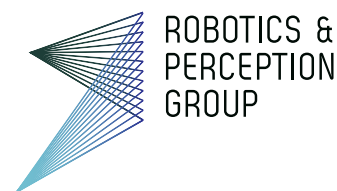




University of Zurich
Department of Informatics



Hans Muster

RPG Thesis Template

Semester Thesis

Robotics and Perception Group
University of Zurich

Supervision

First Supervisor
Second Supervisor

January 2013

Contents

Abstract	iii
Nomenclature	v
1 Introduction	1
1.1 Related Work	2
2 Approach	4
2.1 PTAM method	4
2.1.1 Place Recognition	4
2.1.2 Real Pose Finder	7
2.1.3 Real Pose Finder Alternative	8
2.2 Using ferns	11
2.2.1 Seminaive Bayesian Approach	12
2.2.2 Training	13
3 To be removed	15
3.1 Headings	15
3.2 References	15
3.3 Writing Equations	15
3.4 Including Graphics	16
3.5 Including Code in your Document	16
4 Experiments	17
5 Discussion	18
5.1 Future Work	18
A Something	19

Abstract

Compress the introduction in a few key sentences. No more than half a page.

Nomenclature

Notation

J	Jacobian
H	Hessian
T _{WB}	coordinate transformation from B to W
R _{WB}	orientation of B with respect to W
${}_W\mathbf{t}_{WB}$	translation of B with respect to W , expressed in coordinate system W

Scalars are written in lower case letters (a), vectors in lower case bold letters (**a**) and matrices in upper case bold letters (**A**).

Acronyms and Abbreviations

RPG	Robotics and Perception Group
DoF	Degree of Freedom
IMU	Inertial Measurement Unit
MAV	Micro Aerial Vehicle
ROS	Robot Operating System

Chapter 1

Introduction

Real-time monocular Visual Odometry (VO) algorithms can be used to estimate the 6 DoF pose of a camera relative to its surroundings. This is attractive for applications such as mobile robotics (mostly aerial vehicles where not much power is available) and Augmented Reality (AR) because cameras are small and self-contained and therefore easy to attach to autonomous robots or AR displays. Further, they are cheap, and are now often pre-integrated into mobile computing devices such as PDAs, phones and laptops.

SVO (Semi-direct Bisual Odometry) [3] is a very fast VO algorithm able to run at more than 300 frames per second on a consumer laptop. It builds a map based on keyframes and salient points. Most monocular VO are feature-based where scale and rotation invariant descriptors (SIFT, SURF...) are extracted and matched in order to recover the motion from frame to frame while finally refining the pose with reprojection error minimization with the map. SVO uses a different approach by using direct methods. Instead of matching descriptors, it uses intensity gradient to minimize the error between patches around detected salient points over the frame to frame transformation. Finally, it uses Bundle Adjustment to align with the map and avoid or minimize derive.

The main problem with most existing monocular VO implementations (including SVO) is a lack of robustness. Rapid camera motions, occlusion, and motion blur (phenomena which are common in all but the most constrained experimental settings) can often cause tracking to fail. While this is inconvenient with any tracking system, tracking failure is particularly problematic for VO systems: not only is camera pose lost, but the estimated map could become corrupted as well.

This problem is accentuated during a fast agile maneuver (e.g., a flip) and so a good relocalization is important when these are intended to be performed. The envisaged relocalization scheme proceeds as follows:

- In a training stage, the vehicle explores the environment where the relocalization is supposed to occur. During this stage, an appropriate repre-

sensation of the scene is created.

- The vehicle executes an agile maneuver during which vision-based tracking is no feasible.
- During the actual relocalization phase, the 6 DoF pose in the built map must be estimated.

1.1 Related Work

Place Recognition

Klein and Murray presents in [5] the relocalization method used in PTAM [4]. PTAM is a VO algorithm based on keyframes that are used during the relocalization. The relocalization method consists of two steps. First, given the current frame, the most similar keyframe is retrieved, and its known pose is used as a baseline. As measure of similarity the difference between subsampled, blurred and zero-mean images is used. This measure is a cross correlation. The small blurry images are stored every time there is a new keyframe and the small blurry image of a new frame is computed during the relocalization to be compared with the keyframes.

