

Bachelor's Thesis

Designing and Implementing a Rephotography Application for iOS

Rasmus Diederichsen

First Supervisor: Prof. Dr. Oliver Vornberger

Second Supervisor: Dr. Thomas Wiemann

Department of Computer Science
Department of Cognitive Science

CONTENTS

1	INTRODUCTION	2
1.1	Overview	2
1.2	Previous Approaches	4
1.2.1	Mobile Applications	4
1.2.2	Computational Re-Photography	4
1.3	Goals of this thesis	6
2	CAMERA GEOMETRY	8
2.1	Camera Models	8
2.2	Epipolar Geometry	10
2.3	Essential Matrix Estimation	12
2.3.1	The 8-Point Algorithm	12

INTRODUCTION

This chapter will introduce the notion of rephotography, elaborate on the process of how to make such a photograph and survey existing approaches to simplify it. These include two applications for mobile operating systems which will be briefly discussed. Furthermore, a summary of more sophisticated work by MIT researchers will be given, leading to the problem statement and the goal of this work.

1.1 OVERVIEW

Rephotography or repeat photography denotes the retrieval of the precise viewpoint used for taking a—possibly historic—photograph and capturing another image from the same spot, ideally with the same camera parameters. This allows for documentation and visualisation of changes which the scene has undergone between the two or more captures. For instance when documenting urban development, one can present progress of construction, restoration efforts or changes in the surroundings in a visually striking manner, e.g. by blending the photographs together. Figures 1.1 and 1.2 show examples.

When done manually, the photographer must attempt to find the original viewpoint usually by visual inspection of the original image and trying to match the current camera parameters—camera position, camera rotation, focal length, possibly principal point—to the original. The procedure is often carried out by placing the camera on a tripod and comparing a printout of the original image with what can be seen through the viewfinder or the camera screen. The number of parameters to match as well as the difficulty to estimate them purely from comparing two-dimensional images makes the process error-prone and tedious. Visual acuity and experience of the photographer thus place limits on the accuracy with which the camera pose of the reference image can be reconstructed. Some corrections can be done by post-processing the images and warping the rephotograph with a homography to better match the original.

At the time of writing, few computerised aids are available to the photographer (see below). The advancement of mobile phones and tablet computers with integrated cameras and larger screens presents the opportunity to develop applications which can assist in this endeavour, moving away from the traditional trial-and-error approach. On current digital cameras ¹ this is impossible due to their closed infrastructure not permitting running user programs.

¹ At the time of writing, no commercial manufacturer produces a camera with user-modifiable firm- or software. A project at Stanford (Adams et al., 2010) was discontinued Levoy (2009)



Figure 1.1: Residenzschloss Dresden, destroyed during World War II, © Sergey Larenkov, printed with permission



Figure 1.2: Frauenkirche Dresden, destroyed during World War II, © Sergey Larenkov, printed with permission

1.2 PREVIOUS APPROACHES

1.2.1 *Mobile Applications*

Two applications have been developed to assist a photographer in taking rephotographs. For smartphone operating systems, *rePhoto*² and *Timera*³ exist, both available for Android and iOS devices. These applications support the user by placing a transparent version of the original image over the current camera image, allowing for easier alignment. The captured rephotograph is then presented together with the original image in a blend (c.f. ??) which can be customized in *Timera*.

What is characteristic about both of these applications is that the user must still determine on their own how to actually move the camera. An overlay simplifies the procedure, eliminating some of the inaccuracy introduced into the manual approach by the necessity to move the eyes from printout to camera, but it is still the user's responsibility to determine the necessary motion between the current camera position and the goal position (that of the original image).

1.2.2 *Computational Re-Photography*

A more sophisticated automated approach was presented in by MIT researchers in 2010. Bae et al. (2010) found in preliminary studies that neither a side-by-side view as would be used in the manual approach, nor a linear blend provided by the above applications result in accurate rephotographs.

