

The Battle for Filter Supremacy: A Comparative Study of the Multi-State Constraint Kalman Filter and the Sliding Window Filter

Lee Clement, Valentin Peretroukhin, Jacob Lambert, and Jonathan Kelly

Abstract—Accurate and consistent egomotion estimation is an **important** part of autonomous navigation. For this task, the combination of visual and inertial sensors is an inexpensive, compact, and complementary hardware suite that can be used on many **ground and aerial vehicles**. In this work, we compare two modern approaches to egomotion estimation: the Multi-State Constraint Kalman Filter (MSCKF) and the Sliding Window Filter (SWF). Both filters use an Inertial Measurement Unit (IMU) to estimate the motion of a vehicle and then correct this estimate with observations of salient features from a monocular camera. While the SWF estimates feature positions as part of the filter state itself, the MSCKF optimizes feature positions in a separate procedure without including them in the filter state. We present experimental characterizations and comparisons of the MSCKF and SWF on data from a moving hand-held sensor rig, as well as several traverses from the KITTI dataset. In particular, we compare the accuracy and consistency of the two filters, and analyze the effect of feature track length and feature density on the performance of each filter. In general, our results show the SWF to be more accurate and less sensitive to tuning parameters than the MSCKF. However, the MSCKF is significantly more computationally efficient, has good consistency properties, and improves in accuracy as more features are tracked.

I. INTRODUCTION

The combination of visual and inertial sensors is a powerful tool for autonomous navigation in unknown environments. Indeed, cameras and inertial measurement units (IMUs) are complementary in several respects. Since an IMU **directly** measures linear accelerations and rotational velocities, these values must be integrated to arrive at a new pose estimate. However, the noise inherent in the IMU's measurements is included in the integration as well, and consequently the pose estimates can drift unbounded over time. The addition of a camera is an excellent way to bound this cumulative drift error because the camera's signal-to-noise ratio is highest when the camera is moving slowly. On the other hand, cameras are not robust to motion blur induced by rapid motions. In these cases, IMU data can be relied upon more heavily **in** estimating egomotion.

The question, then, is how best to fuse measurements from these two sensor types to arrive at an accurate estimate of a vehicle's motion over time. This problem is often complicated by the absence of a known map of features from which the camera can generate measurements. Any solution must therefore solve a Simultaneous Localization and Mapping (SLAM) problem, although the importance

All authors are at the Institute for Aerospace Studies, University of Toronto, Canada {lee.clement, v.peretroukhin, jacob.lambert} @mail.utoronto.ca, jkelly@utias.utoronto.ca

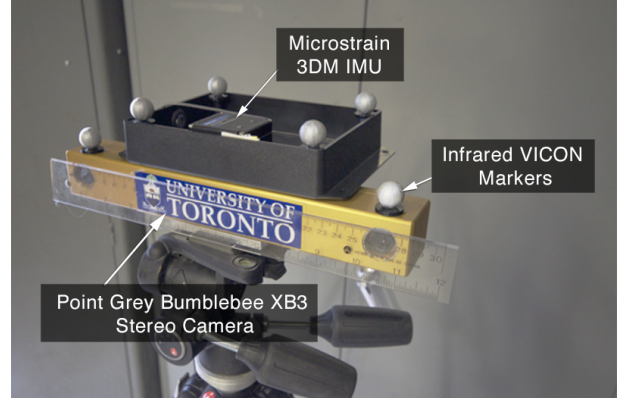


Fig. 1. The hand-held sensor head used in our experiments with the “Starry Night” dataset. The IMU reports translational and rotational velocities, while the stereo camera observes point features. We artificially blinded the stereo camera by using measurements from the left camera only.

placed on the mapping component may vary from algorithm to algorithm.

In this work we characterize, compare, and contrast the performance of two modern solutions to the visual-inertial **SLAM** problem, the Sliding Window Filter (SWF) and the Multi-State Constraint Kalman Filter (MSCKF) [1], [2], on data from a moving hand-held sensor rig, as well as several traverses from the KITTI dataset [3]. The most similar work to ours is that of Leutenegger et al. [4], which compares the accuracy of the MSCKF to a keyframe-based SWF on datasets consisting of relatively planar motion through urban and indoor environments. In contrast to [4], our SWF optimizes over a constant number of timesteps rather than keyframes. We also conduct a more extensive characterization of the sensitivity of the MSCKF to certain **parameters** and compare both algorithms using data from a hand-held sensor rig that mimics the more arbitrary motion of a micro aerial vehicle (MAV).

