

Consistency Analysis for Sliding-Window Visual Odometry

Tue-Cuong Dong-Si and Anastasios I. Mourikis

Dept. of Electrical Engineering, University of California, Riverside

E-mail: tdongsi@ee.ucr.edu, mourikis@ee.ucr.edu

Abstract—In this paper we focus on the problem of *visual odometry*, i.e., the task of tracking the pose of a moving platform using visual measurements. In recent years, several VO algorithms have been proposed that employ nonlinear minimization in a sliding window of poses for this task. Through the use of iterative re-linearization, these methods are capable of successfully addressing the nonlinearity of the measurement models, and have become the de-facto standard for high-precision VO. In this work, we conduct an analysis of the properties of marginalization, which is the process through which older states are removed from the sliding window. This analysis shows that the standard way of marginalizing older poses results in an erroneous change in the rank of the measurements’ information matrix, and leads to underestimation of the uncertainty of the state estimates. Based on the analytical results, we also propose a simple modification of the way in which the measurement Jacobians are computed. This modification avoids the above problem, and results in an algorithm with superior accuracy, as demonstrated in both simulation tests and real-world experiments.

I. INTRODUCTION AND RELATED WORK

Accurate pose tracking is an essential requirement in a number of systems, ranging from large-scale autonomous vehicles to small hand-held devices. If a system operates in an environment where reliable GPS reception is possible, then it can use the GPS signals, possibly in conjunction with additional sensors such as an inertial measurement unit (IMU), to track its position. However, in GPS-denied environments, different sensing modalities must be employed. Among the possible choices (e.g., cameras, laser rangefinders, sonars) cameras stand out due to their low cost, size and power consumption, as well as their widespread availability on mobile devices. Recent advances in the performance of vision sensors and computing hardware have made vision-based algorithms a more attractive option. Motivated by these reasons, in this paper we focus on the task of vision-based motion estimation.

Several different methods for localization using cameras have been proposed. For instance, when landmarks or fiducial points with known positions are available, these can be used to estimate the camera’s position with respect to a known frame (e.g., [1]–[3]). However, when cameras operate in unknown and un-instrumented environments, these methods are not applicable. In these settings, one can employ simultaneous localization and mapping (SLAM)-type methods (see, e.g., [4]–[6] and references therein), which jointly estimate the camera trajectory and features’ positions, or visual odometry (VO)-type methods [7]–[12], which track the pose of the camera only. In this work, our interest is

in the latter type of methods. We are not interested in the problem of loop-closure detection, or in producing a map of the environment. Instead, we focus on the task of performing VO in real time, using the measurements from a monocular or stereo camera *only*.

Several methods exist that perform VO by estimating the camera displacement using the images recorded at two consecutive time steps (e.g., [12]–[14]). In these methods, when a single camera is used, additional sensors, scene information, or a statistical motion model must be employed to infer scale. By using only consecutive images, these methods attain low computational cost, but this often comes at the expense of accuracy. Due to the nonlinear nature of the camera measurement models and the existence of outliers in the image data, such methods may also be prone to failure.

At the other end of the spectrum, the “golden standard” method for vision-based estimation is bundle adjustment (e.g., [15] and references therein). In bundle adjustment the entire history of camera states and feature positions is jointly optimized using nonlinear minimization methods. This can result in high precision, but bundle-adjustment methods cannot operate in real-time in large-scale environments, as their computational complexity continuously increases over time. Incremental implementation of the nonlinear minimization is possible [16], but even in this case the computational cost increases in time, and eventually becomes unsuitable for real-time applications.

As a compromise between bundle adjustment and using pairwise displacement estimates, methods that perform optimization over a *sliding window* of states have recently gained popularity (see, e.g., [9]–[11], [17]–[20] and references therein). These techniques remove older states (features and/or camera poses) from the actively estimated state vector, and carry out iterative minimization to produce estimates for the most recent states. The use of a sliding window of more than two camera poses increases the accuracy and stability of these algorithms. At the same time, the removal of older states means that these methods have a *bounded* computational cost, which makes them suitable for real-time implementation. In fact, by changing the size of the sliding window we can adaptively control the computational requirements, which is an important characteristic of sliding-window algorithms. In this paper, we focus on the properties of these methods.

In addition to differences in the visual front end (i.e., the algorithms used for feature extraction and matching) the main difference between the various sliding-window

algorithms is the way in which older states are removed. It is well-known that the theoretically “correct” way of removing states from the state vector is the process of *marginalization* [18], [19], [21], [22]. When marginalization is carried out the uncertainty of the discarded states is properly modelled in the estimator’s equations, which is a key requirement for precise estimation. However, several successful methods do not follow this approach, and simply fix the values of the states that are removed from the state vector, using them to “bootstrap” the trajectory [9], [10], [17], [20]. In certain cases this is done purely for simplification of the algorithms and to improve computational efficiency. However, in [9] and [10] it is reported that this is done to reduce the estimation error. This fact, which appears to be counterintuitive at first, suggests that the “standard” way of carrying out marginalization may produce inaccuracies.

We note that in our previous work [19] we have shown that when a fixed-lag smoother is employed for tracking the motion of a vehicle using a camera and an IMU, the standard marginalization approach results in *inconsistency*¹. Specifically, due to the marginalization process, two different estimates of the same states are used in computing certain Jacobian matrices in the estimator. In [19] we showed that this causes an infusion of information along directions of the state space where no actual information is provided by the measurements (the un-observable directions). This “artificial” information causes the estimates’ covariance to be underestimated, and results in inconsistency. Moreover, since the accuracy of different states is misrepresented, the estimates’ accuracy is also reduced.

Motivated by the observations of [9] and [10] on sliding-window VO, and the results of [19] on fixed lag smoothing using visual and inertial measurements, in this paper we carry out an analysis of the effects of marginalization in VO. Our results show that, similarly to the case of fixed-lag smoothing, even when *only* camera measurements are used for estimation, the same infusion of “artificial information” takes place. This degrades both the consistency and the accuracy of the trajectory estimates. Additionally, building on this analysis, we present a simple solution to the problem. This solution consists of ensuring that only one estimate of any given state variable is used in computing Jacobian matrices. The resulting algorithm is shown to perform better than competing approaches, in both simulation results and real-world experiments.

II. SLIDING-WINDOW VISUAL ODOMETRY

In this section, we present the “standard” algorithm for sliding window visual odometry [18], [19]. We start by discussing bundle adjustment, which serves to introduce the notation and will also be useful for our derivations in Section III.

¹A recursive estimator is termed consistent when the state estimation errors are zero mean, and their covariance equals the one reported by the estimator [23, Section 5.4].