Other methods can be used for image retrieval, for example using bag of words. Nistér and Stewenius [7] proposes to use a tree structure to store words in order to handle much larger vocabulary or have a much faster retrieval. Every node of the tree would have k child nodes which are the clustering results of k -means. The tree would be built by recursive k -means.

This structure is expensive to build because k -means is very resource consuming. During the online process, new words can be appended to the final leaves.

Özuysal et al. [8] proposes a simplified random forest classifier which relates image patches to objects. It is simplified because instead of using a tree structure they use a linear structure applying all the binary tests to the patch. The result of the tests is a binary descriptor, the list of binary tests is called Fern. Every object is trained with multiple random warps of the known view to introduce information from possible different views of the object. In the end every object can be represented with many binary descriptors and every descriptor should output a probability distribution of possible objects represented. Evaluating multiple Ferns and joining the yielded distributions the final classification is achieved.

Pose Estimation

During the second step of the relocalization of PTAM, the transformation from the retrieved frame is calculated. This transformation will be finally appended

to the known keyframe pose. To do so, an image alignment algorithm, Efficient Second-Order Minimization method (ESM) [2], is employed. ESM is a Gauss-Newton gradient descent algorithm which can be used with different image warp functions, it is a Lucas-Kanade [1] algorithm that uses Second-order functions; therefore, results in a faster convergence.

Geometric methods are typically used to find the transformation from the found keyframe using the classic pipeline of salient points detection, feature extraction and matching. The 5pt algorithm can then be used to find the 6 DoF transformation or the 3pt algorithm if depth is known.

Joint Place and Pose estimation

One approach to solve the relocalization problem was proposed by Williams [11]. In their implementation, they use Random Forest classifiers to characterize a salient object in space. To do so, the classifier needs to be trained with as many as possible representations of the object (multiple views). Therefore, the first time an object is found, multiple warps of the patch are used to initialize its presence in the classifier. On later encounters with the object, the classifier is incrementally trained with additional data. During the relocalization phase salient points are classified using the trained classifier and the 3pt algorithm is used to recover the 6 DoF position. This method is memory expensive and requires GPU power to generate the patch warps.

Shotton et al. [9] also propose a method using random forests. RGB-D data is used to train the classifier. In this case, all the information is encoded in the classifier so no previous data storing or computing (salient point detection, descriptor extraction, etc...) is needed. The classifier is trained to an individual RGB-D pixel, and an RGB-D pixel query will output a probability distribution over the position in \mathbb{R}^3 . This can be applied to all pixels of a frame or to a sparse subset selection of them. Ideally, the camera pose can be inferred from only three pixels, but as the output of the classifier can be very noisy, a second step is applied. From the output from many pixels an energy function is minimized using preemptive RANSAC in order to find a pose that agrees with most of the distributions.

To train this method a very complete dataset of RGB-D images with 6 DoF poses from the environment associated to them is needed. That makes it difficult to be used with SLAM problems where the map gets populated incrementally. An online training method should be developed.

Chapter 2

Approach

Two kind of approaches will be proposed. On one side, local approaches based on the PTAM implementation, where two steps are performed. The first step has been named *Place Recognition* and the second *Real Pose Recognition*. Multiple methods will be proposed to solve the second step. Then, on the other side, a global approach will be proposed. In this case machine learning methods (*ferns*) will be used to recognize points in space.

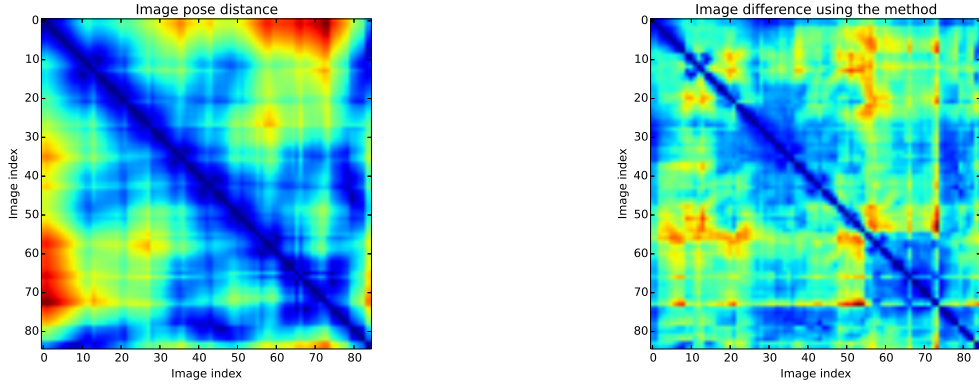
2.1 PTAM method

PTAM is a VO algorithm based on keyframes and so the relocalization method proposed is based on keyframes as well. Every keyframe has associated with it a camera pose that will be used to relocalize. During the relocalization there are two steps involved. We have called the first step *Place Recognition* and the second *Real Pose Finder*.

