

Bachelor Project



**Czech
Technical
University
in Prague**

F3

**Faculty of Electrical Engineering
Department of Cybernetics**

Two-Body Structure from Motion

Petr Hrubý

Supervisor: doc. Ing. Tomáš Pajdla, Ph.D.

Field of study: Open Informatics

Subfield: Computer and Informatic Science

May 2019

I. Personal and study details

Student's name: **Hrubý Petr** Personal ID number: **466292**
Faculty / Institute: **Faculty of Electrical Engineering**
Department / Institute: **Department of Cybernetics**
Study program: **Open Informatics**
Branch of study: **Computer and Information Science**

II. Bachelor's thesis details

Bachelor's thesis title in English:

Two-Body Structure from Motion

Bachelor's thesis title in Czech:

Trojdimenzionální rekonstrukce dvou pohybujících se těles

Guidelines:

1. Review the state of the art in Structure from Motion (SfM) as well as in Multi-Body SfM [1-7].
2. Propose a method for reconstructing scenes consisting of a single moving object in front of a static background captured by multiple image sequences. Consider the scenario when each image sequence captures the object and the background in a static configuration.
3. Implement the method as an extension of the COLMAP [2] SfM pipeline.
4. Demonstrate and evaluate the method on real data.

Bibliography / sources:

- [1] J Schonberger and J-M Frahm. Structure-from-Motion Revisited. CVPR 2016.
- [2] COLMAP <https://colmap.github.io/>
- [3] K Schindler and D Suter. Two-view multibody structure-and-motion with outliers through model selection. IEEE Transactions on Pattern Analysis and Machine Intelligence, 28(6):983–995, 2006.
- [4] K Schindler, D Suter, H Wang. A model-selection framework for multibody structure-and-motion of image sequences. International Journal of Computer Vision, 79(2):159–177, 2008.
- [5] X Xu, L-F Cheong, Z Li. Motion Segmentation by Exploiting Complementary Geometric Models. CVPR 2018
- [6] L Magri, A Fusiello. Robust multiple model fitting with preference analysis and low-rank approximation. BMVC 2015
- [7] E Maset, F Arrigoni, A Fusiello. Practical and efficient multi-view matching. ICCV 2017.

Name and workplace of bachelor's thesis supervisor:

doc. Ing. Tomáš Pajdla, Ph.D., Applied Algebra and Geometry, CIIRC

Name and workplace of second bachelor's thesis supervisor or consultant:

Date of bachelor's thesis assignment: **04.10.2018** Deadline for bachelor thesis submission: **24.05.2019**

Assignment valid until: **30.09.2020**

doc. Ing. Tomáš Pajdla, Ph.D.
Supervisor's signature

doc. Ing. Tomáš Svoboda, Ph.D.
Head of department's signature

prof. Ing. Pavel Ripka, CSc.
Dean's signature

III. Assignment receipt

The student acknowledges that the bachelor's thesis is an individual work. The student must produce his thesis without the assistance of others, with the exception of provided consultations. Within the bachelor's thesis, the author must state the names of consultants and include a list of references.

Date of assignment receipt

Student's signature

Acknowledgements

This research was supported by IMPACT CZ.02.1.01/0.0/0.0/15_003/0000468 EU Structural and Investment Funds, Operational Programme Research, Development and Education. I would like to express my thanks to my supervisor doc. Ing. Tomáš Pajdla, Ph.D. for his guidance and his valuable advice. I would like to thank Ing. Michal Polic, Ing. Stanislav Steidl and Federica Arrigoni, Ph.D. for their advice. I would like to thank my family for their support.

Declaration

I declare that the presented work was developed independently and that I have listed all sources of information used within it in accordance with the methodical instructions for observing the ethical principles in the preparation of university theses.

Prague, 24. May 2019

Abstract

Multi-Body Structure from Motion (MB-SfM) is a problem of 3D reconstruction of objects in a scene, which is not static. The MBSfM problem has not yet been solved in general. We propose a solution to a relaxed version of the problem, where the images depict several static configurations of the scene and the number of the objects is limited to two.

We show in this thesis that this relaxed MBSfM has a practical usage because in some cases the reconstruction of an object can be stabilized, using a background. When the object is placed on the background and multiple sequences of different configurations of the scene are taken, the object can be reconstructed and segmented from the background with the method proposed in this thesis.

Keywords: structure from motion , multi-body, two-body, 3D reconstruction, SfM, MBSfM, dynamic scene

Supervisor: doc. Ing. Tomáš Pajdla, Ph.D.
Aplikovaná algebra a geometrie CIIRC,
Jugoslávských partyzánů 1580/3,
Praha 6

Abstrakt

Trojdimenzionální rekonstrukce více pohybujících se těles (MBSfM) je problém rekonstrukce scény, která není statická. Tento problém zatím nemá obecné řešení. V této práci navrhuje řešení relaxované verze tohoto problému, kde je scéna zachycena v několika statických konfiguracích a počet objektů je omezen na 2.

Dále v této práci ukážeme, že řešení tohoto relaxovaného problému má praktické využití, protože v některých případech může být trojdimenzionální rekonstrukce objektu stabilizována přítomností pozadí ve scéně. Pokud je objekt umístěn na pozadí a tato scéna je vyfotografována v několika statických konfiguracích, může být tento objekt zrekonstruován metodou popsanou v této práci.

Klíčová slova: struktura z pohybu, dvě tělesa, 3D rekonstrukce, SfM, MBSfM, dynamická scéna

Překlad názvu: Trojdimenzionální rekonstrukce dvou pohybujících se těles

Contents

1 Introduction	1	3.1.4 Bundle adjustment	10
1.1 Motivation	1	3.2 Motion segmentation	11
1.2 Thesis structure	2	3.2.1 Algebraic methods	11
1.3 Problem formulation	2	3.2.2 Random sampling methods	12
1.4 Contributions	3	3.2.3 Other methods	13
2 Frequently used concepts	5	3.3 Reconstruction of dynamic scenes	13
2.1 Representation of rotation in 3D	5	4 Proposed solution	15
2.1.1 Norm of the rotation	6	4.1 Reconstruction of the takes	15
2.1.2 Distance between two rotations	6	4.2 Object motion calculation	18
2.2 Camera representation	6	4.3 Bringing the motions into the same coordinate system	20
2.3 PnP	8	4.3.1 Description of the basis transformation	21
3 State of the art	9	4.3.2 Clustering of the basis transformations	24
3.1 Structure from Motion	9	4.3.3 Verification of the basis transformations	27
3.1.1 Feature extraction and matching	9	4.3.4 Transformation of the motions	29
3.1.2 Initial reconstruction	10	4.3.5 Zero motions removing	31
3.1.3 Incremental reconstruction	10	4.4 Motion clustering	33

4.4.1 Inlier recognition	35	5.1 Overview of used objects and backgrounds	59
4.5 Motion verification	36	5.2 Qualitative results of the method	61
4.5.1 Condition for consistent cycles	37	5.3 Comparison to the single body SfM	65
4.5.2 Chordal completion	38	5.3.1 Planar objects	68
4.5.3 Clustering of the motion clusters	44	5.3.2 Repetitive objects	69
4.6 Track building and segmentation	44	5.4 Review of approaches to the Bundle Adjustment	70
4.6.1 Track building	45	6 Future work	73
4.6.2 Track segmentation	45	7 Conclusion	75
4.7 Merging of the reconstructions .	47	Bibliography	77
4.7.1 Order of the merging	47		
4.7.2 Merging of the points	47		
4.7.3 Calculation of the object motion from the transformations	49		
4.7.4 Merging of the cameras	51		
4.8 Bundle adjustment	56		
4.9 Filtering of the reconstruction . .	58		
5 Experiments	59		

Figures

1.1 An example of the input scene; $k = 3, n = 5,$ $t_1 = 1, t_2 = 2, t_3 = 2, t_4 = 3, t_5 = 3,$ the point is connected with the camera if the camera observes the point	3
4.1 Example of different results of the sequential PnP, points with the same color are observed by the same camera. The cameras have been registered in the order green, red, yellow.	17
4.2 Example of a camera after the sequential PnP registration onto the anchor take.	18
4.3 Euler vectors (a) and translations (b) calculated from a dataset consisting of 8 takes. The situation is before the transformation into the same coordinate system. Clusters of motions are observable.	21
4.4 Euler vectors of the rotations between coordinate systems of three reconstructions. The clusters of the vectors are observable.	23
4.5 Euler vectors (a), (b) and translations (c), (d) calculated from a dataset consisting of 8 takes. Images (a), (c) depict the situation before the transformation, images (b), (d) depict the situation after the transformation. The object was always rotated around the vertical axis which is observable in (b). . . .	31
4.6 Examples of results of the spectral clustering of the clusters. Only rotations are depicted.	32
4.7 The Euler vectors (a), (b) and the translations (c), (d) of the object motion. Images (a), (c) depict the situation before the removing, images (b), (d) depict the situation after the removing.	33
4.8 An example of a graph of the clusters. The labels of the vertices are the numbers of the initial and the final takes of the clusters they represent. The valid cycles are $((1, 2), (2, 3), (3, 4), (4, 5), (5, 6), (1, 6)),$ $((1, 2), (2, 6), (1, 6)).$ The invalid cycles are $((1, 2), (2, 3), (2, 6)),$ $((1, 6), (2, 6), (5, 6))$	39
4.9 An example of a cycle after the chordal completion. Red edges are the chords. Red arrows lead from the chord to the 3-cycle which is its predecessor. Blue arrows lead from the 3-cycle to the edges which are its predecessors.	40
4.10 Reconstructions of a single take after the segmentation of the points. Points which do not belong to background nor to the object are not shown.	46
4.11 Points merged from 8 reconstructions of takes.	49
5.1 Images taken from the datasets, which are used for reconstruction. Part 1.	60

5.2 Images taken from the datasets, which are used for reconstruction. Part 2.	61
5.3 Reconstructed datasets, Part 1..	62
5.4 Reconstructed datasets, Part 2..	63
5.5 Reconstructed datasets, Part 3..	64
5.6 Reconstructed datasets, Part 4..	65
5.7 Example of an image from the single body dataset without the background.....	66
5.8 Comparison of the models reconstructed with and without the background. The models on the left side have been reconstructed with the background.....	67
5.9 Comparison of the models of the "Catalog" object reconstructed with and without the background. The model (a) has been reconstructed with the background.....	68
5.10 Examples of repetitive patterns on "Daliborka" (a) and "Vatican" (b)..	69
5.11 Comparison of the reconstruction of the back side of the "Vatican" object. The model (a) has been reconstructed with the background.	69
5.12 Model of the "Lego" object. ...	70

Tables

5.1 Properties of the datasets.	59
5.2 Results of the reconstructions... ..	62
5.3 Number of images from which the object is reconstructed	66
5.4 Number of points from the object	66
5.5 Median reprojection error E_O ..	68
5.6 Overall median reprojection error E	71
5.7 Median reprojection error E_O of the object.....	71
5.8 Median reprojection error E_B of the background	72



Chapter 1

Introduction



1.1 Motivation

Structure from Motion [1] is a problem of reconstruction of a 3D model of a scene from a set of images which depict this scene. The case, where the scene is static, is well understood and implemented in many pipelines, such as [2] [3] [4].

Multi-Body Structure from Motion (MBSfM) is the case of a dynamic scene, where the objects may move between the images. This problem has on the other hand not yet been solved in general. We have relaxed the problem in such way that the number of the moving objects is limited to 2 and the images depict several static configurations of the scene. The objects may move between the different configurations.

The motivation behind this task is that a hypothesis exists, that the presence of a background in a scene may improve the quality of a 3D reconstruction of an object. In this case, one of the objects is the background and the other one is the object which we want to reconstruct.

The task described in this section would be able to reconstruct the object together with the background and to separate the background from the object. Our intention was to perform a couple of experiments, which would prove the

functionality of the method and show for which objects the method improves the results of the reconstruction.

1.2 Thesis structure

In Chapter 1 the problem is formulated. In Chapter 2 some important concepts, such as representation of rotation in 3D or representation and registering of the cameras, are introduced. In Chapter 3 the State of the Art in the Structure from Motion and in motion segmentation is reviewed. A solution to the problem, which is introduced in Section 1.3, is proposed In Chapter 4. This solution is evaluated on real data in Chapter 5. In the same chapter we compare the models reconstructed by our method with the models reconstructed by the State of the Art single body pipeline.

1.3 Problem formulation

The scene consists of a single object on a static background, and it has been captured by multiple image sequences (takes). Each of the takes captures the scene in a static configuration and the object moves between the takes. The task is to reconstruct the object as well as the background and to segment the points belonging to the object and to the background.

take a sequence of images of a static configuration of the scene

k number of takes

$(A_{i,j}, \vec{b}_{i,j})$ motion of the object between the configurations captured by the takes i, j

X_B^i a point from the background in the configuration captured by take i , the background is static, therefore $\forall i \forall j : X_B^i = X_B^j$

X_O^i a point from the object in the configuration captured by take i , the object moves according to $(A_{i,j}, \vec{b}_{i,j})$, so $\forall j : X_O^j = A_{i,j} X_O^i + \vec{b}_{i,j}$

$\{I_i\}_{i=1}^n$ a sequence of images of the scene

t_i a take to which the image I_i belongs

Chapter 2

Frequently used concepts

2.1 Representation of rotation in 3D

Different representations of rotation in space exist, each of which has its advantages and disadvantages. Transformations between the representations can be found in [5].

Rotation matrix. As the rotation of vectors in \mathbb{R}^3 is a linear transformation, it can be represented by a matrix $R \in \mathbb{M}^{3,3}$. R is a full rank orthogonal matrix with determinant equal to 1. Representation by the matrix is suitable for operations with the rotations, as the vector \vec{x} can be rotated easily performing $R\vec{x}$. Compound rotation can as well be easily obtained by multiplication of the rotation matrices. This representation is however not suitable for clustering, storing and viewing, as 9 numbers are necessary to determine the matrix.

Euler vector. (or angle-axis representation) is a vector $\vec{e} \in \mathbb{R}^3$. The direction of the vector is the axis of the rotation and the norm of the vector is the angle of the rotation in radians. This is the minimal representation of the rotation and is therefore suitable for clustering and storing of the rotations, as well as for the viewing of the vectors because they belong to the 3D space. It is however not suitable for operations with the rotations as they are more complicated with the Euler vectors than with the rotation matrices.

■ 2.1.1 Norm of the rotation

The norm of the rotation is the angle of the rotation. It can, therefore, be computed as the norm $\|\vec{e}\|$ of the Euler vector which represents the rotation.

■ 2.1.2 Distance between two rotations

There are different ways to compute the distance between two rotations R_1 , R_2 . We compute it as the norm of the rotation $R_3 = R_2^{-1}R_1$. In order to do that the rotation R_3 has to be converted to the Euler vector form.

■ 2.2 Camera representation

A pinhole camera model [1] is a model of a camera consisting of the camera centre \vec{c} and the projective plane π onto which the points are projected.

f focal length of the camera; the distance between the centre \vec{c} and the projective plane π

$(\vec{\sigma}, \alpha)$ the image coordinate system; $\vec{\sigma}$ is the origin of the projective plane (therefore $\vec{\sigma} \in \pi$), α is the basis of the plane, for standard cameras α is orthonormal but it is not a condition

(\vec{c}, β) the camera coordinate system; the centre of the system is the centre of the camera \vec{c} , the first two vectors from the basis β are identical with the basis α and the last one is the vector $\vec{\sigma} - \vec{c}$.

(\vec{O}, δ) the world coordinate system; \vec{O} is a

X_δ a 3D-point in the world coordinate system

X_δ the point X_δ in the camera coordinate system

p a projection ray; a line connecting the point X with the camera centre C

\vec{x} a 2D point in the image coordinate system which is the projection of the point X_δ

The projection of a point X_δ in the world coordinate system is realized in two steps. At first, the point is transformed to the camera coordinate system, then the projection is found as an intersection of the ray p with the projective plane π . The projection is then represented in the image coordinate system.

Transformation of the point X_δ in the world coordinate system to the point X_β in the camera coordinate system is an affine transformation $\mathbb{R}^3 \rightarrow \mathbb{R}^3$, it can, therefore, be performed by multiplication of a matrix $P \in \mathbb{M}^{3,4}$ by homogeneous coordinates of the point X . The matrix P is the camera projection matrix and it can represent the camera.

$$X_\beta = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = P \begin{bmatrix} X_\delta \\ 1 \end{bmatrix} \quad (2.1)$$

The form of the basis of the camera coordinate system β ensures that the 3D-points on the projective plane π have the third coordinate equal to 1 and the first two coordinates are equal to the representation of the point in the image system $(\vec{\sigma}, \alpha)$, so the coordinates are the homogeneous coordinates of the representation of the 2D-point in the image system.

The intersection of the ray p with the plane π can be obtained by multiplication of the point X_β by a scalar σ such, that the point σX_β is on the projective plane π . Because of the property of the system $(\vec{\sigma}, \beta)$ is the scalar σ equal to z :

$$\vec{x} = \begin{bmatrix} \frac{x_1}{x_3} \\ \frac{x_2}{x_3} \\ 1 \end{bmatrix} \quad (2.2)$$

The transformation of the points from the coordinate system can be decomposed into the transformation to the coordinate system with the orthonormal basis whose last vector is orthogonal to the projective plane and with the center in the point \vec{c} , the camera matrix \vec{c} can, therefore, be decomposed into:

$$P = KR \left[I \mid -\vec{c} \right] \quad (2.3)$$

, where K is the camera calibration matrix [1], R is the rotation of the camera and \vec{c} is the centre of the camera.

2.3 PnP

PnP is a procedure which registers a camera onto a set of 3D points with n correspondences between the 3D points and the 2D features on the image [1].

The simplest type of the PnP algorithm is P6P which requires 6 points to compute the matrix P . The camera matrix P is stacked into a vector \vec{P} which is used as the vector of unknowns. Each correspondence between a point and a feature generates two linear equations which are derived from equations (2.1), (2.2), one equation for each element of the feature \vec{x} .

