilde{m}_1 &= 0, \end{aligned} \quad (2.19)$$

where $[t]_{\times}$ represents the cross product as a skew-symmetric matrix, defined as

$$[t]_{\times} = \begin{bmatrix} 0 & -t_3 & t_2 \\ t_3 & 0 & -t_1 \\ -t_2 & t_1 & 0 \end{bmatrix}. \quad (2.20)$$

The scale factors λ_2 and λ_1 were dropped for being irrelevant here.

In order to find the fundamental matrix, and consequently the transformation from one view to another, having a set of matched image points, $m_{i1} = [u_{i1} \ v_{i1}]^T$ and $m_{i2} = [u_{i2} \ v_{i2}]^T$, it's useful to take in the notation

$$u_i f = 0, \quad (2.21)$$

where $u_i = [u_{i1}u_{i2}, v_{i1}u_{i2}, u_{i2}, u_{i1}v_{i2}, v_{i1}v_{i2}, v_{i2}, u_{i1}, v_{i1}, 1]$, $f = [F_{11}, F_{12}, F_{13}, F_{21}, F_{22}, F_{23}, F_{31}, F_{32}, F_{33}]^T$ and F_{ij} is the element of F at row i and column j , to express the epipolar equation (2.13). For n sets of points the expression becomes

$$Uf = 0, \quad (2.22)$$

where $U = [u_1, \dots, u_n]^T$.

This linear homogeneous equation, and the rank-2 constraint over F , that restrains the matrix to 7 degrees of freedom (since F is also defined up a scalar factor), will permit its unique identification.

The estimation of F may be done using only 7 point matches, or by using 8 or more if enough data points are available. In the latter case, the solution produced is in general unique, and there are several techniques to obtain it. In Zhang 1996's review on the issue [5], the conclusion was that linear techniques are usually sensitive to noise and not very stable, because they ignore the constraints of F and the minimization criterion is not physically meaningful. However, the results could be improved by using normalized data points instead of pixel coordinates.

Nevertheless, non-linear optimization techniques seem to yield better results to the estimation problem of F . Three nonlinear algorithms are available: (i) the distances between points and their corresponding epipolar lines, (ii) the gradient-weighted epipolar errors, and (iii) the distances between points and the reprojections of their corresponding points reconstructed in space. The third one is the most time-consuming method, thus not recommended. The first seems to give the worst results out of the three. Therefore, the second algorithm has been proposed to give the best results in the least amount of computational time.

Note that outliers, such as a wrongly matched pairs of points between the two views of the scene, or points with large location errors, could severely affect the precision of the estimation of F . The reason for this is that all methods are least-square techniques that assume the noise which corrupts the data has zero mean. Hence, more robust techniques have been proposed that are less sensitive to outliers, the most popular being M-estimators, and the least-median-of-squares (LMedS). The LMedS method solves a nonlinear minimization problem yielding the smallest value for the median of squared residuals computed over the entire data set, and proves to be very robust both against false matches, and to erroneous locations, whereas M-estimators is less robust to false matches. LMedS must be solved by a search in the space of possible estimates generated from the data, which is too time-consuming. Thus, as a pragmatic approach is to only test a random sample of the data.

Once the fundamental matrix is obtained through one of those methods, assuming the intrinsic parameters are known, which is the case for this project, where the setup is only one camera, the essential

matrix may be obtained by

$$E = K^T F K, \quad (2.23)$$

where $K = K_1 = K_2$ is the intrinsic parameters matrix. From here, it's possible to retrieve the camera matrix, $P = [R \mid t]$ from the first to the second view, since $E = [t]_{\times} R$.

E may be written as $E = SR$, where S is a skew-symmetric matrix, which in turn can be written as $S = kUZU^T$, where U is an orthogonal matrix, k is the scale and Z , a skew-symmetric matrix, is defined as

$$Z = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}. \quad (2.24)$$

Up to sign, $Z = \text{diag}(1, 1, 0)W$, where

$$W = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (2.25)$$

is an orthogonal matrix. Thus, up to scale, $S = U \text{diag}(1, 1, 0)WU^T$ and

$$E = SR = U \text{diag}(1, 1, 0) (WU^T R) = U \text{diag}(1, 1, 0)V^T, \quad (2.26)$$

is the SVD of E . Having said this, there are two possible factorizations of the essential matrix (ignoring signs),

$$S = UZU^T, \quad R = UWV^T \quad \text{or} \quad UW^TV^T. \quad (2.27)$$

Since $St = [t]_{\times}t = 0$, it follows that $t = U(0, 0, 1)^T = u_3$, the last column of U . The sign of E , and consequently, of t cannot be determined, so there are 4 possible solutions for the camera matrix,

$$P = [UWV^T + u_3] \quad \text{or} \quad [UWV^T - u_3] \quad \text{or} \quad [UW^TV^T + u_3] \quad \text{or} \quad [UW^TV^T - u_3], \quad (2.28)$$

that are represented by (a), (b), (c) and (d), respectively, on Figure 2.3.

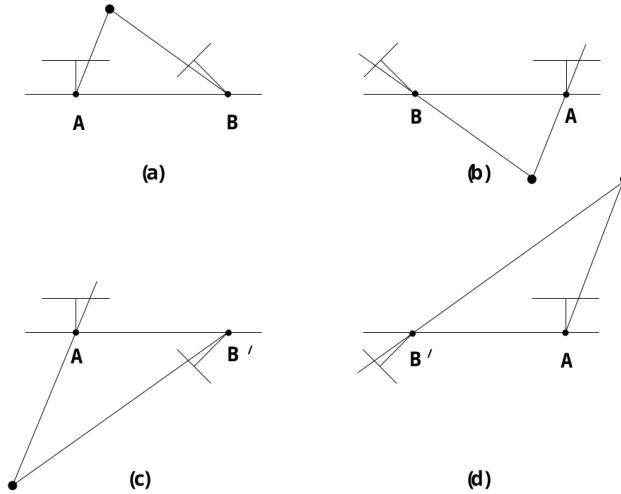


Figure 2.3: Four possible solutions retrieved from E . A and B are the first and second camera center positions when looking at the same scene. Between the right side, (a) and (c), and the left side, (b) and (d), figures, there is a translation direction inversion. Between top, (a) and (b), and bottom, (c) and (d), figures, there is a rotation of 180 degrees around the baseline. [6]

It is shown that the real world point is only in front of both cameras in one of the four solutions, thus, it is enough to use a single point to decide which of the different solutions produces the correct camera matrix. [6]

2.4 Feature Detection and Matching

Apart from methods to determine image transformations from one view of a scene to another, as described previously, it is also necessary to have solid detection and trust-worthy correspondence of the points in each view that are going to be used for computing that transformation.