2.1.1 Place Recognition

During this step, the algorithm tries to find the keyframe image most similar to the last acquired image. The pose associated with the most similar keyframe is used as an initial rough estimation of the current pose. The similarity score should be resistant to view point because the new acquired image will, most probably, never be taken from the same pose as any of the keyframes. Also it should be fast to compute.

The used similarity score is the Cross Correlation between images meaning the sum of the squared error between two zero-mean images. To make computation faster both images are resized become 40×30 . Then, to make the images more resistant to view point changes it is blurred with a 3×3 Gaussian kernel with $\sigma = 2.5$. The resulting stored image is a resized, blurred and zero-mean



(a) Image to image real distance between the keyframe pose

(b) Image to image distance approximated using the Cross Correlation value

Figure 2.1

image called *small-blurry-image*.

During the normal map building pipeline this image is computed and stored every time a new keyframe is added to the map. And then, during the relocalization, the sum of squared difference between every stored *small-blurry-image* and the last acquired frame is computed to find the most similar keyframe and then use its pose as an initial estimation of the current camera pose.

Method validation

To evaluate the method the real distance between two frames is going to be compared with the Cross Correlation value described above. Far away frames should be dissimilar and close by frames should be more similar and have a lower CC value. In figure 2.1a can be seen the pair distances between image using the real pose while in figure 2.1b there is the approximated pair distances using the CC value as distance. It can be seen that they have a similar distribution. Also the correlation of 0.4337 shows that one explains the other in most cases.

Finally, to show that this method can be used, for every image the most CC similar image was taken being it real K closest image. Ideally the most CC similar image should always be the closest image. In figure 2.2 there is the count of occurrences of each K . It can be seen that most images resolve to the first or second closest image using this method.

In figure 2.3 of what would the algorithm return for a given query image for which there is no pose information available. Figure 2.3a is the image seen by the camera at the moment of the relocalization and 2.3b is what the system found to be the most similar, closest, image using the CC procedure explained

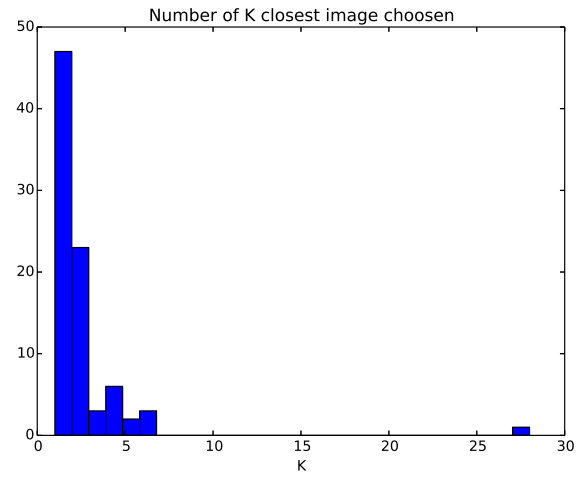


Figure 2.2: Distribution of closest chosen images using CC



(a) Query image



(b) Found closest image

Figure 2.3

previously. In this case the two image's camera pose are not the same, but they should be similar enough in order to find the difference in posterior steps.

2.1.2 Real Pose Finder

The second step of the relocalization algorithm from PTAM tries to refine the pose of the found to be most similar keyframe to explain the current pose of the camera. During this step, in the implementation from PTAM, only rotations are corrected. An image alignment through optimization algorithm (ESM) ?? is used to find the $SE(2)$ transformation between the two and then another minimization is performed to find a transformation in the world frame.

