

Degenerate Motions in Multicamera Cluster SLAM with Non-overlapping Fields of View

Michael J. Tribou^{a,*}, David W. L. Wang^b, Steven L. Waslander^a

^a*Department of Mechanical and Mechatronics Engineering, University of Waterloo, 200 University Avenue West, Waterloo, ON, Canada, N2L 3G1.*

^b*Department of Electrical and Computer Engineering, University of Waterloo, 200 University Avenue West, Waterloo, ON, Canada, N2L 3G1.*

Abstract

An analysis of the relative motion and point feature model configurations leading to solution degeneracy is presented, for the case of a Simultaneous Localization and Mapping system using multicamera clusters with non-overlapping fields-of-view. The SLAM optimization system seeks to minimize image space reprojection error and is formulated for a cluster containing any number of component cameras, observing any number of point features over two keyframes. The measurement Jacobian is transformed to expose a reduced-dimension representation such that the degeneracy of the system can be determined by the rank of a dense submatrix. A set of relative motions sufficient for degeneracy are identified for certain cluster configurations, independent of target model geometry. Furthermore, it is shown that increasing the number of cameras within the cluster and observing features across different cameras over the two keyframes reduces the size of the degenerate motion sets significantly.

Keywords: SLAM, Computer vision, Multicamera cluster, Non-overlapping FOV, Degeneracy analysis, Critical motions

1. Introduction

Precise robotic motion and manipulation tasks with respect to unknown target environments and objects require an accurate, real-time measurement of the relative position and orientation of the robot and target. Multicamera systems are often employed for robotic pose and target model estimation, as each camera is an inexpensive, light-weight, and passive device capable of collecting a large amount of environment information at high rates. Many researchers across different fields have investigated the use of cameras for the purpose of estimating motion and scene structure. As a result, many techniques using a variety of camera types and configurations have been detailed in the literature.

A camera cluster is composed of any number of simple perspective cameras mounted rigidly with respect to each other, as shown in Figure 1, including configurations in which their fields-of-view (FOV) are spatially

disjoint [1]. This arrangement makes effective use of the camera sensors to cover a large combined FOV with high resolution, and in general, is able to overcome the limitations of other camera configurations, such as scale and translation-rotation motion ambiguities [2]. Additionally, by arranging the cameras to look in many directions, the pose estimation is made more robust since when certain cameras do not see any point features suitable for tracking, the other cameras in the cluster can maintain the localization. In this scenario, camera arrangements with a smaller collective FOV may become lost causing the tracking operation to fail.

In order for any pose estimation system to operate successfully, the current state must be uniquely recoverable given the measurable outputs up to, and including the current time step. In the context of a multicamera cluster relative pose system, this means that the image measurements must contain sufficient information to recover the cluster motion and the target model parameters, including the proper global scale metric. Furthermore, the solution must be unique since convergence to a different configuration, which may also agree with the measurements, would likely result in failure of perception and control operations.

*Corresponding author at: University of Waterloo, E3X-4118 – 200 University Avenue West, Waterloo, ON, Canada, N2L 3G1. Telephone: +1-519-635-8971

Email addresses: mtribou@uwaterloo.ca (Michael J. Tribou), dwang@uwaterloo.ca (David W. L. Wang), stevenw@uwaterloo.ca (Steven L. Waslander)



Figure 1: An example camera cluster in which the three component cameras are rigidly-fixed with respect to each other.

When the multicamera cluster is configured such that there is little or no spatial FOV overlap between the component cameras, the sensitivity of the image measurements to the global scale of the reconstructed model is low, particularly around specific motion profiles known as critical motions [3]. When the relative motion of the cluster is at or near critical, the global scale of the solution is extremely difficult, if not impossible, to recover accurately. In the presence of measurement noise, the solution will converge to an incorrect scale value.

This work investigates the degenerate configurations when estimating the Simultaneous Localization and Mapping (SLAM) [4] system states for a calibrated multicamera cluster over two keyframes while observing a set of point features in each camera and using an iterative optimization or recursive filter-based approach, minimizing the image space reprojection error of point feature measurements. This includes Bundle Adjustment (BA) [5] schemes as well as recursive filters such as an extended Kalman filter [6]. The main contribution is the identification of configurations of motion and target model structure leading to non-unique SLAM solutions.

Determining the system configurations leading to solution degeneracy is closely related to the concept of observability in control systems. In the study of observability for nonlinear systems, the local weak observability of the system can be determined by calculating the observability rank condition about any point in the state space [7]. This involves checking the column rank of a matrix containing the partial derivatives with respect to the system states, for increasing orders of Lie derivatives of the measurement model with respect to the system dynamics. When the matrix has full column rank, the system is locally weakly observable about that point.

For a SLAM system using only the visual measure-

ments from the cluster cameras and a non-stationary target, the system does not have a model of the dynamics for the relative motion and therefore, only the zeroth-order Lie derivatives are non-zero. In this case, evaluating the observability rank condition is equivalent to checking the rank of the measurement Jacobian matrix, as will be done here in the degeneracy analysis in Section 4. If the system were to contain a model of the relative motion dynamics, and the extra information that comes with it, the higher-order Lie derivatives of the measurement model would contain non-zero terms and the added matrix rows would only increase the likelihood that the matrix has full column rank at any point in the state space. However, in this analysis, no such assumptions about the relative motion dynamics are made and the degenerate configurations arising from only using image measurements for a set of point features over two keyframes are identified.

The remainder of this paper is arranged as follows: Section 2 contains a review of the previous analyses for degenerate configurations of the multicamera cluster relative pose system; Section 3 presents the multicamera cluster SLAM system; the degenerate configurations of the pose estimation system are identified in Section 4; and finally, conclusions are drawn in Section 5.

2. Related Work

Previous analyses identifying cluster motions leading to degenerate system solutions have assumed that the five degrees of freedom describing relative orientation and translation direction of the cluster are known using the well-studied single camera ego-motion estimation techniques (e.g. [5]). These include the work of Kim *et al.* [8], and Clipp *et al.* [3] for camera clusters with two component cameras, as well as that of the authors [9] for clusters with three component cameras. Of interest are the conditions when the image measurements from the camera cluster are able to allow for estimation of the final degree of freedom, corresponding to the translation magnitude and therefore, global system scale. The analyses show that when each point feature is seen by only one of the two cameras at both keyframes, the global scale of the solution is recoverable only when the relative translational and rotational motion are both non-zero, and does not result in the optical centres of each camera moving in concentric arcs on circles with a common centre at the intersection of the baselines at each keyframe [3]. When a third non-collinear camera is added to the cluster, the set of degenerate motions is reduced to those which result in all the three cameras moving in parallel [9].

Analyses of degeneracies of the full SLAM solution for multicamera clusters have focused on those associated with solving the generalized camera relative pose problem, either linearly using the Generalized Essential Matrix (GEM) [2], or aligning imaging rays in space for minimal cases of camera poses and points [10]. Sturm [11], Stewenius *et al.* [10], and Mouragnon *et al.* [12] discuss some degenerate cases, but Kim and Kanade [13] provide the most complete analysis. They identify the following degenerate configurations for generalized cameras using the seventeen point method [2]:

1. All of the observation rays pass through one common point before and after the camera motion.
2. The camera centres are on a line before and after the motion.
3. Each corresponding ray pair passes through the same local point in the general camera frame before and after the motion.

For a camera cluster with non-overlapping FOV, it is possible that each component camera observes its own mutually exclusive set of feature points over the two keyframes. In this case, the system satisfies condition 3 and the solution to the seventeen point algorithm is always degenerate. However, it is known from previous results that in certain configurations, other solution methods are able to recover an accurate estimate of the motion and structure. Consequently, the seventeen point algorithm does not always recover a solution when one exists. This problem was noticed by Li *et al.* [14], who have since modified the algorithm for use with non-overlapping clusters, but the subsequent degeneracy analysis has not been carried out. More importantly, the degenerate configurations are specific to the linear method of estimation. In this work, the minimization of image-space reprojection error is considered and the configurations for which an optimization of this type will fail are identified in the subsequent analysis.

3. Multicamera Cluster Pose Estimation

3.1. Projective Geometry

The projective space \mathbb{P}^n (refer to Appendix A for a brief introduction) provides a convenient way of representing the camera measurement system in terms of homogeneous transformations and points [5]. It will sometimes be necessary to move between the respective real and projective space representation of points, and the following promotion and demotion operators are defined. The projective promotion operator $\tilde{\rho} : \mathbb{R}^n \rightarrow \mathbb{P}^n$

maps a point \mathbf{x} in the real vector space to its representation in the projective space,

$$\tilde{\rho}(\mathbf{x}) = \begin{bmatrix} \mathbf{x}^\top & 1 \end{bmatrix}^\top. \quad (1)$$

The projective demotion operator $\pi_n : \mathbb{P}^n \rightarrow \mathbb{R}^n$ maps a point $\tilde{\mathbf{x}}$ in the projective space back to the corresponding point \mathbf{x} the real vector space,

$$\begin{aligned} \mathbf{x} &= \pi_n(\tilde{\mathbf{x}}) \\ &= \begin{cases} \text{undefined} & \text{if } x_{n+1} = 0 \\ \left[\frac{\tilde{x}_1}{\tilde{x}_{n+1}} \quad \frac{\tilde{x}_2}{\tilde{x}_{n+1}} \quad \dots \quad \frac{\tilde{x}_n}{\tilde{x}_{n+1}} \right]^\top & \text{if } x_{n+1} \neq 0. \end{cases} \end{aligned} \quad (2)$$

Note that the result of this operator is undefined for points at infinity.

In this work, unless it is ambiguous from the context, the promotion and demotion operators will be implied by the vector notation. The homogeneous coordinates for a given vector $\mathbf{x} \in \mathbb{R}^n$ will simply be written as $\tilde{\mathbf{x}} \in \mathbb{P}^n$, but implicitly, $\tilde{\mathbf{x}} \equiv \tilde{\rho}(\mathbf{x})$, and likewise, $\mathbf{x} \equiv \pi_n(\tilde{\mathbf{x}})$ assuming $\tilde{x}_{n+1} \neq 0$.

3.2. Pin-hole Camera Model

An individual component camera within the cluster is modelled as a simple pin-hole imaging device, which maps 3D points onto a 2D plane called the image plane [15]. An example is shown in Figure 2. A 3D point $\tilde{\mathbf{p}}^{C_i} = [x^{C_i} \ y^{C_i} \ z^{C_i} \ 1]^\top$, represented in the projective space \mathbb{P}^3 , and expressed with respect to the i^{th} camera coordinate frame, C_i , is projected onto the image plane I_i . The intersection of the point feature ray $\tilde{\mathbf{p}}^{C_i}$, through the optical centre, \mathbf{o}_i , with the image plane occurs at the point, $[u \ v]^\top \in \mathbb{R}^2$. It is assumed that each camera has been intrinsically calibrated using one of the many existing offline techniques [5], such that the measurements are made to match the structure shown.