A. Bundle adjustment

We consider the case where a monocular or stereo camera moves in space, observing unknown features. The camera state vector at time-step i , \mathbf{c}_i , consists of the sensor orientation and position with respect to a global frame of reference:

$$\mathbf{c}_i = \begin{bmatrix} \mathbf{q}_{C_i} \\ \mathbf{p}_{C_i} \end{bmatrix} \quad (1)$$

where we have employed a unit-quaternion description of orientation, \mathbf{q}_{C_i} [24]. Assuming calibrated cameras, the observation of feature j at time-step i is described by the perspective camera model:

$$\mathbf{z}_{ij} = \mathbf{h}(\mathbf{C}(\mathbf{q}_{C_i})(\mathbf{p}_{L_j} - \mathbf{p}_{C_i})) + \mathbf{n}_{ij} \quad (2)$$

where \mathbf{p}_{L_j} is the 3D feature position vector, $\mathbf{C}(\mathbf{q}_{C_i})$ is the rotation matrix corresponding to \mathbf{q}_{C_i} (i.e., the rotation matrix from the global frame to the camera frame at time-step i), $\mathbf{h}(\cdot)$ is the perspective measurement function: $\mathbf{h}(\mathbf{p}) = [p_1/p_3 \quad p_2/p_3]^T$, and finally $\mathbf{n}_{ij} \sim \mathcal{N}(\mathbf{0}_{2 \times 1}, \mathbf{R}_{ij})$ is the measurement noise vector, modeled as a zero-mean Gaussian variable with covariance matrix \mathbf{R}_{ij} . If a stereo camera is used, the state vector in (1) represents the pose of one of the two sensors (e.g., the left one), and we obtain two measurements of the form (2), one from each camera. When performing bundle adjustment at time-step k , we simultaneously estimate the entire history of camera poses, $\mathbf{c}_{0:k} = \{\mathbf{c}_0, \dots, \mathbf{c}_k\}$, as well as the positions of all features observed by the camera, $\mathbf{l}_{1:n} = \{\mathbf{p}_{L_1}, \dots, \mathbf{p}_{L_n}\}$. We denote the state vector containing all these quantities by \mathbf{x}_k . The optimal solution is computed by maximizing the following pdf:

$$p(\mathbf{x}_k | \mathbf{z}_{0:k}) = p(\mathbf{x}_k) \prod_{(i,j) \in \mathcal{S}_a(k)} p(\mathbf{z}_{ij} | \mathbf{c}_i, \mathbf{p}_{L_j}) \quad (3)$$

where the set $\mathcal{S}_a(k)$ contains the pairs of indices (i, j) that describe all the feature observations through time k , $p(\mathbf{z}_{ij} | \mathbf{c}_i, \mathbf{p}_{L_j})$ is the Gaussian measurement-likelihood pdf, and $p(\mathbf{x}_k)$ is a pdf describing the prior information for the state vector. For instance, this may express our knowledge of the first camera state, constraints on the global scale in the case of monocular VO, etc. For clarity of presentation, we assume here that this prior information is modelled as a probabilistic constraint, $\mathbf{f}(\mathbf{x}_k) \sim \mathcal{N}(\hat{\mathbf{f}}, \mathbf{R}_p)$. The maximization of the above pdf is equivalent to the minimization of the following cost function:

$$c(\mathbf{x}_k) = \frac{1}{2} \|\mathbf{f}(\mathbf{x}_k) - \hat{\mathbf{f}}\|_{\mathbf{R}_p} + \frac{1}{2} \sum_{(i,j) \in \mathcal{S}_a(k)} \|\mathbf{z}_{ij} - \mathbf{h}(\mathbf{c}_i, \mathbf{p}_{L_j})\|_{\mathbf{R}_{ij}}$$

where we have employed the notation $\|\mathbf{a}\|_{\mathbf{M}} = \mathbf{a}^T \mathbf{M}^{-1} \mathbf{a}$.

$c(\mathbf{x}_k)$ is a nonlinear cost function, which can be minimized using iterative Gauss-Newton minimization [15]. At the ℓ -th iteration of this method, a correction, $\Delta \mathbf{x}^{(\ell)}$, to the current estimate, $\mathbf{x}_k^{(\ell)}$, is computed by solving the linear system:

$$\mathbf{A}^{(\ell)} \Delta \mathbf{x}^{(\ell)} = -\mathbf{b}^{(\ell)} \quad (4)$$

where

$$\mathbf{A}^{(\ell)} = \mathbf{F}^T \mathbf{R}_p^{-1} \mathbf{F} + \sum_{(i,j) \in \mathcal{S}_a(k)} \mathbf{H}_{ij}^{(\ell)T} \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}^{(\ell)} \quad (5)$$

$$\begin{aligned} \mathbf{b}^{(\ell)} = & \mathbf{F}^T \mathbf{R}_p^{-1} (\mathbf{f}(\mathbf{x}_k^{(\ell)}) - \hat{\mathbf{f}}) \\ & - \sum_{(i,j) \in \mathcal{S}_a(k)} \mathbf{H}_{ij}^{(\ell)T} \mathbf{R}_{ij}^{-1} (\mathbf{z}_{ij} - \mathbf{h}(\mathbf{c}_i^{(\ell)}, \mathbf{p}_{L_j}^{(\ell)})) \end{aligned} \quad (6)$$

In the above expressions, \mathbf{F} is the Jacobian of the function $\mathbf{f}(\mathbf{x}_k)$ with respect to \mathbf{x}_k , and $\mathbf{H}_{ij}^{(\ell)}$ is the Jacobian of the measurement function $\mathbf{h}(\mathbf{c}_i, \mathbf{p}_{L_j})$ with respect to \mathbf{x}_k , evaluated at $\mathbf{x}_k^{(\ell)}$. Since the measurement model involves only one camera pose and one feature, $\mathbf{H}_{ij}^{(\ell)}$ has the following sparse structure:

$$\mathbf{H}_{ij}^{(\ell)} = \begin{bmatrix} \mathbf{0} & \dots & \mathbf{H}_{L_{ij}}(\mathbf{x}_k^{(\ell)}) & \dots & \mathbf{H}_{C_{ij}}(\mathbf{x}_k^{(\ell)}) & \dots & \mathbf{0} \end{bmatrix} \quad (7)$$

where $\mathbf{H}_{L_{ij}}$ and $\mathbf{H}_{C_{ij}}$ are the Jacobians with respect to the feature position and the camera pose, respectively:

$$\begin{aligned} \mathbf{H}_{C_{ij}}(\mathbf{x}_k) &= \mathbf{\Gamma}_{ij} \mathbf{C}(\mathbf{q}_{C_i}) \begin{bmatrix} [(\mathbf{p}_{L_j} - \mathbf{p}_{C_i}) \times] \mathbf{C}^T(\mathbf{q}_{C_i}) & -\mathbf{I}_3 \end{bmatrix} \\ \mathbf{H}_{L_{ij}}(\mathbf{x}_k) &= \mathbf{\Gamma}_{ij} \mathbf{C}(\mathbf{q}_{C_i}) \quad (8) \\ \mathbf{\Gamma}_{ij} &= \left. \frac{\partial \mathbf{h}(\mathbf{p})}{\partial \mathbf{p}} \right|_{\mathbf{p}=\mathbf{C}(\mathbf{q}_{C_i})(\mathbf{p}_{L_j} - \mathbf{p}_{C_i})} \end{aligned}$$

where \mathbf{I}_3 denotes the 3×3 identity matrix, and $[\mathbf{a} \times]$ is the skew-symmetric matrix associated with the vector \mathbf{a} . After solving (4) the correction is applied to the state, and the process is repeated until convergence.