Two main approaches exist to the problem of feature detection and matching. The first method finds features in one image and matches these features in the other image(s). The second method, which is more suitable for points in the near space, independently detects features in all images and subsequently matches them on the basis of local correspondences.

It might not be appropriate to extract features from an image having high detail (i.e. high spatial bandwidth) on the finest stable scale possible, because when matching those points with the view of the same scene at a different (coarser) scale those details may no longer be detectable. Therefore, it's important to extract features at a variety of scales to ensure their invariance.

Besides scale, preserving rotation and orientation invariance is also a concern. One way to deal with this problem is to define a descriptor. This is a scaled and oriented patch around the chosen point with the local orientation and scale, from which the dominant orientation can be extracted to guarantee its invariance. [7]

There exist a number of algorithms that can do the job. None of these algorithms, however, is optimal for all possible images, as each of them is optimized for particular computer vision applications, with performance depending heavily on the environmental properties (illumination, image quality, contrast, ...). A comprehensive survey on the state of the art of feature detection and description [8] (by Salahat E. and Qasameh M. in 2017), proposes that ideal features should have the following properties:

- Distinctiveness: the gradient variations surrounding the point should be sufficiently large;
- Locality: to avoid obstruction and deformation between two views of a scene;
- Quantity: enough features to describe the view's content;
- Accuracy: features should be accurate (what is accurate?) enough to be detected independently of image scale, shape or pixel location;
- Efficiency: be detected fast enough for real-time systems;
- Repeatability: a high percentage of features should be detected in both views of the overlapping regions in the scene;
- Invariance: deformative effects on the features, due to scaling, rotation or translation of the scene's view, are minimized;
- Robustness: features should be less sensitive to deformations due to noise, blur, compression, and so on.

The review paper also presents an overview on the recent algorithms proposed on this matter, comparing them in terms of the metrics defined above. The analysis in that article suggests that MSER (Maximally stable extremal regions) and SIFT (Scale Invariant Feature Transform) algorithms enhance performance on computational complexity, accuracy and execution time. From the scale and rotation

invariant algorithms, SURF (Speeded Up Robust Features), proved to be faster than SIFT, although not as robust.

MSER² regions are areas of uniform intensity outlined by contrasting backgrounds. They are constructed by trying multiple thresholds on the input image. The ones that remain unchanged in shape over the range of thresholds are the potential features to be selected as areas of interest. The centroids of these areas can subsequently be used to create a feature descriptor for the matching step.

2.4.1 SIFT - Scale Invariant Feature Transform

The SIFT algorithm consists of a local feature detector and local histogram-based descriptor. At first, SIFT takes the original image, and generates progressively blurred images through a Gaussian filter as $L(u, v, \sigma) = G(u, v, \sigma) * I(u, v)$, where (u, v) are the pixel coordinates, L is the blurred image, I is the original image, G is the Gaussian blur, and finally σ is the amount of blur. Then, it re-sizes the original image to half its size, and generates the blurred images again, and repeats. Each new re-sizing of the previous image is scaled by an octave, and each octave can be blurred several times. This process ensures the scale invariance of potential features.

The set of blurred images obtained, then goes through the Difference of Gaussians (DoG) calculation, to find points of interest, or features, later on. To that end, the images from each octave (scale) are subtracted from each other pairwise, as illustrated in Figure 2.4. Mathematically, what happens is obtaining $D(u, v, \sigma) = L(u, v, k\sigma) - L(u, v, \sigma)$, where k is a constant multiplicative (scaling) factor.

Then, the search for local maxima/minima on the resulting images takes place. Each point is compared to its eight neighbors in the current image and nine neighbors in the image of one scale above and below, as Figure 2.5 shows. The points are only selected if they are larger than all of its neighbors or smaller than all of them.

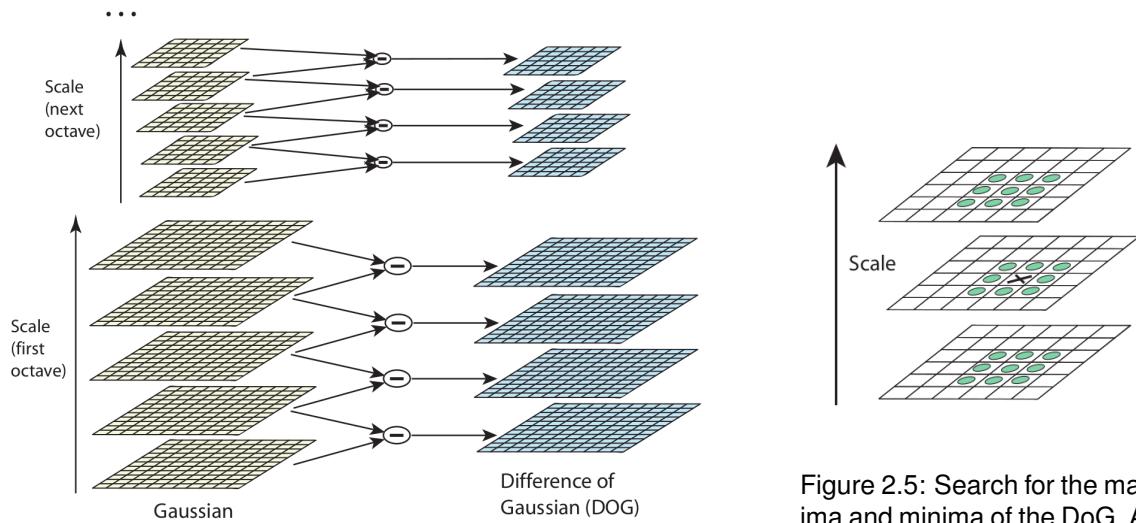


Figure 2.4: Difference of Gaussians (DoG) by octave. To find points of interest (features), the Gaussian-blurred images at each scale (= octave) are subtracted pairwise. This process is computationally more efficient than computing the Laplacian of the Gaussians. [9]

Figure 2.5: Search for the maxima and minima of the DoG. After the DoG, each point is compared with its eight neighbours, and with the 9 corresponding neighbours of the images of the scale above and below. [9]

The location of this maxima/minima is an approximation since it almost never lies on the exact pixel but somewhere between pixels. Thus, the next step is to search for the feature's exact location by using

²J. Matas, O. Chum, M. Urban and T. Pajdla, "Robust Wide Baseline Stereo from Maximally Stable Extremal Regions," 2002.