Because the matrix P has 12 elements, at least 12 linearly independent equations are required to find it. These equations can be generated from 6 correspondences. More advanced versions of PnP exist which require fewer correspondences. If the camera is calibrated, only 3 correspondences are required to compute the pose.

Chapter 3

State of the art

3.1 Structure from Motion

Reconstruction of static scenes is a well-understood problem. Two main approaches to it exist, namely incremental SfM and global SfM. There are different SfM pipelines, such as Bundler [2], COLMAP [3] or OpenMVG[4]. Among these pipelines, we have chosen COLMAP as the base of our method. The most common approach is the incremental SfM, which consists of the following steps: feature extraction and matching, the initial pose calculation, camera registration, point triangulation, and bundle adjustment.

3.1.1 Feature extraction and matching

At first, features, such as SIFT [6], are found in each of the images. Corresponding features in a pair of images are matched using nearest neighbor search on the descriptors of the features. In the exhaustive matching, the matches are found between all pairs of images, therefore the matching has quadratic complexity. Other variants of matching exist, which have linear complexity and therefore are better suited for large datasets. Transitive matches (tracks) may be used in order to increase the number of the tracks. Maset, Arrigoni and Fussielo [7] use eigendecomposition to maximize consistency of the tracks.

improves only the limited number of cameras, is performed, so the time complexity of the reconstruction is better.

■ 3.2 Motion segmentation

In the motion segmentation problem, motion tracks between two or more images are given, and the task is to split the tracks into groups that move together. The problem is strongly connected to the Multi-body Structure from Motion. If all motion groups were correctly determined, the MBSfM problem would be reduced to the static scene reconstruction. On the other hand, if the transformations of the cameras towards each of the motion groups were known, it would be easy to determine which point belongs to which group. The problem is that the segmentation and the models have to be found simultaneously. Different approaches to the problem exist.

■ 3.2.1 Algebraic methods

One group of the algebraic methods is based on matrix factorization. These methods describe the motion segmentation as a subspace clustering problem. Among these methods belong for example [10] and [11], which perform the factorization via SVD.

R. Vidal et. al. [12] have proposed The Multibody Epipolar Geometry, where the multibody fundamental matrix is a sequence of fundamental matrices. A point which is inlier to this model fulfills the epipolar constraint for at least one of the fundamental matrices. The Multibody epipolar constraint leads to a set of polynomial equations, which is solved using GPCA. This solution however cannot be applied if erroneous matches occur. Rao et. al. [13] have proposed an algebraic method which is able to deal with the erroneous matches. The method is based on the merging of the points such that the coding length of the model is minimized using Agglomerative Lossy Compression. Rao et. al. have proposed a Hybrid Perspective Constraint [14], which can describe both planar and perspective motions.

Sometimes it is difficult to perform the segmentation, as degenerate situations exist where two independent motions can be explained with a single fundamental matrix.

■ 3.2.3 Other methods

Another motion segmentation methods exist, which belong neither to the algebraic group nor to the random sampling one. These methods are based for example on the optical flow [24] or on an iterative approach inspired by EM [25]. After the rise of the neural networks, semantic segmentation methods have also been developed [26].

Growing homographies is an approach similar to the random sampling. From each match, an affine transformation is computed, which is subsequently refined to the form of a homography using inliers to the transformation. F. Šrajcar [27] uses the Planar Homology concept to merge the homographies which have arisen from the same motion.

■ 3.3 Reconstruction of dynamic scenes

The most of the motion segmentation methods concentrate on the segmentation between 2 images. K Schindler, D Suter and H Wang [28] have proposed a method for n-view MBSfM. They draw random samples for each pair of subsequent images and after filtering of the worst models they link the models to create chains, which describe the motion in the whole sequence. The chains are selected with the maximum likelihood criterion. The method is more suitable for videos than for general image sets. Method [13] works with a matrix of tracks, so it is able to work with n views.

Little has been done in the actual reconstruction of the dynamic scenes (MBSfM), so the problem still remains unsolved in general. YASFm [27] is a pipeline, which performs a motion segmentation and is, therefore, able to reconstruct multiple objects from the dynamic scenes. However, it is shown in [29], that YASFm is not able to reconstruct the scenes described in Section 1.3.

Chapter 4

Proposed solution

4.1 Reconstruction of the takes

At first, the matches are found between each pair of images. The option `multiple_objects` is selected, so correct matches between points from both the object and the background are retained during the verification step. Each take is reconstructed individually in a standard Structure from Motion pipeline COLMAP [3]. The take from which the reconstruction has been obtained is the **anchor take** of the reconstruction. After a sparse point cloud is obtained from the anchor take, cameras from other takes are registered onto the point cloud via sequential PnP.

Sequential PnP is a procedure, where the camera is registered onto the points using a RANSAC PnP, and the outliers to the found camera pose are reused for further registration. This procedure is repeated until no pose with satisfactory support is found. Therefore it is possible that multiple camera poses are registered from a single camera.

$\text{RANSAC_PnP}(\mathcal{M}, r)$ is a procedure that performs registration of a camera using a set of 2D to 3D correspondences \mathcal{M} , the number of iterations of the RANSAC algorithm is r , the output of the procedure is the registered pose P . $\text{Find_Inliers}(P, \mathcal{M})$ outputs a subset $\mathcal{I} \subseteq \mathcal{M}$ of inliers to the pose P .

Algorithm 1: Sequential PnP

input : \mathcal{M} a finite set of correspondences between 2D features and 3D points
 m minimal number of inliers for a pose
 r number of the iterations of RANSAC

output : \mathcal{P} a finite set of camera matrices

$\mathcal{P} \leftarrow \emptyset$;

while $|\mathcal{M}| \geq m$ **do**

$P = \text{RANSAC_PnP}(\mathcal{M}, r)$;

$\mathcal{I} \leftarrow \text{Find_Inliers}(P, \mathcal{M})$;

if $|\mathcal{I}| \geq m$ **then**

$\mathcal{P} \leftarrow \mathcal{P} \cup \{P\}$;

$\mathcal{M} \leftarrow \mathcal{M} \setminus \mathcal{I}$;

else

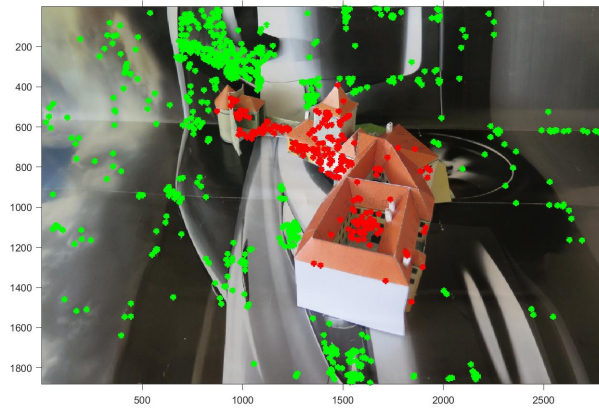
break;

end

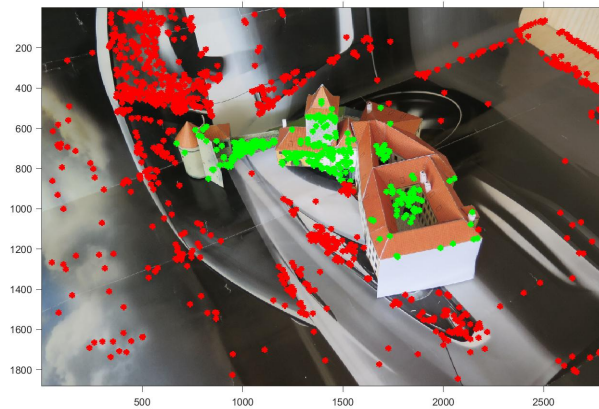
end

The idea behind this is that cameras from takes different from the anchor take have different poses towards the static background and towards the object, therefore in the ideal case the sequential PnP should find one pose towards the background and another one towards the object. The points belonging to the object as well as those belonging to the background can be recognized as the inliers to these poses.

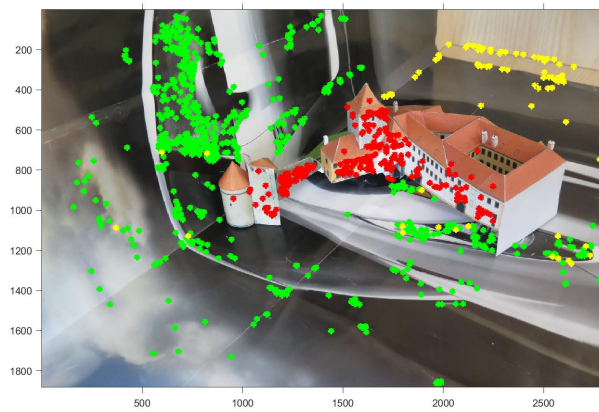
The motion of the object between the anchor take and the take to which the image belongs (the second take) can also be computed from the poses of the camera towards the background and the object which have been found using the sequential PnP.



(a)



(b)



(c)

Figure 4.1: Example of different results of the sequential PnP, points with the same color are observed by the same camera. The cameras have been registered in the order green, red, yellow.

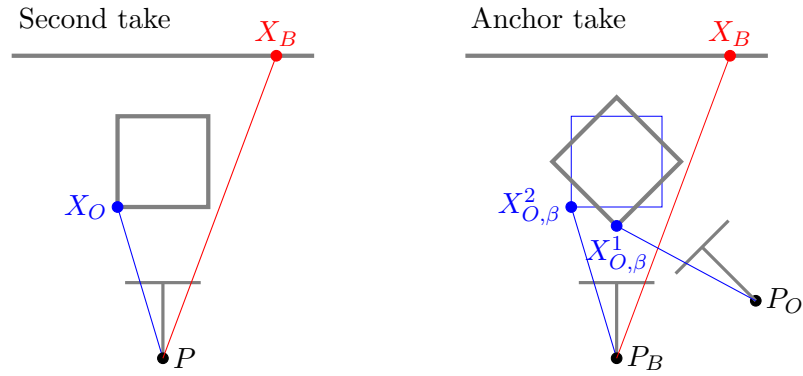


Figure 4.2: Example of a camera after the sequential PnP registration onto the anchor take.

4.2 Object motion calculation

X_O - arbitrary 3D point from the object

P - a camera from the second take which observes the point X_O

\vec{x} - projection of the point X_O on the camera P

$X_{O,\beta}^1$ - reconstruction of the point P in the anchor take reconstruction

β - basis of the anchor take reconstruction

A - matrix of rotation of the object between the anchor take and the second take

\vec{b} - translation of the object between the anchor take and the second take

$X_{O,\beta}^2$ - new position of the point $X_{O,\beta}^1$ after movement of the object from the anchor take configuration to the second take configuration; for every point it holds true: $X_{O,\beta}^2 = AX_{O,\beta}^1 + \vec{b}$

P_B - camera matrix of the camera P registered on the background points of the anchor take reconstruction

P_O - camera matrix of the camera P registered on the object points of the anchor take reconstruction

If the camera P from the second take is registered on the object points from the reconstruction of the anchor take, it projects the object points from the anchor take onto the same 2D points onto which they were projected in the configuration of the second take. Particularly this means that \vec{x} is a projection of the point $X_{O,\beta}^1$ onto the camera P_O

If the camera P from the second take is registered on the background points from the reconstruction of the anchor take, its position is the same as it was in the second take. This means that if the object points were transformed to

the position in which they were in the second take, they would be projected onto the same 2D points onto which they were projected in the second take. This implies that \vec{x} is a projection of the point $X_{O,\beta}^2$ onto the camera P_B .

The following equation must therefore hold true:

$$\sigma_1 P_O \begin{bmatrix} X_{O,\beta}^1 \\ 1 \end{bmatrix} = \sigma_2 P_B \begin{bmatrix} X_{O,\beta}^2 \\ 1 \end{bmatrix} \quad (4.1)$$

We can introduce new symbols for rotation, center and camera calibration matrix for cameras P_O and P_B and write this equation as

$$\sigma_1 K \begin{bmatrix} R_O & -R_O \vec{c}_O \end{bmatrix} \begin{bmatrix} X_{O,\beta}^1 \\ 1 \end{bmatrix} = \sigma_2 K \begin{bmatrix} R_B & -R_B \vec{c}_B \end{bmatrix} \begin{bmatrix} AX_{O,\beta}^1 + \vec{b} \\ 1 \end{bmatrix} \quad (4.2)$$

After elimination of K and rewriting the equation into the polynomial form

$$\sigma_1 R_O X_{O,\beta}^1 - \sigma_1 R_O \vec{c}_O = \sigma_2 R_B A X_{O,\beta}^1 + \sigma_2 R_B \vec{b} - \sigma_2 R_B \vec{c}_B \quad (4.3)$$

Because the equation has to hold true for all $X_{O,\beta}^1$, it holds true also for $X_{O,\beta}^1 = 0$, which implies that the following two equations are also valid:

$$-\sigma_1 R_O \vec{c}_O = \sigma_2 R_B \vec{b} - \sigma_2 R_B \vec{c}_B \quad (4.4)$$

$$\sigma_1 R_O = \sigma_2 R_B A \quad (4.5)$$

We can easily find the relative rotation A from the latter equation as

$$A = \frac{\sigma_1}{\sigma_2} R_B^{-1} R_O$$

Because A , R_B^{-1} and R_O are all rotation matrices, $\sigma_1 = \sigma_2$ and therefore we can write

$$A = R_B^{-1} R_O \quad (4.6)$$

$$-R_O \vec{c}_O = R_B \vec{b} - R_B \vec{c}_B \quad (4.7)$$

We can find translation \vec{b} as

$$\vec{b} = \vec{c}_B - R_B^{-1} R_O \vec{c}_O = \vec{c}_B - A \vec{c}_O \quad (4.8)$$

In Figure 4.1 it is shown that in the sequential PnP order of the objects towards whose points the poses are calculated is arbitrary. That means that apart from the true motions, the inverted ones (motion of the background if the object is considered static) can be calculated. Also, we can see in Figure

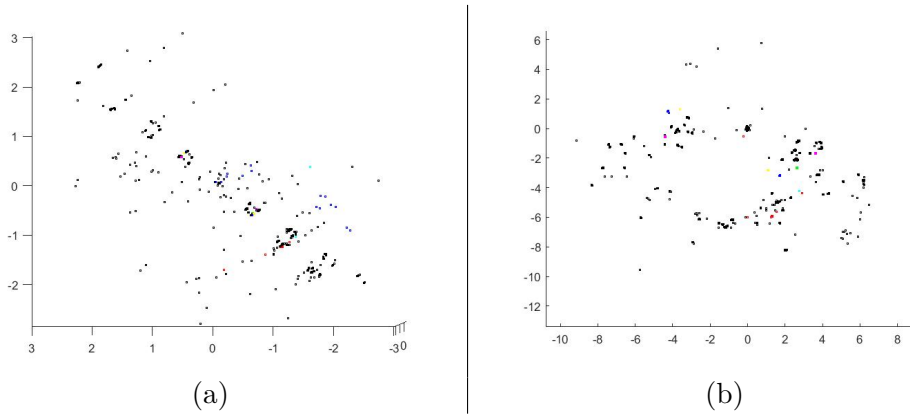


Figure 4.3: Euler vectors (a) and translations (b) calculated from a dataset consisting of 8 takes. The situation is before the transformation into the same coordinate system. Clusters of motions are observable.

4.3.1 Description of the basis transformation

β' - origin coordinate system

β - target coordinate system

$X'_{\beta'}$ - arbitrary 3D point in the origin coordinate system

X_{β} - point $X'_{\beta'}$ in the target coordinate system

$P_{\beta'}$ - camera matrix in the origin coordinate system

P_{β} - camera matrix in the target coordinate system

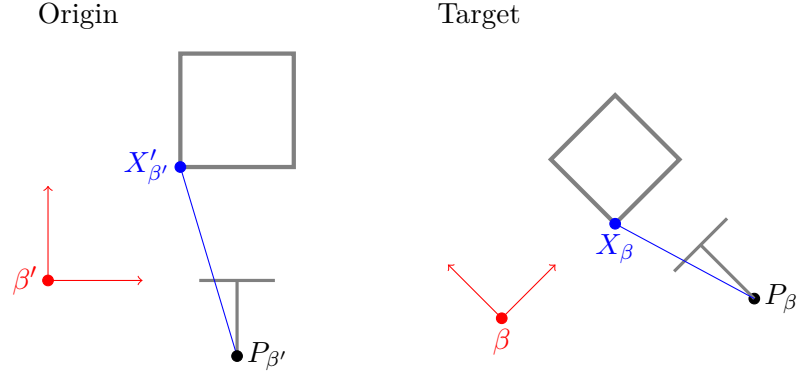
B - change of basis matrix from β' to β , because both bases are orthonormal, it has form of $B = \sigma R$ where σ is a scalar and R is a rotation matrix.