The camera projection matrix, κ_i , maps the point in \mathbb{P}^3 into \mathbb{P}^2 on the image plane. It is assumed in this work, without loss of generality, that the projection matrices for all of the cameras have the form,

$$\kappa_i = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}. \quad (4)$$

The point $\tilde{\mathbf{p}}^{C_i}$ is projected into \mathbb{P}^2 on the image plane,

$$\tilde{\mathbf{p}}^{I_i} = \kappa_i \tilde{\mathbf{p}}^{C_i}, \quad (5)$$

then subsequently mapped to the actual image plane coordinates in \mathbb{R}^2 through the demotion function π_2 for

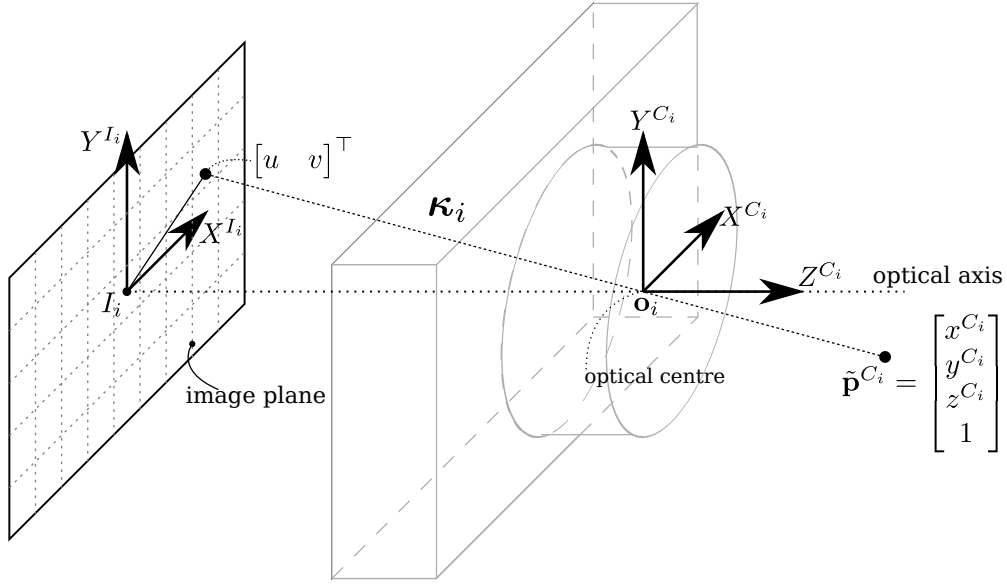


Figure 2: A simple pin-hole camera measurement model is used to relate the camera frame coordinates to the camera image plane coordinates for a feature point.

\mathbb{P}^2 ,

$$\pi_2(\tilde{\mathbf{p}}^{I_i}) = \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} -x^{C_i} \\ z^{C_i} \\ -y^{C_i} \\ z^{C_i} \end{bmatrix}, \quad z^{C_i} \neq 0. \quad (6)$$

Each camera is assumed to have an FOV strictly less than 180 degrees and therefore, is only able to observe points in front of the lens so every point is constrained to have a positive z-axis coordinate,

$$z^{C_i} > 0, \quad (7)$$

which satisfies (6).

3.3. Calibrated Multicamera Cluster

Collectively, the calibrated camera cluster is modelled as a set of n_c component pin-hole cameras with known relative coordinate transformations between each camera coordinate frame. Accordingly, a point $\tilde{\mathbf{p}}^{C_h}$ in the camera frame C_h , can be transformed into any other camera frame C_i by,

$$\tilde{\mathbf{p}}^{C_i} = \mathbf{T}_{C_h}^{C_i} \tilde{\mathbf{p}}^{C_h} \quad (8)$$

where $\mathbf{T}_{C_h}^{C_i} \in SE(3)$, $\forall i, h \in \{1, \dots, n_c\}$, is a homogeneous transformation matrix in $SE(3)$ [16]. Without loss of generality, the coordinate frame for the camera cluster is chosen to coincide with the first camera frame, C_1 . The transformation from camera h to the cluster

frame can be written in shortened form as $\mathbf{T}_{C_h} \equiv \mathbf{T}_{C_h}^{C_1}$, where the cluster frame C_1 is implied when the superscript is neglected. The transformation is shown in Figure 3.

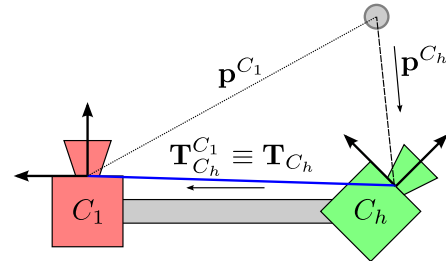


Figure 3: The relative position and orientation of each camera is known relative to the cluster frame, C_1 and therefore, the position of points in any camera frame can be found with respect to the cluster frame using the known transformation, \mathbf{T}_{C_h} .

3.4. Point Feature Target Object Model

The tracked target object or environment, henceforth referred to simply as the target, is a rigid body which contains a set of visible point features. A point feature is a visually distinguishable point on the tracked physical target that corresponds to a unique 3D position in a local target coordinate frame M , and is measurable in

a set of camera images through a relative motion sequence. Image measurements of these point features are extracted from the images using image processing techniques, including feature extraction algorithms like the FAST corner detector [17, 18], the Scale-Invariant Feature Transform (SIFT) [19], or Speeded-Up Robust Features (SURF) [20].

The target model point features are organized into n_k keyframes, each a six degree of freedom pose with respect to the target model reference frame M , along with the n_c images from the cluster cameras captured at that location, as in [21] for a single camera. The coordinate frame of camera h at keyframe k is denoted $C_h K_k$.

Since the relative position and orientation of each camera within the cluster is fixed at all times, the k^{th} keyframe pose is parameterized by the single homogeneous transformation for the cluster coordinate frame at the keyframe, $C_1 K_k$, with respect to the target model reference frame, M , resulting in $\mathbf{T}_{C_1 K_k}^M \in SE(3)$. The C_1 and M frames are applied universally in this keyframe pose definition, and therefore, the transformation will be written simply as $\mathbf{T}_{K_k} \equiv \mathbf{T}_{C_1 K_k}^M$. The pose of camera h at keyframe k is easily found as,

$$\mathbf{T}_{C_h K_k}^M = \mathbf{T}_{K_k} \mathbf{T}_{C_h}. \quad (9)$$

The position of the j^{th} point feature is parameterized by the azimuth and altitude angles of the vector from the origin of the anchor camera coordinate frame through the feature, $\boldsymbol{\mu}_j = [\phi_j, \theta_j]^T$ where $\phi_j, \theta_j \in (-\frac{\pi}{2}, \frac{\pi}{2})$. The depth along this bearing to the point feature, is the value $s_j \in \mathbb{R}^+$. The bearing angles are used to form the unit vector in the camera coordinate frame at the first keyframe,

$$\hat{\mathbf{p}}_j^{h,1} = \begin{bmatrix} \sin \phi_j \cos \theta_j \\ -\sin \theta_j \\ \cos \phi_j \cos \theta_j \end{bmatrix}, \quad (10)$$

and the point feature position is along this bearing at the distance s_j ,

$$\mathbf{p}_j^{h,1} = s_j \hat{\mathbf{p}}_j^{h,1}. \quad (11)$$

An example system with a camera cluster composed of $n_c = 2$ cameras is shown in Figure 4. The cameras in this example are arranged back-to-back with the optical axes looking outwards along the green axes of the associated coordinate frames. The j^{th} point feature is anchored in the second camera at the first keyframe, $C_2 K_1$, and its position with respect to this coordinate frame is represented as $\mathbf{p}_j^{C_2 K_1}$.

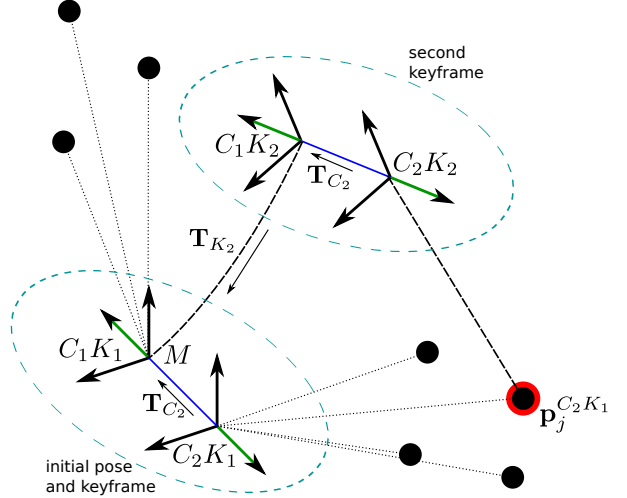


Figure 4: An example target object model with two keyframes for a two-camera back-to-back cluster. The cameras look outwards with the green arrows showing the optical axes. The point feature j is anchored, and therefore, positioned within the $C_2 K_1$ coordinate frame. The relative pose of camera 2, \mathbf{T}_{C_2} , is known from calibration, but the relative pose of keyframe 2, \mathbf{T}_{K_2} , as well as the position of the point features must be estimated.

The parameters representing the poses of the keyframes, together with the positions of the point features observed within them, compose the target model, as well as the full system state. These parameters are estimated using the point feature image measurements within the cluster cameras.

3.5. Multicamera Cluster SLAM System

This work considers the motion and structure estimation for a cluster of n_c cameras observing a set of n_f point features over two keyframes. The cameras within the cluster are arranged with little or no overlap in their FOV where each point feature in the target model is visible in only one camera at the first keyframe. Without loss of generality, the target model frame is chosen to coincide with the pose of the first keyframe $M \equiv C_1 K_1$. This results in the keyframe transformation becoming the identity,

$$\mathbf{T}_{K_1} = \mathbf{I}_{4 \times 4}. \quad (12)$$

Several further assumptions about the system are made to facilitate the analysis in the subsequent sections:

Assumption 3.1. Each point feature is observed and measured by only one of the component cameras at the

first keyframe. The position of the point feature is expressed with respect to the coordinate frame for that camera at the first keyframe. The coordinate frame in which the point feature is parameterized is referred to as the anchor keyframe and camera coordinate frame.

Assumption 3.2. Each point feature is observed and measured by one or more of the component cameras at the second keyframe. At least one of the observations is by a camera for which its motion is not collinear with the initial bearing to the point feature at the first keyframe. The camera and keyframe in which the observation occurs is called the observing keyframe and camera coordinate frame.

Assumption 3.3. The point feature positions and keyframe poses are arranged such that if a camera observes a point feature, the feature position expressed in the observing camera coordinate frame has a finite positive non-zero z-axis component, $0 < z < \infty$.

For Assumption 3.1, the function $h : \{1, \dots, n_f\} \rightarrow \{1, \dots, n_c\}$ maps the point feature index to the anchor camera index. As a result, the anchor camera for the j^{th} point feature is camera $h(j)$. In the following, when it is obvious from the context, the anchor camera index will be written in the shortened-form by dropping the argument, $h \equiv h(j)$.