Image alignment

World frame interpretation

This transformation is not easily interpreted in $SE(3)$, a translation in pixels can mean different things in the world frame depending on the distance to the object. Also, in $SE(3)$ there are three angular DoF, not one. The found transformation can be interpreted in multiple ways in the world frame. In the PTAM implementation it is assumed that the translation will only involve rotations, and so, the next step is the mapping from $SE(2)$ to $SO(3)$.

A Gauss-Newton minimization algorithm is used to find the $SO(3)$ model that modifies the image in the same way as the found $SE(2)$ model. The minimization is over the $SO(3)$ parameters ξ and the error is expressed as in 2.1

$$\delta_i(\xi) = T_{SE(2)} u_i - \pi(T_{SO(3)}(\xi) p_i) \quad \text{where} \quad u_i = \pi(p_i) \quad (2.1)$$

where u_i is a pixel position (during the implementation $u_0 = (5, 0)$ and $u_1 = (-5, 0)$ are used). Also the Jacobian of δ is needed during the minimization process.

$$\frac{\partial \delta(\xi)}{\partial \xi} = - \frac{\partial \pi(b)}{\partial b} \Big|_{b=p} \frac{\partial T_{SO(3)}(\xi)}{\partial \xi} \Big|_{\xi=0} \quad p \quad (2.2)$$

with

$$\frac{\partial \pi(b)}{\partial b} \Big|_{b=p} = \frac{f}{z} \begin{bmatrix} 1 & 0 & -\frac{x}{z} \\ 0 & 1 & -\frac{y}{z} \end{bmatrix} \quad \text{where} \quad p = \begin{bmatrix} x \\ y \\ z \end{bmatrix} \quad \text{and} \quad f = \text{focal length}$$

(2.3)

$$\frac{\partial T_{SO(3)}(\xi)}{\partial \xi_k} = G_k \quad \text{where } G = SO(3) \text{ Generator} \quad (2.4)$$

Method validation

To visualize the results of the different optimization procedures described above, the found transformation has been applied to the image. In the first step, the $SE(2)$ transformation between two image is computed and the transformed image on every step is one of the used variables. In figure 2.4 can be seen the final transformation found from 2.3b to 2.3a. In figure 2.6 the error can be visualized, there it can be seen that translation and rotation are well corrected but there still is a misalignment caused mostly by a change on scale which is not taken into account during the alignment.



Figure 2.4: $SE(2)$ transformed image

On the other side, during the second minimization where there $SO(3)$ translation is found, no image is actually involved, only two pixel coordinates. To generate a visualization of the translation the calibration of the camera is needed. Every pixel is unprotected form the image into the image plane, then the transformation is applied to it and finally it is projected back to the image. The result of the described procedure can be seen in figure 2.5.

2.1.3 Real Pose Finder Alternative

During the mapping of an area the VO algorithm finds landmarks in the world frame which are associated with detected featured points in keyframes (i.e. every featured point in an image is related to a $3D$ position in the world frame). Given a new image, some extracted featured points can be related to a keyframe using the description-matching fashion which at the same time are related to world positions. From this information the full 6 DoF translation $SE(3)$ can be



Figure 2.5: $SO(3)$ transformed image

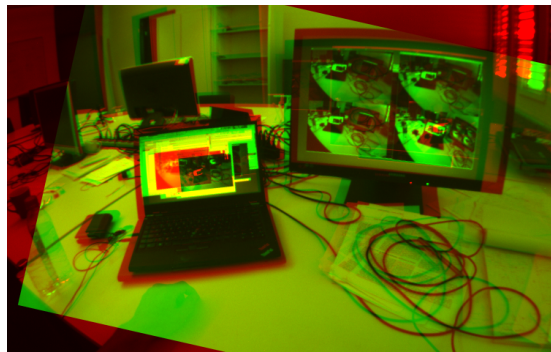


Figure 2.6: $SE(2)$ transformation error visualization

computed using the three point algorithm.

The 3pt algorithm and its implementation is described in [6]. In this case the described *Central absolute pose* is dealt with.