\vec{o}'_{β} - origin of the β' coordinate system in the β coordinate system

Then every point in the origin coordinate system can be transformed to the target system as follows:

$$X_{\beta} = BX'_{\beta'} + \vec{o}'_{\beta} = \sigma RX'_{\beta'} + \vec{o}'_{\beta} \quad (4.9)$$

\vec{x} - projection of the 3D point onto the camera



The 3D point is projected onto the same point on the camera in both reconstructions, so the following holds true:

$$\vec{x} = \sigma_1 P_\beta \begin{bmatrix} X_\beta \\ 1 \end{bmatrix} = \sigma_2 P_{\beta'} \begin{bmatrix} X'_{\beta'} \\ 1 \end{bmatrix} \quad (4.10)$$

For some scalars σ_1, σ_2 . We assume that both pictures have been taken with the same camera K . We introduce $\alpha = \frac{\sigma_1}{\sigma_2}$, R_β, \vec{c}_β as rotation and center of matrix P_β and $R_{\beta'}, \vec{c}_{\beta'}$ as rotation and center of matrix $P_{\beta'}$ and write

$$\alpha K \begin{bmatrix} R_\beta & -R_\beta \vec{c}_\beta \end{bmatrix} \begin{bmatrix} \sigma R X'_{\beta'} + \vec{\sigma}'_\beta \\ 1 \end{bmatrix} = K \begin{bmatrix} R_{\beta'} & -R_{\beta'} \vec{c}_{\beta'} \end{bmatrix} \begin{bmatrix} X'_{\beta'} \\ 1 \end{bmatrix} \quad (4.11)$$

$$\alpha (R_\beta (\sigma R X'_{\beta'} + \vec{\sigma}'_\beta) - R_\beta \vec{c}_\beta) = R_{\beta'} X'_{\beta'} - R_{\beta'} \vec{c}_{\beta'} \quad (4.12)$$

$$\alpha \sigma R_\beta R X'_{\beta'} + \alpha R_\beta \vec{\sigma}'_\beta - \alpha R_\beta \vec{c}_\beta = R_{\beta'} X'_{\beta'} - R_{\beta'} \vec{c}_{\beta'} \quad (4.13)$$

This equation has to hold true for every $X'_{\beta'}$, so following two equations must be valid:

$$\alpha \sigma R_\beta R = R_{\beta'} \quad (4.14)$$

$$\alpha R_\beta \vec{\sigma}'_\beta - \alpha R_\beta \vec{c}_\beta = -R_{\beta'} \vec{c}_{\beta'} \quad (4.15)$$

Because R_β, R and $R_{\beta'}$ are all rotation matrices, in order for the first equation to be valid, $\alpha \sigma$ must be equal to 1, so

$$\alpha = \frac{1}{\sigma} \quad (4.16)$$

We can therefore write:

$$R = R_\beta^{-1} R_{\beta'} \quad (4.17)$$

Rotation between coordinate systems can be found with one camera pair.

The latter equation can be multiplied by σ :

$$R_\beta \vec{\sigma}'_\beta - R_\beta \vec{c}_\beta = -\sigma R_{\beta'} \vec{c}_{\beta'} \quad (4.18)$$

We introduce new symbols $\vec{t}_\beta = -R_\beta \vec{c}_\beta$ and $\vec{t}'_{\beta'} = -R_{\beta'} \vec{c}_{\beta'}$ and write:

$$R_\beta \vec{o}'_\beta + \vec{t}_\beta = \sigma \vec{t}'_{\beta'} \quad (4.19)$$

$$R_\beta \vec{o}'_\beta - \sigma \vec{t}'_{\beta'} = \vec{t}_\beta \quad (4.20)$$

This equation can be rewritten in a element-wise form:

$$\begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix} \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} - \sigma \begin{bmatrix} t'_1 \\ t'_2 \\ t'_3 \end{bmatrix} = - \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} \quad (4.21)$$

If we know more such camera pairs, this can be written as a following system of linear equations:

$$\begin{bmatrix} R_{11} & R_{12} & R_{13} & -t'_1 \\ R_{21} & R_{22} & R_{23} & -t'_2 \\ R_{31} & R_{32} & R_{33} & -t'_3 \\ & \dots & & \end{bmatrix} \begin{bmatrix} o_1 \\ o_2 \\ o_3 \\ \sigma \end{bmatrix} = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \dots \end{bmatrix} \quad (4.22)$$

Translation and scale between the coordinate systems can be found with at least two camera pairs.

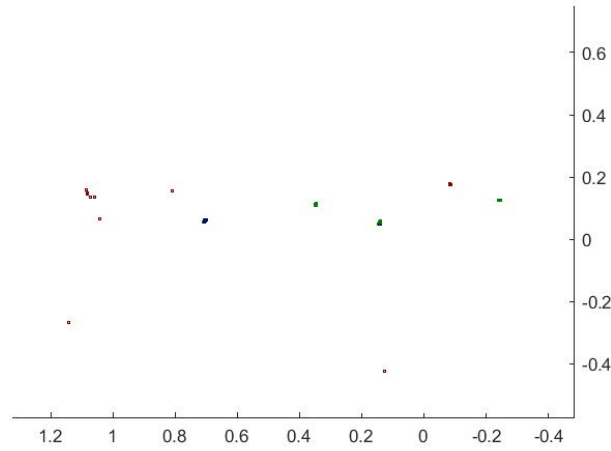


Figure 4.4: Euler vectors of the rotations between coordinate systems of three reconstructions. The clusters of the vectors are observable.

finds inliers to the transformation $(R_{a,b}, \vec{\sigma}_{a,b}, \sigma_{a,b})$ from the set \mathcal{S} according to the procedure described in this section.

Algorithm 2: Cluster bases

```

input :  $\mathcal{S}$  a finite set of camera pairs
          $a$  the origin take
          $b$  the target take
          $r$  number of the iterations of RANSAC
output :  $\mathcal{C}$  a finite set of finite sets (clusters) of camera pairs
 $\mathcal{S} \leftarrow \text{Select\_Pairs}(\mathcal{M}, a, b)$ ;
 $\mathcal{C} \leftarrow \emptyset$ ;
while  $|\mathcal{S}| > 1$  do
     $\mathcal{I}_{best} = \emptyset$ ;
    for  $i \leftarrow 1$  to  $r$  do
         $(P_{a,1}, P_{b,1}) \leftarrow \text{Random\_Pair}(\mathcal{S})$ ;
         $(P_{a,2}, P_{b,2}) \leftarrow \text{Random\_Pair}(\mathcal{S})$ ;
         $R_{a,b} \leftarrow \text{Find\_Rotation}((P_{a,1}, P_{b,1}))$ ;
         $(\vec{\sigma}_{a,b}, \sigma_{a,b}) \leftarrow \text{Find\_TS}((P_{a,1}, P_{b,1}), (P_{a,2}, P_{b,2}), R_{a,b,1})$ ;
         $\mathcal{I} \leftarrow \text{Find\_Inliers\_Basis}(\mathcal{S}, R_{a,b}, \vec{\sigma}_{a,b}, \sigma_{a,b})$ ;
        if  $|\mathcal{I}| > |\mathcal{I}_{best}|$  then
             $\mathcal{I}_{best} \leftarrow \mathcal{I}$ ;
        end
    end
     $\mathcal{C} = \mathcal{C} \cup \{\mathcal{I}_{best}\}$ ;
    if  $\mathcal{I}_{best} = \emptyset$  then
        break;
    end
end

```

■ Inlier recognition

In order to use the RANSAC algorithm, a procedure for inlier recognition has to be chosen. The randomly selected hypothesis consists of two camera pairs $(P_{a,1}, P_{b,1}), (P_{a,2}, P_{b,2})$. Each camera can be represented according to Section 2.2 as $P_{t,i} = K_{t,i} * [R_{t,i} - R_{t,i} \vec{c}_{t,i}]$. We want to check whether a pair $(P_{a,3}, P_{b,3})$ is consistent with this hypothesis. In order to be so, both rotation and translation has to be consistent.

We can see from equation (4.17), that rotation can be checked with one camera pair. We therefore calculate $R_{a,b,3} = R_{b,3}^{-1} R_{a,3}$ and then we find

distance D between $R_{a,b,3}$ and $R_{a,b,1} = R_{b,1}^{-1}R_{a,1}$ using a procedure in Section 2.1.2. This distance is in radians and therefore scale-invariant and thus the rotational consistency can be checked by comparison of the distance D with a fixed value, which can be a parameter of the pipeline. If the rotation is inconsistent, the camera pair is rejected. Otherwise, we follow to checking of the translational consistency.

According to the equation (4.22) we compute the origin $\vec{\sigma}_{a,b}$ and scaling $\sigma_{a,b}$ from camera pairs $(P_{a,1}, P_{b,1})$ and $(P_{a,2}, P_{b,2})$. If the transformation is consistent with the pair $(P_{a,2}, P_{b,2})$, it will transform centre $\vec{c}_{a,3}$ of the camera $P_{a,3}$ to the centre $\vec{c}_{b,3}$ of the camera $P_{b,3}$. So we transform the point $\vec{c}_{a,3}$ according to equation (4.9) as $\vec{c}'_{a,3} = \sigma_{a,b}R_1\vec{c}_{a,3} + \vec{\sigma}_{a,b}$ and we compare it with $\vec{c}_{b,3}$. Translation, however, is not scale invariant, so we use the apical angle concept. We compute the center of mass \vec{m} of the points from the reconstruction of the take b which are visible by the camera $P_{b,3}$ and then we find the angle between vectors $\vec{c}'_{a,3} - \vec{m}$ and $\vec{c}_{b,3} - \vec{m}$. If the angle is smaller than a fixed value (which can again be a parameter of the pipeline), the pair $(P_{a,3}, P_{b,3})$ is accepted as an inlier, otherwise it is rejected.

Find_Centre(P) returns the centre of the camera P , Angle(\vec{x}, \vec{y}) returns the angle between the vectors \vec{x}, \vec{y} . Point_Of_Mass(P) returns the point of mass of the points observed by the camera P .

Algorithm 3: Find inliers basis

input : \mathcal{S} a finite set of camera pairs
 $(R, \vec{\sigma}, \sigma)$ the hypothetical basis transformation
 t_{rot} rotational threshold
 t_{trans} translational threshold
output : \mathcal{I} a finite set of inliers $\mathcal{I} \subseteq \mathcal{S}$
 $\mathcal{I} \leftarrow \emptyset$;
foreach $(P_{a,3}, P_{b,3}) \in \mathcal{S}$ **do**
 $R_3 \leftarrow \text{Find_Rotation}((P_{a,3}, P_{b,3}))$;
 if $\text{Dist}(R, R_3) \leq t_{rot}$ **then**
 $\vec{c}_{a,3} \leftarrow \text{Find_Centre}(P_{a,3})$;
 $\vec{c}'_{a,3} \leftarrow \sigma R \vec{c}_{a,3} + \vec{\sigma}$;
 $\vec{c}_{b,3} \leftarrow \text{Find_Centre}(P_{b,3})$;
 $\vec{m} \leftarrow \text{Point_Of_Mass}(P_{b,3})$;
 if $\text{Angle}(\vec{c}'_{a,3} - \vec{m}, \vec{c}_{b,3} - \vec{m}) \leq t_{trans}$ **then**
 $\mathcal{I} \leftarrow \mathcal{I} \cup \{(P_{a,3}, P_{b,3})\}$
 end
 end
end

4.3.3 Verification of the basis transformations

We have found clusters of the camera pairs and their corresponding transformations. But we still need to group the transformations between different coordinate systems, as we want to transform all motions to the same coordinate system and we need that the transformations are either all towards the background or all towards the object. In order to do so, we use the verification via cycles of the transformations of length three.

We have three takes a, b, c together with their basis transformations and the corresponding clusters of camera pairs. If a vector in the first coordinate system is transformed to the second coordinate system, then to the third coordinate system, and then back to the first coordinate system, the resulting vector is equal to the original vector. This forms a constraint for coordinate change matrices which change the coordinates in a cycle.

$R_{a,b}$ - rotation between bases of takes a, b
 $\sigma_{a,b}$ - scale between first and second basis
 $\vec{\sigma}_{a,b}$ - origin of the first basis in the second coordinate system
 $R_{b,c}$ - rotation between second and third basis
 $\sigma_{b,c}$ - scale between second and third basis
 $\vec{\sigma}_{b,c}$ - origin of the second basis in the third coordinate system
 $R_{c,a}$ - rotation between third and first basis
 $\sigma_{c,a}$ - scale between third and first basis
 $\vec{\sigma}_{c,a}$ - origin of the third basis in the first coordinate system

x_a - arbitrary point in the first coordinate system
 x_b - point x_a in the second coordinate system
 x_c - point x_a in the third coordinate system

Transformations between the points:

$$x_b = \sigma_{a,b} R_{a,b} x_a + \vec{\sigma}_{a,b} \quad (4.23)$$

$$x_c = \sigma_{b,c} R_{b,c} x_b + \vec{\sigma}_{b,c} \quad (4.24)$$

$$\begin{aligned}
 x_a &= \sigma_{c,a} R_{c,a} x_c + \vec{\sigma}_{c,a} = \sigma_{c,a} R_{c,a} (\sigma_{b,c} R_{b,c} x_b + \vec{\sigma}_{b,c}) + \vec{\sigma}_{c,a} \\
 x_a &= \sigma_{c,a} R_{c,a} (\sigma_{b,c} R_{b,c} (\sigma_{a,b} R_{a,b} x_a + \vec{\sigma}_{a,b}) + \vec{\sigma}_{b,c}) + \vec{\sigma}_{c,a} \\
 x_a &= \sigma_{c,a} R_{c,a} \sigma_{b,c} R_{b,c} (\sigma_{a,b} R_{a,b} x_a + \vec{\sigma}_{a,b}) + \sigma_{c,a} R_{c,a} \vec{\sigma}_{b,c} + \vec{\sigma}_{c,a}
 \end{aligned}$$

4.3.4 Transformation of the motions

We have verified the clusters via cycle consistency and divided them into connected components. For each connected component we create another graph but this time the vertices represent the takes and there is an edge between the vertices a, b if a cluster exists in the connected component which relates bases of reconstructions of takes a, b . The edge has a weight of $\frac{1}{R_{cy}}$ where R_{cy} is the rotational consistency of the most consistent cycle among those which contain a cluster relating takes a, b . Such cluster and its corresponding basis transformation is assigned to the edge.

To obtain transformations between all takes, it is sufficient to have a spanning tree of the second graph and transformations between takes which are not connected can be computed transitively. In this case, the minimum spanning tree of the second graph is used.

We have selected the central take as the one with the highest degree in the spanning tree of the second graph and we want to transform all computed motions to the coordinate system of the central take. In order to do so, we need to know how the computed rotation and translation of the object behaves under the change of the basis.

β' - origin basis

β - target basis

B - change of basis matrix from β' to β , because both bases are orthonormal, it has form of $B = \sigma R$ where σ is a scalar and R is a rotation matrix.

\vec{o}_{β} - origin of the β' coordinate system in the β coordinate system

$X_{\beta'}^1$ - arbitrary 3D point from the object in the origin coordinate system

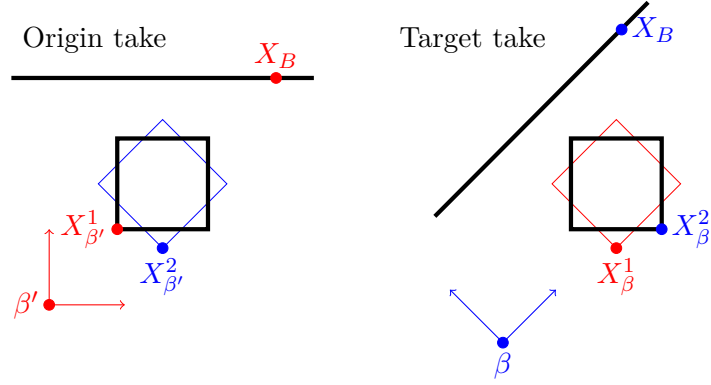
$X_{\beta'}^2$ - new position of the point $X_{\beta'}^1$ after moving the object to a position in the second take; point is in the origin coordinate system

$A_{\beta'}, \vec{b}_{\beta'}$ - rotation and translation of the object in the origin coordinate system

X_{β}^1 - point $X_{\beta'}^1$ in the target coordinate system

X_{β}^2 - point $X_{\beta'}^2$ in the target coordinate system

$A_{\beta}, \vec{b}_{\beta}$ - rotation and translation of the object in the target coordinate system



It holds true:

$$X_{\beta'}^2 = A_{\beta'} X_{\beta'}^1 + \vec{b}_{\beta'} \quad (4.26)$$

$$X_{\beta}^2 = A_{\beta} X_{\beta}^1 + \vec{b}_{\beta} \quad (4.27)$$

$$X_{\beta}^1 = \sigma R X_{\beta'}^1 + \vec{\sigma}'_{\beta} \quad (4.28)$$

$$X_{\beta}^2 = \sigma R X_{\beta'}^2 + \vec{\sigma}'_{\beta} \quad (4.29)$$

We can see that both second and fourth equations are equal to X_{β}^2 , we can therefore write:

$$\begin{aligned} A_{\beta} X_{\beta}^1 + \vec{b}_{\beta} &= \sigma R X_{\beta'}^2 + \vec{\sigma}'_{\beta} \\ A_{\beta} (\sigma R X_{\beta'}^1 + \vec{\sigma}'_{\beta}) + \vec{b}_{\beta} &= \sigma R (A_{\beta'} X_{\beta'}^1 + \vec{b}_{\beta'}) + \vec{\sigma}'_{\beta} \\ A_{\beta} (\sigma R) X_{\beta'}^1 + A_{\beta} \vec{\sigma}'_{\beta} + \vec{b}_{\beta} &= \sigma R A_{\beta'} X_{\beta'}^1 + \sigma R \vec{b}_{\beta'} + \vec{\sigma}'_{\beta} \end{aligned}$$

The equation must be valid for every $X_{\beta'}^1$, which means:

$$A_{\beta} (\sigma R) = \sigma R A_{\beta'} \quad (4.30)$$

$$A_{\beta} \vec{\sigma}'_{\beta} + \vec{b}_{\beta} = \sigma R \vec{b}_{\beta'} + \vec{\sigma}'_{\beta} \quad (4.31)$$

We can express relative rotation in the target basis from the first equation as:

$$A_{\beta} = R A_{\beta'} R^{-1} \quad (4.32)$$

We can see that transformed relative rotation can be obtained only with relative rotation between the coordinate systems, we do not need to find the translation nor the scale between the coordinate systems.

We can express relative translation in the target basis from the second equation as:

$$\vec{b}_{\beta} = \sigma R \vec{b}_{\beta'} + \vec{\sigma}'_{\beta} - R A_{\beta'} R^{-1} \vec{\sigma}'_{\beta} \quad (4.33)$$

If there is a direct connection between the origin and the central takes in the spanning tree, we can directly transform the translation and the rotation according to the equations (4.32), (4.33) using the transformation assigned to the edge connecting the takes. If there is no such connection, we can follow the path from the origin take to the central take in the spanning tree and transform the rotation and the translation sequentially.