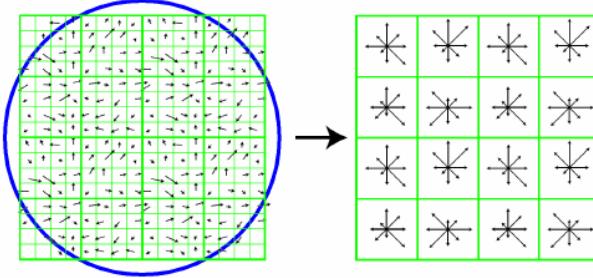


Figure 2.6: SIFT feature descriptor. A window of 16×16 squares is centered around the point of interest. The blue circle represents the Gaussian window. In each 4×4 sub-region, a histogram of orientations is determined, represented on the image on the right. In total there are $8 \times 16 = 128$ values to describe the feature. [9]

a Taylor expansion (up to the quadratic terms) of the scale-space function $D(u, v, \sigma)$, shifted from the origin, according to

$$D(x) = D + \frac{\partial D^T}{\partial x} x + \frac{1}{2} x^T \frac{\partial^2 D}{\partial x^2} x. \quad (2.29)$$

D and its derivatives are determined at each point by using $x = (u, v, \sigma)^T$ as an offset.

Apart from removing low-contrast sub-pixels just determined, it's also important to remove edges since the DoG has a strong response to them (delta functions). An unwanted peak in the DoG will have a large gradient perpendicular to the edge but a small gradient along the edge direction, and because of that property it can be easily detected and eliminated. This is accomplished by computing a Hessian matrix on the point.

Finally, to have a rotation invariant feature, it's necessary to create an oriented patch around the point. For each image, $L(u, v)$, at a certain scale, the gradient magnitude, m and orientation, θ , are determined as

$$\begin{aligned} m(u, v) &= \sqrt{(L(u+1, v) - L(u-1, v))^2 + (L(u, v+1) - L(u, v-1))^2} \\ \theta(u, v) &= \tan^{-1}((L(u, v+1) - L(u, v-1)) / (L(u+1, v) - L(u-1, v))) \end{aligned} \quad (2.30)$$

This information is used to create a histogram of orientations. Peaks on the histogram correspond to dominant directions of local gradients. Local peaks that are within 80% of the highest peak are used to create more features with that orientation.

For the descriptor of the feature, a 16×16 window around the point of interest is set. For each 4×4 region, the orientations are put into an 8 bin histogram (each bin with a 45 degree range). This can be seen through Figure 2.6. The length of each arrow corresponds to the sum of the gradient magnitudes near that direction within the region. Doing this for the 16 regions, there will be 128 values describing the feature. [9]

2.4.2 SURF - Speedup Robust Features

SURF starts from an integral image, which is an intermediate representation of the original image. The entry of an integral image $I_\Sigma(x)$ at a location $x = (u, v)$ represents the sum of all pixels in the input image I of a rectangular region formed by the point x and the origin as

$$I_\Sigma(x) = \sum_{i=0}^{i \leq x} \sum_{j=0}^{j \leq y} I(i, j), \quad (2.31)$$

allowing a faster computation with box-type convolution filters.

The applied box filter is an approximation of Gaussian second-order derivatives. Instead, of having to apply the Gaussian filter iteratively to the output image and to different scales as seen on SIFT, this can be done directly with the box filter by up-scaling it and altering its masks in parallel. Figure 2.7 makes a comparison between the Gaussian second-order derivatives, Laplacians L_{yy} and L_{xy} , and a box filter applied on a point on the y and xy-direction, D_{yy} and D_{xy} .

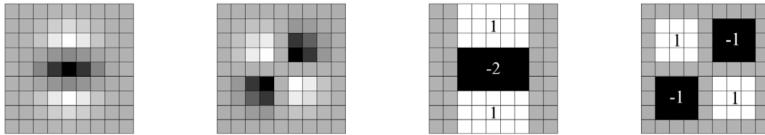


Figure 2.7: Comparison between Gaussian second order derivatives (from the left, the first image in y-direction - L_{yy} - and second image in xy-direction - L_{xy}) and the corresponding images (from the left, the third image in y-direction - D_{yy} - and fourth image in xy-direction - D_{xy}) using a box filter. [10]

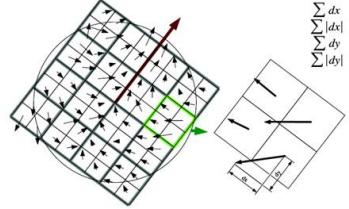


Figure 2.8: SURF feature descriptor. The sampling window is rotated towards the dominant orientation. The vector $(\sum d_x, \sum d_y, \sum |d_x|, \sum |d_y|)$ is generated per each region through the Haar-wavelet responses in x and y direction. [11]

The Hessian matrix at point x and scale σ is defined as

$$\mathcal{H}(x, \sigma) = \begin{bmatrix} L_{xx}(x, \sigma) & L_{xy}(x, \sigma) \\ L_{xy}(x, \sigma) & L_{yy}(x, \sigma) \end{bmatrix}, \quad (2.32)$$

where $L_{xx}(x, \sigma)$ is the convolution of the Gaussian second-order derivative $\frac{\partial^2}{\partial x^2}g(\sigma)$ with the image I at point x , and similarly for $L_{xy}(x, \sigma)$ and $L_{yy}(x, \sigma)$, which are weighted on the calculation of the determinant by $w = \frac{|L_{xy}(1.2)|_F |D_{yy}(9)|_F}{|L_{yy}(1.2)|_F |D_{xy}(9)|_F} \approx 0.9$. The determinant is then defined as

$$\det(\mathcal{H}_{approx}) = D_{xx}D_{yy} - (wD_{xy})^2, \quad (2.33)$$

and is used to find the local change around a point. Points where the determinant is maximal are chosen as features.

SURF's feature descriptor was designed to be scale and rotation invariant. The descriptor is sampled over a window that is proportional to the image scale, so that when a scaled version of the feature in another image is found, the descriptor will be sampled relatively, satisfying the scale-invariance requirement.

To obtain rotation invariance, the sampling window is rotated towards the dominant orientation. The window size is $20s$, where s is the scale at which the point was detected, and is divided into a quadratic grid with 4×4 square sub-regions. For each region, the Haar-wavelet responses are computed in the x and y directions. The descriptor is the sum of the x responses over its four quadrants, sum of the absolute values of the x responses and similarly for y, hence it may be defined as $(\sum d_x, \sum d_y, \sum |d_x|, \sum |d_y|)$. Having 4 values per region and 16 regions, yields a total of 64 values to describe the feature. Figure 2.8 helps at visualizing the process.

The dominant orientation, here lightly and carelessly explained, is captured through the calculation of Haar-wavelet responses over a circular region around the point of interest with a radius of $6s$. Those responses are then summed within a sliding orientation window covering an angle of $\frac{\pi}{3}$. [10] [11] [12]

2.4.3 Matching step

Regarding the matching of feature descriptors, one of the most common search algorithms used is the Nearest Neighbor Search. Fast Approximate Nearest Neighbour (FLANN)³ provides a library for performing this kind of searches in high-dimensional spaces. It contains a collection of algorithms that the authors found to work best for nearest neighbor search and a system for automatically choosing the best algorithm and optimum parameters depending on the dataset. Pyramid Match Kernel⁴ also seems to be successful approach to the matching problem. It forms a partial matching between two sets of feature vectors with linear time complexity.