In this setup, the relevant parameters of a historic image's camera are reconstructed, including the focal length, principal point and the six degrees of freedom in camera pose. This subsection will give a high-level overview, while a more in-depth discussion of the relevant concepts is deferred until ??.

The software runs on a laptop connected to a digital camera. After reconstructing the scene in 3D by use of two images captured by the user, they are then directed by the software to the desired viewpoint, the user does not have to find it by themselves. On the screen, they are shown the current camera image alongside two arrows indicating in which direction to move—one for movement in the sensor plane and one for movement along the optical axis.

Bae et al. (2010) identify five primary obstacles in viewpoint reconstruction of a historic photograph.

1. The necessary camera motion has six degrees of freedom—three for translation and three for rotation—which are challenging for the user to adjust simultaneously, as changing one parameter will often necessitate adjustments for the others to improve the matching. Furthermore, the number of degrees of freedom makes it difficult to communicate to the user how they must move the camera.

² <http://projectrephoto.com/>

³ <http://www.timera.com/Explore>

2. Computing relative translation between two cameras from corresponding image points is possible only up to an unknown scale (see ??), meaning it is impossible to determine e.g. if an object viewed by the camera is small and close or large and further away. This poses the problem of how to determine if the user is close to the desired viewpoint and whether or not they have come closer or moved further away over iterations.
3. Relative pose estimation from corresponding points fails when the motion between the two cameras approaches zero, which is the ultimate goal one wishes to achieve. When naïvely comparing the current camera image to the reference photograph, the estimate for relative rotation and translation would become increasingly unreliable as the camera approaches the original viewpoint.
4. Automated computation of relative camera pose will rely on feature detection to find correspondences. However, historical images will often be vastly different from the current scene. Not only may the scene itself have changed considerably, but also the historical image—having been taken by a historical camera—may differ in contrast, sharpness and colours. Feature detectors may not be able to reliably find correspondences when comparing an old with a new photograph.
5. The calibration data—most importantly, focal length and principal point—of the historical camera are often unknown. The calibration data is needed for relative pose computation.

Initially, after loading a historical image, the user is instructed to take two photographs of the scene with a reasonably wide baseline (about 20°). One of them, termed *first frame* is supposed to be taken from some distance from the original viewpoint, the *second frame* should be the user's best eyeballed approximation of it. The wide baseline allows for a more reliable 3D-reconstruction of the scene used to tackle problems 2. and 3.

SIFT features (Lowe, 2004) are computed and matched between the two images. Given these correspondences, 3D coordinates of the points can be computed. A selection of these is reprojected into the second frame after which the user identifies six or more points in the historical photograph corresponding to these points in the second frame. This allows estimating extrinsic and intrinsic camera parameters of the historical camera by running an optimisation algorithm on an initial estimate for relative rotation and translation between first frame and reference image as well as sensor skew, focal length and principal point of the historical camera (problem 5.). The principal point's initial guess is found again with help of the user who identifies three sets of parallel lines in the historical image (see Hartley and Zisserman, 2004, chapter 8.8).

The result is that the location of the reference camera relative to the first camera is known. During the homing process, the current camera

frame is compared to the first frame (not the reference frame, avoiding problem 3.), which avoids degeneracy due to the wide baseline. Given the locations of the reference camera and the current frame's camera, each relative to the first frame, one can compute the location of the reference relative to the current frame and thus guide the user in the right direction. Hence, the reference photograph is not needed anymore after this initial step, circumventing problem 4.

During homing, the current camera frame is warped according to the necessary rotation before being shown to the user, allowing them to focus only on the translation (problem 1.). This is possible since for rephotography dealing with structures usually at some distance, the rotation will be small, otherwise the warped image would be unusable.

A remaining problem (2.) is that the scale of the necessary translation is unknown, so that only the direction is known. This poses the question of how to determine whether the user has come closer to the goal or not. It may be feasible to find the original viewpoint nonetheless, if it was possible to determine at least when the user reaches it, but this impossible without further information. On top of that, it would make for a better user experience if also the distance to the goal could be communicated.