Similarly for Assumption 3.2, the j^{th} point feature is observed and measured by $n_o(j) \in \mathbb{N}^+$ cameras at the second keyframe. The indices of the observing cameras are found using the function $i : \{1, \dots, n_f\} \times \{1, \dots, n_o(j)\} \rightarrow \{1, \dots, n_c\}$, such that the k^{th} observation of the j^{th} point feature at the second keyframe is measured by camera $i(j, k)$. Once again, this will be shortened to exclude the feature index j and observation index k when it is implied by the context, $i \equiv i(k) \equiv i(j, k)$.

The motion of the camera cluster between the two keyframes is parameterized by six values describing the relative translation and orientation of the first keyframe with respect to the second keyframe. The translation parameters, t_x, t_y, t_z , form the relative translation vector, $\mathbf{t}_K = [t_x, t_y, t_z]^T$, and the rotation parameters, $\omega_x, \omega_y, \omega_z$, form the relative rotation matrix, $\mathbf{R}_K \in SO(3)$. Together, the rotation and translation form the transformation, $\mathbf{T}_K \in SE(3)$.

The resulting state vector, $\mathbf{x} \in \mathbb{R}^n$, where $n = 6 + 3n_f$, is composed of the parameters for the n_f point features, along with the relative translation and orientation states

for the cluster motion between the keyframes,

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix}, \quad (13)$$

where

$$\mathbf{x}_1 = [s_1, \dots, s_{n_f}]^T \in \mathbb{R}_+^{n_f}, \quad (14)$$

are the radial distances to the point features,

$$\mathbf{x}_2 = [\mathbf{t}_K^T, \boldsymbol{\omega}_K^T]^T \in \mathbb{R}^6, \quad (15)$$

are the relative position and orientation of the first keyframe with respect to the second keyframe and,

$$\mathbf{x}_3 = [\boldsymbol{\mu}_1^T, \boldsymbol{\mu}_2^T, \dots, \boldsymbol{\mu}_{n_f}^T]^T \in \mathbb{R}^{2n_f}, \quad (16)$$

are the bearings to the point features in their respective anchor camera coordinate frames. This state order has been specifically chosen in order to facilitate the analysis of the degeneracies of the solution presented in Section 4.

3.6. Camera Cluster Measurement Model

The system measurement vector, $\mathbf{z} \in \mathbb{R}^m$, is formed by stacking the image plane coordinates of all of the point feature observations in all cameras at both keyframes, where $m = 2(n_f + m_o)$ with

$$m_o = \sum_{j=1}^{n_f} n_o(j), \quad (17)$$

as the total number of observations of all of the point features at the second keyframe.

The measurement model, relating the observed point feature locations in the camera image planes, to the system states, can be written as a series of coordinate transformations. Suppose that the j^{th} point feature, anchored in the coordinate frame $C_h K_1$, is measured by camera i at $C_i K_2$. An example of this chain of transformations is shown for the simple back-to-back two-camera cluster system in Figure 4. In this particular case, the point feature j is anchored in $C_2 K_1$ and observed in $C_2 K_2$.

The point feature position parameters give the location of the j^{th} feature in its anchor keyframe and camera frame $C_h K_1$, resulting in $\mathbf{p}_j^{h,1}$. This point feature is first transformed into the target model coordinate frame by,

$$\tilde{\mathbf{p}}_j^M = \mathbf{T}_{K_1} \mathbf{T}_{C_h} \tilde{\mathbf{p}}_j^{h,1} \quad (18)$$

$$= \mathbf{T}_{C_h} \tilde{\mathbf{p}}_j^{h,1}, \quad (19)$$

which is the transformation provided by the known cluster calibration.

The point feature position, with position estimate expressed in the target model reference frame, is transformed into the coordinate frame of the observing keyframe and camera $C_i K_\ell$ using the relative keyframe pose transformation, \mathbf{T}_{K_ℓ} , and the cluster calibration,

$$\tilde{\mathbf{p}}_j^{i,\ell} = \begin{bmatrix} x_j^{i,\ell} & y_j^{i,\ell} & z_j^{i,\ell} & 1 \end{bmatrix}^\top \quad (20)$$

$$= (\mathbf{T}_{C_i})^{-1} (\mathbf{T}_{K_\ell})^{-1} \tilde{\mathbf{p}}_j^M \quad (21)$$

$$= (\mathbf{T}_{C_i})^{-1} (\mathbf{T}_{K_\ell})^{-1} \mathbf{T}_{C_h} \tilde{\mathbf{p}}_j^{h,1}. \quad (22)$$

Finally, the point is projected into \mathbb{P}^2 and onto the image plane of camera C_i using the corresponding projection matrix, κ_i ,

$$\tilde{\mathbf{u}}_j^{i,\ell} = \begin{bmatrix} u_x & u_y & u_z \end{bmatrix}^\top \quad (23)$$

$$= \kappa_i \tilde{\mathbf{p}}_j^{i,\ell} \quad (24)$$

$$= \begin{bmatrix} -x_j^{i,\ell} \\ -y_j^{i,\ell} \\ z_j^{i,\ell} \end{bmatrix}, \quad (25)$$

which leads to the resulting measurement vector $\mathbf{z}_j^{i,\ell} \in \mathbb{R}^2$ and mapping $\mathbf{g}_j^{i,\ell} : \mathbb{R}^n \rightarrow \mathbb{R}^2$ for the observation of point feature j in camera i at keyframe ℓ ,

$$\mathbf{z}_j^{i,\ell} = \mathbf{g}_j^{i,\ell}(\mathbf{x}) = \pi_2(\tilde{\mathbf{u}}_j^{i,\ell}). \quad (26)$$

Each of the intermediate transformations in (22) can be represented by a rotation matrix and translation vector,

$$\mathbf{T}_{C_h} = \begin{bmatrix} \mathcal{R}_{C_h} & \mathbf{t}_{C_h} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \quad (27)$$

$$\mathbf{T}_{K_\ell} = \begin{bmatrix} \mathcal{R}_{K_\ell} & \mathbf{t}_{K_\ell} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} = \begin{cases} \begin{bmatrix} \mathbf{I}_{4 \times 4} & \\ & \end{bmatrix}, & \ell = 1 \\ \begin{bmatrix} \mathcal{R}_K^\top & -\mathcal{R}_K^\top \mathbf{t}_K \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}, & \ell = 2 \end{cases} \quad (28)$$

$$\mathbf{T}_{C_i} = \begin{bmatrix} \mathcal{R}_{C_i} & \mathbf{t}_{C_i} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}. \quad (29)$$

When (27)–(29) are substituted into (22) along with (11), the coordinates of the point feature position in \mathbb{R}^3 become,

$$\mathbf{p}_j^{i,\ell} = s_j \mathcal{R}_{C_i}^\top \mathcal{R}_{K_\ell}^\top \mathcal{R}_{C_h} \hat{\mathbf{p}}_j^{h,1} - \mathcal{R}_{C_i}^\top \mathcal{R}_{K_\ell}^\top \mathbf{t}_{K_\ell} + \mathcal{R}_{C_i}^\top \mathcal{R}_{K_\ell}^\top \mathbf{t}_{C_h} - \mathcal{R}_{C_i}^\top \mathbf{t}_{C_i} \quad (30)$$

$$= \mathcal{R}_{C_i}^\top (s_j \mathcal{R}_{K_\ell}^\top \mathcal{R}_{C_h} \hat{\mathbf{p}}_j^{h,1} - \mathcal{R}_{K_\ell}^\top \mathbf{t}_{K_\ell} + \mathcal{R}_{K_\ell}^\top \mathbf{t}_{C_h} - \mathbf{t}_{C_i}) \quad (31)$$

$$= \mathcal{R}_{C_i}^\top \mathbf{q}_j^{i,\ell}, \quad (32)$$

where

$$\mathbf{q}_j^{i,\ell} = s_j \hat{\mathbf{a}}_{j,\ell} + \mathbf{b}_\ell + \mathbf{c}_{h,\ell} + \mathbf{d}_i \quad (33)$$

with

$$\hat{\mathbf{a}}_{j,\ell} = \mathcal{R}_{K_\ell}^\top \mathcal{R}_{C_h} \hat{\mathbf{p}}_j^{h,1} \quad (34)$$

$$\mathbf{b}_\ell = -\mathcal{R}_{K_\ell}^\top \mathbf{t}_{K_\ell} \quad (35)$$

$$\mathbf{c}_{h,\ell} = \mathcal{R}_{K_\ell}^\top \mathbf{t}_{C_h} \quad (36)$$

$$\mathbf{d}_i = -\mathbf{t}_{C_i} \quad (37)$$

and

$$\mathcal{R}_{C_i} = [\hat{\mathbf{n}}_{i,x} \quad \hat{\mathbf{n}}_{i,y} \quad \hat{\mathbf{n}}_{i,z}] \in SO(3), \quad (38)$$

where $\hat{\mathbf{n}}_{i,x}$, $\hat{\mathbf{n}}_{i,y}$, and $\hat{\mathbf{n}}_{i,z}$ are the orthonormal basis vectors for the observing camera i frame with respect to the camera 1 coordinate frame. An example system consisting of the cameras observing point features over two keyframes is shown in Figure 5 with the intermediate variables labelled.

The set of camera observation vectors for point feature j is defined as the displacements between the anchor camera coordinate frame at the first keyframe, $C_h K_1$, and the centres of each of the observing cameras at the second keyframe, $C_{i(k)} K_2$, $\forall k \in \{1, \dots, n_o(j)\}$. The set of vectors are,

$$V = \{\mathbf{v}_{\alpha,\beta} \in \mathbb{R}^3 | \alpha, \beta \in \mathbb{N}^+, \alpha \leq n_f \text{ and } \beta \leq n_o(\alpha)\}. \quad (39)$$

Therefore, if it is included in the set V , the camera observation vector can be written as,

$$\mathbf{v}_{\alpha,\beta} = \mathbf{b}_2 + \mathbf{c}_{h(\alpha),2} + \mathbf{d}_{i(\beta)}. \quad (40)$$

and is illustrated in the example system shown in Figure 6.

The image coordinates of each individual point feature observation measurements are then compiled into a vector for point feature j containing all of the individual observations of that feature at both keyframes,

$$\mathbf{z}_j = \begin{bmatrix} \mathbf{z}_j^{i(1),2} \\ \vdots \\ \mathbf{z}_j^{i(n_o),2} \\ \mathbf{z}_j^{h,1} \end{bmatrix} \in \mathbb{R}^{2+2n_o}, \quad (41)$$

where $n_o \equiv n_o(j)$ is the number of observations of the j^{th} point feature at the second keyframe.