B. Marginalization of old states

By exploiting the sparse structure of $\mathbf{A}^{(\ell)}$, we can speed up the solution of the linear system (4) considerably. However, as the camera continuously moves and observes new features, the size of the state vector \mathbf{x}_k constantly increases. Therefore, in order to obtain a real-time algorithm with bounded computational complexity, we marginalize out older states.

We consider the following scenario: The moving camera observes features during the time interval $[0, k]$, and bundle adjustment is carried out at time-step k . Then, the states $\mathbf{x}_m = \{\mathbf{c}_0, \dots, \mathbf{c}_{m-1}, \mathbf{p}_{L_1}, \dots, \mathbf{p}_{L_q}\}$ (i.e., the m oldest camera poses and the q oldest landmarks, which we can no longer observe) are marginalized out, and only the states $\mathbf{x}_r = \{\mathbf{c}_m, \dots, \mathbf{c}_k, \mathbf{p}_{L_{q+1}}, \dots, \mathbf{p}_{L_n}\}$ remain active in the sliding window. After marginalization, the states \mathbf{x}_m and all the measurements that involve them are discarded. We will use \mathcal{S}_m to denote the set of indices (i, j) describing all the camera observations that involve either marginalized camera poses or marginalized landmarks, or both. These measurements provide information that is useful for the estimation of the remaining states, and this information should *not* be completely discarded. To express this information, we maintain in memory a vector $\mathbf{b}_p(k)$ and a matrix $\mathbf{A}_p(k)$, which are defined as follows:

$$\mathbf{b}_p(k) = \mathbf{b}_{mr}(k) - \mathbf{A}_{rm}(k) \mathbf{A}_{mm}(k)^{-1} \mathbf{b}_{mm}(k) \quad (9)$$

$$\mathbf{A}_p(k) = \mathbf{A}_{rr}(k) - \mathbf{A}_{rm}(k) \mathbf{A}_{mm}(k)^{-1} \mathbf{A}_{mr}(k) \quad (10)$$

where

$$\mathbf{b}_m(k) = \begin{bmatrix} \mathbf{b}_{mm}(k) \\ \mathbf{b}_{mr}(k) \end{bmatrix} \quad (11)$$

$$\begin{aligned} &= \mathbf{F}^T \mathbf{R}_p^{-1} (\mathbf{f}(\hat{\mathbf{x}}_k(k)) - \hat{\mathbf{f}}) \\ & - \sum_{(i,j) \in \mathcal{S}_m} \mathbf{H}_{ij}^T(k) \mathbf{R}_{ij}^{-1} (\mathbf{z}_{ij} - \mathbf{h}(\hat{\mathbf{c}}_i(k), \hat{\mathbf{p}}_{L_j}(k))) \end{aligned} \quad (12)$$

$$\mathbf{A}_m(k) = \begin{bmatrix} \mathbf{A}_{mm}(k) & \mathbf{A}_{mr}(k) \\ \mathbf{A}_{rm}(k) & \mathbf{A}_{rr}(k) \end{bmatrix} \quad (13)$$

$$= \mathbf{F}^T \mathbf{R}_p^{-1} \mathbf{F} + \sum_{(i,j) \in \mathcal{S}_m} \mathbf{H}_{ij}^{(\ell)}(k) \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}^{(\ell)}(k) \quad (14)$$

In the above, the size of the matrix partitions is defined by the length of the vectors \mathbf{x}_m and \mathbf{x}_r , and all quantities are evaluated using the state estimate $\hat{\mathbf{x}}_k(k)$ (i.e., the estimate of \mathbf{x}_k computed using bundle adjustment at time-step k). We point out that the matrix \mathbf{A}_m represents the information contained in the prior and the discarded measurements, and $\mathbf{A}_p(k)$ is its Schur complement. Thus, as desired, $\mathbf{A}_p(k)$ represents all the information that the prior and the discarded measurements provide for estimating \mathbf{x}_r .

Proceeding further, as the camera keeps moving and observing features in the time interval $[k+1, k']$, the sliding window of states is augmented by the new camera and landmark states $\mathbf{x}_n = \{\mathbf{c}_{k+1}, \dots, \mathbf{c}_{k'}, \mathbf{p}_{L_{n+1}}, \dots, \mathbf{p}_{L_{n'}}\}$. Now, at time-step k' , the sliding window contains the states \mathbf{x}_r and \mathbf{x}_n . To obtain an estimate for the active state vector we once again employ iterative minimization of an appropriate cost function [19]. Similarly to the previous case, at the ℓ -th iteration the correction to the active states $\{\mathbf{x}_r, \mathbf{x}_n\}$ is computed by solving the linear system $\mathbf{A}^{(\ell)} \Delta \mathbf{x} = -\mathbf{b}^{(\ell)}$, with:

$$\begin{aligned} \mathbf{b}^{(\ell)} = & \mathbf{\Pi}_r^T \mathbf{b}_p(k) + \mathbf{\Pi}_r^T \mathbf{A}_p(k) (\mathbf{x}_r^{(\ell)} - \hat{\mathbf{x}}_r(k)) \\ & - \sum_{(i,j) \in \mathcal{S}_a(k')} \mathbf{H}_{ij}^{(\ell)T} \mathbf{R}_{ij}^{-1} (\mathbf{z}_{ij} - \mathbf{h}(\mathbf{c}_i^{(\ell)}, \mathbf{p}_{L_j}^{(\ell)})) \end{aligned} \quad (15)$$

$$\mathbf{A}^{(\ell)} = \mathbf{\Pi}_r^T \mathbf{A}_p(k) \mathbf{\Pi}_r + \sum_{(i,j) \in \mathcal{S}_a(k')} \mathbf{H}_{ij}^{(\ell)T} \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}^{(\ell)} \quad (16)$$

where the set $\mathcal{S}_a(k')$ contains the (i, j) indices corresponding to all the active measurements at time-step k' (i.e., all measurements involving states in \mathbf{x}_r and \mathbf{x}_n), and $\mathbf{\Pi}_r = [\mathbf{I}_{\dim \mathbf{x}_r} \ \mathbf{0}]$. After the iterations have converged, we can again marginalize out some older states if desired, and proceed in the same way.

III. ESTIMATOR CONSISTENCY

This section presents the main results of this paper, which show the effects of marginalization on the estimator's consistency. Specifically, we prove that due to the marginalization the rank of the information matrix associated with the feature measurements is erroneously increased.

We start by considering what this information matrix would be if we had *not* performed marginalization, and instead carried out bundle adjustment for the entire trajectory in the time interval $[0, k']$. In that case, the matrix describing

the information given by the measurements for the camera poses and feature positions would be given by:

$$\mathbf{J}_{\text{ba}}(k') = \sum_{(i,j) \in \mathcal{S}} \mathbf{H}_{ij}^T(k') \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}(k') \quad (17)$$

where $\mathcal{S} = \mathcal{S}_a(k') \cup \mathcal{S}_m$ is the set describing all the available measurements in $[0, k']$. The above expression is the sum of the information contribution of each of these measurements.