A key observation in this regard is that the actual scale of the translation is irrelevant, it is sufficient that there be a way to make the scale consistent accross iterations. That is, it is unnecessary to know whether the goal is a specific distance away, if one can ensure that the translations computed one after the other can be somehow meaningfully compared. For this, Bae et. al observe that when triangulating 3D coordinates from corresponding points, their computed distance from the camera (the first frame) is inversely proportional to the distance between the cameras. Therefore, in each iteration, the scale of the world is computed by triangulating correspondences between the first and current frames. The scale is compared to the scale computed in the initial step for the first and second frames. Scaling the current translation vector by the ratio of the two scales makes the length consistent across iterations and decreasing with the distance to the goal.

The results of this method appear to be very successful, but two main drawbacks exist.

- The prototype is not very convenient, as it requires a (laptop) computer and a digital camera to carry out which is impractical for spontaneous rephotography.
- The application is not available to the public, neither in source nor binary form. It is therefore impossible to adapt it for more mobility.

1.3 GOALS OF THIS THESIS

This work's objective can thus be summarised as follows.

1. Implement in a prototypal fashion the process from (Bae et al., 2010) for a mobile operating system so it can be run on a smartphone or tablet and direct the user in approximate real-time.
2. Evaluate the approach and attempt to reproduce the results.

For a proof-of-concept application, several simplifying assumptions are made. Firstly, it is assumed that the “historic” photograph is captured with the same camera as the application is run on and that the camera is calibrated. Secondly, no strong visual differences between the reference and current scenes are assumed so that the reference image is accessible to the same feature detection algorithm without the user manually labelling correspondences.

The application targets iOS 8 and current hardware, as image processing is computationally intense, and has been tested on an iPad Air 2.

2

CAMERA GEOMETRY

This chapter will introduce the geometry of image projection, largely following (Hartley and Zisserman, 2004, chapters 6,7) and the geometry of two views (the epipolar geometry, Ma et al. (2003, chapter 5.1)) and how it can be used to recover relative camera position from two images of the same scene. It largely follows.

2.1 CAMERA MODELS

Given a camera C whose center of projection is the origin and a point X in camera-centric coordinates, the central projection of X onto C 's image plane is depicted in a side-view in Figure 2.1. The image plane is virtually moved to the front of the camera, otherwise the image would be mirrored at the principal point as in real cameras. Let f be the focal length, which is the distance of the image plane to the optical centre. If $X = (X, Y, Z)$, then $x = (f\frac{X}{Z}, f\frac{Y}{Z}, f)$ by use of the intercept theorem, with $(fX, fY)^T$ being the coordinates in the image plane.

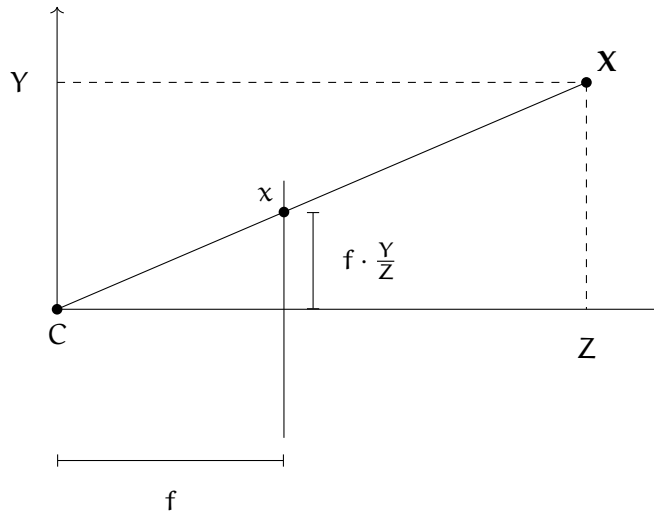


Figure 2.1: Central projection for a pinhole camera

Homogeneous vectors are the elements of projective geometry. They can be obtained from cartesian coordinates by appending a 1-element. All projective entities which differ only by a scalar factor are equivalent, one writes $x \sim y$ if $x = \lambda y, \lambda \neq 0$. This has the added effect that points at infinity can be represented by vectors whose last coordinate is zero.