The full system measurement vector is composed of the observations of the n_f point features at both keyframes,

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_{n_f} \end{bmatrix} \in \mathbb{R}^m. \quad (42)$$

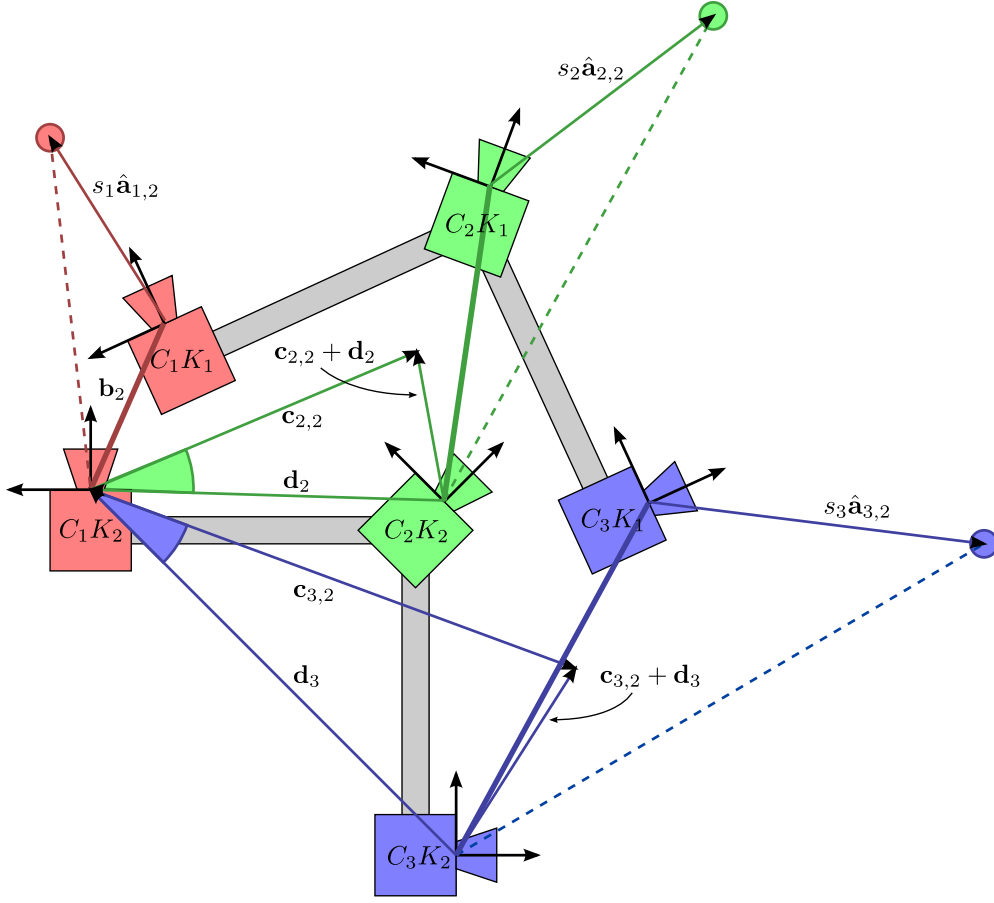


Figure 5: An example three-camera cluster observing point features over two keyframes with intermediate vectors labelled. Vectors are parameterized with reference to keyframe C_1K_2 .

4. Degeneracy Analysis

4.1. Solution Degeneracies

Typical optimization methods attempt to minimize a nonlinear cost function, $c : \mathbb{R}^n \rightarrow \mathbb{R}$, and determine the optimal state vector estimate, $\check{\mathbf{x}}^* \in \mathbb{R}^n$, such that,

$$\check{\mathbf{x}}^* = \arg \min_{\mathbf{x}} c(\mathbf{x}). \quad (43)$$

The optimization proceeds iteratively, starting with an initial state estimate, $\check{\mathbf{x}}_0 \in \mathbb{R}^n$. Each iteration seeks to update the current state estimate, $\check{\mathbf{x}}_k$, with a vector $\delta_k \in \mathbb{R}^n$,

$$\check{\mathbf{x}}_{k+1} = \check{\mathbf{x}}_k + \delta_k, \quad (44)$$

such that the sequence $\{\check{\mathbf{x}}_0, \check{\mathbf{x}}_1, \check{\mathbf{x}}_2, \dots\} \rightarrow \check{\mathbf{x}}^*$. In this analysis, a cost function relating to sum of squared measurement error is assumed,

$$c(\check{\mathbf{x}}_k) = \frac{1}{2} \bar{\mathbf{z}}_k^\top \bar{\mathbf{z}}_k, \quad (45)$$

where $\bar{\mathbf{z}}_k = \mathbf{z} - \mathbf{g}(\check{\mathbf{x}}_k) \in \mathbb{R}^m$ is the measurement error vector at iteration k . A commonly-used method for BA is the Levenberg-Marquardt (LM) method [5], although other optimization methods may be used such as Gauss-Newton, gradient descent, or Newton step. Each of these optimization methods can operate using the sum of squared reprojection error and the parameter update δ_k is defined as the solution to,

$$\mathbf{N} \delta_k = \mathbf{J}^\top \bar{\mathbf{z}}_k \quad (46)$$

where \mathbf{N} is the normal matrix, which varies by optimization method, and \mathbf{J} is the measurement Jacobian such that,

$$\mathbf{J} = \left. \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\check{\mathbf{x}}_k} \in \mathbb{R}^{m \times n}. \quad (47)$$

Solving for δ_k , the solution becomes,

$$\delta_k = \mathbf{N}^{-1} \mathbf{J}^\top \bar{\mathbf{z}}_k, \quad (48)$$

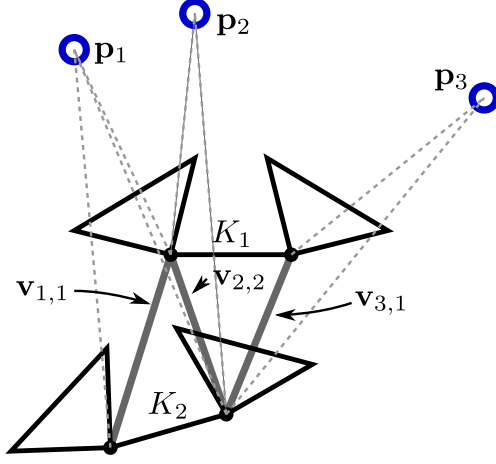


Figure 6: An example two-camera cluster observing three features. Each feature is observed in one camera at the first keyframe, but some are seen by different cameras at the second keyframe. The camera observation vectors, $\mathbf{v}_{1,1}$, $\mathbf{v}_{2,2}$, $\mathbf{v}_{3,1}$ link the cameras which see the same feature at the different keyframes.

where $\bar{\mathbf{z}}_k = \mathbf{z} - \mathbf{g}(\check{\mathbf{x}}_k)$ is the measurement error. A unique δ_k can be found as long as the matrix \mathbf{N} is invertible and the Jacobian has full rank. Therefore, the system in (46) is degenerate when,

$$\text{rank}(\mathbf{J}) < n, \quad (49)$$

and the solution is under-constrained. The cases where the system becomes degenerate are the focus of this work and are investigated in the next section.

4.2. Identification of Degenerate Configurations

This section identifies the configurations of the camera cluster geometry, relative motion, and target model structure, for which the Jacobian \mathbf{J} falls below full column rank. It is shown that for the assumptions stated previously, the $m \times n$ measurement Jacobian matrix \mathbf{J} has full rank if and only if a $m_o \times 6$ matrix, \mathbf{M}_2 has full rank.

To determine the rank of the Jacobian matrix, the structure of the sub-blocks formed for the point feature observations is investigated. Each point feature j is observed by only one camera at the first keyframe. This

observation adds two rows to the Jacobian,

$$\mathbf{J}_j^{h,1} = \left. \frac{\partial \mathbf{g}_j^{h,1}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\check{\mathbf{x}}} \quad (50)$$

$$= \begin{bmatrix} \mathbf{0}_{2 \times n_f + 6} & \mathbf{0}_{2 \times 2(j-1)} & \left. \frac{\partial \mathbf{g}_j^{h,1}(\mathbf{x})}{\partial \boldsymbol{\mu}_j} \right|_{\mathbf{x}=\check{\mathbf{x}}} & \mathbf{0}_{2 \times 2(n_f-j)} \end{bmatrix}, \quad (51)$$

where the only non-zero elements are in the 2×2 block relating the measurement coordinates to the point feature bearing states,

$$\left. \frac{\partial \mathbf{g}_j^{h,1}(\mathbf{x})}{\partial \boldsymbol{\mu}_j} \right|_{\mathbf{x}=\check{\mathbf{x}}} = \begin{bmatrix} -\tan(\phi_j)^2 - 1 & 0 \\ \frac{\tan(\phi_j) \tan(\theta_j)}{\cos(\phi_j)} & \frac{1}{\cos(\phi_j) \cos(\theta_j)^2} \end{bmatrix}. \quad (52)$$

The non-zero element structure of these rows are shown in Figure 7.

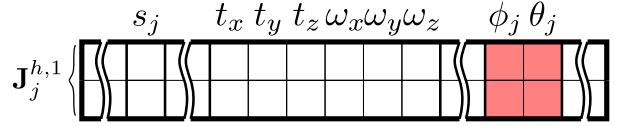


Figure 7: Structure of the Jacobian rows for the measurement of the point feature j in its anchor camera at the first keyframe. Non-zero elements are shown as the shaded cells. The columns not related to feature j contain only zeros and have been removed for conciseness.

Each point feature j is also observed and measured by at least one camera at the second keyframe. The Jacobian matrix rows corresponding to the k^{th} observation of point feature j in the second keyframe with camera $i(k)$ are the partial derivatives of the measurement equation (26) with respect to the system states,

$$\mathbf{J}_j^{i(k),2} = \left. \frac{\partial \mathbf{g}_j^{i(k),2}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\check{\mathbf{x}}}, \quad (53)$$

$$= \left(\frac{\partial \pi_2(\tilde{\mathbf{u}}_j^{i(k),2})}{\partial \tilde{\mathbf{u}}_j^{i(k),2}} \right) \left(\frac{\partial \tilde{\mathbf{u}}_j^{i(k),2}}{\partial \mathbf{x}} \right) \bigg|_{\mathbf{x}=\check{\mathbf{x}}}, \quad (54)$$

using the chain rule. Dropping the implied $\mathbf{x} = \check{\mathbf{x}}$, the first term evaluates to,

$$\frac{\partial \pi_2(\tilde{\mathbf{u}}_j^{i(k),2})}{\partial \tilde{\mathbf{u}}_j^{i(k),2}} = \frac{1}{(u_z)^2} \begin{bmatrix} u_z & 0 & -u_x \\ 0 & u_z & -u_y \end{bmatrix} \quad (55)$$

$$= \frac{1}{(z_j^{i(k),2})^2} \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \end{bmatrix} [\tilde{\mathbf{u}}_j^{i(k),2}]_{\mathbf{x}}, \quad (56)$$

where $[\mathbf{a}]_{\times}$ is the skew-symmetric matrix such that $[\mathbf{a}]_{\times} \mathbf{b} = \mathbf{a} \times \mathbf{b}$, $\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^3$.