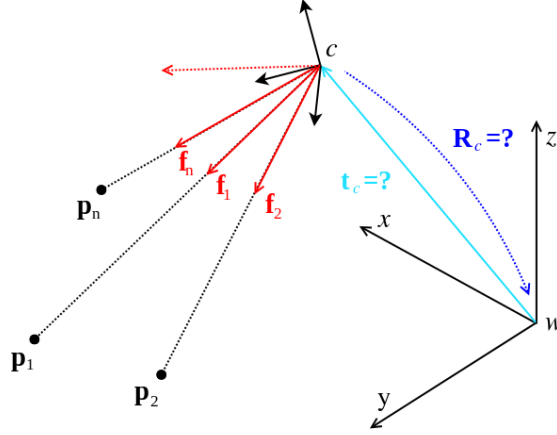


Figure 2.7: From camera frame bearings (or image pixels if the camera calibration is available) and fix frame points the transformation between the two frames can be computed using the 3pt algorithm. Image from [6]

First, descriptors from every featured point are extracted, both SIFT and SURF can be used updating a configuration file. Second, a brute force KNN matching is performed from all the extracted descriptors from on image and all the descriptors of the second image. The first, and second most similar descriptor are retrieved. Then, only good matches are kept, that is, using the matching technique described by Lowe, only matches with a descriptor ratio between the first and the second closest match of 0.8 or less are kept, only discriminant matches are used.

Finally, the 3pt algorithm is fed with the pixel positions from the query image and the landmarks from the other. Because there are still outliers after the described simple filtering, this process is run in RANSAC framework.

Method validation

Figures 2.8 and 2.9 are an example of the described above. The pose outputted was used to correctly relocalize and so it was correct.

The first part, where descriptors are extracted, matches and filtered, can be seen in figure 2.8. It can be seen that most matches are correct after this first filtering.

The inliers found during the RANSAC process can be seen in 2.9, it has to be kept in mind that the found transformation is not from image to image

but from world, using world points, to frame. These world points can not be visualized and its precision could lead to not being included to the inliers group.

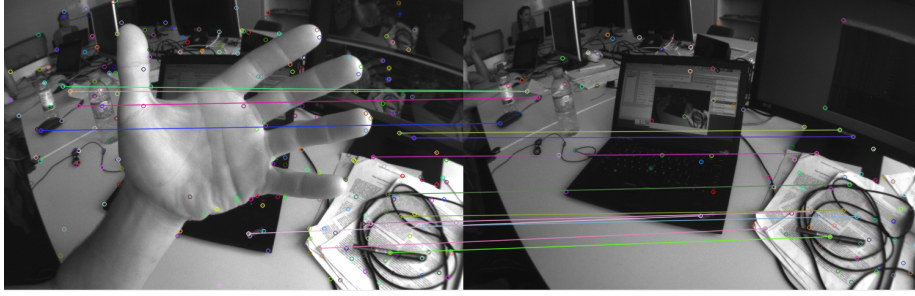


Figure 2.8: Accepted matches using SIFT

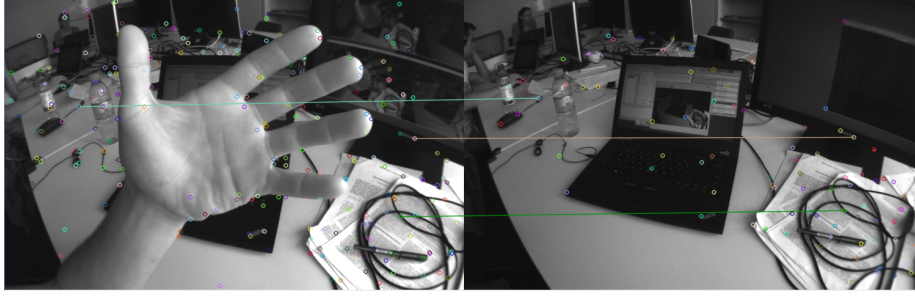


Figure 2.9: Inliers from RANSAC used to calculate the pose with the three-point algorithm

2.2 Using ferns

As said previously with the three-point algorithm it is possible to recover the 6 DoF of the camera pose from the relation from pixel coordinates and points in space. In this case machine learning techniques are used to model this relationship. In the classifier scheme, an object in space is a class and multiple views seen from the camera should all be classified as this class.