Let us now return to the scenario described in the preceding section, i.e., marginalization of the states \mathbf{x}_m at time-step k , and a new estimation step at time-step k' . In this case the estimator uses *the same* measurements, and thus the information matrix must contain a summation of the same number of terms as above. However, the important distinction is that the information contribution of the measurements that were discarded upon marginalization was evaluated at time-step k , and expressed by the matrix $\mathbf{A}_p(k)$ (see (10) and (16)). Thus, the information matrix for the entire history of states in $[0, k']$ is given by:

$$\mathbf{J}_{\text{mar}}(k') = \sum_{(i,j) \in \mathcal{S}_m} \mathbf{H}_{ij}^T(k) \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}(k) + \sum_{(i,j) \in \mathcal{S}_a(k')} \mathbf{H}_{ij}^T(k') \mathbf{R}_{ij}^{-1} \mathbf{H}_{ij}(k') \quad (18)$$

Comparing (17) and (18) we clearly see that, since $\mathcal{S} = \mathcal{S}_a(k') \cup \mathcal{S}_m$, the *only* difference between these two information matrices are the state estimates used for computing the Jacobians. Apart from that, the structure of the matrices is the same in both cases. Yet, perhaps surprisingly, the mere fact that the Jacobians are evaluated using different state estimates causes the rank of these two matrices to differ. Specifically, in Section III-A, we prove that

$$\text{rank}(\mathbf{J}_{\text{mar}}(k')) = \text{rank}(\mathbf{J}_{\text{ba}}(k')) + 3 \quad (19)$$

In other words, when marginalization takes place, the estimator *appears* to have more information (i.e., information along more directions of the state space) than when we perform bundle adjustment. Clearly, this increase is incorrect, since the estimators use the same measurements in both cases, and thus have access to the same information.

Since the sliding-window VO estimator believes it has more information, it underestimates the uncertainty of the state estimates it produces, i.e., it becomes *inconsistent*. Moreover, this over-confidence in the estimates' accuracy is not uniform: as discussed in Section III-B, the estimator mistakenly "believes" that the global orientation is observable. It therefore has undue confidence in its orientation estimates, which ultimately degrades the accuracy of the computed state estimates, as corrections are not properly applied to all states. This reasoning shows that the artificial increase in the rank of the information matrix leads both to inconsistency *and* to suboptimal estimates. This is corroborated by the experimental results of Section V.

It is worth pointing out that the results of our analysis are not restricted to the particular choice of states to marginalize. Even though we here focus on sliding-window estimation, where the oldest states are always marginalized, the same increase in the rank of the information matrix would occur if one chooses the states to marginalize in a different way

(for instance, a common approach is to keep a selection of keyframes evenly spread in time and/or space).

A. Proof of (19)

We now prove (19) for the case in which the visual measurements are recorded by a stereo pair of cameras. Due to limited space some of the intermediate results are omitted, and the interested reader is referred to [25] for the full details.

We start by noting that $\mathbf{J}_{\text{ba}}(k')$ and $\mathbf{J}_{\text{mar}}(k')$ can be written as follows:

$$\mathbf{J}_{\text{ba}}(k') = \mathbf{H}^T(k') \mathbf{Diag}(\mathbf{R}_{ij}^{-1}) \mathbf{H}(k') \quad (20)$$

$$\mathbf{J}_{\text{mar}}(k') = \mathbf{H}^T(k, k') \mathbf{Diag}(\mathbf{R}_{ij}^{-1}) \mathbf{H}(k, k') \quad (21)$$

where $\mathbf{Diag}(\cdot)$ denotes a block diagonal matrix, $\mathbf{H}(k')$ is a matrix with block rows $\mathbf{H}_{ij}(k')$, for all $(i, j) \in \mathcal{S}$, while $\mathbf{H}(k, k')$ is a matrix with block rows $\mathbf{H}_{ij}(k)$, for $(i, j) \in \mathcal{S}_m$, and $\mathbf{H}_{ij}(k')$, for $(i, j) \in \mathcal{S}_a(k')$.

We see that, similarly to $\mathbf{J}_{\text{ba}}(k')$ and $\mathbf{J}_{\text{mar}}(k')$, the matrices $\mathbf{H}(k')$ and $\mathbf{H}(k, k')$ have the same structure, and the only difference lies in the state estimates at which the Jacobians are evaluated. In the matrix $\mathbf{H}(k, k')$ the Jacobians of all measurements that involve marginalized states (the matrices \mathbf{H}_{ij} where $(i, j) \in \mathcal{S}_m$) are evaluated using the state estimates available at time-step k , while all other Jacobians are evaluated using the estimates available at time-step k' . On the other hand, in the matrix $\mathbf{H}(k')$, *all* Jacobians are evaluated using the latest state estimates at time-step k' .

Proceeding further, we note that since $\mathbf{Diag}(\mathbf{R}_{ij}^{-1})$ is a full-rank matrix, we have that $\text{rank}(\mathbf{J}_{\text{ba}}(k')) = \text{rank}(\mathbf{H}(k'))$, and $\text{rank}(\mathbf{J}_{\text{mar}}(k')) = \text{rank}(\mathbf{H}(k, k'))$. Thus, to prove (19), it suffices to show that

$$\text{rank}(\mathbf{H}(k, k')) = \text{rank}(\mathbf{H}(k')) + 3 \quad (22)$$

To this end, we utilize the structure of the measurement Jacobians (see (8)) to factorize the matrix $\mathbf{H}(k')$ as:

$$\mathbf{H}(k') = \mathbf{D}(k') \bar{\mathbf{H}}(k') \quad (23)$$

where $\mathbf{D}(k') = \mathbf{Diag}(\mathbf{H}_{L_{ij}}(k'))$, $(i, j) \in \mathcal{S}$, and $\bar{\mathbf{H}}(k')$ is a matrix with block rows given by

$$\bar{\mathbf{H}}_{ij}(k') = [\mathbf{0} \quad \dots \quad \bar{\mathbf{H}}_{C_{ij}}(k') \quad \dots \quad \mathbf{I}_3 \quad \dots \quad \mathbf{0}] \quad (24)$$

$$\bar{\mathbf{H}}_{C_{ij}}(k') = [[(\hat{\mathbf{p}}_{L_j}(k') - \hat{\mathbf{p}}_{C_i}(k')) \times] \mathbf{C}^T(\hat{\mathbf{q}}_{C_i}(k')) \quad -\mathbf{I}_3]$$

A similar factorization can be obtained for $\mathbf{H}(k, k') = \mathbf{D}(k, k') \bar{\mathbf{H}}(k, k')$. We can now employ the following result, which allows us to compute the rank of the product of two matrices [26, 4.5.1]:

$$\text{rank}(\mathbf{AB}) = \text{rank}(\mathbf{B}) - \dim(\mathcal{N}(\mathbf{A}) \cap \mathcal{R}(\mathbf{B})) \quad (25)$$

where \mathcal{N} and \mathcal{R} denote the null space and the range of a matrix, respectively. An important intermediate result, whose proof is given in [25], is the following:

Lemma 1: When stereo camera measurements are used,

$$\dim(\mathcal{N}(\mathbf{D}(k')) \cap \mathcal{R}(\bar{\mathbf{H}}(k'))) = 0 \quad (26)$$

$$\dim(\mathcal{N}(\mathbf{D}(k, k')) \cap \mathcal{R}(\bar{\mathbf{H}}(k, k'))) = 0 \quad (27)$$

Thus, by combining the results of Lemma 1, along with the decomposition in (23) and the property (25), we see that $\text{rank}(\mathbf{H}(k')) = \text{rank}(\bar{\mathbf{H}}(k'))$ and $\text{rank}(\mathbf{H}(k, k')) = \text{rank}(\bar{\mathbf{H}}(k, k'))$. Thus, to prove (22) it suffices to show that

$$\text{rank}(\bar{\mathbf{H}}(k, k')) = \text{rank}(\bar{\mathbf{H}}(k')) + 3 \quad (28)$$

To prove this result, we first apply elementary row and column operations to the matrix $\bar{\mathbf{H}}(k')$, in order to compute its rank. The intermediate steps are straightforward but the matrices involved are quite large, and therefore the proof is detailed in [25]. There, it is shown that when at least three non-collinear features are available, it is:

$$\text{rank}(\bar{\mathbf{H}}(k')) = 3n' + 6k' \quad (29)$$

Note that the state vector consists of n' landmarks and $k' + 1$ camera poses, which means that the matrix $\bar{\mathbf{H}}(k')$ has $3n' + 6k' + 6$ columns, equal to the number of state variables. Thus, the above result shows that $\bar{\mathbf{H}}(k')$ is rank deficient by 6.

Proceeding with the proof of (28), we apply similar elementary row and column operations to the matrix $\bar{\mathbf{H}}(k, k')$ to compute its rank. However, recall that in this matrix some of the Jacobians are evaluated using the state estimates at time step k , and some using the state estimates at time-step k' . As a result, for certain state variables, *two different* estimates appear in the equations. This means that certain cancellations that occurred when applying elementary row and column operations in $\bar{\mathbf{H}}(k')$ do *not* happen when these operations are applied to $\bar{\mathbf{H}}(k, k')$. As a result, the rank of the matrix is increased. Specifically, in [25] it is shown that:

$$\text{rank}(\bar{\mathbf{H}}(k, k')) = 3n' + 6k' + 3 \quad (30)$$

This result, in conjunction with (29), completes the proof.

B. Physical interpretation

Equation (19) shows that when marginalization takes place the estimator erroneously believes to have information along three more directions of the state space. To identify these directions, we can examine the nullspaces of the matrices $\mathbf{J}_{\text{mar}}(k')$ and $\mathbf{J}_{\text{ba}}(k')$. First, note that in the preceding section it was shown that $\text{rank}(\mathbf{J}_{\text{ba}}(k')) = \text{rank}(\bar{\mathbf{H}}(k')) = 3n' + 6k'$. Since $\mathbf{J}_{\text{ba}}(k')$ is a $(3n' + 6k' + 6) \times (3n' + 6k' + 6)$ matrix, this result means that $\mathbf{J}_{\text{ba}}(k')$ has a nullspace of dimension 6. To obtain a basis for this nullspace, we define the $(3n' + 6k' + 6) \times 6$ matrix

$$\mathbf{N}(\hat{\mathbf{x}}_{\mathbf{m}}(k'), \hat{\mathbf{x}}_{\mathbf{r}}(k'), \hat{\mathbf{x}}_{\mathbf{n}}(k')) = \begin{bmatrix} \mathbf{I}_3 & -[\hat{\mathbf{p}}_{L_1}(k') \times] \\ \vdots & \vdots \\ \mathbf{I}_3 & -[\hat{\mathbf{p}}_{L_{n'}}(k') \times] \\ \mathbf{0}_{3 \times 3} & \mathbf{C}(\hat{\mathbf{q}}_{C_0}(k')) \\ \mathbf{I}_3 & -[\hat{\mathbf{p}}_{C_0}(k') \times] \\ \vdots & \vdots \\ \mathbf{0}_{3 \times 3} & \mathbf{C}(\hat{\mathbf{q}}_{C_{k'}}(k')) \\ \mathbf{I}_3 & -[\hat{\mathbf{p}}_{C_{k'}}(k') \times] \end{bmatrix} \quad (31)$$

It is easy to verify that the following property holds (here we are assuming a state variable ordering of all features followed

by all camera poses):

$$\mathbf{J}_{\text{ba}}(k') \cdot \mathbf{N}(\hat{\mathbf{x}}_{\mathbf{m}}(k'), \hat{\mathbf{x}}_{\mathbf{r}}(k'), \hat{\mathbf{x}}_{\mathbf{n}}(k')) = \mathbf{0}$$

which means that the columns of the matrix \mathbf{N} (which are linearly independent) form a basis for the nullspace of $\mathbf{J}_{\text{ba}}(k')$. The nullspace of the information matrix $\mathbf{J}_{\text{ba}}(k')$ describes changes in the state that cannot be detected using the available measurements (i.e., the unobservable subspace). Close examination of the columns of \mathbf{N} reveals that the first block column corresponds to global translations of the entire state vector, while the second corresponds to global rotations. This should come as no surprise, since by using only measurements of unknown features only the *relative* camera motion can be determined, and not the global pose.

Let us now examine the situation when marginalization takes place. In this case, based on the results of the preceding section, we see that $\text{rank}(\mathbf{J}_{\text{mar}}(k')) = 3n' + 6k' + 3$, which in turn means that the nullspace of $\mathbf{J}_{\text{mar}}(k')$ is only of dimension three. Multiplying $\mathbf{J}_{\text{mar}}(k')$ with the first block column (first three columns) of the matrix \mathbf{N} yields a zero matrix, and thus we conclude that the nullspace of $\mathbf{J}_{\text{mar}}(k')$ is spanned by the first three columns of \mathbf{N} . We thus see that the directions that correspond to the global orientation are “missing” from the nullspace of the information matrix. In other words, the sliding-window VO estimator incorrectly “believes” that the global orientation is observable. As discussed in Section III, this results in inconsistent estimates, and a degradation of accuracy.