Substituting (56) back into (54) and recognizing that,

$$\tilde{\mathbf{u}}_j^{i(k),2} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathcal{R}_{C_{i(k)}}^{\top} \mathbf{q}_j^{i(k),2}, \quad (57)$$

the Jacobian rows can be written as,

$$\mathbf{J}_j^{i(k),2} = \frac{1}{(z_j^{i(k),2})^2} \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \end{bmatrix} [\tilde{\mathbf{u}}_j^{i(k),2}]_{\times} \frac{\partial \tilde{\mathbf{u}}_j^{i(k),2}}{\partial \mathbf{x}}, \quad (58)$$

$$= \frac{1}{(z_j^{i(k),2})^2} \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \mathcal{R}_{C_{i(k)}}^{\top} [\mathbf{q}_j^{i(k),2}]_{\times} \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{x}}, \quad (59)$$

$$= \frac{1}{(z_j^{i(k),2})^2} \begin{bmatrix} -(\hat{\mathbf{n}}_{i(k),y} \times \mathbf{q}_j^{i(k),2})^{\top} \\ (\hat{\mathbf{n}}_{i(k),x} \times \mathbf{q}_j^{i(k),2})^{\top} \end{bmatrix} \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{x}}, \quad (60)$$

where the partial derivatives of the point feature position at the second keyframe with respect to the system states are written,

$$\frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{x}} = \begin{bmatrix} \mathbf{0}_{3 \times (j-1)} & \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial s_j} & \mathbf{0}_{3 \times (n_f-j)} & \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{t}_K} & \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \boldsymbol{\omega}_K} \\ & & \mathbf{0}_{3 \times 2(j-1)} & \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \boldsymbol{\mu}_j} & \mathbf{0}_{3 \times 2(n_f-j)} \end{bmatrix} \quad (61)$$

with the position change with respect to the radial distance,

$$\frac{\partial \mathbf{q}_j^{i(k),2}}{\partial s_j} = \hat{\mathbf{a}}_{j,2}, \quad (62)$$

the translation between the keyframes,

$$\frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{t}_K} = \mathbf{I}_{3 \times 3}, \quad (63)$$

the rotation between the keyframes,

$$\frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \boldsymbol{\omega}_K} = -[s_j \hat{\mathbf{a}}_{j,2} + \mathbf{c}_{h,2}]_{\times}, \quad (64)$$

and the initial bearings to the point feature,

$$\frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \boldsymbol{\mu}_j} = s_j \mathcal{R}_K \mathcal{R}_{C_h} \begin{bmatrix} \cos \phi_j \cos \theta_j & -\sin \phi_j \sin \theta_j \\ 0 & -\cos \theta_j \\ -\sin \phi_j \cos \theta_j & -\cos \phi_j \sin \theta_j \end{bmatrix}. \quad (65)$$

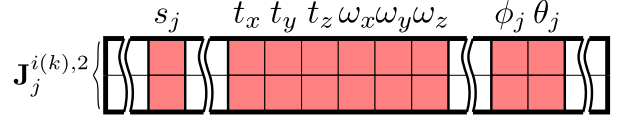


Figure 8: Structure of the Jacobian rows for an observation of the point feature j at the second keyframe.

The structure of the Jacobian rows associated with the observations of point feature j at the second keyframe is shown in Figure 8.

The full measurement Jacobian is formed by stacking all of the observations of all of the point features at both keyframes,

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_1^{i(1,1),2} \\ \vdots \\ \mathbf{J}_1^{i(1,n_o(1)),2} \\ \mathbf{J}_1^{h(1),1} \\ \vdots \\ \mathbf{J}_{n_f}^{i(n_f,1),2} \\ \vdots \\ \mathbf{J}_{n_f}^{i(n_f,n_o(n_f)),2} \\ \mathbf{J}_{n_f}^{h(n_f),1} \end{bmatrix}. \quad (66)$$

The configurations for which this measurement Jacobian possesses full rank can be identified by checking the rank of a reduced-dimension matrix, as shown in the following Lemma.

Lemma 4.1. *For a multicamera cluster SLAM system satisfying Assumptions 3.1, 3.2, and 3.3, the rank of the measurement Jacobian matrix \mathbf{J} in (66) is full if and only if the rank of the matrix,*

$$\mathbf{M}_2 = \begin{bmatrix} \begin{bmatrix} -\hat{\mathbf{a}}_{1,2}^{\top} [\mathbf{v}_{1,1}]_{\times} & \hat{\mathbf{a}}_{1,2}^{\top} [\mathbf{v}_{1,1}]_{\times} [\mathbf{w}_1]_{\times} \\ \vdots & \vdots \\ -\hat{\mathbf{a}}_{1,2}^{\top} [\mathbf{v}_{1,n_o(1)}]_{\times} & \hat{\mathbf{a}}_{1,2}^{\top} [\mathbf{v}_{1,n_o(1)}]_{\times} [\mathbf{w}_1]_{\times} \end{bmatrix} \\ \vdots \\ \begin{bmatrix} -\hat{\mathbf{a}}_{n_f,2}^{\top} [\mathbf{v}_{n_f,1}]_{\times} & \hat{\mathbf{a}}_{n_f,2}^{\top} [\mathbf{v}_{n_f,1}]_{\times} [\mathbf{w}_{n_f}]_{\times} \\ \vdots & \vdots \\ -\hat{\mathbf{a}}_{n_f,2}^{\top} [\mathbf{v}_{n_f,n_o(n_f)}]_{\times} & \hat{\mathbf{a}}_{n_f,2}^{\top} [\mathbf{v}_{n_f,n_o(n_f)}]_{\times} [\mathbf{w}_{n_f}]_{\times} \end{bmatrix} \end{bmatrix}, \quad (67)$$

where

$$\mathbf{w}_j = s_j \hat{\mathbf{a}}_{j,2} + \mathbf{c}_{h(j),2},$$

is full.

Proof. The strategy is to first show that the columns of \mathbf{J} corresponding to the point feature positions $(s_j, \boldsymbol{\mu}_j)$ have full rank. Consequently, the only way for the Jacobian to have less than full rank is when the columns corresponding to the keyframe motion $(\mathbf{t}_k, \mathbf{R}_k)$ have rank less than six.

By Assumption 3.3, the position of the j^{th} point feature in its anchor camera and keyframe ensures that $\cos(\phi_j) > 0$ and $\cos(\theta_j) > 0$. As a result, the block (52) always has a rank of 2 since the determinant is non-zero,

$$\det \left(\frac{\partial \mathbf{g}_j^{h,1}(\mathbf{x})}{\partial \boldsymbol{\mu}_j} \bigg|_{\mathbf{x}=\tilde{\mathbf{x}}} \right) = \frac{-1}{\cos(\phi_j)^3 \cos(\theta_j)^2} \neq 0. \quad (68)$$

Therefore, it is possible to diagonalize the sub-block using elementary row and column operations without changing the rank of the matrix. After diagonalization, the new matrix rows, $\mathbf{K}_j^{h,1}$, have the structure shown in Figure 9. As a result, the columns corresponding to the bearing states $\boldsymbol{\mu}_j$ have full rank for all of the point features.

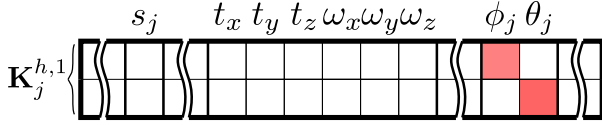


Figure 9: Structure of the Jacobian rows for the measurement of the point feature j in the first keyframe after diagonalizing the bearing sub-block.

Additionally, the columns associated with the point feature radial depth parameter s_j for an observation at the second keyframe contain only zeros when,

$$\begin{bmatrix} -\hat{\mathbf{n}}_{i(k),y}^\top \\ \hat{\mathbf{n}}_{i(k),x}^\top \end{bmatrix} [\mathbf{v}_{j,k}]_\times \hat{\mathbf{a}}_{j,2} = \mathbf{0}_{2 \times 1}, \quad (69)$$

the displacement between the anchor and observing camera coordinate frames is collinear with the initial bearing to the point feature in the anchor camera frame. In this case, there is no information about the depth of the feature within this measurement since the triangulation baseline has zero length. However, by Assumptions 3.2 and 3.3, there exists an observation $k \in \{1, \dots, n_o(j)\}$ such that,

$$\begin{bmatrix} s_x \\ s_y \end{bmatrix} = \begin{bmatrix} -\hat{\mathbf{n}}_{i(k),y}^\top \\ \hat{\mathbf{n}}_{i(k),x}^\top \end{bmatrix} [\mathbf{v}_{j,k}]_\times \hat{\mathbf{a}}_{j,2} \neq \mathbf{0}_{2 \times 1}, \quad (70)$$

and therefore, at least one non-zero element in the column. The matrix rows $\mathbf{J}_j^{i(k),2}$ are manipulated using the

row operations matrix,

$$\mathbf{O}_j^{i(k),2} = \begin{cases} \begin{bmatrix} \hat{\mathbf{n}}_{i(k),x}^\top [\mathbf{q}_j^{i(k),2}]_\times \hat{\mathbf{a}}_{j,2} & 0 \\ -\hat{\mathbf{n}}_{i(k),x}^\top [\mathbf{q}_j^{i(k),2}]_\times \hat{\mathbf{a}}_{j,2} & -\hat{\mathbf{n}}_{i(k),y}^\top [\mathbf{q}_j^{i(k),2}]_\times \hat{\mathbf{a}}_{j,2} \end{bmatrix}, & \text{if } s_x, s_y \neq 0 \\ \begin{bmatrix} 0 & 1 \\ \hat{\mathbf{n}}_{i(k),x}^\top [\mathbf{q}_j^{i(k),2}]_\times \hat{\mathbf{a}}_{j,2} & 0 \end{bmatrix}, & \text{if } s_x = 0, s_y \neq 0 \\ \begin{bmatrix} 1 & 0 \\ 0 & -\hat{\mathbf{n}}_{i(k),y}^\top [\mathbf{q}_j^{i(k),2}]_\times \hat{\mathbf{a}}_{j,2} \end{bmatrix}, & \text{if } s_x \neq 0, s_y = 0 \\ \mathbf{I}_{2 \times 2}, & \text{if } s_x, s_y = 0, \end{cases} \quad (71)$$

in order to achieve the desired structure $\mathbf{K}_j^{i(k),2}$ for the Jacobian rows,

$$\mathbf{K}_j^{i(k),2} = \mathbf{z}_j^{i(k),2} \mathbf{O}_j^{i(k),2} \mathbf{J}_j^{i(k),2}. \quad (72)$$

which is shown in Figure 10.

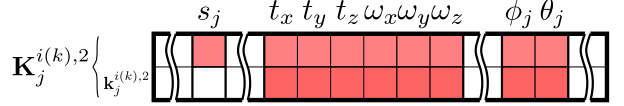


Figure 10: Element structure of the modified Jacobian rows associated with measuring point feature j at the second keyframe.