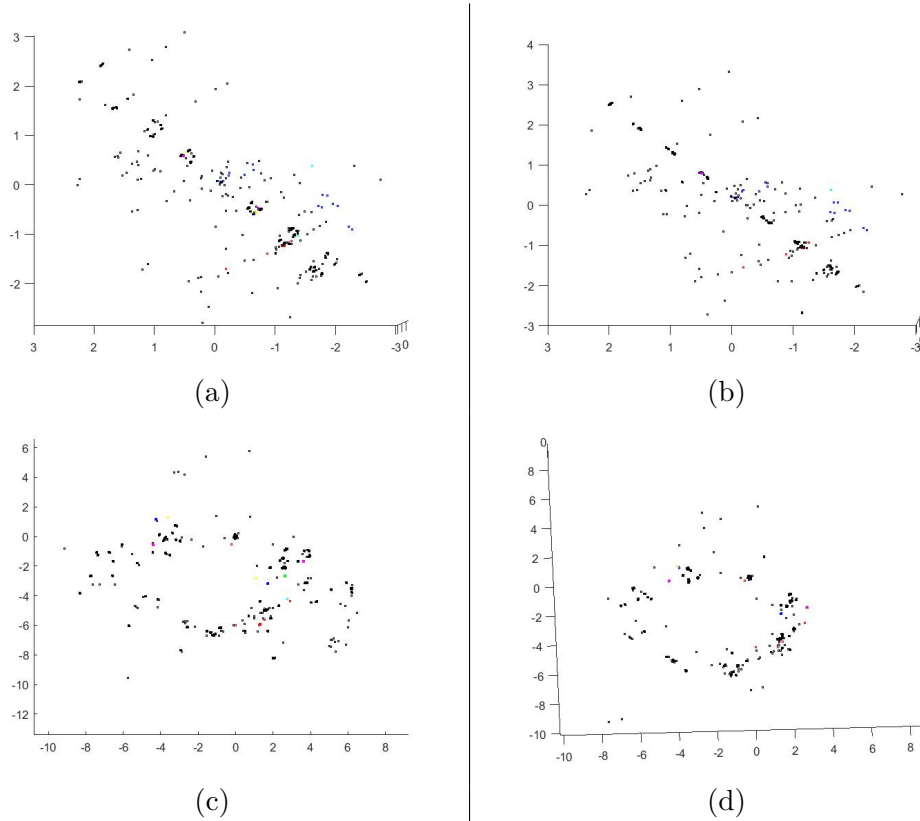


Figure 4.5: Euler vectors (a), (b) and translations (c), (d) calculated from a dataset consisting of 8 takes. Images (a), (c) depict the situation before the transformation, images (b), (d) depict the situation after the transformation. The object was always rotated around the vertical axis which is observable in (b).

4.3.5 Zero motions removing

Due to the multiple registration of a camera towards the points (Figure 4.1) of the same object, zero motions arise among the calculated motions of the object which is observable in Figures 4.7(b), 4.7(d). However these motions can make a problem in a following cyclic verification of the motions (Section 4.5.1), so we need to dispose of them. In order to do that we use spectral clustering to discover clusters of the zero motions.

We use multi-view spectral clustering [30], [23] by the kernel addition to get clusters of points which are distinguished by both rotation and translation. The clustering using the RANSAC would not be sufficient because in some cases the cluster of the zero motions appears to have too high extent. Figure 4.7(b) The rotational distance between two motions is computed according to Section 2.1.2 and the translational distance between the points is computed using the Euclidean distance between the translations. The distances D are converted into affinities $A = \frac{1}{D}$. The normalized Laplacian [30] is used in order to eliminate the scaling of the translations. The final kernel is computed as a sum of both Laplacians. Because the number of all motions is too high to perform eigendecomposition effectively, the motions are divided according to the origin and final takes of the motions. The number of clusters k is equal to the number of eigenvalues which are lower or equal to 1 while the maximum eigenvalue is 2 because of the normalization. The k-means algorithm is used to obtain the clusters from the first k eigenvectors.

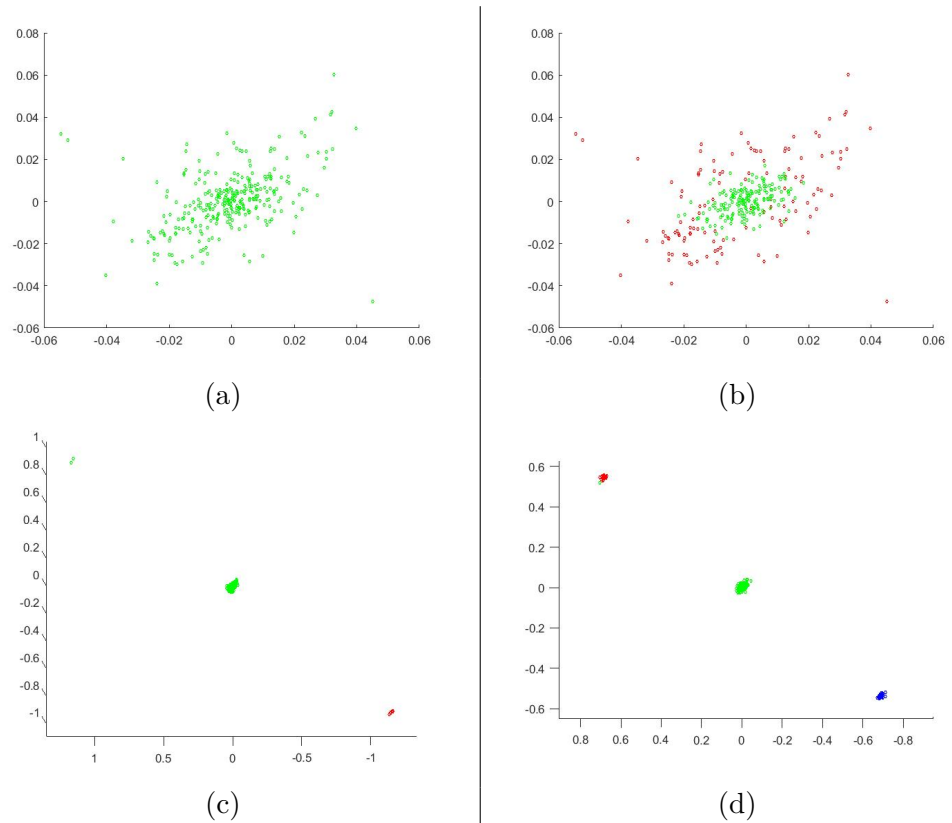


Figure 4.6: Examples of results of the spectral clustering of the clusters. Only rotations are depicted.

All clusters are then tested on the distance from zero. For rotation, the norm (Section 2.1.1) of the median Euler vector among the motions in the cluster is

used and the median translation has to be lower than a fixed fraction of the median distance from cameras from the central take to the points which they observe. If both rotation and translation are near enough to zero, the whole cluster is removed. According to Figure 4.7 this procedure is able to remove the cluster of the zero motions. If the zero cluster is split like in Figure 4.6(b), both zero clusters are removed.

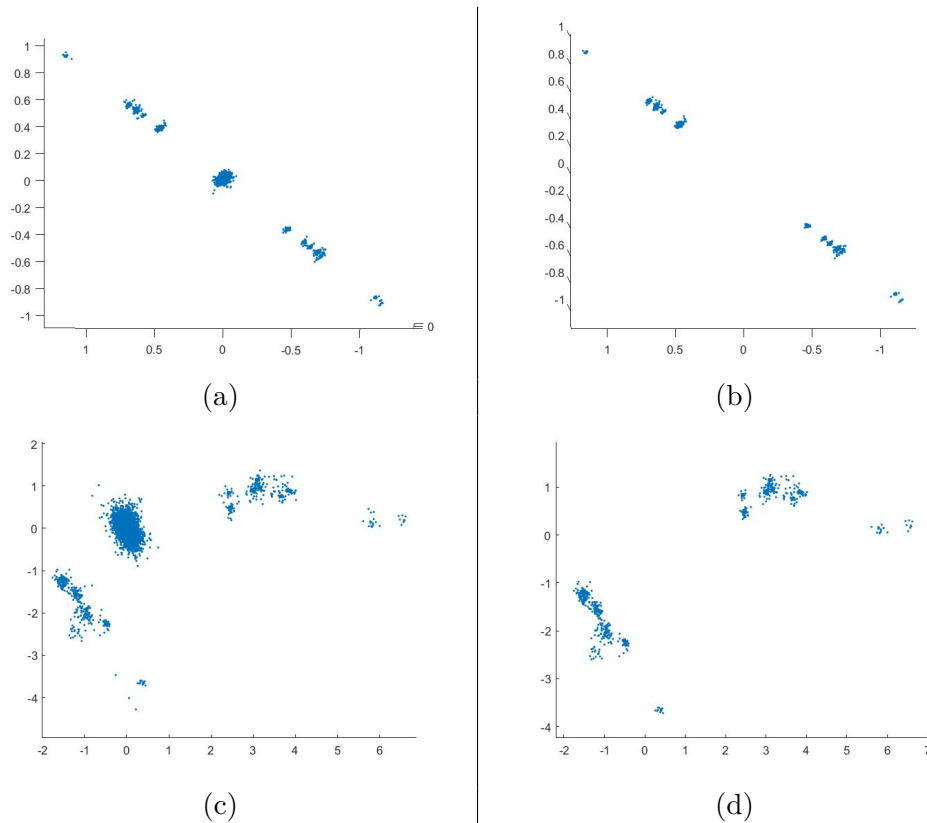


Figure 4.7: The Euler vectors (a), (b) and the translations (c), (d) of the object motion. Images (a), (c) depict the situation before the removing, images (b), (d) depict the situation after the removing.

4.4 Motion clustering

We have transformed the motions into the same coordinate system and removed the zero motions. We can, therefore, follow to clustering the motions. The main purpose of the clustering is that although the zero motions have been eliminated, there are motions, which have been computed from pairs where the first camera has been registered towards the background points and the second one has been registered towards the object points, as well as motions which have been computed from pairs where the sequential PnP

performed in the reversed order and pairs where one or two cameras have been registered wrongly. We want to distinguish these groups of cameras from each other in order to distinguish the points which the cameras observe.

Similarly to the clustering of the basis transformations, we use RANSAC but unlike that, one camera pair is sufficient to compute the motion according to equations (4.6), (4.8), so this time the hypothesis consists of one motion only. Because of a small total number of hypotheses, they do not have to be drawn randomly but all hypotheses can be checked.

The motions are divided into subsets $S_{a,b}$ according to their original and target takes $a, b, a < b$ and they are clustered separately. For each motion A_1, \vec{b}_1 in the subset $S_{a,b}$, its support is computed as the number of inliers to this motion from the set $S_{a,b}$. The motion with the highest support is selected and together with its inliers it builds a new cluster. Outliers to the motion are reused in the next iteration and the procedure is repeated until all motions are in some track.

The clusters are then represented by a single Euler vector and a single translation vector. Both are medians of the motions in the cluster.

Algorithm 4: Cluster motions

```

input :  $\mathcal{S}$  a finite set of transformed camera pairs
          $a$  the origin take
          $b$  the target take
output :  $\mathcal{C}$  a finite set of finite sets (clusters) of camera pairs
 $\mathcal{S} \leftarrow \text{Select\_Pairs}(\mathcal{M}, a, b)$ ;
 $\mathcal{C} \leftarrow \emptyset$ ;
while  $|\mathcal{S}| > 0$  do
   $\mathcal{I}_{best} = \emptyset$ ;
  foreach  $(P_{a,1}, P_{b,1}) \in \mathcal{S}$  do
     $R_{a,b} \leftarrow \text{Find\_Rotation}((P_{a,1}, P_{b,1}))$ ;
     $\vec{t}_{a,b} \leftarrow \text{Find\_Translation}((P_{a,1}, P_{b,1}))$ ;
     $\mathcal{I} \leftarrow \text{Find\_Inliers\_Motions}(\mathcal{S}, R_{a,b}, \vec{t}_{a,b})$ ;
    if  $|\mathcal{I}| > |\mathcal{I}_{best}|$  then
       $\mathcal{I}_{best} \leftarrow \mathcal{I}$ ;
    end
  end
   $\mathcal{C} = \mathcal{C} \cup \{\mathcal{I}_{best}\}$ ;
  if  $\mathcal{I}_{best} = \emptyset$  then
    break;
  end
end

```

4.4.1 Inlier recognition

As well as at the clustering of the transformations of the coordinate systems, a procedure for inlier recognition has to be chosen. The hypothesis consists of one motion (A_1, b_1) which has been computed from a camera pair $(P_{1,1}, P_{2,1})$. We want to check whether a motion (A_2, b_2) is consistent with this hypothesis. In order to be so, both the rotation A and the translation b have to be consistent.

We find distance D_r between the rotations A_1 and A_2 using a procedure in Section 2.1.2. This distance is in radians and therefore scale-invariant and thus the rotational consistency can be checked by comparison of the distance D_r with a fixed value, which can be a parameter of the pipeline. If the rotation is inconsistent, the camera pair is rejected. Otherwise, we follow to the checking of the translational consistency.

4.5.1 Condition for consistent cycles

If we move the object from the original position in the first take to a position of the second take, then from a position in the second take to a position in the third take and then from the position in the third take back to the original position in the first take, a position of an arbitrary point on the object will be equal to its original position. From this a condition for cycles of length 3 can be derived.

X_1 - Original position of a point.

X_2 - Point X_1 after its movement from the first to the second take

X_3 - Point X_2 after its movement from the second to the third take

A_{12} - Rotation from the first to the second take

b_{12} - Translation from the first to the second take

A_{23} - Rotation from the second to the third take

b_{23} - Translation from the second to the third take

A_{31} - Rotation from the third to the first take

b_{31} - Translation from the third to the first take

We assume that all points, rotations and translations are in the same coordinate system. Relationships between the points are as follows:

$$X_2 = A_{12}X_1 + b_{12}$$

$$X_3 = A_{23}X_2 + b_{23}$$

$$X_1 = A_{31}X_3 + b_{31}$$

It can therefore be written

$$X_1 = A_{31}(A_{23}X_2 + b_{23}) + b_{31} = A_{31}(A_{23}(A_{12}X_1 + b_{12}) + b_{23}) + b_{31} = A_{31}A_{23}(A_{12}X_1 + b_{12}) + A_{31}b_{23} + b_{31} = A_{31}A_{23}A_{12}X_1 + A_{31}A_{23}b_{12} + A_{31}b_{23} + b_{31}$$

From this equation the constraint for consistency of the cycles of the length 3 can be derived as:

$$A_{31}A_{23}A_{12} = I \quad (4.34)$$

$$A_{31}A_{23}b_{12} + A_{31}b_{23} + b_{31} = 0 \quad (4.35)$$

If the length of the cycle is n , we have a sequence of takes $\{t_i\}_{i=1}^{n+1}$ where $t_{n+1} = t_1$ in order for the takes to form a cycle. The corresponding sequence of translations and rotations is $\{A_{t_i, t_{i+1}}, b_{t_i, t_{i+1}}\}_{i=1}^n$, where $A_{a,b}$ is a rotation from take a to take b and $b_{a,b}$ is a translation from take a to take b . The point X_{t_1} is gradually transformed as follows:

$$X_{t_{i+1}} = A_{t_i, t_{i+1}}X_{t_i} + b_{t_i, t_{i+1}} \quad (4.36)$$

graph. We want to find all chordless cycles in this original graph. In order to do that we fill these cycles with additional chords until they are fully connected and then we can discover the original cycle and check its consistency.

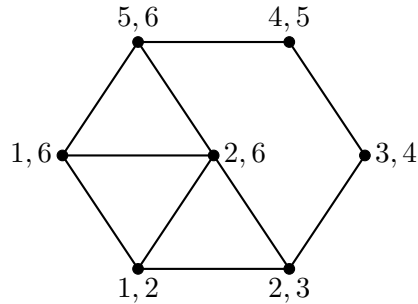


Figure 4.8: An example of a graph of the clusters. The labels of the vertices are the numbers of the initial and the final takes of the clusters they represent. The valid cycles are $((1,2), (2,3), (3,4), (4,5), (5,6), (1,6))$, $((1,2), (2,6), (1,6))$. The invalid cycles are $((1,2), (2,3), (2,6))$, $((1,6), (2,6), (5,6))$

A triplet of vertices can be a subset of a cycle if it does not contain multiple vertices corresponding to the same pair of takes, as in the cyclic transformation each take can be visited only once except for the first one which is the last one as well. The triplets satisfying this condition are inserted in a queue. Then until the queue is empty or no new chord can be added and no cycle can be verified remove the triplets from the queue.

Once a triplet is removed, the count of the existing edges between the vertices in the triplet is detected. The only possible values are 0, 1, 2, 3. If the count is 0 or 1, the triplet is returned back to the queue, as nothing can be done with it yet. If the count is 2, the missing chord is added to the graph and the triplet is assigned to the chord as its predecessor. The two other edges are assigned to the triplet as its predecessors. If the count is 3, the original cycle can be discovered and verified, all three edges are predecessors of the triplet.

The procedure is not optimal in the terms of number of operations as some cycles can be discovered more times, as the chordal completion algorithm has a greedy nature and is not optimal.

`Add_Edges(V, \mathcal{E})` adds edges from \mathcal{E} to a graph represented by an adjacency matrix V , `Valid(\mathcal{E})` returns true if edges from a set \mathcal{E} can build a valid cycle.

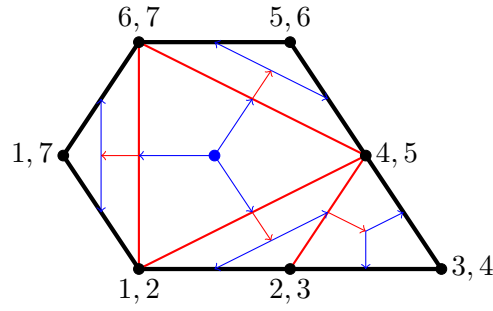


Figure 4.9: An example of a cycle after the chordal completion. Red edges are the chords. Red arrows lead from the chord to the 3-cycle which is its predecessor. Blue arrows lead from the 3-cycle to the edges which are its predecessors.

$\text{Check_3}((R_i, t_i), (R_j, t_j), (R_k, t_k))$ checks the cycle of length 3 according to equations (4.34), (4.35). $\text{Check_N}(C)$ checks the cycle C of length N according to equations (4.37), (4.38).