IV. IMPROVING THE ESTIMATOR’S PERFORMANCE

In this section we describe a simple modification of the standard sliding-window VO algorithm that prevents the increase of the rank of the measurement information matrix, and improves the performance of the estimator. As shown in Section III-A, the erroneous increase in the rank of the information matrix is caused by the fact that two different estimates of certain states appear in the measurement Jacobians. Specifically, these are the camera poses and/or landmarks in $\mathbf{x}_{\mathbf{r}}$ that are “connected” to marginalized states via measurements in the set \mathcal{S}_m . For example, let us assume that the camera pose \mathbf{c}_{i^*} is one of the camera poses that remains in the sliding window after marginalization at time step k . Moreover, consider that the feature $\mathbf{p}_{L_{j^*}}$ was observed from this camera pose and was marginalized at time step k . Then the information matrix $\mathbf{H}_{i^*j^*}^T(k) \mathbf{R}_{i^*j^*}^{-1} \mathbf{H}_{i^*j^*}(k)$ will appear in the summation in (10), and will be used to compute the matrix $\mathbf{A}_{\mathbf{p}}(k)$. Later on, when we perform iterative estimation at time step k' , the camera pose \mathbf{c}_{i^*} is still in the sliding window, but now the Jacobians that involve this pose are evaluated using the estimate $\hat{\mathbf{c}}_{i^*}(k')$. Thus, two different estimates of \mathbf{c}_{i^*} are used for Jacobian computations.

To avoid this problem, a simple solution is to change the state estimates that are used for Jacobian computations. Specifically, when an active state is connected to already marginalized states via measurements (e.g., \mathbf{c}_{i^*} in the above example), then we use the estimate that was available at the time of marginalization (e.g., $\hat{\mathbf{c}}_{i^*}(k)$) for all subsequent

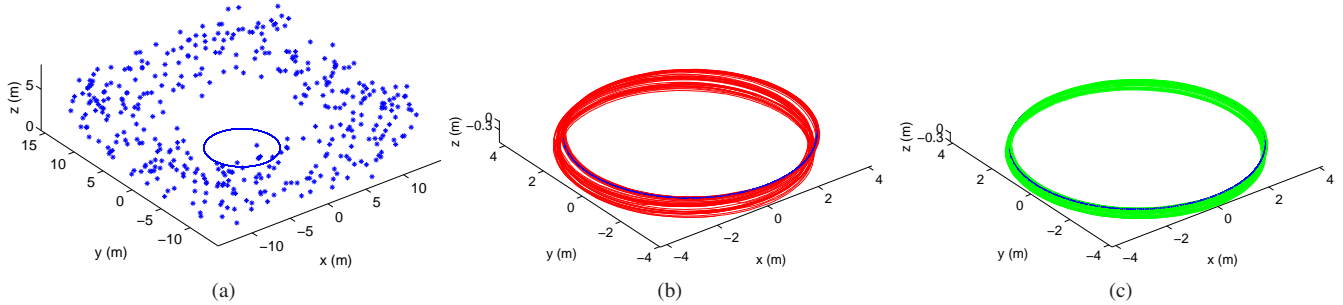


Fig. 1: Simulation setup and example results for stereo-based VO. (a) The simulation environment and the true trajectory (circle) (b) The trajectory estimated by S-VO, in one example Monte-Carlo trial (red) (c) The trajectory estimated by M-VO, in the same trial (green).

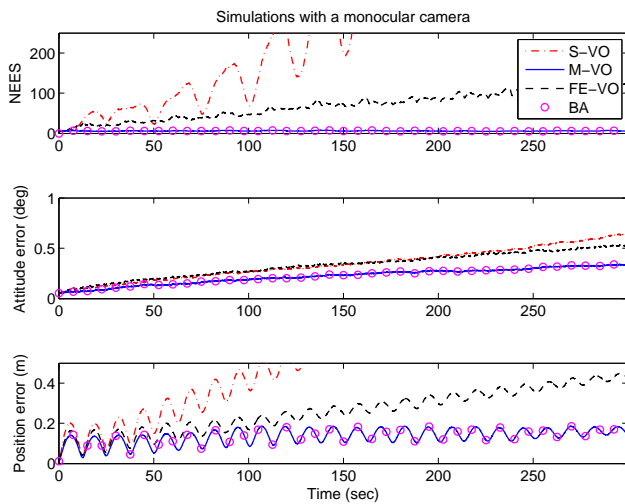


Fig. 2: Simulation results for VO using a monocular camera. From top to bottom: (a) The average value of the camera-pose NEES over time (b) The RMS errors of the camera attitude over time (c) The RMS errors of the camera position over time.

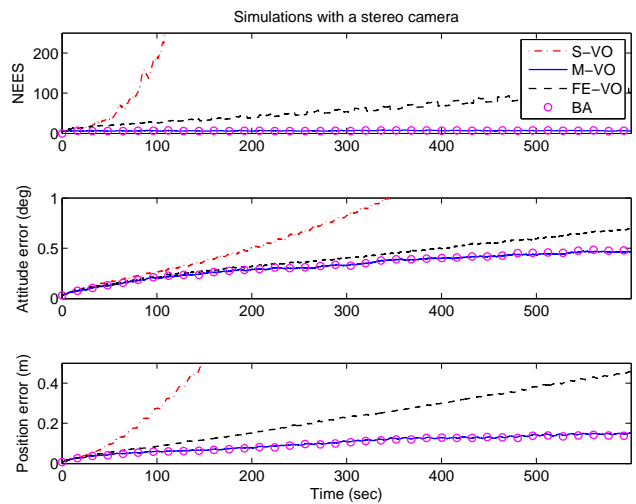


Fig. 3: Simulation results for VO using a stereo camera. From top to bottom: (a) The average value of the camera-pose NEES over time (b) The RMS errors of the camera attitude over time (c) The RMS errors of the camera position over time.

Jacobian computations. In this way only a single estimate of this state appears in the information matrix, and the increase in rank is averted. We stress that we use these “older” estimates when computing the Jacobians *only*, and we still allow the actual estimates to be updated normally in the Gauss-Newton iterations. Clearly, the use of “older” estimates for computing Jacobians will inevitably lead to larger linearization errors. However, as indicated by the results presented in the next section, the effect of this loss of linearization accuracy is not significant, while avoiding the creation of fictitious information leads to significantly improved precision.

V. RESULTS

A. Simulation results

In this section, we present simulation results that demonstrate the performance of the modified sliding-window VO algorithm presented in Section IV. In our simulation setup, we consider a camera (stereo or monocular) that moves along

a circular trajectory of radius 4 m in a $24 \times 24 \times 5$ m room with 600 visual point features randomly placed near the walls (as shown in Fig. 1a). The camera moves with constant velocity and angular velocity of 2 m/s and 0.5 rad/s, respectively. The camera frame rate is 10 Hz and 5 Hz for monocular and stereo camera, respectively, the field of view is 45° , the focal length is 500 pixels, and the measurement noise standard deviation is 1 pixel. For the simulations with a stereo camera, we set the stereo baseline equal to 0.12 m. The features can be observed for up to 30 consecutive camera poses. Therefore, in the sliding-window VO we choose a sliding window containing the 40 latest camera poses and the landmarks seen in these poses. In these simulations, we compare the performance of (i) the modified sliding-window VO algorithm (termed M-VO) presented in Section IV, (ii) sliding-window VO with the standard linearization approach (termed S-VO), (iii) sliding-window VO with fixed estimates for the previous states, similarly to [10] (termed FE-VO), and finally, (iv) the standard bundle adjustment (termed BA)

that estimates the entire history of the camera poses (i.e., no marginalization).