A *fern* is a descriptor made from a set of binary tests 2.5, made in a similarly to random forests but flattened. When used as a classifier every possible evaluation of a *fern* will contain a posterior provability distribution for every class, like random forests' final leaves do.

$$f_j = \begin{cases} 1, & \text{if } I(d_j, 1) < I(d_j, 2) \\ 0, & \text{otherwise} \end{cases} \quad (2.5)$$

The result of the tests is encoded into the bits of an unsigned integer 2.1. If a *fern* has S tests then its representation can have $K = 2^S$ possible values which are then used as index of its posterior distribution.

Listing 2.1: Fern evaluation

```
uint8_t fern = 0;
for (size_t j = 0; j < S; ++j)
{
    //Shift bits
    fern <<= 1;
    if (I(d_j,1) < I(d_j,2))
    {
        //Change last bit
        fern++;
    }
}
```

One fern is usually not descriptive enough to correctly classify, in [8] it is claimed that with 50 ferns and $S = 11$ a problem with 200 different classes is tractable. In memory, it would involve $50 \times 2^{11} \times 200 = 20480000$ elements to be stored in memory. If these are stored as *float* then 78 MB are needed, which is tractable. It should be noticed that the problem does not scale well on S , but at the same time, it is the most critic parameter.

2.2.1 Seminaive Bayesian Approach

Let $c_i, i = 1, \dots, H$ be a set of classes and $f_j, j = 1, \dots, N$ be a set of binary tests 2.5 calculated over an image patch that is being classified. Formally, we are looking for

$$\hat{c}_i = \underset{c_i}{\operatorname{argmax}} \quad P(C = c_i | f_1, f_2, \dots, f_N) \quad (2.6)$$

then, Bayes' formula tells

$$P(C = c_i | f_1, f_2, \dots, f_N) = \frac{P(f_1, f_2, \dots, f_N | C = c_i) P(C = c_i)}{P(f_1, f_2, \dots, f_N)} \quad (2.7)$$

Assuming a uniform prior and since the denominator is a scaling factor the problem is reduced to

$$\hat{c}_i = \underset{c_i}{\operatorname{argmax}} \quad P(f_1, f_2, \dots, f_N | C = c_i) \quad (2.8)$$

Since this simple tests are very simple many of them are needed ($N \approx 300$), the storage of joint probabilities for every outcome is no feasible because 2^N entries

for each class would be needed. One way to avoid this storage is to assume total independence between tests, this way the class provability distributions could be calculated like

$$P(f_1, f_2, \dots, f_N | C = c_i) = \prod_{j=1}^N P(f_j | C = c_i) \quad (2.9)$$

but this ignores the possible correlation between tests. With ferns a trade off is achieved, tests are grouped in sets F and then the conditional provability becomes

$$P(f_1, f_2, \dots, f_N | C = c_i) = \prod_{k=1}^M P(F_k | C = c_i) \quad (2.10)$$

It can be seen here why S is important, an increment in the number of tests in a *fern* means an increment in the modelled correlation.

2.2.2 Training

During the training step many different views of a world point are needed to correctly model it. In the original work [8] only one image is used to train and to generate more possible views of the objects multiple (about 10,000) randomly generated warps are applied to it, all this warped images are used for training. Those warps include affine transformations, noise addition and smoothing with a Gaussian kernel (with all the parameters randomly picked from a uniform distribution).

In the studied case here, the algorithm can already provide multiple real views of the object (around 5-10), but still on each patch many randomly generated warps are applied and used to train (usually 100). In this case, the warps are simplified to include only image rotation and scale.

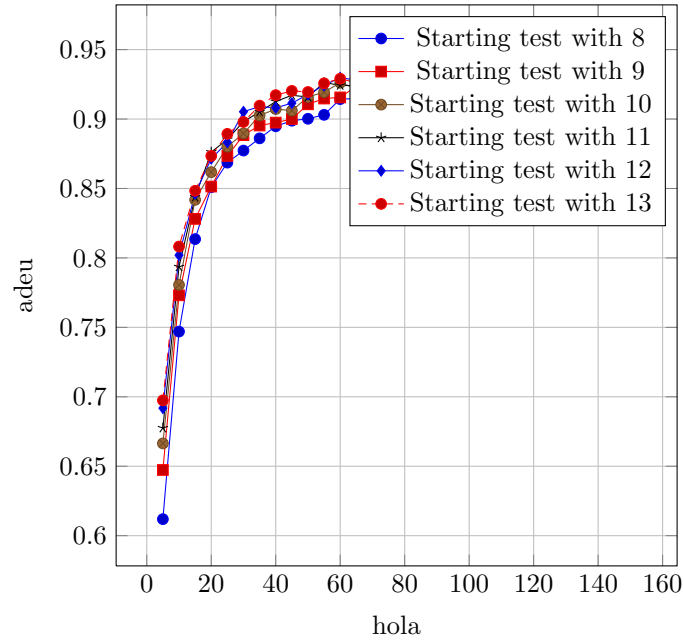
During the training, every warped patch is evaluated with every *fern*. The count of a class evaluation on each *fern* is performed along with the count of per class patches used. Those values are used during the classification step to calculate the provability distributions