Algorithm 5: Init graph

input : n number of the motion clusters
 $\{a_i, b_i\}_{i=1}^n$ a sequence of takes of the clusters, $a_i < b_i$
output : G adjacency matrix of the initial graph $V \in M^{n,n}$
 $G \leftarrow 0^{n,n}$;
for $i \leftarrow 1$ **to** n **do**
 for $j \leftarrow i + 1$ **to** n **do**
 if $a_i = a_j$ **and** $b_i = b_j$ **then**
 $G_{i,j} \leftarrow -2$;
 $G_{j,i} \leftarrow -2$;
 else if $a_i = a_j$ **or** $a_i = b_j$ **or** $b_i = a_j$ **or** $b_i = b_j$ **then**
 $G_{i,j} \leftarrow 0$;
 $G_{j,i} \leftarrow 0$;
 else
 $G_{i,j} \leftarrow -1$;
 $G_{j,i} \leftarrow -1$;

Algorithm 6: Chordal completion

```

input :  $n$  number of the motion clusters
           $\{R_i, t_i\}_{i=1}^n$  a sequence of motions belonging to the clusters
           $\{a_i, b_i\}_{i=1}^n$  a sequence of takes of the clusters,  $a_i < b_i$ 
output :  $V$  adjacency matrix of the verified graph  $V \in M^{n,n}$ 
 $V \leftarrow 0^{n,n}$ ;
 $G \leftarrow \text{Init\_Graph}(n, \{a_i, b_i\}_{i=1}^n)$ ;
 $pos \leftarrow 1$ ;
 $T \leftarrow ()$ ;
 $Q \leftarrow$  empty queue;
for  $i \leftarrow 1$  do
  for  $j \leftarrow i + 1$  do
    for  $k \leftarrow j + 1$  do
      if  $G_{i,j} \neq -2$  and  $G_{j,k} \neq -2$  and  $G_{i,k} \neq -2$  then
         $T_{pos} \leftarrow (i, j, k)$ ;
         $Q.\text{enqueue}(pos)$ ;
         $pos \leftarrow pos + 1$ ;
      end
    end
  end
end
while  $Q$  not empty do
   $pos \leftarrow Q.\text{dequeue}$ ;
   $(i, j, k) \leftarrow T_{pos}$  if  $G_{i,j} \geq 0$  and  $G_{j,k} \geq 0$  and  $G_{i,k} \geq 0$  then
     $\mathcal{E} \leftarrow \text{Check\_Cycle}((i, j, k), G, \{R_i, t_i\}_{i=1}^n, T)$ ;
     $V \leftarrow \text{Add\_Edges}(V, \mathcal{E})$ ;
  else if  $G_{i,j} \geq 0$  and  $G_{j,k} \geq 0$  then
     $G_{i,k} \leftarrow pos$ ;
  else if  $G_{i,j} \geq 0$  and  $G_{i,k} \geq 0$  then
     $G_{j,k} \leftarrow pos$ ;
  else if  $G_{j,k} \geq 0$  and  $G_{i,k} \geq 0$  then
     $G_{i,j} \leftarrow pos$ ;
  else
     $Q.\text{enqueue}(pos)$ ;
  end
end

```

■ **Cycle discovering**

If the 3-cycle consists only of original edges, it can be directly checked on the consistency. Otherwise, the original edges are discovered via DFS. The edges from the cycle are inserted into the stack. Once an edge is removed

from the stack, it is controlled whether it is an original edge or an additional chord. If it is an original edge, it is saved to a set of edges. If it is an additional chord, it has a predecessor cycle. The other two edges from this predecessor cycle are also inserted to the stack. The search finishes when the stack is empty and the set of edges contains all edges which belong to the cycle.

Algorithm 7: Check cycles

input : (i, j, k) ids of the clusters in the cycle
 G adjacency matrix of a graph
 $\{R_i, t_i\}_{i=1}^n$ a sequence of motions of the clusters
 T a sequence of triplets of clusters

output : \mathcal{E} set of the verified edges
 $\mathcal{E} \leftarrow \emptyset$;

if $G_{i,j} + G_{j,k} + G_{i,k} = 0$ **then**
 | **if** $Check_3((R_i, t_i), (R_j, t_j), (R_k, t_k))$ **then**
 | | $\mathcal{E} \leftarrow \{ (i, j), (j, k), (i, k) \}$;

else
 | $Q \leftarrow$ empty queue;
 | $Q.enqueue((i, j))$;
 | $Q.enqueue((j, k))$;
 | $Q.enqueue((i, k))$;
 | **while** Q not empty **do**
 | | $(i, j) \leftarrow Q.dequeue$;
 | | **if** $G_{i,j} = 0$ **then**
 | | | $\mathcal{E} \leftarrow \mathcal{E} \cup \{(i, j)\}$;
 | | **else**
 | | | $pred \leftarrow G_{i,j}$;
 | | | $(a, b, c) \leftarrow T_{pred}$;
 | | | **if** $i \neq a$ or $j \neq b$ **then**
 | | | | $\mathcal{E} \leftarrow \mathcal{E} \cup \{(a, b)\}$;
 | | | **if** $i \neq a$ or $j \neq c$ **then**
 | | | | $\mathcal{E} \leftarrow \mathcal{E} \cup \{(a, c)\}$;
 | | | **if** $i \neq b$ or $j \neq c$ **then**
 | | | | $\mathcal{E} \leftarrow \mathcal{E} \cup \{(b, c)\}$;

| **if** not $Valid(\mathcal{E})$ **then**
 | | $\mathcal{E} \leftarrow \emptyset$;

| **else**
 | | $C \leftarrow Build_Cycle(\mathcal{E}, \{R_i, t_i\}_{i=1}^n)$;
 | | **if** not $Check_N(C)$ **then**
 | | | $\mathcal{E} \leftarrow \emptyset$;

These edges have to be tested whether they can build a valid cycle and eventually sorted to build one. A valid cycle is a sequence of motions $\{A_{t_i, t_{i+1}}, b_{t_i, t_{i+1}}\}_{i=1}^n$ where $t_1 = t_{i+1}$ and no other takes in the sequence $\{t_i\}_{i=1}^{n+1}$ repeat. The reason for discarding the cycles with repeating takes is that a sequence of forward motions $(A_{x,y}, b_{x,y}), (A_{y,z}, b_{y,z})$ followed by a sequence of backward motions $(A'_{z,y}, b'_{z,y}), (A'_{y,x}, b'_{y,x})$ could be a discovered cycle. This cycle would pass the condition from equation (4.37), but if it did, the forward and backward motions would be connected, which should not happen. The cycle is however not valid, as the take y repeats in it. An example of a cycle which is not valid is the cycle (1, 2), (2, 3), (2, 6) in Figure 4.8.

We have a set of edges. Each edge has two adjacent vertices, an initial one and a final one. The task is to build a valid cycle from these edges. At first, an arbitrary edge is selected and added to the sequence. Among the remaining edges, such edge is selected that one of its adjacent vertices is the final edge of the sequence. If the vertex is the final vertex of the edge, the edge has to be reversed and the motions belonging to the vertices are reversed, too. After all the edges have been inserted into the sequence, we test whether it is a valid cycle and if so, the conditions (4.38), (4.37) are used to check the consistency of the cycle. This time however the thresholds are not fixed but they depend linearly on the length of the cycle because longer cycles have higher uncertainty.

Algorithm 8: Build cycle

```

input :  $\mathcal{E}$  set of edges
           $\{R_i, t_i\}_{i=1}^n$  a sequence of motions of the clusters
output :  $C$  sequence of motions which builds a valid cycle
 $C \leftarrow ()$ ;
 $last \leftarrow -1$ ;
 $pos \leftarrow 1$ ;
while  $\mathcal{E} \neq \emptyset$  do
  foreach  $(a, b) \in \mathcal{E}$  do
    if  $last = -1$  or  $last = a$  then
       $C_{pos} \leftarrow (R_a, t_a)$ ;
       $pos \leftarrow pos + 1$ ;
       $last \leftarrow b$ ;
       $\mathcal{E} \leftarrow \mathcal{E} \setminus (a, b)$ ;
    else if  $last = b$  then
       $C_{pos} \leftarrow (R_a^{-1}, -R_a^{-1}t_a)$   $pos \leftarrow pos + 1$ ;
       $last \leftarrow a$ ;
       $\mathcal{E} \leftarrow \mathcal{E} \setminus (a, b)$ ;

```

■ 4.6.1 Track building

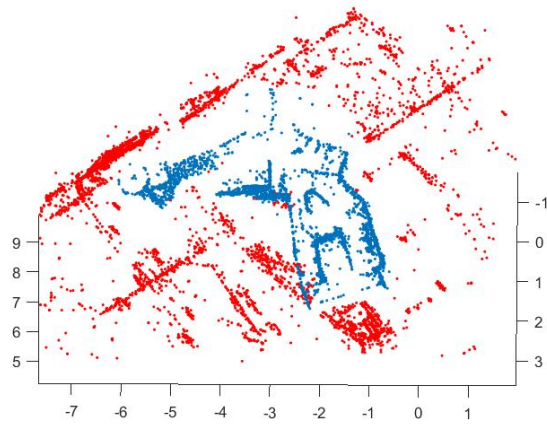
A graph is built whose vertices are the 3D points from all reconstructions. Two 3D points X, Y from different reconstructions a, b are connected by an edge if there is a 2D feature in some image, which observes the 3D point X in the reconstruction a and the 3D point Y in the reconstruction b . The weight of this edge is equal to the number of such 2D features. Each connected component of the graph is a tentative track. It is however possible due to mismatches, that points reconstructed from different real points are connected in the graph and therefore should occur in the same track.

A track is inconsistent if it contains multiple 3D points from the same reconstruction. If a track is inconsistent, it is split via the spectral clustering. The Laplacian of the graph is used directly without changes. The initial number of clusters is 2 and it is incremented until the k-means clustering of the eigenvectors of the Laplacian gives only consistent tracks.

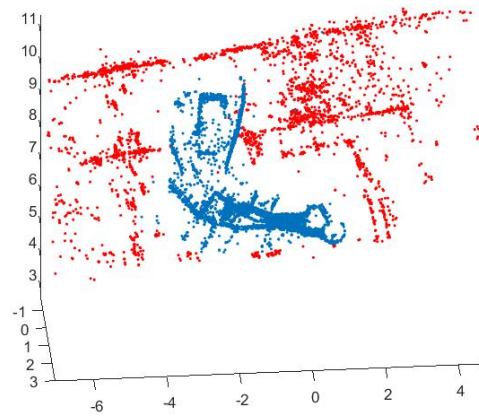
■ 4.6.2 Track segmentation

The motions in the clusters verified according to equations (4.37), (4.38) have been calculated from two cameras. The first camera was registered towards the background points and the second one towards the object points. After the cameras were registered, no further points were triangulated. Therefore for each camera pair (P_a, P_b) in the cluster the points observed by the first camera P_a should belong to the background and the points observed by the second camera P_b should belong to the object.

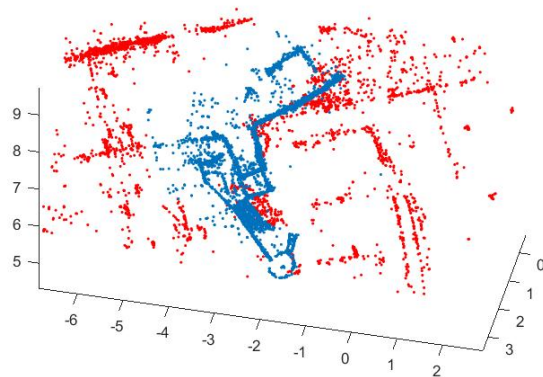
The cameras which are on the first position some pair in the cluster build the group G_a , the cameras on the second position build the group G_b . For each track, a score is introduced with 0 as the initial value. If the track is observed by a camera from the group G_a , its score is increased. If the track is observed by a camera from the group G_b , its score is decreased. The tracks with a positive score are assigned to the background, the tracks with a negative score are assigned to the object and the tracks whose score is zero are ignored.



(a)



(b)



(c)

Figure 4.10: Reconstructions of a single take after the segmentation of the points. Points which do not belong to background nor to the object are not shown.

4.7 Merging of the reconstructions

4.7.1 Order of the merging

The number of common points from the object between the reconstructions of takes i, j is $n_{i,j}^o$, the number of common points from the background between the reconstructions i, j is $n_{i,j}^b$. A graph is built whose vertices are the takes. The weight of the edge between the vertices i, j is $w_{i,j} = \min(n_{i,j}^o, n_{i,j}^b)$. If the weight should be zero, the edge does not exist. The first two takes in the order are the takes a, b which are adjacent to the edge with the highest weight. One of these takes is selected as the reference take. These two vertices are collapsed afterward. Weight of the edges between this collapsed vertex $\{a, b\}$ and an another vertex c is $w_{\{a,b\},c} = w_{a,c} + w_{b,c}$.

The take which is connected to the collapsed vertex with the edge with the highest weight is selected. This take is added as the next one in the order and it is merged with the collapsed vertex in the same way as the vertex was originally collapsed. This is repeated until all vertices are collapsed or there is no vertex connected to the collapsed vertex.

The order is determined greedily, however it is not necessary for each reconstruction to have common points with the reconstruction of the reference take.

4.7.2 Merging of the points

The final coordinate system is the system of the reference take r . The other reconstructions are transformed to this system in the given order. For each take two transformations are necessary, one for the object points and another one for the background points. $T_{r,s}^O$ is a set of all tracks from the object which contain points in both the reference take and in the take s which is to be transformed. P_s^O is a sequence of points from the take s which are contained in a track from $T_{r,s}^O$. $\{p_i^s, p_i^r\}_{i=1}^n$ is a sequence of pairs of corresponding points from $T_{r,s}^O$. The first point is from the take s , the second point is from the take r . The task is to find the scale, rotation and translation which transforms each p_i^s to its corresponding p_i^r as $p_i^r = \sigma R p_i^s + \vec{t}$. The transformation is

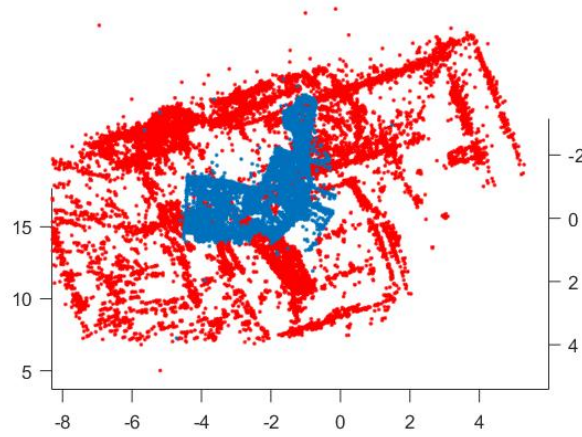


Figure 4.11: Points merged from 8 reconstructions of takes.

4.7.3 Calculation of the object motion from the transformations

s the take from whose coordinate system we transform

r the reference take to which we transform

$(\sigma_{s,r}, R_{s,r}^B, \vec{\sigma}_{s,r}^B)$ transformation of points of the background from the take s to the reference take r

$(\sigma_{s,r}, R_{s,r}^O, \vec{\sigma}_{s,r}^O)$ transformation of points of the object from the take s to the reference take r

X_O^s a point from the reconstruction s from the object

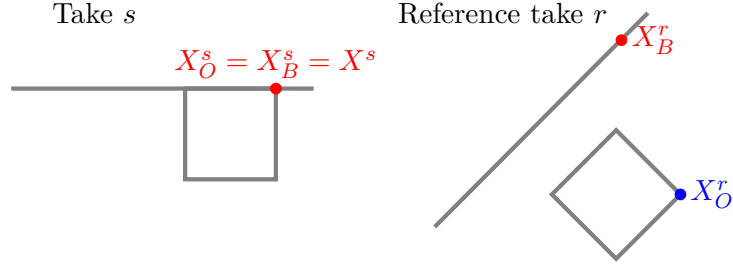
X_B^s a point from the reconstruction s from the background which is on the same position as the point X_O^s : $X_O^s = X_B^s$

X_O^r the point X_O^s transformed by $(\sigma_{s,r}, R_{s,r}^O, \vec{\sigma}_{s,r}^O)$

X_B^r the point X_B^s transformed by $(\sigma_{s,r}, R_{s,r}^B, \vec{\sigma}_{s,r}^B)$

$(A_{r,s}, \vec{b}_{r,s})$ motion of the object from the take r to the take s

$(A_{s,r}, \vec{b}_{s,r})$ motion of the object from the take s to the take r



The points X_O^s , X_B^s are on the same position, they can be replaced by a single point X^s , therefore we can write:

$$X_B^r = \sigma_{s,r} R_{s,r}^B X^s + \vec{\partial}_{s,r}^B \quad (4.44)$$

$$X_O^r = \sigma_{s,r} R_{s,r}^O X^s + \vec{\partial}_{s,r}^O \quad (4.45)$$

If the object was transformed from the position in the take r to the position in the take s , it would be on the same position where X_B^r is.

$$X_B^r = A_{r,s} X_O^r + \vec{b}_{r,s} \quad (4.46)$$

$$A_{r,s} (\sigma_{s,r} R_{s,r}^O X + \vec{\partial}_{s,r}^O) + \vec{b}_{r,s} = \sigma_{s,r} R_{s,r}^B X + \vec{\partial}_{s,r}^B \quad (4.47)$$

This equation has to hold true for all X , therefore also for the zero vector, therefore:

$$A_{r,s} \vec{\partial}_{s,r}^O + \vec{b}_{r,s} = \vec{\partial}_{s,r}^B \quad (4.48)$$

$$A_{r,s} \sigma_{s,r} R_{s,r}^O = \sigma_{s,r} R_{s,r}^B \quad (4.49)$$

As the scaling $\sigma_{s,r}$ is the same for the transformation of the object points and the background points, $A_{r,s}$ can be found from equation (4.49) as:

$$A_{r,s} = R_{s,r}^B (R_{s,r}^O)^{-1} \quad (4.50)$$

$\vec{b}_{r,s}$ can be found from equations (4.48) and (4.50) as:

$$\vec{b}_{r,s} = \vec{\partial}_{s,r}^B - R_{s,r}^B (R_{s,r}^O)^{-1} \vec{\partial}_{s,r}^O \quad (4.51)$$