2.5 Representation of eye orientations in 3 dimensions

The eye is a rotating body (no translations), and thus any orientation can be described by a unique series of three rotations around each of the axis defined in three dimensional space (roll axis, x, pitch axis, y, and yaw axis, z, respectively):

$$R_x(\psi) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \psi & -\sin \psi \\ 0 & \sin \psi & \cos \psi \end{bmatrix}, R_y(\phi) = \begin{bmatrix} \cos \phi & 0 & \sin \phi \\ 0 & 1 & 0 \\ -\sin \phi & 0 & \cos \phi \end{bmatrix}, R_z(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (2.34)$$

where each column vector of the matrices from (2.34) appears counter-clockwise when the axis about which they occur points toward the observer. An arbitrary rotation may then be defined as some multiplication (i.e. a serial order) of those three, noting that the order by which they are multiplied matters (rotations are non-commutative).

Instead of using multiple rotations done after each other, Euler's theorem states that the orientation of the rotating body can also be parametrized by a single rotation with an angle, ρ , about an axis in 3D, $\hat{n} = (n_1, n_2, n_3)$. That rotation is denoted as $R(\hat{n}, \rho)$ and is given by

$$R(\hat{n}, \rho) = \begin{pmatrix} \cos \rho + n_1^2(1 - \cos \rho) & n_1 n_2(1 - \cos \rho) - n_3 \sin \rho & n_1 n_3(1 - \cos \rho) + n_2 \sin \rho \\ n_1 n_2(1 - \cos \rho) + n_3 \sin \rho & \cos \rho + n_2^2(1 - \cos \rho) & n_2 n_3(1 - \cos \rho) + n_2 \sin \rho \\ n_1 n_3(1 - \cos \rho) - n_2 \sin \rho & n_2 n_3(1 - \cos \rho) + n_1 \sin \rho & \cos \rho + n_3^2(1 - \cos \rho) \end{pmatrix}. \quad (2.35)$$

2.5.1 Rotation axis and angle

There are multiple ways by which to obtain the axis and angle of the rotation matrix. The following one uses the fact that a vector, u , parallel to the rotation axis, is necessarily an eigenvector of the rotation matrix with the eigenvalue $\lambda = 1$, as described by

$$Ru = u. \quad (2.36)$$

Rewriting (2.36) as $(R - I)u = 0$, it can be shown that

$$0 = R^T 0 + 0 = R^T(R - I)u + (R - I)u = (R^T R - R^T + R - I)u = (I - R^T + R - I)u = (R - R^T)u. \quad (2.37)$$

³M. Muja and D. Lowe, "Fast approximate nearest neighbors with automatic algorithm configuration," 2009.

⁴K. Grauman and T. Darrell, "The Pyramid Match Kernel: Efficient Learning with Sets of Features", 2005

Because $(R - R^T)$ is a skew-symmetric matrix, u may be chosen such that $[u]_{\times} = (R - R^T)$ (this notation is described in (2.20)), and (2.37) then becomes a cross-product of vector u with itself,

$$(R - R^T) u = [u]_{\times} u = u \times u = 0, \quad (2.38)$$

obtaining

$$u = \begin{pmatrix} R_{32} - R_{23} \\ R_{13} - R_{31} \\ R_{21} - R_{12} \end{pmatrix}, \quad (2.39)$$

as the rotation axis. This does not work if R is symmetric, to do so, it is necessary to diagonalize the matrix and find the eigenvector which corresponds to the eigenvalue of 1.

The rotation angle can be obtained through $\|u\| = 2 \sin \rho$, or through a more direct method by computing the trace of $R(\hat{n}, \rho)$, defined as

$$\text{Tr } R(\hat{n}, \rho) = 1 + 2 \cos \rho. \quad (2.40)$$

2.5.2 Rotation vector in head-fixed coordinates

There are other ways to represent a rotation that may be more appropriate for this project, such as the next one, which comes from Neuroscience. Here, there is a rotation vector, r , directed along the rotation axis, \hat{k} , that has its length varies with the amount of rotation, ρ , around it. A description of this vector is given by $r = \tan(\frac{\rho}{2})\hat{k}$, that can be obtained through the previous rotation matrix as

$$\begin{aligned} \cos \rho &= 1/2 \cdot (R_{11} + R_{22} + R_{33} - 1) \\ k_x &= \sin \theta (R_{32} - R_{23}) / 2 \\ k_y &= \sin \theta (R_{13} - R_{31}) / 2 \\ k_z &= \sin \theta (R_{21} - R_{12}) / 2. \end{aligned} \quad (2.41)$$

Note that in this format, the units are given in half-radians.

2.5.3 Quaternions

Quaternions are 4-dimensional complex algebraic objects that are related to a rotation around an axis, \hat{n} , just like the representation before, by an angle ρ as follows,

$$q = \cos(\rho/2) + i \sin(\rho/2)\hat{n} \equiv q_0 + \mathbf{i} \cdot \mathbf{q}, \quad (2.42)$$

where q_0 and \mathbf{q} are the scalar and vectorial parts of the quaternion, respectively, and \mathbf{i} is the complex 3D vector. For the description of eye movements, only the vectorial part is necessary. Quaternions are related to rotation vectors by $r = \mathbf{q}/q_0$. [13]

2.5.4 Coordinate system

The three-dimensional coordinate systems typically used in Neuroscience and Computer Vision use different conventions, as shown on Figure 2.9. For the remainder of this section, the neuroscience convention will be used to explain the next concepts.

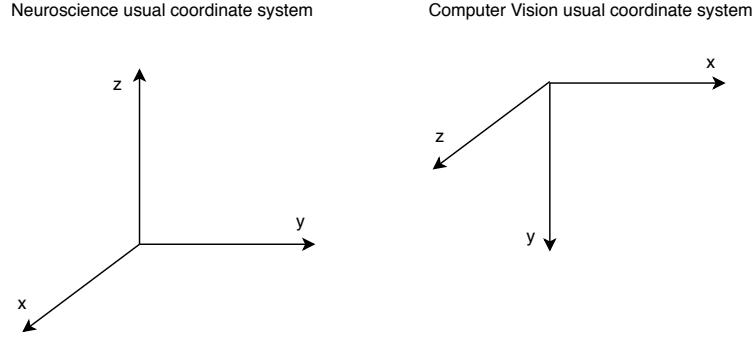


Figure 2.9: The neuroscience and computer vision 3 dimensional coordinate systems normally used are different; both use the right-hand rule. In Neuroscience, torsion is along the frontal visual x-axis, vertical angles along the inter-aural y-axis, and horizontal angles along the bottom-up z-axis.