$$p_{k,c_i} = \frac{N_{k,c_i} + N_r}{N_{c_i} + K \times N_r} \quad (2.11)$$

where N_{k,c_i} is the number views of class c_i that evaluates to k and N_{c_i} is the number of views of class c_i to train. A regularization term $N_r = 1$ is used to avoid provabilities evaluating to zero.

Method Validation

To validate the classifier a small test has been set up. One image is taken and from it 30 random patches are used to train the classifier. Then, 100 random warps are applied to those patches and classified. This process is reproduced 100 times.



hola

Chapter 3

To be removed

Describe the main steps in your algorithm. An illustration is always helpful.

Here are some L^AT_EX tips:

3.1 Headings

Your report can be structured using several different types of headings. Use the commands `\chapter{.}`, `\section{.}`, `\subsection{.}`, and `\subsubsection{.}`. Use the asterisk symbol `*` to suppress numbering of a certain heading if necessary, for example, `\section*{.}`.

3.2 References

References to literature are included using the command `\cite{.}`. For example [4, 10]. Your references must be entered in the file `bibliography.bib`. Making changes or adding new references in the bibliography file can be done manually or by using specialized software such as *JabRef* which is free of charge.

Cross-referencing within the text is easily done using `\label{.}` and `\ref{.}`. For example, this paragraph is part of chapter 2; more specifically on page 15.

3.3 Writing Equations

The most common way to include equations is using the `equation` environment. Use `\eqref{.}` to reference an equation, e.g. (3.1).

$$\begin{aligned} C(\mathbf{x}) &= \frac{1}{2} \sum_{i \in \mathcal{I}} \sum_{k \in \mathcal{K}_i} \mathbf{e}_{i,k}(\mathbf{x})^T \mathbf{W}_{i,k} \mathbf{e}_{i,k}(\mathbf{x}) \\ \hat{\mathbf{x}}^{LS} &= \operatorname{argmin}_{\mathbf{x}} C(\mathbf{x}), \end{aligned} \tag{3.1}$$

$$T_i = \begin{bmatrix} \mathbf{R}_i & \mathbf{p}_i \\ 0 & 1 \end{bmatrix} \quad \text{with } \mathbf{R}_i \in SO(3), \mathbf{p} \in \mathbb{R}^3. \quad (3.2)$$

3.4 Including Graphics

The easiest way to include figures in your document is to use pdf figures if you use `pdflatex` to compile. Figure 3.1 was created with the use of the open source program `ipe`.

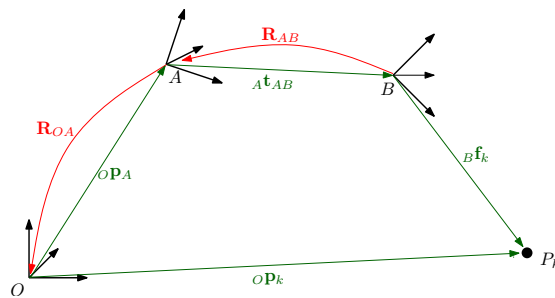


Figure 3.1: Example of a figure.

3.5 Including Code in your Document

You may include samples from your Matlab code using the `lstlistings` environment, for example

Listing 3.1: Matlab Example

```
% Evaluate y = 2x
for i = 1:length(x)

    y(i) = 2*x(i);

end
```

Listing 3.2: C++ Example

```
% sum all elements in a list
int sum=0;
for (list<int>::iterator it=mylist.begin(); it!=mylist.end(); ++it)
    sum += *it;
```

Chapter 4

Experiments

Provide numerical results, plots and timings. Interpret the data.