The second row of the matrix $\mathbf{K}_j^{i(k),2}$, labelled $\mathbf{k}_j^{i(k),2}$, becomes,

$$\mathbf{k}_j^{i(k),2} = [\mathbf{0} \quad k_{s_j} \quad \mathbf{0} \quad \mathbf{k}_{t_k} \quad \mathbf{k}_{\omega_k} \quad \mathbf{0} \quad \mathbf{k}_{\boldsymbol{\mu}_j} \quad \mathbf{0}] \quad (73)$$

$$= \frac{1}{z_j^{i(k),2}} (\hat{\mathbf{n}}_{i(k),z}^\top \mathbf{q}_j^{i(k),2}) (\mathbf{q}_j^{i(k),2} \times \hat{\mathbf{a}}_{j,2})^\top \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{x}} \quad (74)$$

$$= ([\mathbf{b}_2 + \mathbf{c}_{h,2} + \mathbf{d}_i]_\times \hat{\mathbf{a}}_{j,2})^\top \frac{\partial \mathbf{q}_j^{i(k),2}}{\partial \mathbf{x}}. \quad (75)$$

where the element associated with the radial distance parameter is zero,

$$k_{s_j} = ([\mathbf{b}_2 + \mathbf{c}_{h,2} + \mathbf{d}_i]_\times \hat{\mathbf{a}}_{j,2})^\top \hat{\mathbf{a}}_{j,2} \quad (76)$$

$$= (\mathbf{b}_2 + \mathbf{c}_{h,2} + \mathbf{d}_i)^\top [\hat{\mathbf{a}}_{j,2}]_\times \hat{\mathbf{a}}_{j,2} \quad (77)$$

$$= 0, \quad (78)$$

and the columns for the keyframe motion states are now,

$$\mathbf{k}_{t_k} = -\hat{\mathbf{a}}_{j,2}^\top [\mathbf{b}_2 + \mathbf{c}_{h,2} + \mathbf{d}_i]_\times \quad (79)$$

$$= -\hat{\mathbf{a}}_{j,2}^\top [\mathbf{v}_{j,k}]_\times, \quad (80)$$

and,

$$\mathbf{k}_{\omega_k} = \hat{\mathbf{a}}_{j,2}^\top [\mathbf{b}_2 + \mathbf{c}_{h,2} + \mathbf{d}_i]_\times [s_j \hat{\mathbf{a}}_{j,2} + \mathbf{c}_{h,2}]_\times \quad (81)$$

$$= \hat{\mathbf{a}}_{j,2}^\top [\mathbf{v}_{j,k}]_\times [\mathbf{w}_j]_\times. \quad (82)$$

All of the modified Jacobian matrix rows for the point feature j observations at both keyframes are then compiled into a single block,

$$\mathbf{K}_j = \begin{bmatrix} \mathbf{K}_j^{i(1),2} \\ \vdots \\ \mathbf{K}_j^{i(n_o),2} \\ \mathbf{K}_j^{h,1} \end{bmatrix} \quad (83)$$

which maintains the same rank as the original Jacobian block for the point feature,

$$\mathbf{J}_j = \left. \frac{\partial \mathbf{g}_j(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\hat{\mathbf{x}}}, \quad (84)$$

since the manipulations are performed by full rank elementary row and column operations matrices. The resulting matrix block has the structure shown in Figure 11.

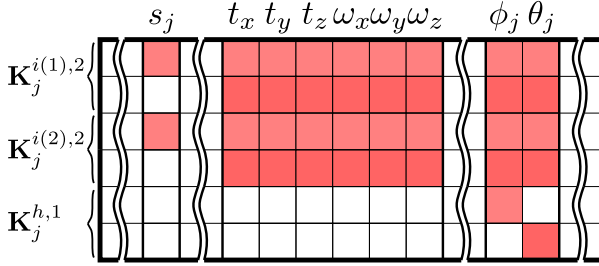


Figure 11: Structure of the manipulated Jacobian block for point feature j stacking all of the observations at both keyframes.

It can be shown that the odd-numbered rows of \mathbf{K}_j may always be written as a linear combination of the resulting even-numbered rows. Additionally, elementary row operations can eliminate all but the last elements in each of the columns associated with the point feature bearing states, ϕ_j and θ_j . Finally, the non-zero element in the s_j column can be used to eliminate the remaining non-zero elements in the first row. Subsequently, this new matrix, \mathbf{L}_j has the same rank as the original Jacobian block \mathbf{J}_j for this point feature and the structure is shown in Figure 12.

The matrix \mathbf{M} is formed by stacking and reordering all of the \mathbf{L}_j matrices for $j = 1 \dots n_f$ and has the same rank as the original Jacobian \mathbf{J} . The structure of matrix \mathbf{M} is shown in Figure 13.

The matrix \mathbf{M} is a block-diagonal matrix composed of three sub-matrices:

- $\mathbf{M}_1 \in \mathbb{R}^{n_f \times n_f}$ is diagonal,

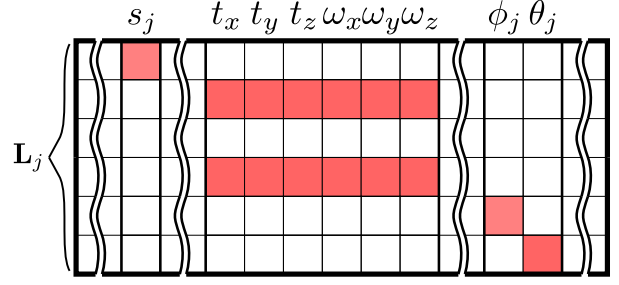


Figure 12: Structure of the Jacobian block \mathbf{L}_j for the point feature j observations, resulting from manipulations to the original Jacobian.

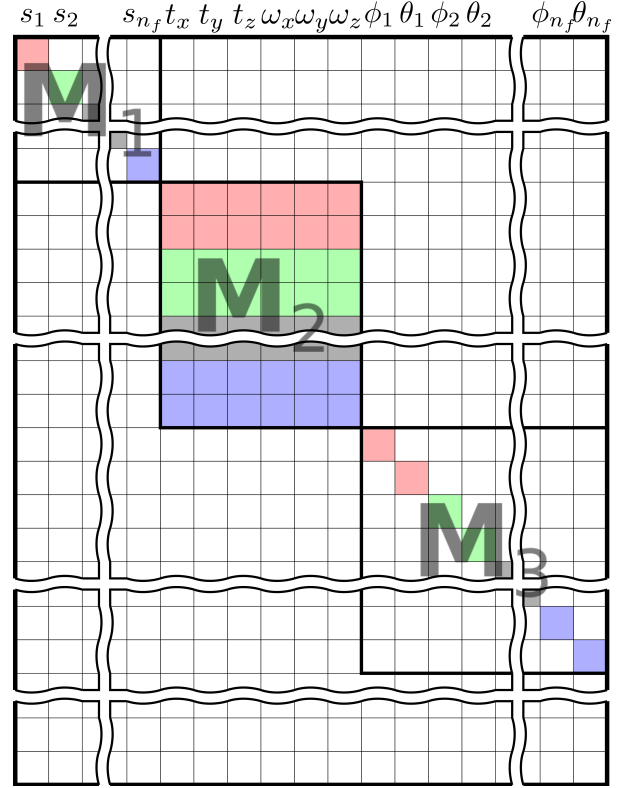


Figure 13: Structure of the matrix \mathbf{M} , resulting from stacking and reordering the rows of the matrices \mathbf{L}_j . The sub-blocks of \mathbf{M} are shown and all must be full rank for \mathbf{M} , and therefore \mathbf{J} , to be full rank.

- $\mathbf{M}_2 \in \mathbb{R}^{m_o \times 6}$,
- $\mathbf{M}_3 \in \mathbb{R}^{2n_f \times 2n_f}$ is diagonal.

As a result, \mathbf{M} and \mathbf{J} are full rank if and only if all of the following are satisfied,

- $\text{rank}(\mathbf{M}_1) = n_f$,

- $\text{rank}(\mathbf{M}_2) = 6$, and
- $\text{rank}(\mathbf{M}_3) = 2n_f$.

It is clear that both \mathbf{M}_1 and \mathbf{M}_3 are full rank by construction, and therefore, \mathbf{M} and by extension \mathbf{J} are full rank if and only if \mathbf{M}_2 is full rank, which concludes the proof. \square

Therefore, determining if the measurement Jacobian \mathbf{J} is full rank is simplified to checking the rank of the reduced-dimension matrix \mathbf{M}_2 . It follows directly from Lemma 4.1, that the degeneracy of the multicamera cluster SLAM system can be determined by checking the rank of \mathbf{M}_2 .

Corollary 4.2. *For a multicamera cluster SLAM system satisfying Assumptions 3.1, 3.2, and 3.3, the solution is degenerate and under-constrained, if and only if $\text{rank}(\mathbf{M}_2) < 6$.*

4.2.1. Rank of \mathbf{M}_2

The matrix \mathbf{M}_2 is a dense $m_o \times 6$ block with a single row for each observation of the point features at the second keyframe. Each row of \mathbf{M}_2 in (67) specifies the six Plücker coordinates for a line in \mathbb{R}^3 since each set of coordinates satisfies the Grassmann-Plücker relation [22],

$$([\mathbf{v}_{j,k}]_{\times} \hat{\mathbf{a}}_{j,2}) \cdot ([\mathbf{w}_j]_{\times} [\mathbf{v}_{j,k}]_{\times} \hat{\mathbf{a}}_{j,2}) \quad (85)$$

$$= -\mathbf{w}_j \cdot \underbrace{([\mathbf{v}_{j,i(k)}]_{\times} \hat{\mathbf{a}}_{j,2})_{\times} ([\mathbf{v}_{j,i(k)}]_{\times} \hat{\mathbf{a}}_{j,2})}_{= \mathbf{0}_{3 \times 1}} \quad (86)$$

$$= 0. \quad (87)$$

The matrix \mathbf{M}_2 will not have full rank when the m_o sets of coordinates are linearly-dependent. This is similar to the problem of identifying motion singularities for series-parallel mechanisms. However, the current problem is more complex since the common connection points which sometimes allow for the simplification of the singularity condition in mechanisms are not present in the cluster SLAM system.

4.3. Sufficient Conditions for Degeneracy

In this section, the structure of the matrix \mathbf{M}_2 from (67) will be exploited to identify cluster configurations and motions that are sufficient for degeneracy of the solution, independent of the number of point features observed and their constellation geometry.

It is immediately apparent that the system must include six point feature observations at the second keyframe for the matrix \mathbf{M}_2 to possibly have full rank. The

first three columns in \mathbf{M}_2 are a stack of cross products involving the camera observation vectors and the bearings to the point features. When they all have a common collinear vector operand, the resulting row vectors are all coplanar, with the normal defined by the collinear vector operand.

As expected, the system will be degenerate if the cluster consists of only one component camera since the rows will all have a common camera observation vector, consistent with how monocular vision systems are unable to recover the six degrees of freedom motion solution in a SLAM system. Additionally, the system is degenerate when only one point feature is observed by six cameras at the second keyframe since all of the matrix rows will contain the common point feature unit vector, $\hat{\mathbf{a}}_{1,2}$, in the cross product term.