The inverse motion $(A_{s,r}, \vec{b}_{s,r})$ transforms each point back to the original position:

$$Y = A_{s,r} (A_{r,s} Y + \vec{b}_{r,s}) + \vec{b}_{s,r} \quad (4.52)$$

This has to hold for every Y , so:

$$0 = A_{s,r} \vec{b}_{r,s} + \vec{b}_{s,r} \quad (4.53)$$

$$I = A_{s,r}A_{r,s} \quad (4.54)$$

From equations (4.53), (4.54) the inverse motion $(A_{s,r}, \vec{b}_{s,r})$ can be found as:

$$A_{s,r} = A_{r,s}^{-1} = R_{s,r}^O (R_{s,r}^B)^{-1} \quad (4.55)$$

$$\vec{b}_{s,r} = -A_{s,r} \vec{b}_{r,s} = -R_{s,r}^O (R_{s,r}^B)^{-1} (\vec{\sigma}_{s,r}^B - R_{s,r}^B (R_{s,r}^O)^{-1} \vec{\sigma}_{s,r}^O) \quad (4.56)$$

$$\vec{b}_{s,r} = \vec{\sigma}_{s,r}^O - R_{s,r}^O (R_{s,r}^B)^{-1} \vec{\sigma}_{s,r}^B \quad (4.57)$$

4.7.4 Merging of the cameras

- P_s a camera in the reconstruction of the take s
- P_r^B the camera after the transformation to the reference take towards the background points
- P_r^O the camera after the transformation to the reference take towards the object points
- I image from which the camera P_s arises
- t take from which the camera P_s and therefore also the image I arises
- s the reconstruction onto which the camera P has been registered
- d descent of the camera P ; if $s = t$, the camera is in the anchor take, so $d = A$, otherwise the camera has been registered using the sequential PnP. If it was registered towards the background, then $d = B$, if it was registered towards the object points, then $d = O$
- R_s rotation of the camera P_s in the original position
- \vec{c}_s centre of the camera P_s in the original position
- R_r^B rotation of the camera P_r^B in the coordinate system of the reference take r towards the background points
- \vec{c}_r^B centre of the camera P_r^B in the coordinate system of the reference take r towards the background points
- R_r^O rotation of the camera P_r^O in the coordinate system of the reference take r towards the object points
- \vec{c}_r^O centre of the camera P_r^O in the coordinate system of the reference take r towards the object points

4. Proposed solution

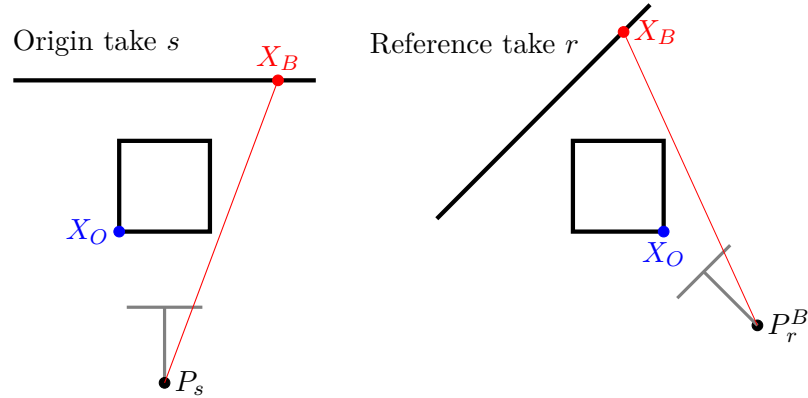
At first, the cameras are transformed to the reference take r to the position towards the background. The procedure of the transformation of the camera P_s to P_r^B however differs depending on t , s and o .

■ $s = r; d \in \{A, B\}$

The camera is already in the desired position, no transformation is necessary.

■ $s \neq r; d \in \{A, B\}$

The camera is in another coordinate system.



The transformation to the reference take is performed according to the equations (4.17), (4.18) as:

$$R_r^B = R_s(R_{s,r}^B)^{-1} \quad (4.58)$$

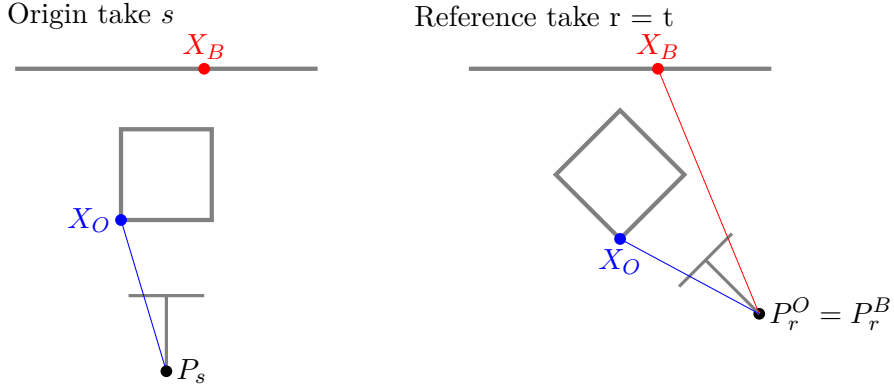
$$R_r^B \vec{c}_r^B = R_r^B \vec{o}_{s,r}^B + \sigma_{s,r} R_s \vec{c}_s \quad (4.59)$$

$$\vec{c}_r^B = \vec{o}_{s,r}^B + \sigma_{s,r} (R_r^B)^{-1} R_s \vec{c}_s \quad (4.60)$$

The transformation $(\sigma_{s,r}, R_{s,r}^B, \vec{o}_{s,r}^B)$ transforms the points from the background in the coordinate system s to the points in the reference take r . The camera in the transformed pose (R_r^B, \vec{c}_r^B) therefore observes the points from the background in the reference take, which is the desired result.

■ $s \neq r; o = O; t = r$

The camera comes from the reference take r but it has been registered onto another reconstruction towards the object points.



The pose of the camera P_r^O observing the object points in the reference take transformed by $(\sigma_{s,r}, R_{s,r}^O, \vec{c}_{s,r}^O)$ can be computed according to the equations (4.17), (4.18):

$$R_r^O = R_s (R_{s,r}^O)^{-1} \quad (4.61)$$

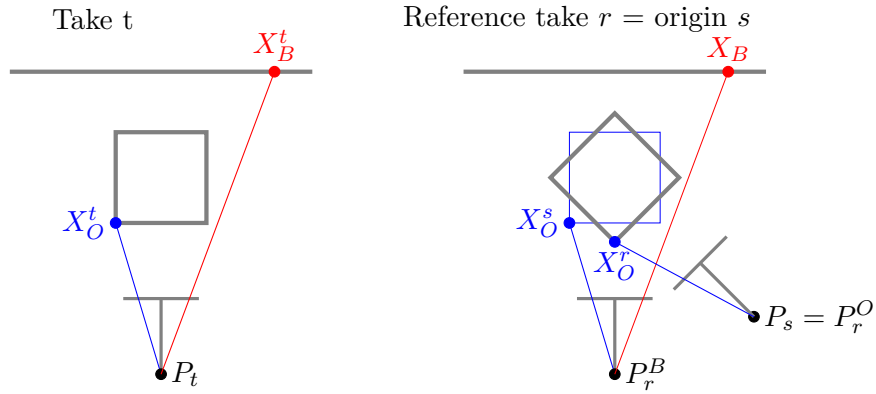
$$R_r^O \vec{c}_r^O = R_r^O \vec{\sigma}_{s,r}^O + \sigma_{s,r} R_s \vec{c}_s \quad (4.62)$$

$$\vec{c}_r^O = \vec{\sigma}_{s,r}^O + \sigma_{s,r} (R_r^O)^{-1} R_s \vec{c}_s \quad (4.63)$$

The pose of the transformed camera is towards the object points. But because the camera arises from the reference take, the pose towards the object and the background in the reference take is the same, so $R_r^B = R_r^O$, $\vec{c}_r^O = \vec{c}_r^B$.

■ $s = r; o = O; t \neq r$

The camera is already in the target coordinate system, but it registered towards the object. It has to be moved to the position where it would observe the background points.



The camera P_s observes a point X on the position X_O^r where it is in the reference take. If the camera was registered towards the background, it would observe the point on its original position X_O^s .

$$X_O^s = A_{r,t} X_O^r + \vec{b}_{r,t} \quad (4.64)$$

Where $A_{r,t}$, $\vec{b}_{r,t}$ are calculated according to equations (4.55), (4.57). Rotation and centre of the camera registered towards the background points can be found according to equations (4.6), (4.7) as:

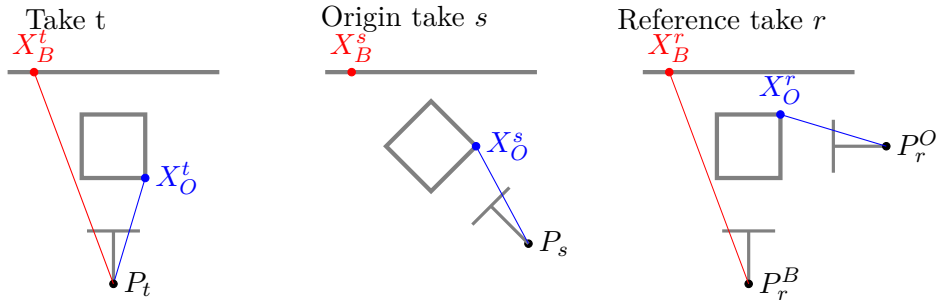
$$R_r^B = R_s A_{r,t}^{-1} = R_s A_{t,r} \quad (4.65)$$

$$R_r^B \vec{c}_r^B = R_r^B \vec{b}_{r,t} + R_s \vec{c}_s \quad (4.66)$$

$$\vec{c}_r^B = \vec{b}_{r,t} + (R_r^B)^{-1} R_s \vec{c}_s \quad (4.67)$$

■ $s \neq r; o = O; t \neq r$

The camera does not come from the reference take, so the transformation $(\sigma_{s,r}, R_{s,r}^O, \vec{c}_{s,r}^O)$ does not bring the desired result. This case can however be solved as a combination of the two previous cases.



If the camera P_s is transformed with $(\sigma_{s,r}, R_{s,r}^O, \vec{c}_{s,r}^O)$ according to equations (4.61), (4.63) to the reference take r , the transformed camera P_r^O observes the object points of the reference take on the same features where the original camera P_s observed the object points in the take s .

$$R_r^O = R_s(R_{s,r}^O)^{-1} \quad (4.68)$$

$$\vec{c}_r^O = \vec{c}_{s,r}^O + \sigma_{s,r}(R_r^O)^{-1}R_s\vec{c}_s \quad (4.69)$$

This converts the problem to the previous one where the task is to transform the camera P_r^O observing the object to the camera P_r^B which observes the background. This can be done according to the equations (4.65), (4.67) where $A_{r,t}$, $\vec{b}_{r,t}$ are calculated according to equations (4.55), (4.57):

$$R_r^B = R_r^O A_{r,t}^{-1} = R_r^O A_{t,r} \quad (4.70)$$

$$\vec{c}_r^B = \vec{b}_{r,t} + (R_r^B)^{-1}R_s\vec{c}_s \quad (4.71)$$

■ Camera averaging

After all the cameras are transformed to the reference take r to the pose towards the points from the background, all cameras which arise from the same image I should have the same pose.

A_I is a set of all transformed cameras P_r^B which arise from the image I . R_I is rotation calculated as the median of the Euler vectors which represent the rotations R_r^B of the cameras from A_I , \vec{c}_I is camera center calculated as the median of the centers \vec{c}_r^B of the cameras from A_I . All cameras from A_I can therefore be merged to one camera P_I with pose (R_I, \vec{c}_I) and whose observations is a union of all observations from all cameras from A_I

If the take t from which the image I arises is the reference take r , the camera P_I observes points from both the background and the object, therefore the merging is finished. But if $t \neq r$, the camera P_I observes only the points from the background. In this case the camera needs to be split into two cameras P_I^B, P_I^O , where P_I^B observes the background and P_I^O observes the object.

P_i^O camera which arises from the image I_i and observes the points from the object; if $t_i = r$, then $P_i^O = P_i^B$, otherwise must according to equations (4.72), (4.74) hold true $R_i^O = R_i^B A_{r,t_i}$, $\vec{c}_i^O = A_{r,t_i}^{-1}(\vec{c}_i^B - \vec{b}_{r,t_i})$

$\{O_i\}_{i=1}^l$ sequence of the observations; an observation O_i is a triplet (x_i, p_i, \vec{f}_i) where x_i is index of the 3D point which is observed, p_i is index of the camera which observes the point and \vec{f}_i is a 2D feature onto which the point is projected;

Parameters of the BA are all 3D points $\{X_i\}_{i=1}^m$, the poses of the cameras $\{P_i^B\}_{i=1}^n$ which observe the background points and the motions $\{(A_{r,t}, \vec{b}_{r,t})\}_{t=1, t \neq r}^k$. Poses of the cameras $\{P_i^O\}_{i=1}^n$ are not parameters as they can be computed from the cameras P_i^B and from the motions $(A_{r,t}, \vec{b}_{r,t})$ and because we want to ensure that the motion of the object from the take t to r computed using equations (4.6), (4.8) is the same if computed from any camera pair arising from the take t .

The task is to find such parameters which minimize the sum of squares of the reprojection errors over all observations

$$\operatorname{argmin} \left(\sum_{i=1}^d \|\vec{f}_i - q_i\|^2 \right) \tag{4.75}$$

Where q_i is the projection of the point X_{x_i} onto the camera $P_{p_i}^B$ if $d_{x_i} = B$ and the projection of the point X_{x_i} onto the camera $P_{p_i}^O$ if $d_{x_i} = O$. The projections are found according to Section 2.2.

Two different iterative approaches to the BA have been tried. The first one is the classical gradient descent which improves all parameters in each iteration. The other one is alternating minimization, where in the first iteration the cameras are fixed and the points are improved and in the next one the points are fixed and the cameras are improved; this is repeated until convergence. These approaches are evaluated in Section 5.4.

4.9 Filtering of the reconstruction

This is an arbitrary step, which may improve the quality of the reconstruction in the case of a high number of misclassified points in the reconstruction.

If a point is misclassified (arises from the background but is assigned to the object or vice versa), it will probably be observed by some of the cameras with a high reprojection error. This method requires a threshold t and the points, whose reprojection error after the projection onto an arbitrary camera is greater than t , are removed from the reconstruction.

This method works correctly only if the cameras are estimated correctly, otherwise it would remove the correctly classified points, as well, which would make the quality of the reconstruction worse. From the same reason a high threshold t is used, which still removes the most of the misclassified points. The filtering may be performed before or after the BA.

Chapter 5

Experiments

5.1 Overview of used objects and backgrounds

In order to prove that our method works well, we have selected several objects and we have created a dataset of images for each of the objects. We have selected appropriate backgrounds for the objects according to [29]. Each dataset is divided into 4 to 8 takes, each of which depicts a static configuration of the object towards the background, as it is described in 1.3.

Object	Background	Takes	Images	Special	Figure
Daliborka	Nissan GTR	8	338	NONE	5.1(a)
Lycan	Nissan GTR	8	140	NONE	5.1(b)
Colosseum	Peugeot 908	8	294	NONE	5.1(c)
Vatican	Peugeot 908	8	270	NONE	5.1(d)
Ganesha	Peugeot 908	4	160	Only translation	5.1(e)
Salt lamp	Seat Ibiza	8	330	Translucent	5.1(f)
Buddha	Peugeot 908	8	300	Shiny	5.1(g)
Pillow	Seat Ibiza	8	278	Repetitive	5.1(h)
Transformer	Audi S8	8	329	NONE	5.2(a)
Ship	Audi S8	8	277	NONE	5.2(b)
Catalog	Seat Ibiza	5	150	Planar	5.2(c)
Lego	Seat Ibiza	4	264	Repetitive, Planar	5.2(d)

Table 5.1: Properties of the datasets.