2.5.5 Donders' and Listing's laws

As described in the beginning of this section, the orientation of a freely rotating rigid body can be quantified by a horizontal angle, θ , a vertical angle, ϕ , and a torsional angle, ψ (around the z, y and x, axis, respectively), relative to a starting position (the primary direction) for which all three angles are zero. Thus, 3 degrees of freedom.

However, as referred before, Donders' law states that the orientation of the eye, when looking in a specific direction at infinity, is always the same; in other words, as described by Donders: "with the head erect and looking at infinity, any gaze direction has a unique torsional angle, regardless the path followed by the eye to get there". Thus, the eye has 2 degrees of freedom for pointing.

Moreover, the 3D orientations of the eye can be uniquely described by head-fixed Cartesian coordinate system of single-axis rotations that bring the eye from the primary position to the current orientation. Described in this way, Listing's law holds that the rotation axis corresponding to all possible eye orientations are confined to the yz-plane, called Listing's plane. The primary position of the eye is an unique position from which any other eye position can be reached through rotation around an axis lying in Listing's plane. This hypothetical rotation axis is defined by a vector $r = (r_x, r_y, r_z)$, where r_x , the torsional component perpendicular to the Listing's plane, is zero. Figure 2.10 illustrates that this law is very precise, by showing 3500 different eye orientations made by a head-restrained monkey, where the left image is the frontal view on the rotation axis, meaning the gazed points coordinates of the horizontal, r_y , and vertical, r_z , components of the axis, and the right-hand plot is a side view, with the coordinates expressed by the horizontal and torsional, r_x , components. A plane is clearly seen around $r_x = 0$ with little deviation ($\sigma_x \approx 0.6$ deg).

It should be noted, however, that the dynamic angular velocity axis, ω , that rotates the eye from one orientation in the Listing's plane to another, e.g. during a saccadic eye movement, is **not** confined to the plane itself. It will typically have a torsional component, $\omega_x \neq 0$, as it depends on the initial eye position, $r_1 = (0, r_{1y}, r_{1z})$ and the saccade difference vector in Listing's Plane, $d_{12} \equiv r_2 - r_1 = (0, d_{12y}, d_{12z})$ that carries the eye to the final eye position, as $\omega_x = r_{1y} \cdot d_{12z} - r_{1z} \cdot d_{12y}$ [3].

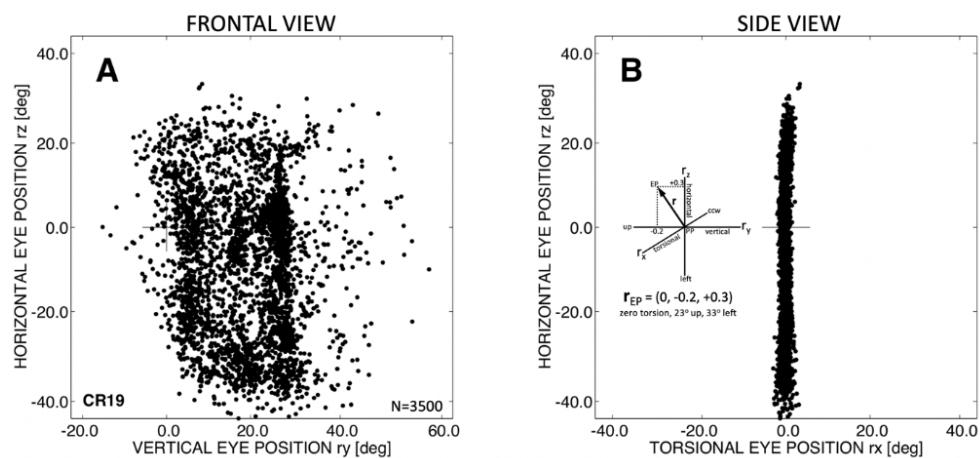


Figure 2.10: Listing's law, illustrated for 3500 eye orientations made by a head-restrained monkey, freely looking around in the laboratory room. Eye fixations are represented by the head-fixed Cartesian components of the 3D rotation axes, expressed in half-radians, and calculated as $\text{angle} = 2 \cdot \text{atan}(\text{component})$. [3]

Chapter 3

Method and Proposal

In this section, the Computer Vision convention of 3D body orientation (referred in section 2.5) is used for the sake of clarity.

The three main aspects that constrain the estimation accuracy of the pose transformation between views and the computational time in which this can be accomplished, are

- the feature detection in both images of the respective views;
- the matching of corresponding features;
- and the computation of the transformation through the set of feature pairs.

Hence, the project will be divided into sub-projects that optimize precision/accuracy and speed for each of these aspects. In this section, the state of the art described previously is taken into consideration, and a comparison between the various techniques outlined for each of the three aspects, as well as, a novel method for computing the transformation that is applicable in the eye model, is proposed. The final aim is to study the benefits/disadvantages of using a camera instead of the current IMU sensor, and to understand how they could be best used in a complementary manner to reach better estimates for the eye's orientation.

Regarding the transformation, two methods were presented in the state of the art, Orthogonal Procrustes Problem and Epipolar Geometry. The first one may be suitable enough for pure rotations. However, as stated in section 1.3, in the current eye model, the camera is not fixed to the center of the eye, so there is also a translation component associated with the transformation. In order to know the exact rotation, this issue could be solved by eliminating the translation, as it was described in subsection 2.2.1. Yet, the camera is being utilized on the project, does not provide depth information. Therefore, to use the Procrustes, it's necessary to project the points gathered from the digital images to a sphere, and for that reason, it doesn't make sense to use the translation correction. In such cases, epipolar geometry would have to be relied on. Epipolar geometry actually demands the existence of a translation component in order to work, since the transformation is taken from the essential matrix, $E = [t] \times R$, that depends on both the translation and the rotation. For that reason, the novel method is proposed to be compared with those two.

The first step to the propounded technique is to find the relationship between the camera orientation and the camera translation. To do so, the diagram of Figure 3.1 is considered. Let the World reference frame W have its origin in the center of rotation of the eye and be aligned with the eye at the rest position. The transformation from the world, W , to the camera's reference frame on the first view of the scene,

C_1 , shown in the figure, is defined by

$${}_{C_1}^W T = \begin{bmatrix} {}_{C_1}^W R & {}_{C_1}^W RB \\ 0 & 1 \end{bmatrix}, \quad (3.1)$$

where $B = [b_x b_y b_z]^T$ is the baseline and ${}_{C_1}^W R$ corresponds to the rotation from the world to the first-view camera's reference frame. Similarly,

$${}_{C_2}^W T = \begin{bmatrix} {}_{C_2}^W R & {}_{C_2}^W RB \\ 0 & 1 \end{bmatrix}, \quad (3.2)$$

is the transformation from the world to the second-view's camera reference frame, C_2 .