Fig. 2 shows results using a monocular camera for VO, and Fig. 3 shows results for the case of stereo. In both figures, the consistency of the estimators is measured by the average normalized estimation error squared (NEES) for the latest camera pose, and accuracy is measured by the root mean squared (RMS) error of the orientation and position. All results are averaged over 50 Monte-Carlo runs. In these plots, we can observe that the M-VO algorithm clearly outperforms S-VO, both in terms of consistency (i.e., NEES) and accuracy (i.e., RMS errors). When using a monocular camera, the average NEES over all Monte-Carlo runs and all timesteps equals 6.090 for M-VO, very close to 6.017, which is the value obtained for BA. The average NEES for S-VO is a staggering 2350. Since the pose error state is of dimension 6, the “ideal” NEES value would equal 6. Similarly, when using a stereo camera, the average NEES equals 6.351, 6.497, and 6435 for the BA, M-VO, and S-VO algorithms, respectively. Most importantly, in both monocular and stereo VO, we see that the performance of the new M-VO algorithm is *almost indistinguishable* from that of the “golden standard” of BA, even though its computational cost is orders of magnitude smaller.

For illustration purposes, two sample estimated trajectories, computed by the S-VO and M-VO algorithms respectively, are also provided in Fig. 1. This plot clearly demonstrates that the M-VO method can yield results that are substantially more accurate than those of the standard method. In particular, as shown in Fig. 1, the position errors in the vertical axis are significantly larger for S-VO.

We now turn our attention to the FE-VO method. This method shares similar characteristics to our M-VO, in the sense that it uses “older” estimates of previous states in computing Jacobians. This improves the consistency of the method, as shown in Figs. 2 and 3. Specifically, FE-VO has an average NEES of 68.68 and 52.48, respectively, in the monocular and stereo cases. In terms of RMS errors, FE-VO also outperforms S-VO significantly. However, both in terms of NEES and RMS errors, the FE-VO approach performs worse than the proposed M-VO. This can be explained by the fact that in FE-VO the estimates of the older states are fixed and not updated, which degrades accuracy. In contrast, in the M-VO method the older states are updated normally, thus attaining higher precision.

B. Real-world experiment

The performance of the proposed algorithm is also validated in a real-world setting. For this purpose, we tested the algorithm on Epoch A (Campus) of the New College dataset [27], using stereo images. Only the first epoch is used because when the robot passes through a dark tunnel to a different area, no point features can be detected in the captured images. Since the stereo camera is the sole sensor used in this experiment, it is impossible for the algorithm to recover the camera pose. In the section of the dataset that we used, the camera moved for about 7 min, performing three

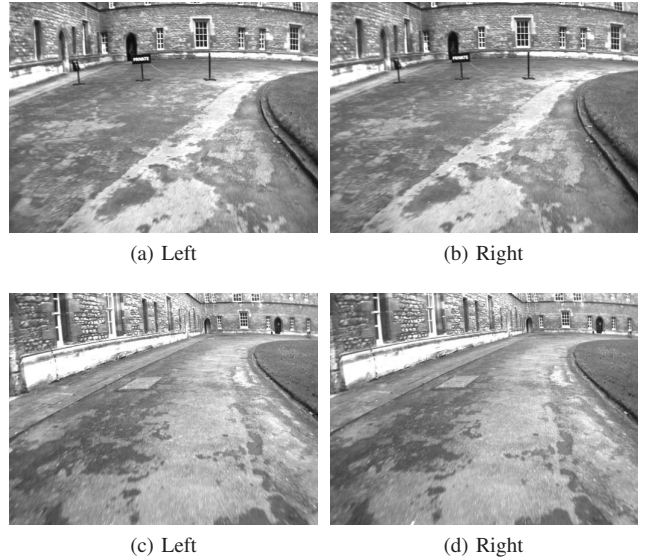


Fig. 4: Sample stereo images from the New College dataset.

loops around the main oval quad in New College, Oxford. The dataset consists of more than 15000 images of resolution 512×384 pixels, captured by a PointGrey Bumblebee stereo rig at 20 Hz. Features are extracted using the Harris corner detector [28], and matched using normalized cross-correlation.

In Fig. 5, the trajectory estimates of the S-VO and the M-VO are shown in red and blue, respectively. Unfortunately, ground truth is not available for the New College dataset, but the camera was driven so that the trajectory in each loop was identical. Compared with previous results on the New College dataset [29], the results obtained by the S-VO and M-VO are similar to the best estimates. It should be noted that only ego-motion is estimated from stereo images, and we do not address loop closure, compared to [29]. By inspection of the trajectory estimates, we can deduce that the position errors of the S-VO are larger than M-VO, both in the x - y plane and along the z -axis. In particular, the side views of the trajectories in Fig. 5b and the elevation estimate plot in Fig. 5c show that the accuracy of the two VO methods is very different in the vertical direction. Since the camera moves on flat ground, its elevation (z coordinate) should remain approximately constant throughout its motion. The results show that by using the prior linearization points to preserve the observability properties of the estimator, we achieve better overall estimation accuracy.

VI. CONCLUSIONS

In this paper, we presented an analysis of the properties of sliding-window minimization for visual odometry. Estimators that employ this approach attain bounded computational cost by marginalizing out older states, so as to maintain an approximately constant number of states active at all times. By analyzing the details of the Jacobian computations needed for the marginalization equations, we have proven that the standard linearization method will introduce erroneous

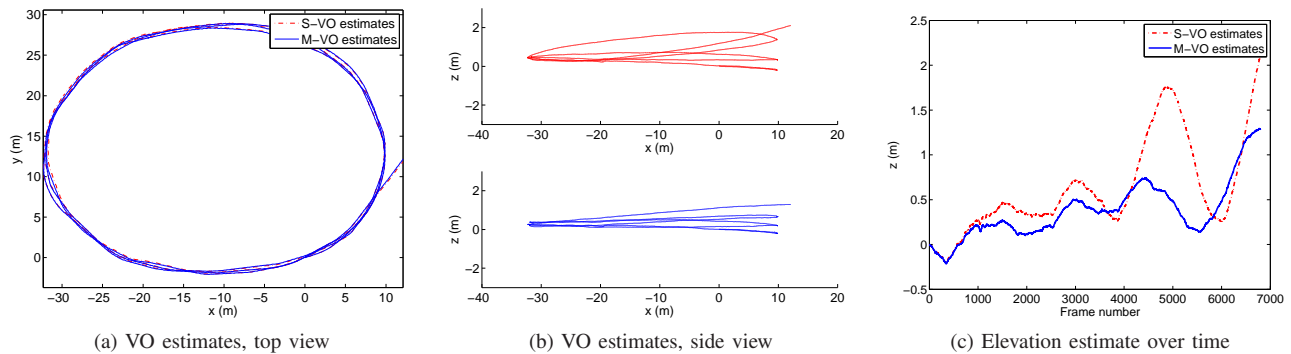


Fig. 5: Real-world results for VO using a stereo camera.

information into the estimator, resulting in inconsistency. Based on this analysis, we proposed a modified linearization scheme, to prevent the infusion of artificial information, and improve estimation performance. Our simulation tests and real-world experiments demonstrated that this modified sliding-window VO estimator outperforms competing methods, both in terms of accuracy and consistency.