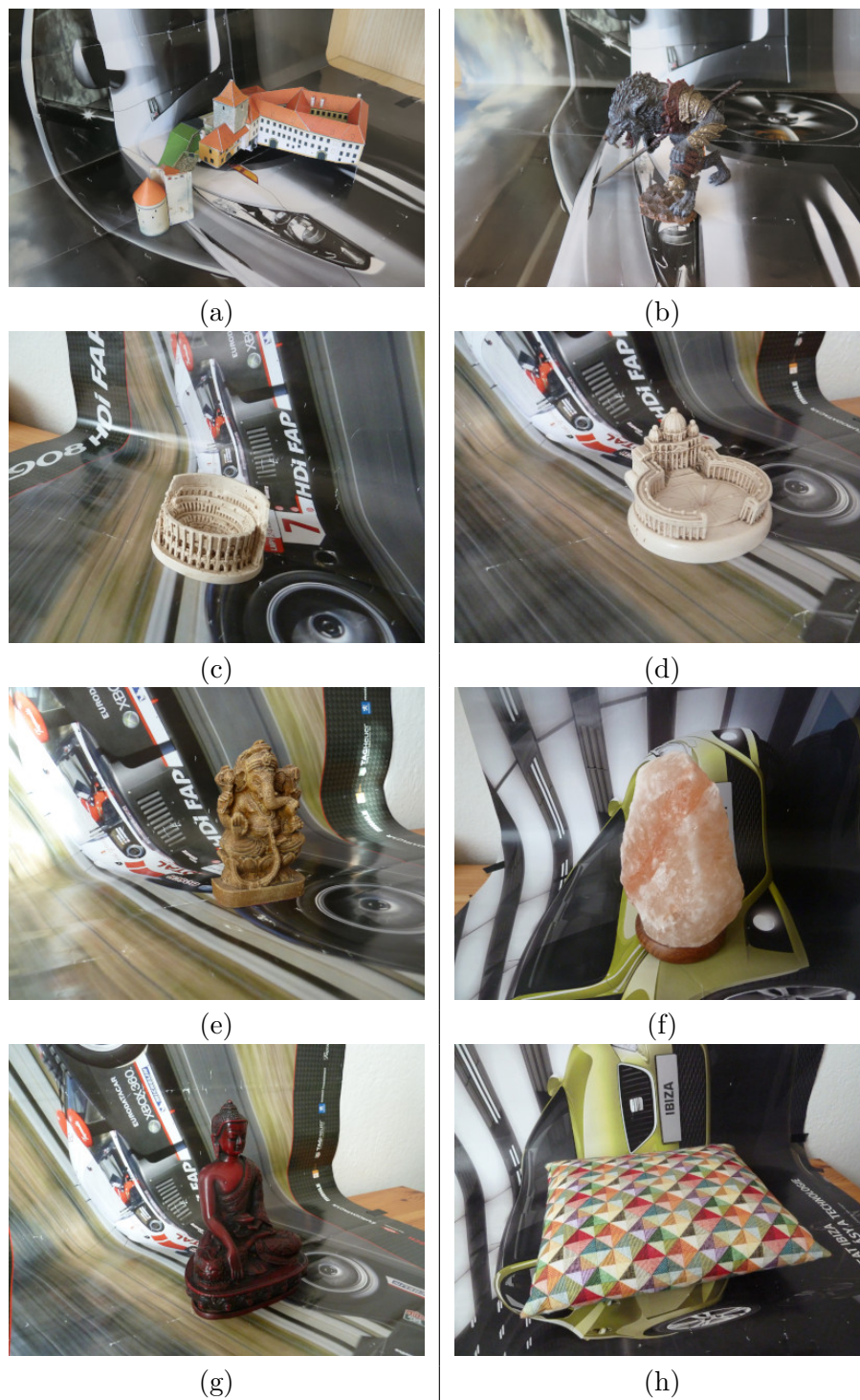


Figure 5.1: Images taken from the datasets, which are used for reconstruction. Part 1.

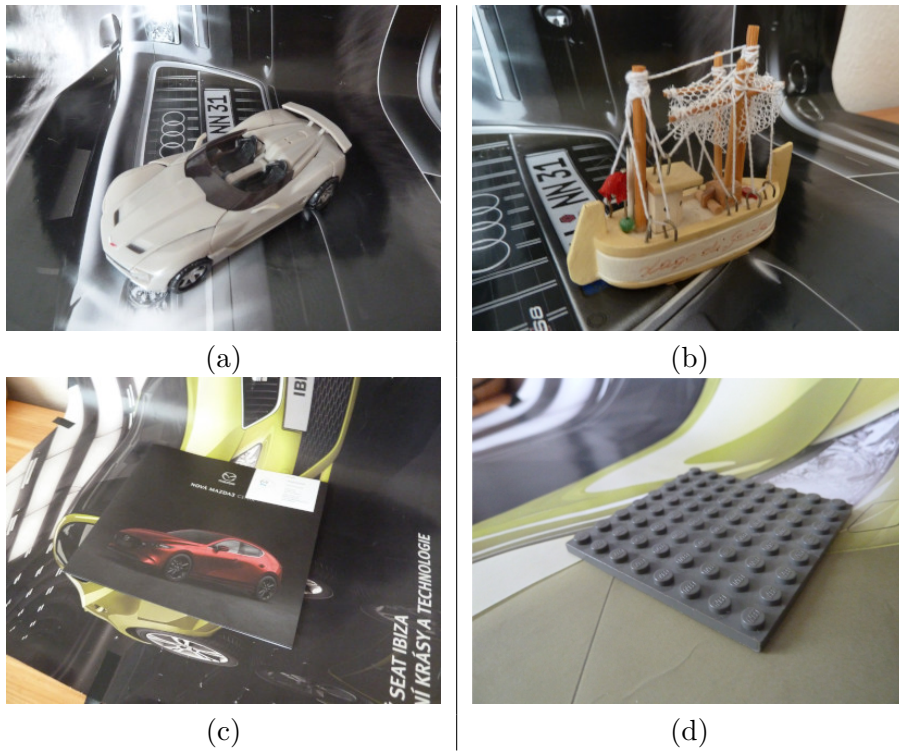


Figure 5.2: Images taken from the datasets, which are used for reconstruction. Part 2.

5.2 Qualitative results of the method

Except for the Buddha dataset, the datasets have been successfully reconstructed by our method. Table 5.2 and figures 5.3, 5.4 depict the results of the reconstructions. The blue points belong to the object and the red ones to the background. According to Section 4.6.2 the only criterion to label the background and the object is the order of the sequential PnP registration. The registration can be arbitrary, so the labels of the background and the object can be swapped, which actually happened at the datasets 5.4(a), 5.4(c), 5.4(f).

The table shows whether the background and the object have been swapped as well as numbers of the points from the object and from the background. If the labels of the background and the object have been swapped, Points Obj. is the number of points from the true object, which has been labeled as the background.

The reconstruction of the Buddha dataset has failed probably due to the lack of features on one side of the object so the sequential PnP could not register towards the object and the cyclic check could not be performed.

Dataset	Figure	Points Obj.	Points Bck.	Swapped
Daliborka	5.3(b)	15904	16553	NO
Lycan	5.3(d)	10286	8138	NO
Colosseum	5.4(b)	23535	25990	NO
Vatican	5.4(d)	17175	15282	NO
Ganesha	5.4(f)	19199	10528	YES
Salt lamp	5.4(h)	2322	24235	NO
Pillow	5.5(b)	92166	18335	YES
Transformer	5.5(d)	2070	12606	NO
Ship	5.5(f)	4665	9879	NO
Catalog	5.5(h)	5047	8421	YES
Lego	5.6(b)	354	7765	NO

Table 5.2: Results of the reconstructions.

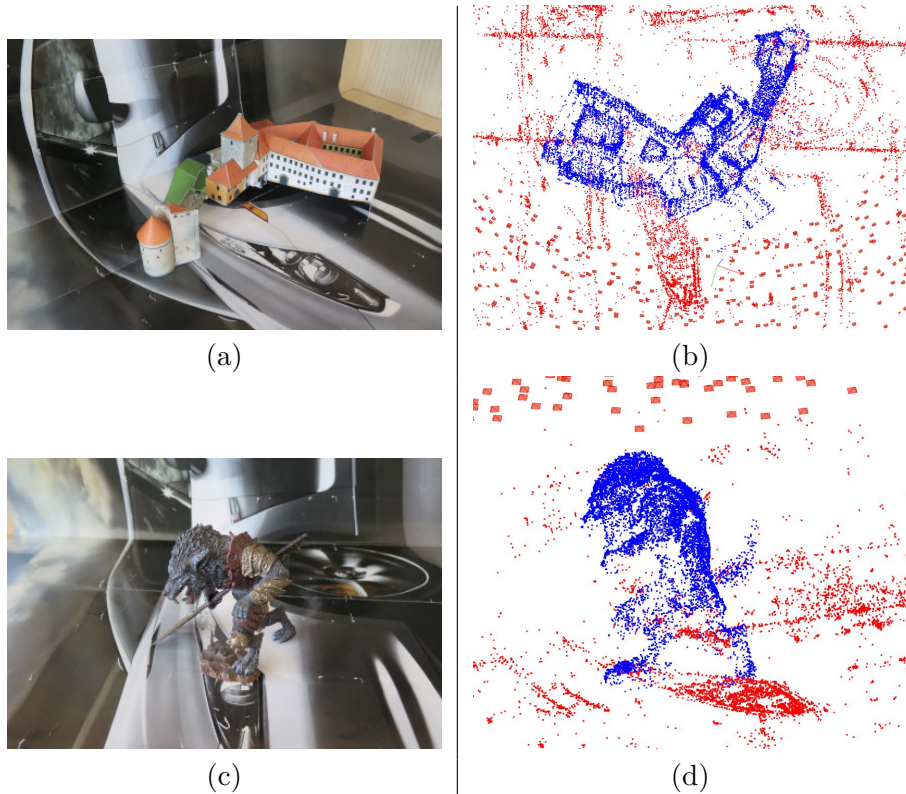
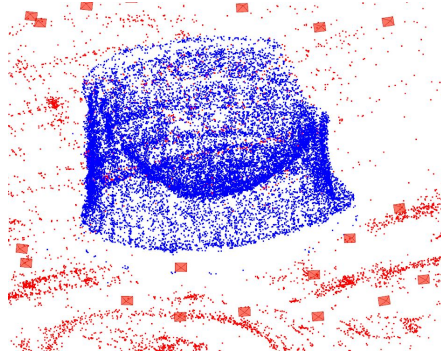


Figure 5.3: Reconstructed datasets, Part 1.



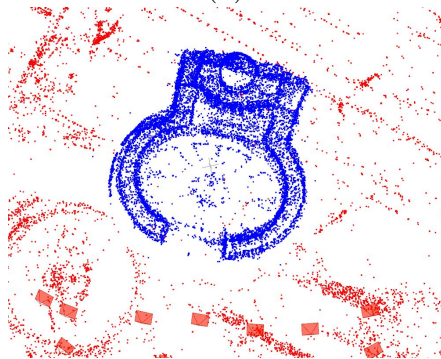
(a)



(b)



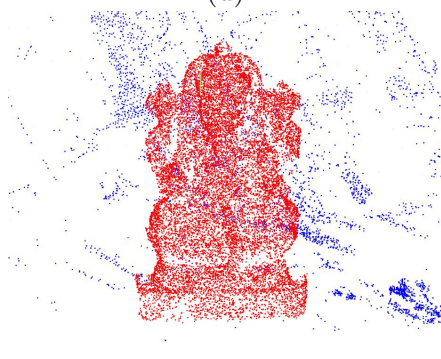
(c)



(d)



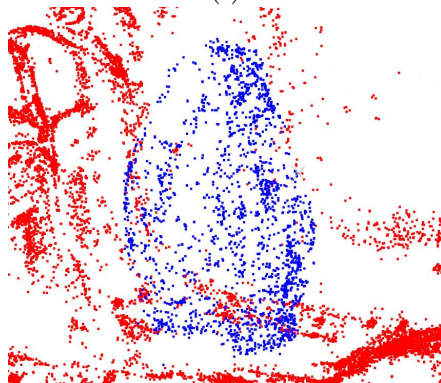
(e)



(f)



(g)



(h)

Figure 5.4: Reconstructed datasets, Part 2.

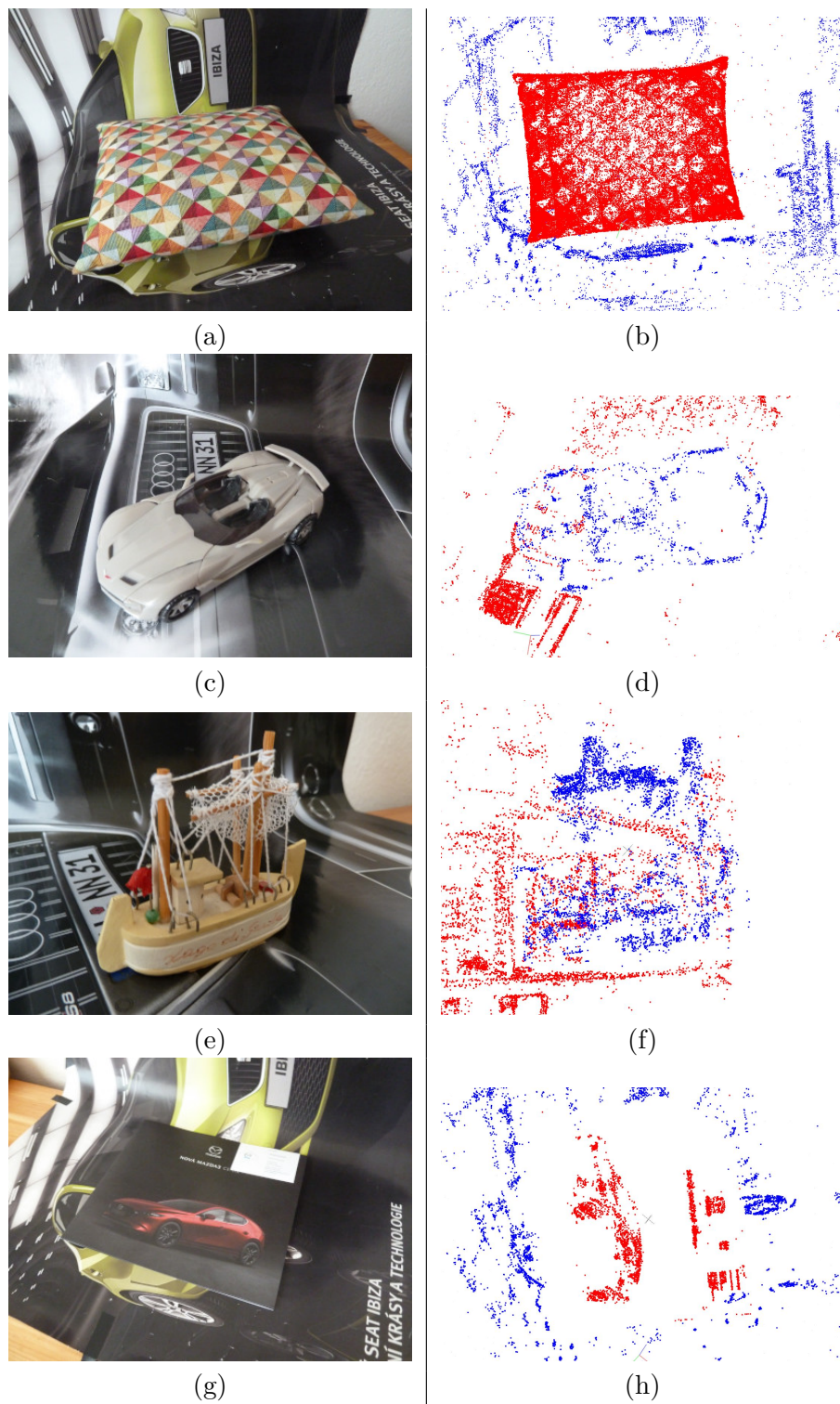


Figure 5.5: Reconstructed datasets, Part 3.

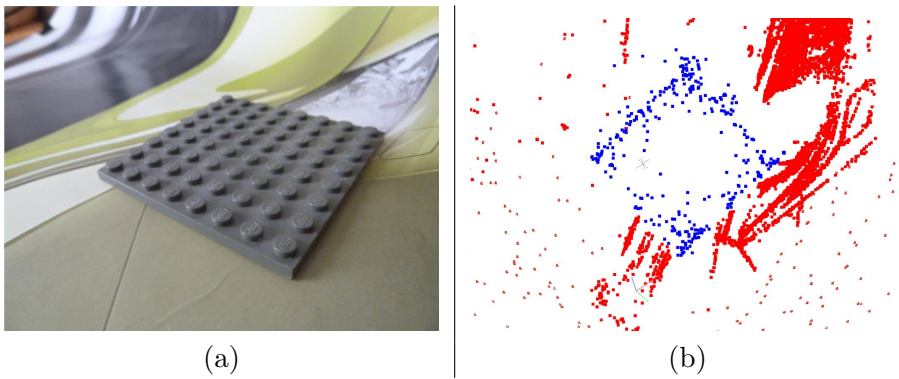


Figure 5.6: Reconstructed datasets, Part 4.

Quality of the models which have been successfully reconstructed differs. The main reason for the lower quality of the reconstruction is the lower density of features, which is especially apparent at the Transformer dataset depicted in Figure 5.2(a). An interesting result is the reconstruction of a planar object "Catalog" which is depicted in Figure 5.2(b). The object "Lego" in Figure 5.2(b) could not be reconstructed with the default settings of COLMAP, it was however reconstructed with a lowered threshold for minimal inlier ratio in the PnP registration, which allowed registration of the images towards the object.

5.3 Comparison to the single body SfM

Apart from the datasets which depict the object on some background, we have reconstructed some of the objects without any background in the classical single body SfM COLMAP [3]. The motivation was to determine whether and eventually in which cases can the presence of the background actually help with the reconstruction of the object. The numbers of the images in the datasets are given in Table 5.3. The models reconstructed with the background have been adjusted using the alternating minimization, which is proven to be more suitable in Section 5.4.

For each object, we compare two models, first of which has been reconstructed with the background using our method, while the second one has been reconstructed without the background using a standard SfM pipeline. The parameters which we compare are the number of reconstructed points, the median reprojection error, and the quality of the reconstruction (e.g.

Dataset	MBSfM	Single body
Colosseum	294	326
Vatican	270	304
Ganesha	160	187
Salt lamp	330	302
Pillow	278	275
Transformer	329	290
Ship	277	294
Catalog	150	151
Lego	264	307

Table 5.3: Number of images from which the object is reconstructed



Figure 5.7: Example of an image from the single body dataset without the background.

whether all parts of the object have been reconstructed). Table 5.4 shows the numbers of points and Table 5.5 shows the reprojection errors.

Dataset	MBSfM	MBSfM filtered	Single body
Colosseum	23535	22932	49418
Vatican	17175	17094	38833
Ganesha	19199	19192	34539
Salt lamp	2322	2320	10049
Pillow	92166	86351	198728
Transformer	2070	2037	14372
Ship	4665	4636	20571
Catalog	5047	4466	10143
Lego	354	90	-

Table 5.4: Number of points from the object

It is apparent from Table 5.4 that the object reconstructed with the background using our method has usually fewer points than the same object reconstructed without the background. The possible reason for this is that only the points towards which a camera has been registered in the sequential PnP can be recognized to belong either to the background or to the object. Therefore some of the points which have been reconstructed with the background are not present in the final reconstruction. Another possible reason may be that the features on the object are ignored because of stronger features on the background. For some of the objects, the images in the dataset without the background outnumber the images in the dataset with the background (Table 5.3), but the difference is too low to make such a difference in the numbers of the points.