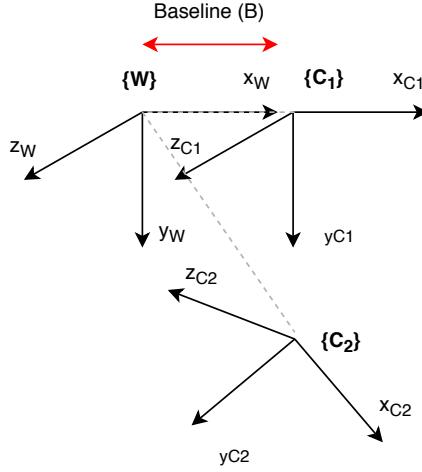


Figure 3.1: Hypothetical coordinate reference frames of the world W , the camera looking at a scene the first view C_1 and the camera looking at the same scene through the second view C_2 after the transformation. Baseline, B is represented as a translation in the x direction only, but in reality it's a vector defined by $[b_x \ b_y \ b_z]^T$. In this case the rotation from W to C_1 is the identity.

Now, the transformation from the first view to the second view of the scene can be translated by ${}_{C_2}^{C_1} T = {}_{C_1}^{C_1} T {}_{C_2}^W T = ({}_{C_1}^W T)^T {}_{C_2}^W T$, which is equivalent to

$${}_{C_2}^{C_1} T = \begin{bmatrix} ({}_{C_1}^W R)^T & -{}_{C_1}^W RB \\ 0 & 1 \end{bmatrix} \begin{bmatrix} {}_{C_2}^W R & {}_{C_2}^W RB \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} ({}_{C_1}^W R)^T {}_{C_2}^W R & ({}_{C_1}^W R)^T {}_{C_2}^W RB - {}_{C_1}^W RB \\ 0 & 1 \end{bmatrix}. \quad (3.3)$$

Thus, the translation component can be characterized as a function of the rotations described by the camera as

$${}_{C_2}^{C_1} t = B(({}_{C_1}^W R)^T {}_{C_2}^W R - {}_{C_1}^W R). \quad (3.4)$$

This constraint can be used to minimize the distances between the matching pairs of image points, initialized by the rotation matrix determined with the Orthogonal Procrustes Problem. This proposed technique will be called Minimization of the Back projection Error, and may be defined as

$$\underset{R, B, z_{1i}}{\text{minimize}} \quad J(m_{1i}, m_{2i}, R, B, z_{1i}), \quad (3.5)$$

where $R = {}_{C_2}^{C_1} R$ is the rotation matrix from the first to the second view, z_{1i} is the depth of the points when

looking at the scene from the first view, which is an unknown quantity, since the camera doesn't gather any direct depth information, $m_{1i} = [u_{1i} \ v_{1i}]^T$ and $m_{2i} = [u_{2i} \ v_{2i}]^T$ are the matched points in the images taken from both views and index i refers to a generic point varying from 1 to N pairs of matched points.

The real world points, M_{1i} and M_{2i} stared at from the first and second view of the scene, respectively, are related by R and t as

$$\begin{aligned} M_{2i} &= RM_{1i} + t \\ M_{1i} &= R^T M_{2i} - R^T t, \end{aligned} \quad (3.6)$$

where $t = {}_{C_2}^{C_1} t$, and can be obtained from the image points through,

$$\begin{aligned} KM_{1i} &= z_{1i}\tilde{m}_{1i} \\ KM_{2i} &= z_{2i}\tilde{m}_{2i}, \end{aligned} \quad (3.7)$$

where K is the intrinsic parameters matrix of the camera. As such, those points may be back projected to retrieve the image points on the complementary image planes as,

$$\begin{aligned} z_{1i}\tilde{m}_{1i}' &= K(R^T K^{-1} z_{2i}\tilde{m}_{2i} - R^T t) \\ z_{2i}\tilde{m}_{2i}' &= K(RK^{-1} z_{1i}\tilde{m}_{1i} + t). \end{aligned} \quad (3.8)$$

where \tilde{m}_{1i}' and \tilde{m}_{2i}' represent the estimation of \tilde{m}_{1i} and \tilde{m}_{2i} . This results in the following values for $m'_{1i} = [u'_{1i} \ v'_{1i}]$ and z_{1i} ,

$$\begin{aligned} u'_{1i} &= \frac{z_{2i}(KR^T K^{-1})_1\tilde{m}_{2i} - (KR^T)_1t}{z_{2i}(KR^T K^{-1})_3\tilde{m}_{2i} - (KR^T)_3t} \\ v'_{1i} &= \frac{z_{2i}(KR^T K^{-1})_2\tilde{m}_{2i} - (KR^T)_2t}{z_{2i}(KR^T K^{-1})_3\tilde{m}_{2i} - (KR^T)_3t} \\ z'_{1i} &= z_{2i}(KR^T K^{-1})_3\tilde{m}_{2i} - (KR^T)_3t, \end{aligned} \quad (3.9)$$

and similarly in,

$$\begin{aligned} u'_{2i} &= \frac{(KRK^{-1})_1z_{1i}\tilde{m}_{1i} + t_1}{(KRK^{-1})_3z_{1i}\tilde{m}_{1i} + t_3} \\ v'_{2i} &= \frac{(KRK^{-1})_2z_{1i}\tilde{m}_{1i} + t_2}{(KRK^{-1})_3z_{1i}\tilde{m}_{1i} + t_3} \\ z_{2i} &= (KRK^{-1})_3z_{1i}\tilde{m}_{1i} + t_3, \end{aligned} \quad (3.10)$$

for $m'_{2i} = [u'_{2i} \ v'_{2i}]$ and z_{2i} values, where $(\cdot)_j$ represents the row j of the 3×3 matrices. Thus, the objective function to minimize can defined by,

$$J(m_{1i}, m_{2i}, R, B, z_{1i}) = \sum_{i=1}^N [(u'_{1i} - u_{1i})^2 + (u'_{2i} - u_{2i})^2 + (v'_{1i} - v_{1i})^2 + (v'_{2i} - v_{2i})^2], \quad (3.11)$$

with the translation, that of course is expressed as seen before by,

$$t = B(R - {}_{C_1}^W R), \quad (3.12)$$

where ${}_{C_1}^W R$ is known because it either refers to the rotation from the world to the camera's starting position (which is known), or to the rotation computed in a previous iteration. And also, with z_{1i} expressed through the depth of the points on the second view,

$$z_{1i} = R^T z_{2i} - R^T t. \quad (3.13)$$

Chapter 4

Preliminary Work

4.1 Specifications

4.1.1 Hardware

The IMU sensor on the current eye model is LP-Research Motion Sensor CAN bus and USB version (LPMS-CU)¹ being used through the USB 2.0 Bus.