When representing the points as homogeneous quantities, the central projection can be expressed by a matrix multiplication. This can be written with homogeneous coordinates as

$$\begin{pmatrix} f\frac{X}{Z} \\ f\frac{Y}{Z} \\ 1 \end{pmatrix} \sim \begin{pmatrix} fX \\ fY \\ Z \end{pmatrix} = \underbrace{\begin{pmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}}_{\text{Projection Matrix of } C} \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix} \quad (2.1)$$

or in short

$$\mathbf{x} = \mathbf{P}\mathbf{X} \quad (2.2)$$

The above situation is a special case wherein the camera centre C defines the origin and the optical and image axes are the coordinate axes. Thus, the rotation and translation of the camera relative to this coordinate system is zero. More generally, there might be a world coordinate frame with different origin and different axes, so that a coordinate transform must be applied to \mathbf{X} before the projection.

Let $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ be a rotation matrix giving the camera's rotation relative to the world frame and $\mathbf{t} \in \mathbb{R}^{3 \times 1}$ its translation such that

$$\mathbf{X}_{\text{cam}} = \mathbf{R}\mathbf{X}_{\text{world}} + \mathbf{t} \quad (2.3)$$

Then the projection of a point \mathbf{X} in world coordinates onto the image plane becomes

$$\mathbf{x} = \mathbf{K} \cdot [\mathbf{R} \mid \mathbf{t}] \mathbf{X} \quad (2.4)$$

Real cameras are not ideal pinhole cameras. Furthermore, it is useful to have the dimension of all values be pixel units. A camera has five intrinsic parameters and can be written in matrix form as

$$\mathbf{K} = \begin{pmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \quad (2.5)$$

where f_x and f_y are the focal lengths in x - and y -directions expressed in pixel units (f_x and f_y are not necessarily identical, if the sensor has non-square pixels), s the sensor skew (the pixels may not be rectangular; their edges may not be perpendicular) which is usually zero, and the coordinates of the principal point (c_x, c_y) with respect to the origin of the image plane which usually placed at the upper left corner. The principal point is the intersection of the image plane with the optical axis.

The intrinsic camera parameters assembled in \mathbf{K} are therefore essential to relate world points to image points which will be important for this application. In theory, these parameters could be obtained from the camera's vendor who knows the precise manufacturing specifications. In practice, only the focal lengths f_x, f_y are known, in most cases only one with the assumption of square pixels. Usually, the principal point is assumed to be at the sensor centre and the pixels are assumed to be rectangular. In practice however, there are variances introduced by different causes such as imprecise manufacturing or physical impacts which may decentre the lens such that the principal point is no longer at the centre.

A further complication is introduced by the camera lens which will often have a non-negligible distortion, most prominently radial distortion as depicted in [Figure 2.2](#). It can be modeled by the application of a distortion factor to the ideal undistorted image coordinates (\tilde{x}, \tilde{y}) .

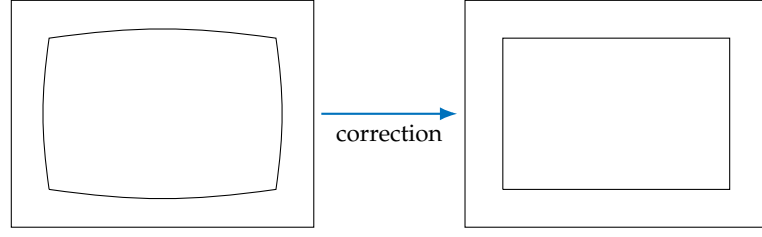


Figure 2.2: Radially distorted image on the left, the corrected image on the right.