When the camera observation vectors are all parallel, the SLAM solution is degenerate. Each camera observation vector can be written as a scalar multiple, $\exists \gamma_{m,n} \in \mathbb{R}$ such that $\mathbf{v}_{m,n} = \gamma_{m,n} \mathbf{v} \in \mathbb{R}^3$. In this case, the matrix \mathbf{M}_2 will have less than full rank since,

$$\text{rank} \begin{pmatrix} -\gamma_{1,1} \hat{\mathbf{a}}_{1,2}^T [\mathbf{v}]_{\times} \\ \vdots \\ -\gamma_{n_f, n_o(n_f)} \hat{\mathbf{a}}_{n_f,2}^T [\mathbf{v}]_{\times} \end{pmatrix} \leq 2 < 3, \quad (88)$$

and $\text{rank}(\mathbf{M}_2) < 6$. This condition includes the previously known two-camera cluster concentric circle degeneracies, since the motion causes the camera centres to move in parallel. Adding more cameras to the cluster reduces the configurations for which the relative motion will lead to the camera observation vectors being parallel. Additionally, when point features are observed across different cameras within the cluster at the two keyframes, it becomes less likely that all of the camera observation vectors are parallel. However, there do exist certain combinations of cluster motions for which the camera observation vectors remain parallel regardless of the feature point locations, and the camera cluster system becomes degenerate. Some example configurations are presented in Figure 14.

When the relative motion of the camera cluster is such that a point feature which was observed in one camera at the first keyframe is observed by a different camera at the second keyframe, there is a camera observation vector between the positions of the two cameras when they observed the particular point feature. This can create a set of camera observation vectors which are non-parallel even when the relative motion is a pure translation. As a result, it is possible that the system is non-degenerate and the solution can be found in this case, depending on the rank of the matrix \mathbf{M}_2 . Observ-

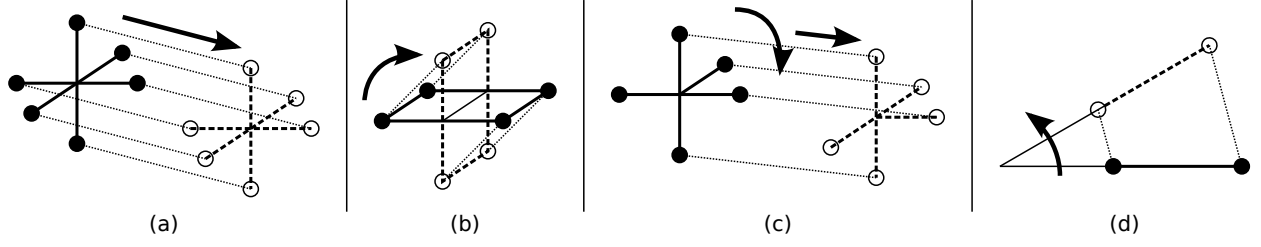


Figure 14: Examples of camera cluster motions sufficient for degeneracy. The black dots are the cameras at the first keyframe connected by solid lines, white dots are the cameras at the second keyframe connected by dashed lines. The dotted lines are the camera observation vectors which are all parallel. Motions include (a) pure translation with no intercamera correspondence, (b) rotation axis in the plane of planar four-camera cluster, (c) 90 degrees rotation with translation, (d) two-camera concentric circles motion.

ing common point features over multiple cameras is an effective way of avoiding the set of camera observation vectors becoming parallel and reducing the set of sufficient motions for system degeneracy.

4.4. Necessary and Sufficient Conditions for Degeneracy

In the case when the motion does not produce parallel camera observation vectors, it is necessary to evaluate the rank of the matrix \mathbf{M}_2 before concluding that the system is non-degenerate. The matrix \mathbf{M}_2 can be regarded as a set of motion constraints on a mechanism where a non-full rank means that the framework is not rigid and the configuration can change without violating the constraints. The full analysis of these singular configurations is beyond the scope of this work, but this section presents a set of example configurations for some general case systems to numerically demonstrate the effect of adding additional cameras and point feature observations on the degenerate configuration set.

Figures 15a and 15b show typical surface meshes for the relative cluster translation, \mathbf{t}_K , leading to \mathbf{M}_2 losing rank for example two and three-camera cluster systems observing six point features with no overlap in the camera observations and non-zero relative rotation. The surface is computed numerically as the locus of zero determinant for the matrix \mathbf{M}_2 using a set of randomly-positioned point features. A similar surface is generated for any non-zero rotation \mathbf{R}_K .

As more point feature observations are added to the system, the number of degeneracies is reduced. With more than six point feature observations, the size of the matrix \mathbf{M}_2 becomes $m_o \times 6$, where $m_o > 6$ and therefore the rank of the matrix cannot be checked by computing the determinant directly. Instead, for the matrix to not be full rank, all of the m_o choose 6 submatrices of size 6×6 formed by the rows of \mathbf{M}_2 must have a zero determinant.

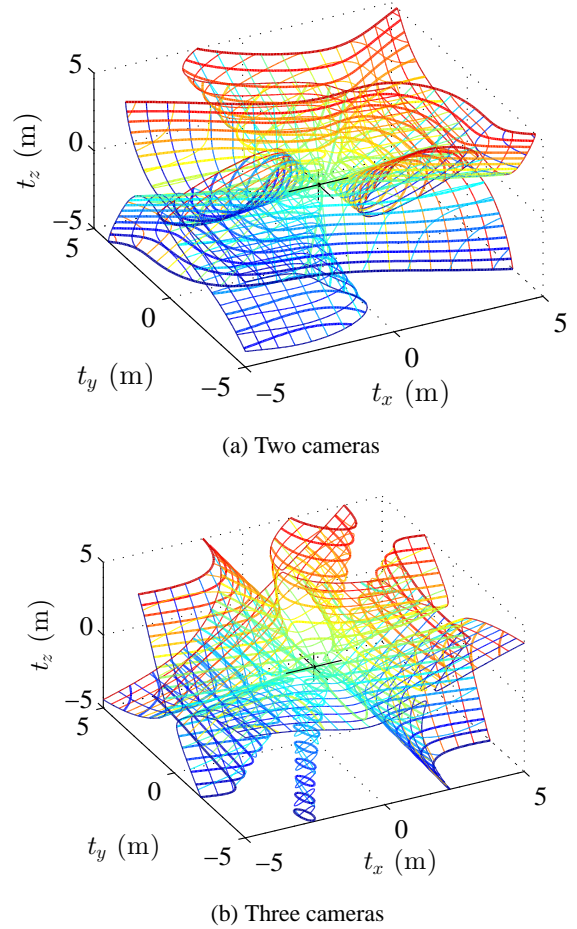


Figure 15: Degenerate \mathbf{t}_K for (a) two and (b) three-camera clusters, observing six points with no overlap and non-zero rotation.

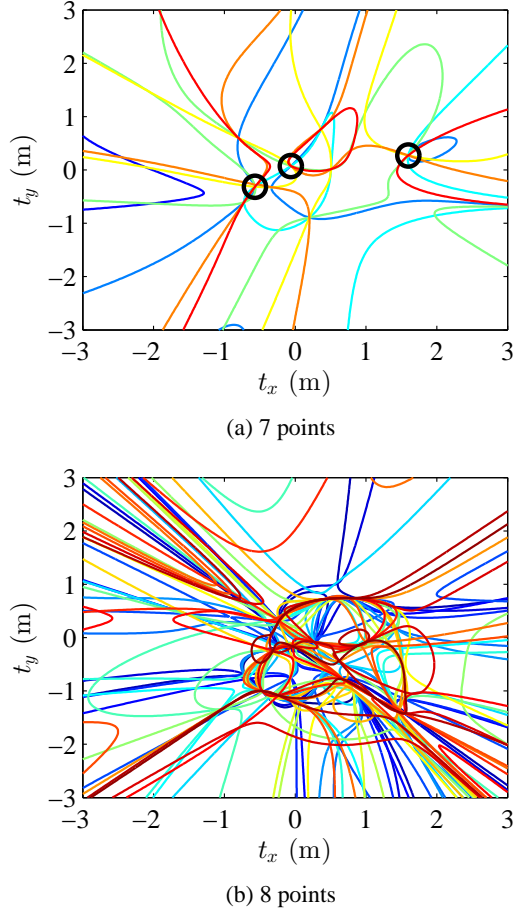


Figure 16: The loci of zero determinants for the m_o choose 6 submatrices of \mathbf{M}_2 at $t_z = 1$ m in the three-camera cluster case for (a) $m_o = 7$ and (b) $m_o = 8$ point features. The degenerate points are the intersections marked with black circles.

Each of the submatrices generates a surface as in Figure 15, and therefore, \mathbf{M}_2 has deficient rank at the \mathbf{t}_k where all of the surfaces intersect.

A two-dimensional cross-section at $t_z = 1$ m is shown in Figure 16 for the three-camera case observing 7 and 8 point features with no overlap and non-zero rotation. The degenerate cluster translations correspond to the points on the graph where all of the curves intersect and are marked as black circles. These intersections are determined numerically using the computed loci for the determinants of the submatrices. If all of the loci intersect with each other within a certain epsilon ball, the location is selected as a degeneracy of \mathbf{M}_2 .

Notice that the number of degenerate motions is reduced from the curves of each colour for a subset of

six feature observations, to a small finite set of points at any given cross-section. While the indicated degenerate positions are subject to numerical precision considerations, these examples are more informative in demonstrating that the system is non-degenerate in almost all configurations.

It is observed that the degenerate points in Figure 16a connect as lines in \mathbb{R}^3 at different slices of t_z . When observing eight points in general position, Figure 16b shows that there are no translations for which the system is degenerate. While not exhaustive, numerical analysis of the singular configurations of \mathbf{M}_2 shows that the set of degenerate motions in the cluster system has been reduced from the previous surface with $m_o = 6$, to a set of lines with $m_o = 7$ and the empty set for $m_o = 8$.

Adding point feature observations to the system is an effective way to reduce the set of degenerate motions for the camera cluster system; however, the sufficient conditions in Section 4.3 remain no matter how many point features are observed. Examples of these degeneracies are shown in Figure 17a for a two-camera cluster and Figure 17b for a three-camera cluster with rotation purely in the camera centre plane [9].

The indicated degeneracy corresponds to the motion causing the camera observation vectors to move in parallel and extends along a line in \mathbb{R}^3 . In order to eliminate these motions, it is necessary to add more cameras to the cluster, observe point features over different cameras, or both, to ensure that not all of the camera observation vectors are parallel.

Multiple cameras observing the same point feature at the different keyframes adds extra camera observation vectors which are less likely to be parallel with the rest of the vector set and produce a full-rank \mathbf{M}_2 . An example two-camera cluster observing eight point features is shown in Figure 18. Camera 2 observes a feature at the second keyframe that was measured by camera 1 initially, and camera 1 observes a feature from camera 2. Significantly, there is no relative rotation between the keyframes of this system, but the rank of \mathbf{M}_2 is full nearly everywhere. This is an improvement on the previous completely non-overlapping cluster configurations which required non-zero rotation to have a full-rank \mathbf{M}_2 .