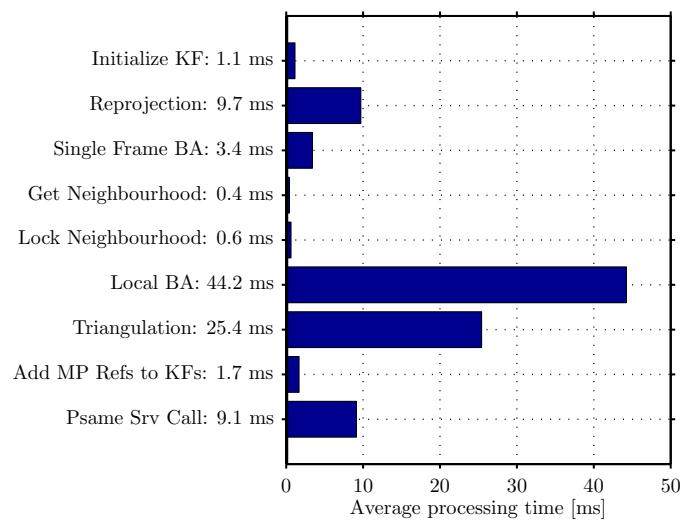


Figure 4.1: Example of a figure.

Chapter 5

Discussion

Explain both, the advantages and limitations of your approach.

5.1 Future Work

How would you extend the work? Can you propose another approach?

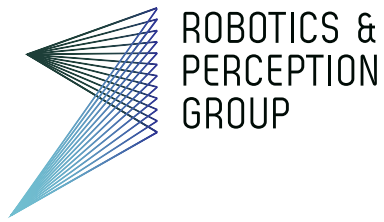
Appendix A

Something

In the appendix you can provide some more data, a tutorial on how to run your code, a detailed proof etc.

Bibliography

- [1] Simon Baker and Iain Matthews. Lucas-Kanade 20 Years On: A Unifying Framework. *International Journal of Computer Vision*, 56(3):221–255, February 2004.
- [2] S Benhimane and E Malis. Homography-based 2D visual servoing. *Robotics and Automation, 2006. ICRA ...*, 2006.
- [3] Christian Forster, Matia Pizzoli, and Davide Scaramuzza. SVO: Fast Semi-Direct Monocular Visual Odometry. *Proc. IEEE Intl. Conf. on Robotics ...*, 2014.
- [4] Georg Klein and David Murray. Parallel Tracking and Mapping for Small AR Workspaces. *IEEE and ACM International Symposium on Mixed and Augmented Reality*, November 2007.
- [5] Georg Klein and David Murray. Improving the agility of keyframe-based SLAM. *Computer Vision ECCV 2008*, 2008.
- [6] Laurent Kneip and Paul Furgale. Opengv: A unified and generalized approach to real-time calibrated geometric vision.
- [7] D Nister and H Stewenius. Scalable recognition with a vocabulary tree. *... Vision and Pattern Recognition, 2006 ...*, 2006.
- [8] M Ozuysal and M Calonder. Fast keypoint recognition using random ferns. *Pattern Analysis and ...*, 2010.
- [9] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene Coordinate Regression Forests for Camera Relocalization in RGB-D Images. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2930–2937, 2013.
- [10] Hauke Strasdat, J M M Montiel, and Andrew J Davison. Real-time Monocular SLAM: Why Filter? *Proc. IEEE International Conference on Robotics and Automation (ICRA)*, pages 2657 – 2664, 2010.
- [11] Brian Williams, Georg Klein, and Ian Reid. Real-time SLAM relocalisation. *... Vision, 2007. ICCV 2007. IEEE 11th ...*, pages 1–8, 2007.



Title of work:

RPG Thesis Template

Thesis type and date:

Semester Thesis, January 2013

Supervision:

First Supervisor

Second Supervisor

Student:

Name:	Hans Muster
E-mail:	muster@student.ethz.ch
Legi-Nr.:	97-906-739

Statement regarding plagiarism:

By signing this statement, I affirm that I have read the information notice on plagiarism, independently produced this paper, and adhered to the general practice of source citation in this subject-area.

Information notice on plagiarism:

http://www.ethz.ch/students/semester/plagiarism_s_en.pdf

Zurich, 22. 5. 2014: _____