The camera handled was an IDS (Image Developing Systems) uEye² model UI-3271LE-C-HQ-VU.

The computer being used to gather the data is an Intel NUC Kit NUC6i7KYK³ running on Windows 10 OS.

4.1.2 Software

OpenCV (Open Source Computer Vision) Library⁴ was used for the programs implemented, using C++. It was mandatory to use the IDS uEye software libraries to handle the cameras and the LPMS software⁵ library LPSensor version 1.3.4 (!) to work with the IMU sensor.

All the work developed regarding software is available at <https://github.com/Mrrvm/IEEC-prework>.

The code developed to capture images from the camera takes an average of 0.12s to execute, including allocating temporary memory, freezing the video to take the picture, passing it to an OpenCV structure necessary later on and freeing the memory allocated.

To calibrate the camera using a 4mm lens, one set of 65 chessboard (of 8x8 squares with 20mm square length) images was taken for model 3271LE-C-HQ-VU. It calibrated with 96% of images of the set, according to the number of chessboards detected.

4.2 Comparison of rotations around each axis between the IMU sensor and the camera

As an initial approach, the Orthogonal Procrustes problem was implemented using the SVD class⁶ from OpenCV, by projecting the image points into a sphere of radius 1.

¹<http://lp-research.com/lpms-cu/>

²https://www.prolinx.co.jp/supplier/IDS/uEye_Manual_En/index.html?hw_modellbezeichnungen.html

³<https://www.intel.com/content/www/us/en/products/docs/boards-kits/nuc/nuc-kit-nuc6i7kyk-features-configurations.html>

⁴<https://opencv.org/>

⁵<https://bitbucket.org/lpresearch/openmat/downloads/>

⁶https://docs.opencv.org/3.4.2/df/df7/classcv_1_1SVD.html

To obtain the matches, first SURF⁷ algorithm was implemented in order to detect the image's points of interest. The initial focus here was to test a reliable computation of a transformation between two images, and not the minimal speed in which this could be possibly done. Using FLANN⁸, the point matches were determined. The minimal distance between both images features is calculated, to further enhance the matching. If the pair of points distances less than double the minimal distance or small arbitrary value (0.02) in the event that the distance is very small, it's defined as a good feature to keep [14].

In order to test this implementation, the structure on Figure 4.1 was used to provoke rotations around each axis, torsional, horizontal and vertical. Because neither the camera nor the sensor are centered at the joint of the structure, there is a translation associated to the rotations. The experiments are described below.

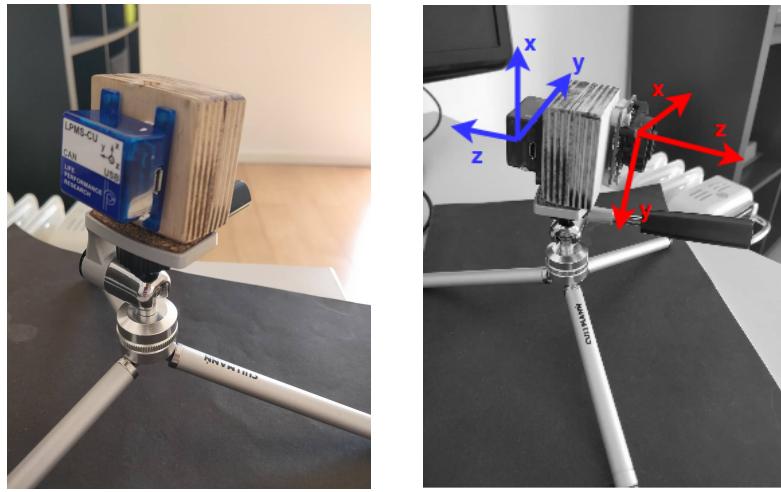


Figure 4.1: Structure used for making the experiments. The image on the right shows the different coordinate systems of the IMU sensor and the camera.

For all tables, θ indicates the absolute angle of rotation expected, θ_{Cam} is the absolute rotation angle obtained by the camera using the Procrustes problem without considering translation and θ_{IMU} is the absolute rotation angle determined by the IMU sensor.

4.2.1 Rotation around the torsional axis

Table 4.1: Rotation angle obtained when rotating the structure around the torsional component. To obtain 90°, a paper strip was placed vertically and the structure was rotated until the strip appeared horizontally on the image.

θ	θ_{Cam}	θ_{IMU}	$n_{matches}$
90°	86.76°	90.9°	7

Table 4.1 shows that, the IMU sensor gave a closer approximation to the expected value. Due to the environment conditions, it wasn't possible to find more than 7 matches between the 2 images (see figure 4.2), which then lead to a poorer estimation through the camera.



Figure 4.2: Matches, represented by colored lines, obtained for rotating 90° around the torsional axis.

Table 4.2: Rotation angle obtained when rotating the structure around the horizontal component. To obtain 30°, the camera was placed at 1 meter away from the grid and the center of the image was shifted by 57cm horizontally. To obtain 5°, the camera was placed at 4 meters away from the grid and the center of the image as shifted by 35cm horizontally.

θ	θ_{Cam}	θ_{IMU}	$n_{matches}$
30°	31.75°	10.63°	11
5°	5.27°	6.12°	27

4.2.2 Rotation around the horizontal axis

In table 4.2, the estimation through the IMU for the 30° rotation is very far from what was expected, and might be caused by the fact that the IMU could not use gravity to determine the orientation since it's horizontal. Figures 4.3 and 4.4 show the matches obtained between the images for the 30° and 5° rotation, respectively.

4.2.3 Rotation around the vertical axis

Table 4.3: Rotation angle obtained when rotating the structure around the vertical component. To obtain 20°, the camera was placed at 1 meter away from the grid and the center of the image was shifted by 35cm vertically. To obtain 5°, the camera was placed at 4 meters away from the grid and the center of the image as shifted by 35cm vertically.

θ	θ_{Cam}	θ_{IMU}	$n_{matches}$
19°	21.37°	16.56°	23
5°	5.41°	5.35°	28

Lastly, table 4.3, shows improvement when using the the camera to determine rotations around the vertical axis. Figures 4.5 and 4.6 show the matches obtained between the images for the 20° and 5° rotation, respectively.

Overall, except for the torsional case where there very few matches, the camera proved to gather better approximations to reality in relation to the sensor. However, further experiments, with a better environment and using a different structure to hold the camera and the sensor must be done, in order to get more accurate results. Because the structure has to be moved by hand, it is very prone to error. The experiments weren't done directly on the eye model, as it still needs some tweaks to support the camera in its structure.