$$\begin{pmatrix} x_d \\ y_d \end{pmatrix} = L(r) \begin{pmatrix} \tilde{x} \\ \tilde{y} \end{pmatrix} \quad (2.6)$$

where L is a nonlinear function of the distance r from the distortion center—usually coincident with the principal point. The function can be approximated as an exponential with a Taylor expansion

$$L(r) = 1 + \sum_{i=1}^k \kappa_i r^i \quad (2.7)$$

for a given k (see [Hartley and Zisserman, 2004](#), chapter 7.4). The intrinsic camera parameters which consist in the entries of K and distortion coefficients κ_i must be determined in order to accurately relate world coordinates to image coordinates. They can be found by calibrating the camera. Different methods exist (e.g [Zhang, 2000](#)) but will not be examined here.

2.2 EPIPOLAR GEOMETRY

Epipolar geometry is the geometry which relates the image points in two views of the same scene. [Figure 2.3](#) shows the basic setup.

We consider a scene viewed by two cameras with optical centres c_1 and c_2 , where c_1 defines the origin, world points $\mathbf{X}_i \in \mathbb{R}^3$, where the subscript denotes the coordinate frame—the first camera, arbitrarily chosen to be the left one, or the second camera—and homogeneous image points \mathbf{x}_i which are the projections of \mathbf{X}_i onto the image planes and thus correspond to the same world point. The cameras are related by a rigid body transform (R, T) , where R is a 3×3 rotation matrix and T the translation between the camera centres. Throughout this work, the direction of coordinate frame transformation will be such that

$$\mathbf{X}_2 = R\mathbf{X}_1 + T. \quad (2.8)$$

It is obvious that the following relation holds

$$\lambda_i \mathbf{x}_i = \mathbf{X}_i, \lambda_i > 0 \quad (2.9)$$

that is, the world point lies on a ray through the optical centre and the image point.

Given the corresponding points \mathbf{x}_i in two images, the ultimate goal is to retrieve the euclidean transform (R, T) .

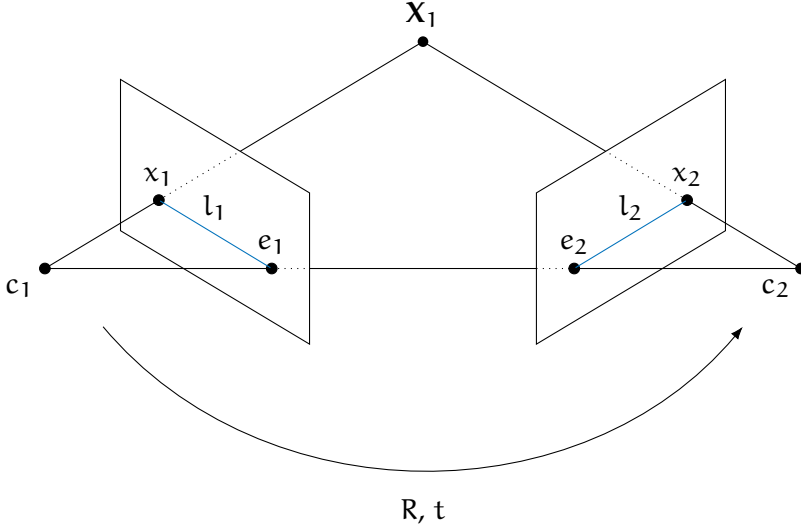


Figure 2.3: Basic epipolar geometry with camera centres c_1, c_2 , image points x_1, x_2 , a world point X_1 , epipoles e_1, e_2 and epipolar lines l_1, l_2

Starting from equation (2.8), one can derive

$$\begin{aligned}
 X_2 &= RX_1 + T \\
 \lambda_2 x_2 &= R\lambda_1 x_1 + T && \text{(by equation (2.9))} \\
 \lambda_2 \hat{T}x_2 &= \hat{T}R\lambda_1 x_1 + \hat{T}T && \hat{T} \in \mathbb{R}^{3 \times 3} \text{ with } \hat{T}x = T \times x \\
 \lambda_2 x_2 \hat{T}x_2 &= x_2 \hat{T}R\lambda_1 x_1 + 0 && T \times T = 0 \\
 \lambda_2 \cdot 0 &= x_2 \hat{T}R\lambda_1 x_1 && \hat{T}x_2 \text{ is perpendicular to } x_2 \\
 0 &= x_2 \underbrace{\hat{T}R}_E x_1 && (2.10)
 \end{aligned}$$

The product $E = \hat{T}R$ is the essential matrix and the constraint it imposes on corresponding image points the essential constraint (see [Ma et al., 2003](#), chapter 5).