ACKNOWLEDGEMENTS

This work was supported by the National Science Foundation (grant no. IIS-1117957), the UC Riverside Bourns College of Engineering, and the Hellman Family Foundation.

REFERENCES

- [1] N. Trawny, A. I. Mourikis, S. I. Roumeliotis, A. E. Johnson, and J. Montgomery, "Vision-aided inertial navigation for pin-point landing using observations of mapped landmarks," *Journal of Field Robotics*, vol. 24, no. 5, pp. 357–378, May 2007.
- [2] R. M. Haralick, C.-N. Lee, K. Ottenberg, and M. Noelle, "Review and analysis of solutions of the three point perspective pose estimation problem," *Intl. Journal of Computer Vision*, vol. 13, no. 3, pp. 331–356, Dec. 1994.
- [3] A. Ansar and K. Daniilidis, "Linear pose estimation from points or lines," *IEEE Transaction on Pattern Analysis and Machine Intelligence*, vol. 25, no. 5, pp. 578–589, May 2003.
- [4] A. J. Davison and D. W. Murray, "Simultaneous localisation and map-building using active vision," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 865–880, Jul. 2002.
- [5] A. Davison, I. Reid, N. Molton, and O. Stasse, "MonoSLAM: Real-time single camera SLAM," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [6] J. Montiel, J. Civera, and A. Davison, "Unified inverse depth parametrization for monocular SLAM," in *Proceedings of Robotics: Science and Systems*, Philadelphia, PA, Aug. 16-19 2006, pp. 81–88.
- [7] D. Nister, O. Naroditsky, and J. Bergen, "Visual odometry," in *Proc. IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*, Washington, DC, June 17-26 2004, pp. 652–659.
- [8] A. Howard, "Real-time stereo visual odometry for autonomous ground vehicles," in *Proc. Intl. Conf. on Intelligent Robots and Systems*, Nice, France, Sept. 22-26 2008, pp. 3946–3952.
- [9] C. Engels, H. Stewenius, and D. Nister, "Bundle adjustment rules," in *Proc. Photogrammetric Computer Vision Conf.*, Bonn, Germany, Sep. 20-22 2006, pp. 266–271.
- [10] D. Nister, O. Naroditsky, and J. Bergen, "Visual odometry for ground vehicle applications," *Journal of Field Robotics*, vol. 23, no. 1, pp. 3–20, Jan. 2006.
- [11] E. Mouragnon, M. Lhuillier, M. Dhome, F. Dekeyser, and P. Sayd, "Real time localization and 3D reconstruction," in *Proc. IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*, New York, NY, June 17-22 2006, pp. 363–370.
- [12] L. H. Matthies, "Dynamic stereo vision," Ph.D. dissertation, School of Computer Science, Carnegie Mellon University, 1989.
- [13] S. I. Roumeliotis, A. E. Johnson, and J. F. Montgomery, "Augmenting inertial navigation with image-based motion estimation," in *Proc. IEEE Intl. Conf. on Robotics and Automation*, Washington D.C, May 2002, pp. 4326–33.
- [14] D. D. Diel, "Stochastic constraints for vision-aided inertial navigation," Master's thesis, Massachusetts Institute of Technology, Jan. 2005.
- [15] B. Triggs, P. McLauchlan, R. Hartley, and Fitzgibbon, "Bundle adjustment – a modern synthesis," in *Vision Algorithms: Theory and Practice*. Springer Verlag, 2000, pp. 298–375.
- [16] M. Kaess, A. Ranganathan, and F. Dellaert, "iSAM: Incremental smoothing and mapping," *IEEE Transactions on Robotics*, vol. 24, no. 6, pp. 1365–1378, 2008.
- [17] K. Konolige, M. Agrawal, and J. Sola, "Large-scale visual odometry for rough terrain," in *Proc. Intl. Symposium on Research in Robotics*, Hiroshima, Japan, Nov. 26-29 2007, pp. 201–212.
- [18] G. Sibley, L. Matthies, and G. Sukhatme, "Sliding window filter with application to planetary landing," *Journal of Field Robotics*, vol. 27, no. 5, pp. 587–608, 2010.
- [19] T.-C. Dong-Si and A. I. Mourikis, "Motion tracking with fixed-lag smoothing: Algorithm and consistency analysis," in *IEEE International Conference on Robotics and Automation*, Shanghai, China, May 9-13 2011, pp. 5655–5662.
- [20] P. McLauchlan, "The variable state dimension filter applied to surface-based structure from motion," School of Electrical Engineering, Information Technology and Mathematics, University of Surrey, UK, Tech. Rep. VSSP-TR-4/99, 1999.
- [21] A. Ranganathan, M. Kaess, and F. Dellaert, "Fast 3D pose estimation with out-of-sequence measurements," in *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems*, San Diego, CA, Oct. 29 - Nov. 2 2007, pp. 2486–2493.
- [22] K. Konolige and M. Agrawal, "FrameSLAM: From bundle adjustment to real-time visual mapping," *IEEE Transactions on Robotics*, vol. 24, no. 5, pp. 1066–1077, oct. 2008.
- [23] Y. Bar-Shalom, X. R. Li, and T. Kirubarajan, *Estimation with Applications to Tracking and Navigation*. John Wiley & Sons, 2001.
- [24] N. Trawny and S. I. Roumeliotis, "Indirect Kalman filter for 3D pose estimation," Dept. of Computer Science & Engineering, University of Minnesota, Minneapolis, MN, Tech. Rep. 2005-002, Mar. 2005.
- [25] T.-C. Dong-Si and A. I. Mourikis, "Consistency analysis for sliding-window visual odometry," Dept. of Electrical Engineering, University of California, Riverside, Tech. Rep., 2011, http://www.ee.ucr.edu/~mourikis/tech_reports/visual_odometry.pdf.
- [26] C. D. Meyer, *Matrix Analysis and Applied Linear Algebra*. SIAM, 2000.
- [27] M. Smith, I. Baldwin, W. Churchill, R. Paul, and P. Newman, "The New College vision and laser data set," *The International Journal of Robotics Research*, vol. 28, no. 5, pp. 595–599, May 2009. [Online]. Available: <http://www.robots.ox.ac.uk/NewCollegeData/>
- [28] C. Harris and M. Stephens, "A combined corner and edge detector," in *Proc. 4th Alvey Vision Conf.*, Manchester, UK, Aug. 31 - Sep. 2 1988, pp. 147–151.
- [29] C. Mei, G. Sibley, M. Cummins, P. Newman, and I. Reid, "RSLAM: A system for large-scale mapping in constant-time using stereo," *International Journal of Computer Vision*, vol. 94, no. 2, pp. 198–214, 2011.