The example systems in this section demonstrate numerically the effect of adding cameras and point feature observations to the camera cluster SLAM solution. For any algorithm, the overall strategy should be to reduce the number of degenerate configurations by increasing the number of point features observed on the target, and then eliminating the remaining sufficient conditions for degeneracy by adding cameras to the cluster, or observ-

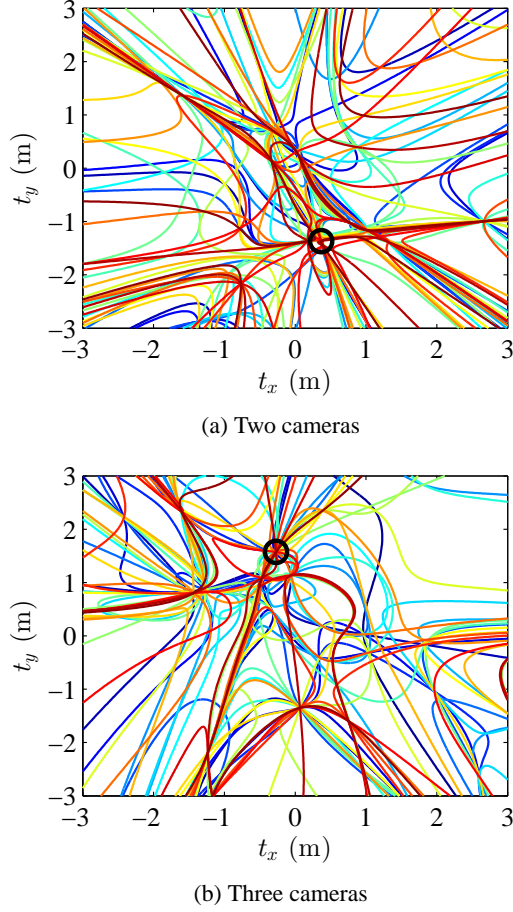


Figure 17: The loci of zero determinants for the 8 choose 6 submatrices of \mathbf{M}_2 at $t_z = 1$ m in the (a) two-camera and (b) three-camera cluster case with rotation axis within the camera centre plane observing eight point features

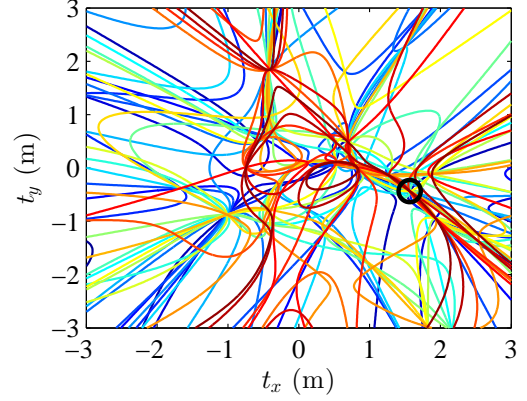


Figure 18: The loci of zero determinants for the 8 choose 6 submatrices of \mathbf{M}_2 at $t_z = 1$ m in the two-camera cluster system with zero rotation between keyframes, but two common features across cameras.

ing point features across cameras, such that it is impossible for all of the camera observation vectors to be parallel through the motion. This will ensure a well-constrained solution to the localization and mapping problem.

5. Conclusions

This work presented a detailed analysis of the degenerate configurations of the calibrated non-overlapping FOV multicamera cluster SLAM problem for an optimization based on minimizing a least-squares cost function with respect to the image-plane reprojection error. The system is reduced to a simple matrix rank test on a matrix consisting of rows of Plücker coordinates for lines in \mathbb{R}^3 . Sufficient configurations for solution degeneracy caused by the relative motion were identified for n_c -camera clusters observing any number of point features over two keyframes. This leads to the novel general conclusion that if all of the camera observation vectors, formed as the displacement between the pairs of cameras observing a particular point feature, are parallel in a common coordinate frame, then the system is degenerate. It is further shown for several example systems that with the addition of more cameras to the cluster, more point feature observations, and observations of the point features across different cameras, the set of degenerate configurations is significantly reduced as it becomes impossible for all of the identified vectors to be parallel and the redundant observations prevent all of the determinants of the submatrices from going to zero concurrently.

Future work will focus on fully characterizing the necessary and sufficient conditions for the system to become degenerate, including the degeneracies related to the geometry of the point feature constellation from the standpoint of geometric algebra techniques [23]. Additionally, the results from this work will be used to generate metrics for deciding when and where to add keyframes in a real-time SLAM system to accurately construct and constrain the generated target model and avoid the degeneracies within the state space caused by measurements from this type of sensor.

Acknowledgements

This work was partially funded by the National Sciences and Engineering Research Council of Canada (NSERC) under Grant No. CRDPJ 397768-10. Partial funding also comes from the NSERC through the Alexander Graham Bell Canada Graduate Scholarship - Doctoral (CGS-D) award.

Appendix A. Projective Geometry

The projective space, \mathbb{P}^n , consists of the real vector space \mathbb{R}^n , with the addition of points at infinity [5]. Only a very brief description of the projective space is presented here and the reader is referred to [5] for a more thorough introduction.

A point in the projective space is represented by the $n + 1$ homogeneous coordinates,

$$\tilde{\mathbf{x}} = [\tilde{x}_1 \quad \tilde{x}_2 \quad \dots \quad \tilde{x}_{n+1}]^T \in \mathbb{P}^n. \quad (\text{A.1})$$

The points at infinity in \mathbb{R}^n are represented by those with coordinate $x_{n+1} = 0$. For finite points in \mathbb{R}^n – when $x_{n+1} \neq 0$ – the coordinates of the corresponding point $\mathbf{x} \in \mathbb{R}^n$ are determined by,

$$\mathbf{x} = [x_1 \quad x_2 \quad \dots \quad x_n]^T \quad (\text{A.2})$$

$$= \left[\frac{\tilde{x}_1}{\tilde{x}_{n+1}} \quad \frac{\tilde{x}_2}{\tilde{x}_{n+1}} \quad \dots \quad \frac{\tilde{x}_n}{\tilde{x}_{n+1}} \right]^T. \quad (\text{A.3})$$

Note that there is no way of mapping a point at infinity back to \mathbb{R}^n since it would require division by zero.

Each point along a ray in the projective space maps to the same point in the real vector space. As a result, the points $\tilde{\mathbf{x}}$ and $\lambda\tilde{\mathbf{x}}$, for $\lambda \in \mathbb{R}$, map to the same point $\mathbf{x} \in \mathbb{R}^n$. Not surprisingly, there is an extra degree of freedom in the projective vectors using $n + 1$ coordinates to represent a n -dimensional space. Finally, it is possible to represent any point $\mathbf{x} \in \mathbb{R}^n$ in the corresponding

projective space \mathbb{P}^n simply by augmenting the coordinates,

$$\tilde{\mathbf{x}} = [\mathbf{x}^T \quad 1]^T. \quad (\text{A.4})$$

The projective spaces allow for projective and coordinate transformations to be represented as linear matrix operations.

References

References

- [1] P. Baker, C. Fermuller, Y. Aloimonos, R. Pless, A spherical eye from multiple cameras (makes better models of the world), in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2001, pp. 576–583.
- [2] R. Pless, Using many cameras as one, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Vol. 2, 2003, pp. II–587–593.
- [3] B. Clipp, J. H. Kim, J. M. Frahm, M. Pollefeys, R. Hartley, Robust 6DOF motion estimation for non-overlapping, multi-camera systems, in: Proceedings of the IEEE Workshop on Applications of Computer Vision (WACV), 2008, pp. 1–8.
- [4] S. Thrun, W. Burgard, D. Fox, Probabilistic robotics, The MIT Press, 2005.
- [5] R. Hartley, A. Zisserman, Multiple view geometry in computer vision, Cambridge University Press, 2003.
- [6] A. J. Davison, I. D. Reid, N. D. Molton, O. Stasse, MonoSLAM: Real-time single camera SLAM, IEEE Transactions on Pattern Analysis and Machine Intelligence 29 (6) (2007) 1052–1067.
- [7] R. Hermann, A. Krener, Nonlinear controllability and observability, IEEE Transactions on Automatic Control 22 (5) (1977) 728–740.
- [8] J. H. Kim, M. J. Chung, B. T. Choi, Recursive estimation of motion and a scene model with a two-camera system of divergent view, Pattern Recognition 43 (6) (2010) 2265–2280.
- [9] M. J. Tribou, S. L. Waslander, D. W. L. Wang, Scale recovery in multicamera cluster SLAM with non-overlapping fields of view, submitted to *Computer Vision and Image Understanding* (Aug. 2013).
- [10] H. Stewenius, D. Nister, M. Oskarsson, K. Astrom, Solutions to minimal generalized relative pose problems, in: Proceedings of the Workshop on Omnidirectional Vision, 2005.
- [11] P. Sturm, Multi-view geometry for general camera models, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Vol. 1, 2005, pp. 206–212.
- [12] E. Mouragnon, M. Lhuillier, M. Dhome, F. Dekeyser, P. Sayd, Generic and real-time structure from motion using local bundle adjustment, Image and Vision Computing 27 (8) (2009) 1178–1193.
- [13] J. S. Kim, T. Kanade, Degeneracy of the linear seventeen-point algorithm for generalized essential matrix, Journal of Mathematical Imaging and Vision 37 (1) (2010) 40–48.
- [14] H. Li, R. Hartley, J. H. Kim, A linear approach to motion estimation using generalized camera models, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2008, pp. 1–8.
- [15] Y. Lu, J. Z. Zhang, Q. M. J. Wu, Z. N. Li, A survey of motion-parallax-based 3-D reconstruction algorithms, IEEE Transactions on Systems, Man, and Cybernetics 34 (4) (2004) 532–548.
- [16] R. M. Murray, Z. Li, S. S. Sastry, A mathematical introduction to robotic manipulation, CRC Press, 1994.

- [17] E. Rosten, T. Drummond, Fusing points and lines for high performance tracking., in: Proceedings of the IEEE International Conference on Computer Vision (ICCV), Vol. 2, 2005, pp. 1508–1511.
- [18] E. Rosten, T. Drummond, Machine learning for high-speed corner detection, in: Proceedings of the European Conference on Computer Vision (ECCV), Vol. 1, 2006, pp. 430–443.
- [19] D. G. Lowe, Object recognition from local scale-invariant features, Proceedings of the IEEE International Conference on Computer Vision (ICCV) 2 (1999) 1150–1157.
- [20] H. Bay, A. Ess, T. Tuytelaars, L. Van Gool, Speeded-up robust features (SURF), Computer Vision and Image Understanding 110 (3) (2008) 346–359.
- [21] G. Klein, D. Murray, Parallel tracking and mapping for small AR workspaces, in: Proceedings of the IEEE and ACM International Symposium on Mixed and Augmented Reality (ISMAR), 2007, pp. 225–234.
- [22] N. L. White, Grassmann-Cayley algebra and robotics, Journal of Intelligent and Robotic Systems 11 (1-2) (1994) 91–107.
- [23] L. Dorst, D. Fontijne, S. Mann, Geometric Algebra for Computer Science: An Object-Oriented Approach to Geometry, Morgan Kaufmann, 2010.