⁷https://docs.opencv.org/3.4/d5/df7/classcv_1_1xfeatures2d_1_1SURF.html

⁸https://docs.opencv.org/2.4/modules/flann/doc/flann_fast_approximate_nearest_neighbor_search.html

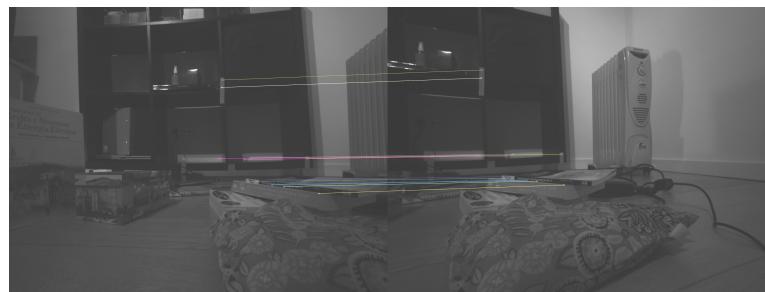


Figure 4.3: Matches, represented by colored lines, obtained for rotating 30° around the horizontal axis.



Figure 4.4: Matches, represented by colored lines, obtained for rotating 5° around the horizontal axis.

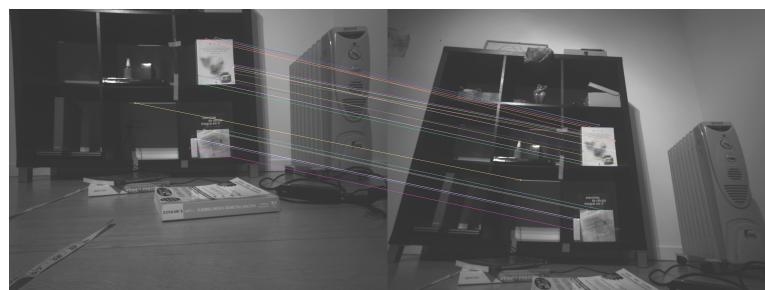


Figure 4.5: Matches, represented by colored lines, obtained for rotating 20° around the vertical axis.



Figure 4.6: Matches, represented by colored lines, obtained for rotating 5° around the vertical axis.

Chapter 5

Thesis Development Plan

As mentioned on the introduction, the IMU sensor has drift associated and since there is complementary parallel work being developed on this research topic, it is necessary to have a good working algorithm to estimate the camera transformation from view to view at reasonable speed. Therefore, out of necessity, this project, as it already has by the preleminary work, starts by determining a good computation method for the transformation between views of the camera.

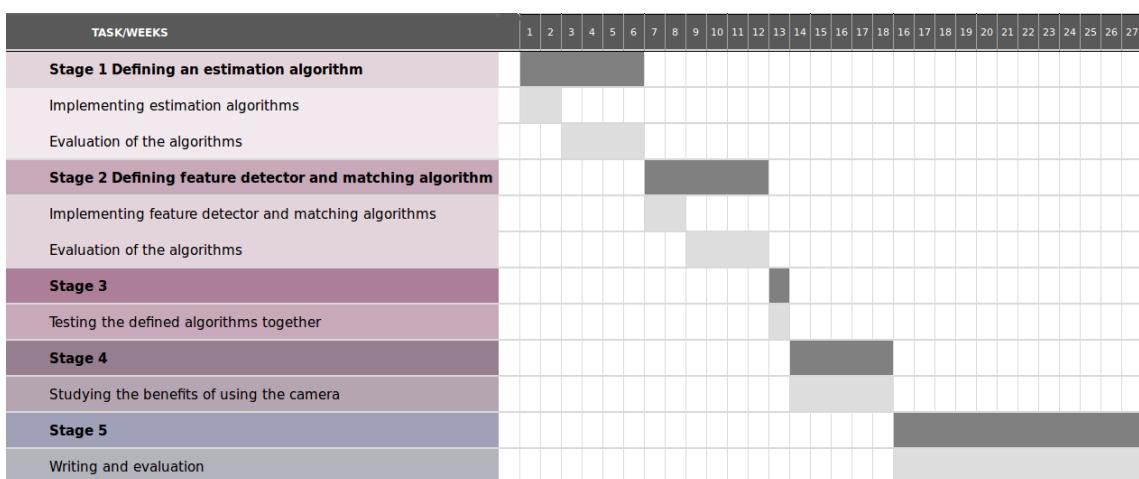


Figure 5.1: Gantt chart of the thesis development. Starts in 18th February 2019.

Keeping that in mind, the development of the thesis can be separated into the following main stages,

1. Defining an estimation algorithm for the camera views transformations;
2. Defining a feature detector and matching algorithm;
3. Testing the conjunction of 1 and 2 over accuracy and speed;
4. Studying the advantages/disadvantages of the camera over the IMU sensor or potential symbiotic relation
5. Evaluation and Writing,

organized into the Gantt chart in Figure 5.1.

Bibliography

- [1] J. van Opstal. http://www.mbfys.ru.nl/~johnvo/OrientWeb/orient_1.html#Abstract. Accessed on December 15th 2018.
- [2] M. Lucas, "Construction and Characterization of a Biomimetic Robotic Eye Model with Three Degrees of Rotational Freedom: A Testbed for Neural Control of Eye Movements," 2017.
- [3] J. van Opstal, "Editorial: 200 Years Franciscus Cornelis Donders," Donders Center for Neuroscience, Dept. of Biophysics, Radboud University Nijmegen, The Netherlands, December 2018.
- [4] J. Gower and G. Dijksterhuis, *Procrustes Problems*. Oxford Statistical Science Series, 2004. Chapter 3 and 4.
- [5] Z. Zhang, "Determining the Epipolar Geometry and its Uncertainty: A review," 1996.
- [6] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2 ed., 2003. Part II: Chapter 9.
- [7] R. Szeliski, *Computer Vision: Algorithms and Applications*. Springer, 2011. Chapter 4.
- [8] E. Salahat and M. Qasaimeh, "Recent Advances in Features Extraction and Description Algorithms: A Comprehensive Survey," 2017.
- [9] D. Lowe, "Distinctive Image Features from Scale-Invariant Keypoints," 2004.
- [10] H. BayTinne, T. TuytelaarsLuc and L. Van Gool, "SURF: Speeded Up Robust Features," 2006.
- [11] H. BayTinne, A. Ess, T. TuytelaarsLuc and L. Van Gool, "Speeded Up Robust Features (SURF)," 2007.
- [12] D. Mistry and A. Banerjee, "Comparison of Feature Detection and Matching Approaches: SIFT and SURF," 2017.
- [13] J. van Opstal, A. Berthoz , *Multisensory Control of Movement: Representaion of the eye position in three dimensions* . Oxford University Press, 1993. Chapter 2.
- [14] OpenCV 3.0.0-dev documentation. https://docs.opencv.org/3.0-beta/doc/tutorials/features2d/feature_flann_matcher/feature_flann_matcher.html#feature-flann-matcher. Accessed on January 2nd 2018.