An intuition for the essential matrix can be obtained from [Figure 2.3](#). Given an image point in one frame x_1 , in attempting to find the point x_2 corresponding to the same world point X_1 , the epipolar geometry restricts the search space to one dimension—the epipolar line of x_1 in the second camera’s image plane. The camera centres and the world point define an epipolar plane. The backprojection of x_1 is the ray through x_1 and the optical centre c_1 . All points on this ray are mapped to the same point on the image plane of c_2 . Depending on how far away the world point X_1 is from c_1 , its image on the second camera’s image plane will vary—but it will be on the intersection of the image plane and the epipolar plane, the epipolar line of x_1 . The line l_1 may be identified with its coimage (the orthogonal complement of its preimage) $l_1 = e_1 \times x_2$ so that

$$\forall x : x \in l_1 \Leftrightarrow x \cdot l_1 = 0. \quad (2.11)$$

The coimage of the epipolar line is the vector perpendicular to the epipolar plane, so every vector in this plane will have an inner product of 0 with this vector. Constrained to vectors in the image planes, this means that all vectors on the epipolar line will have an inner

product of 0 with this vector l_1 . Referring back to equation (2.10), it can be seen that multiplication with E will yield a term $x_2 E = l$ which fulfills

$$l \cdot x_1 = 0, \quad (2.12)$$

which is precisely the relation stated in equation (2.11). The essential matrix thus maps an image point onto its epipolar line in the other image.

2.3 ESSENTIAL MATRIX ESTIMATION

Estimating the essential matrix between two cameras and decomposing it into relative rotation and translation is a necessity in the endeavour to communicate a necessary camera movement to the application's user. The most prominent algorithms in this regard are the 8-Point algorithm due to [Longuet-Higgins \(1987\)](#) and the 5-Point algorithm proposed by [Nistér \(2004\)](#). To illustrate the mathematical tractability of the problem, the former will be presented below.

2.3.1 *The 8-Point Algorithm*

BIBLIOGRAPHY

- Adams, A., Talvala, E.-V., Park, S. H., Jacobs, D. E., Ajdin, B., Gelfand, N., Dolson, J., Vaquero, D., Baek, J., Tico, M., Lensch, H. P. A., Matusik, W., Pulli, K., Horowitz, M., and Levoy, M. (2010). The frankencamera: An experimental platform for computational photography. *ACM Transaction on Graphics*, 29(4):29:1–29:12.
- Bae, S., Durand, F., and Agarwala, A. (2010). Computational re-photography. *ACM Transactions on Graphics*, 29(3).
- Hartley, R. I. and Zisserman, A. (2004). *Multiple View Geometry in Computer Vision*. Cambridge University Press, second edition.
- Levoy, M. (2009). Camera 2.0: New computing platforms for computational photography. <http://graphics.stanford.edu/projects/camera-2.0/>. accessed 2015-07-28.
- Longuet-Higgins, H. C. (1987). Readings in computer vision: Issues, problems, principles, and paradigms. chapter A Computer Algorithm for Reconstructing a Scene from Two Projections, pages 61–62. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110.
- Ma, Y., Soatto, S., Kosecka, J., and Sastry, S. S. (2003). *An Invitation to 3-D Vision: From Images to Geometric Models*. Springer Verlag.
- Nistér, D. (2004). An efficient solution to the five-point relative pose problem. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(6):756–777.
- Zhang, Z. (2000). A flexible new technique for camera calibration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(11):1330–1334.