The most significant difference in the numbers of points is at the objects which have a low density of the features such as the Salt lamp, the Transformer (Figure 5.8(a), 5.8(b)) and the Ship. The difference at the objects which have a higher density of features is less significant and visually unrecognizable. (Figure 5.8(c) 5.8(d))

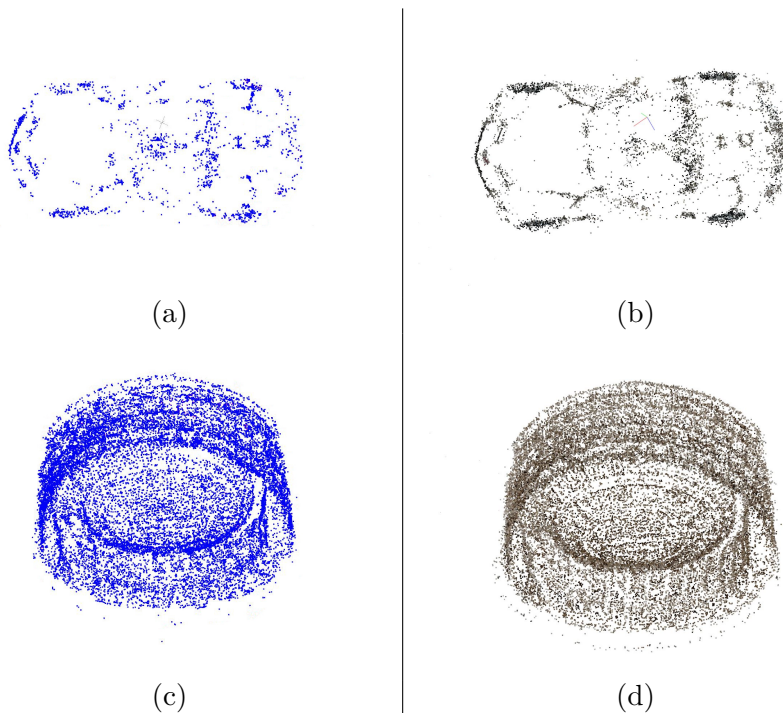


Figure 5.8: Comparison of the models reconstructed with and without the background. The models on the left side have been reconstructed with the background.

Dataset	MBSfM	MBSfM filtered	Single body
Colosseum	3.2662	3.2085	1.4825
Vatican	3.4622	3.4427	1.7924
Ganesha	0.8209	0.8208	1.8522
Salt lamp	4.1629	4.1610	1.8889
Pillow	1.0564	0.9710	2.3420
Transformer	6.3239	6.1637	2.6443
Ship	6.0914	6.0112	1.7392
Catalog	2.0319	2.0035	2.8318
Lego	20.4614	7.2846	-

Table 5.5: Median reprojection error E_O

For the majority of the objects, the reprojection error of the reconstruction with the background is higher than the error of the reconstruction without the background. The exceptions are the Ganesha, the Catalog and the Pillow.

5.3.1 Planar objects

Our method allows a reconstruction of planar objects by placing them onto a non-planar background. The object with the background are then reconstructed as a general object and we can afterward separate the object from the background. This has been demonstrated on the object "Catalog", whose reconstruction can be seen in Figure 5.9(a). However, the single-body pipeline COLMAP managed to reconstruct the dataset, as well (Figure 5.9(b)).

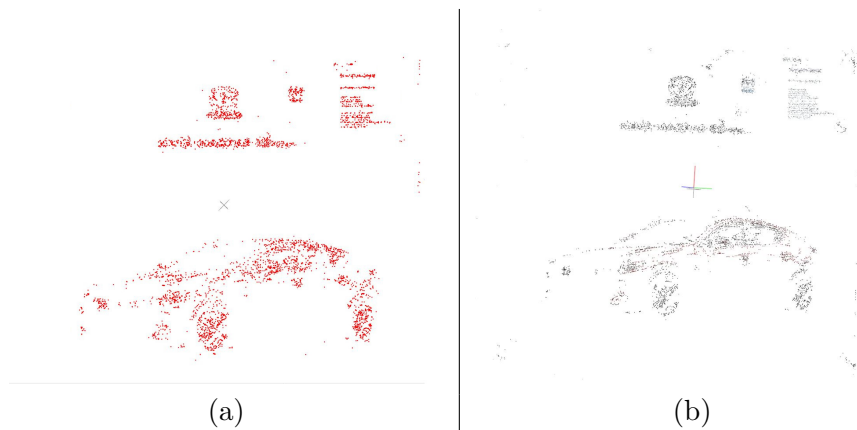


Figure 5.9: Comparison of the models of the "Catalog" object reconstructed with and without the background. The model (a) has been reconstructed with the background.

5.3.2 Repetitive objects

In most cases, the reconstruction without the background has higher a number of reconstructed points and lower reprojection errors at the same time. In some cases, however, our method can perform better. One of these cases are objects with repetitive patterns.

According to Table 5.5, a repetitive object "Pillow" is one of the objects which have lower reprojection error after the reconstruction with the background, than without the background.

Objects "Daliborka" and "Vatican" contain repetitive patterns (Figure 5.10). These patterns could be reconstructed using our method. Figure 5.11 shows that the repetitive back side of the Vatican object has been reconstructed better with the background than without it.

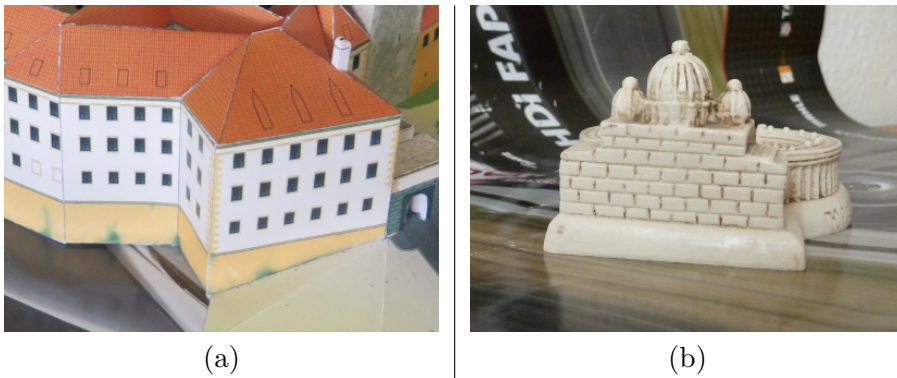


Figure 5.10: Examples of repetitive patterns on "Daliborka" (a) and "Vatican" (b).

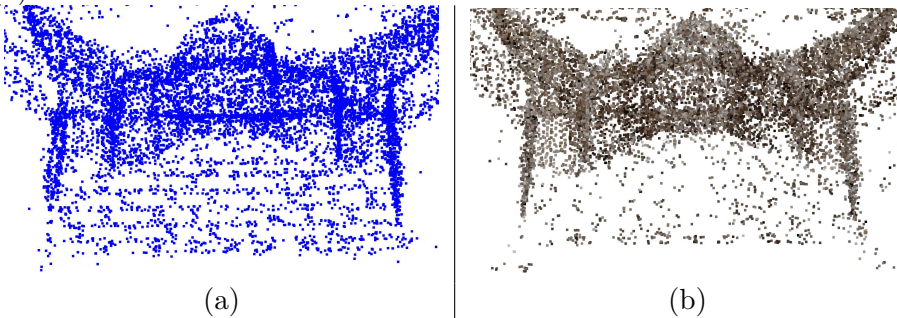


Figure 5.11: Comparison of the reconstruction of the back side of the "Vatican" object. The model (a) has been reconstructed with the background.

The object "Lego" is repetitive and almost planar, and therefore it is difficult to reconstruct. The standard single body pipeline failed to find the initial

pair of cameras and therefore could not reconstruct the object. Our method managed to reconstruct the object, as the relative pose could be found using the points from the background. This model is shown in Figure 5.12

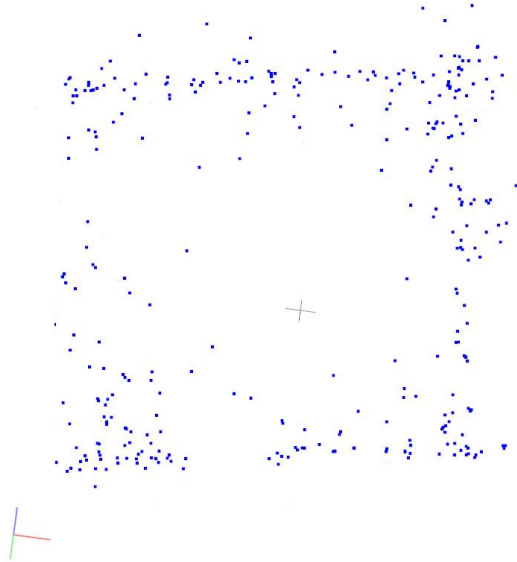


Figure 5.12: Model of the "Lego" object.

5.4 Review of approaches to the Bundle Adjustment

In Section 4.8 we have proposed two approaches to the bundle adjustment, namely the gradient descent and the alternating minimization. We are interested in the median reprojection error E , as well as in the partial reprojection errors E_O, E_B . E_O is calculated only from the observations of points which belong to the object and E_B is calculated only from the observations of the points which belong to the background. If the labels of the background and the object have been swapped, E_O is calculated from the true object, which has been labeled as the background.

For each dataset, we compare the values E, E_O, E_B of the model before the adjustment, after the adjustment using the gradient descent, and after the adjustment using the alternating minimization in order to determine which of the approaches is more appropriate. Table 5.6 shows values for E , Table 5.7 shows E_O and Table 5.8 shows E_B . The best value is bold; if the error after the adjustment is worse, the value is red.

Dataset	Original	Gradient	Alternating
Daliborka	3.4294	1.8704	1.5888
Lycan	2.3397	2.0448	1.3930
Colosseum	5.9503	3.4171	2.6892
Vatican	6.6690	3.2737	2.8050
Ganesha	2.8834	1.4558	1.3600
Salt lamp	8.9735	2.3805	2.4062
Pillow	7.8455	1.4150	1.4340
Transformer	15.9971	2.2689	2.2509
Ship	11.3332	2.8786	2.8458
Catalog	9.8959	4.8998	4.5814
Lego	9.6033	2.0808	2.0773

Table 5.6: Overall median reprojection error E

Dataset	Original	Gradient	Alternating
Daliborka	2.4018	3.3470	2.1918
Lycan	1.6981	3.5872	1.6719
Colosseum	3.3818	5.7818	3.2662
Vatican	3.5849	5.7980	3.4622
Ganesha	1.9734	0.8508	0.8209
Salt lamp	3.4133	7.4074	4.1629
Pillow	6.9275	1.0436	1.0564
Transformer	6.9450	9.2377	6.3239
Ship	6.5760	9.2265	6.0914
Catalog	5.5661	2.3659	2.0319
Lego	50.6459	20.7596	20.4614

Table 5.7: Median reprojection error E_O of the object

According to Table 5.6 the alternating minimization outperforms the gradient descent in the most cases. In addition to that, Table 5.7 shows, that in 7 of 11 cases the partial error E_O is made worse by the gradient descent. The same thing happens only at one dataset in the case of the alternating minimization. Therefore the alternating minimization is more suitable for our purpose.

Dataset	Original	Gradient	Alternating
Daliborka	4.7270	1.2718	1.2373
Lycan	3.3340	1.1679	1.1408
Colosseum	9.4377	2.3787	2.3417
Vatican	10.4341	2.5497	2.5125
Ganesha	7.5913	5.8102	4.4917
Salt lamp	9.5084	2.2928	2.3459
Pillow	13.8701	36.6333	6.7692
Transformer	16.8831	2.0978	2.0942
Ship	13.6002	2.3304	2.3812
Catalog	14.5464	8.1427	9.4299
Lego	9.4550	2.0490	2.0440

Table 5.8: Median reprojection error E_B of the background

Yet another advantage of the alternating minimization is that the reconstruction is naturally stabilized during the bundle adjustment, as in each step either the cameras or the points are fixed.



Chapter 6

Future work

We have solved a relaxed version of the MBSfM problem, where the number of the objects is limited to 2 and the images depict several static configurations of a scene and it is known which image belongs to which scene. This algorithm can be improved in many ways.

In order to dispose of the dependence on the thresholds in the clustering phase, the uncertainties of the cameras can be utilized. In that case, a series of experiments would have to be performed to show which of the options, or the combination of both, would perform better.

The solution can be generalized in such a way, that the number of moving objects would be arbitrary. In that case, it would be challenging to handle the situations where some of the objects remain static between some of the takes, while the other ones move. The algorithm can be improved such that it would be able to assign the images to their takes, so the labels of the images would not be necessary.

In order to improve the quality of the reconstruction, new points could be triangulated from camera pairs which do not belong to the same object. Many points have been reconstructed but they have not been assigned to the object nor to the background. An additional labeling of these points can be performed.

The sequential PnP can be replaced by another motion segmentation method. This algorithm can as well be used in a general version of the MBSfM.



Chapter 7

Conclusion

In this thesis we have proposed a solution to the problem introduced in Section 1.3, where the task is to reconstruct a scene consisting of two objects. This scene is captured by multiple takes, which depict different static configurations of the scene. We have implemented the method as an extension of a COLMAP pipeline [3] and we have demonstrated its functionality on real data.

We have reconstructed several objects using our method, where the background served as the second object. In most cases the reconstruction of an object with a background using our method produces a model, which is worse than the model reconstructed in a single-body pipeline without the background. We have however shown, that for some objects, especially the planar ones or the objects which contain repetitive structures, our method produces better models than the single-body SfM pipeline.



Bibliography

- [1] R. Hartley and A. Zisserman. Multiple view geometry in computer vision. 2003.
- [2] S. M. Seitz N. Snavely and R. Szeliski. Photo tourism: Exploring photo collections in 3d. *SIGGRAPH*, pages 835–846, 2006.
- [3] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [4] Pierre Moulon, Pascal Monasse, Renaud Marlet, and Others. Openmvg. <https://github.com/openMVG/openMVG>.
- [5] Tomáš Pajdla. Elements of geometry for robotics. 2019.
- [6] David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 2004, 2004.
- [7] A. Fusiolo E. Maset, F. Arrigoni. Practical and efficient multi-view matching. *ICCV*, 2017.
- [8] M. A. Fischler and R. C. Bolles. Random sample consensus: A paradigm for model fitting with apphcatlons to image analysis and automated cartography. *ACM 0001-0782/81/0600-0381 00.75*, 1981.
- [9] R. I. Hartley B. Triggs, P. F. McLauchlan and A. W. Fitzgibbon. Bundle adjustment - a modern synthesis. *Proceedings of the International Workshop on Vision Algorithms: Theory and Practice*, 2000.
- [10] T. E. Boult and L. G. Brown. Factorization-based segmentation of motions. *Visual Motion*, 1991.

- [11] J. Costeira and T. Kanade. A multibody factorization method for independently moving objects. *Int. Journal of Computer Vision*, vol. 29, no. 3, 1998.
- [12] S. Soatto R. Vidal, Y. Ma and S. Sastry. Two-view multibody structure from motion. *IJCV*, 68(1):725,, 2006.
- [13] R. Vidal S. Rao, R. Tron and Y. Ma. Motion segmentation in the presence of outlying, incomplete, or corrupted trajectories. *PAMI*, 32(10):18321845,, 2010.
- [14] S. S. Sastry S. R. Rao, A. Y. Yang and Y. Ma. Robust algebraic segmentation of mixed rigid-body and planar motions from two views. *IJCV*, 88(3):425-446, 2010.
- [15] M. Irani and P. Anandan. A unied approach to moving ob ject detection in 2d and 3d scenes. *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 6, pp. 577-589, 1998.
- [16] C. S. Kenney M. Zuliani and B. S. Manjunath. The multiransac algorithm and its application to detect planar homographies. *ICIP*, pages 153-156, 2005.
- [17] K. Schindler and D. Ozden. Two-view multibody structure-and-motion with outliers through model selection. *IEEE Transactions on Pat-tern Analysis and Machine Intelligence*, 28(6):983–995, 2006.
- [18] W. Zhang and J. Kořecká. Nonparametric estimation of multiple structures with outliers. *European Conference on Computer Vision*, volume 4358, pages 60–74, 2006.
- [19] R. Toldo and A. Fusiello. Robust multiple structures estimation with j-linkage. *European Conference on Computer Vision*, volume 5302, pages 537–547, 2008.
- [20] L. Magri and A. Fusiello. T-linkage: a continuous relaxation of j-linkage for multi-model fitting. *Computer Vision and Pattern Recognition*, pages 3954–3961, 2014.
- [21] Jin Yu Trung Thanh Pham, Tat-Jun Chin and David Suter. The random cluster model for robust geometric fitting. *Pattern Analysis and Machine Intelligence*, 36(8):1658– 1671, 2014.
- [22] L. Magri and A. Fusiello. Robust multiple model fitting with preference analysis and low-rank approximation. *BMVC*, 2015.
- [23] Z. Li X. Xu, L. F. Cheong. Motion segmentation by exploiting complementary geometric models. 2018.
- [24] C. Rabe A. Wedel J. Klappstein, T. Vaudrey and R. Klette. Moving object segmentation using optical flow and depth information. *PSIVT*, pages 611-623, 2008.

- [25] J. Ju H. Jung and J. Kim. Rigid motion segmentation using randomized voting. *CVPR*, pages 1210-1217, 2014.
- [26] V. Murino C. Rubino, M. Crocco and A. D. Bue. Semantic multi-body motion segmentation. *WACV*, pages 1145-1152, 2015.
- [27] Filip Šrajer. Image matching for dynamic scenes. <http://cmp.felk.cvut.cz/srajeffil/theses/filip-srajer-diploma-thesis.pdf>, 2016.
- [28] H Wang K Schindler, D Suter. A model-selection framework for multi-body structure-and-motion of image sequences. *International Journal of Computer Vision*, 79(2):159–177, 2008.
- [29] P. Hruby and T. Pajdla. Multi-body structure from motion experiments. 2019.
- [30] Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17 (4), 2007.
- [31] D. S. Blostein K. S. Arun, T. S. Huang. Least-squares fitting of two 3-d point sets. *IEEE transactions on pattern analysis and machine intelligence*, vol. PAMI-9, No. 5. September 1987, 1987.