II. BACKGROUND

Visual-inertial navigation systems (VINS) have been applied broadly in robotics [2], [5]–[10], and there is a considerable body of work covering a wide range of estimation algorithms for the camera-IMU sensor pair. These techniques are often characterized as either *loosely coupled* or *tightly coupled*. In loosely coupled systems, image and IMU measurements are processed individually before being fused into a single estimate, while tightly coupled systems process all information together. The decoupling of inertial

and visual measurements in loosely coupled systems limits computational complexity [4], but at the cost of information: processing camera and IMU measurements separately makes optimal estimation of biases impossible [9]. In this work, we consider tightly coupled algorithms, which are preferable for accurate and consistent visual-inertial navigation.

For both tightly coupled and loosely coupled VINS, a popular estimator is the Extended Kalman filter (EKF) or one of its variants [2], [6], [8], [10], [11], though methods also exist that employ unscented Kalman filters [12], particle filters [13], [14], or batch optimization methods [15], [16].

EKF-SLAM is an efficient recursive algorithm for small-scale, online tasks, but maintaining the entire map as part of the state is a weakness for navigation over long distances. Indeed, the computational complexity of EKF-SLAM scales far too poorly with the number of features to be applied naively to this problem.

Furthermore, EKF-based VINS are inconsistent, that is, the state uncertainties are underestimated. Li and Mourikis [9] showed that the Jacobians in the linearized model of a VINS have different observability properties than the actual nonlinear system, reducing the reliability of such estimators.

Finally, EKF-SLAM is *forgetful*: because the filter state includes only the most recent vehicle pose, a given update step can never modify past poses even if later feature measurements ought to constrain them. By locking in past poses, EKF-SLAM condemns itself to sub-optimally estimating both vehicle motion and feature positions.

The shortcomings of the EKF motivate the use of filters that operate on a subset of the problem. Algorithms operating in constant time with respect to map size are particularly desirable. The SWF analyzed in this work is a modern example of such algorithms, operating in constant time by performing a batch optimization over a set number of states, which is suboptimal but still leads to exceptional accuracy for its cost. This makes the SWF useful for delicate, large scale operations such as planetary landing [5]. The MSCKF [1], [2] can be thought of as a hybrid of EKF-SLAM and the SWF in the sense that it maintains a variable window of poses and applies batch updates using all observations of each landmark. It is also notable for its computational efficiency, achieving accurate, real-time position tracking onboard a smartphone [17]. In Sections III and IV, we discuss both the SWF and MSCKF in detail.

III. SLIDING WINDOW FILTER

The aim of the Sliding Window Filter (SWF) is to estimate a vehicle's motion by optimizing a sliding window of vehicle poses and observed landmarks. The optimization problem in the SWF is typically solved as a non-linear least squares problem using Gauss-Newton (GN) or Levenberg-Marquardt (LM) optimization over a sliding window of K poses $\mathbf{x}_{k \in [1, K]}$ observing M features $f_{j \in [1, M]}$:

$$\mathbf{x} := [\mathbf{x}_0^T \quad \dots \quad \mathbf{x}_K^T \quad \mathbf{p}_G^{f_1 G T} \quad \dots \quad \mathbf{p}_G^{f_M G T}]^T. \quad (1)$$

In our notation, \mathbf{p}_A^{BA} is the vector from the origin of frame \mathcal{F}_A to the origin of frame \mathcal{F}_B expressed in \mathcal{F}_A , \mathbf{C}_{BA} is the

rotation matrix from \mathcal{F}_A to \mathcal{F}_B , and $\mathbf{1}$ is the identity matrix.

The optimization applies an update, $\delta \mathbf{x}^*$, by solving the linear system:

$$(\mathbf{H}^T \mathbf{T}^{-1} \mathbf{H}) \delta \mathbf{x}^* = -\mathbf{H}^T \mathbf{T}^{-1} \mathbf{e}(\bar{\mathbf{x}}). \quad (2)$$

Showing non-zero blocks only, the matrix \mathbf{H} is given by

$$\mathbf{H} = \begin{bmatrix} -\mathbf{H}_{\mathbf{x},1} & \mathbf{1} & -\mathbf{G}_{\mathbf{f},1} \\ \vdots & \vdots & \vdots \\ -\mathbf{H}_{\mathbf{x},K} & \mathbf{1} & -\mathbf{G}_{\mathbf{f},K} \end{bmatrix}, \quad (3)$$

where

$$\mathbf{G}_{\mathbf{x},k} = [\mathbf{G}_{\mathbf{x},k}^1 \quad \dots \quad \mathbf{G}_{\mathbf{x},k}^M]^T, \quad (4)$$

$$\mathbf{G}_{\mathbf{x},k}^j = \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \bigg|_{\bar{\mathbf{p}}_{C_k}^{f_j C_k}} \begin{bmatrix} -\mathbf{C}_{CG} \mathbf{C}_{IG,k} \\ \mathbf{C}_{CG} (\mathbf{C}_{IG,k} (\mathbf{p}_G^{f_j G} - \mathbf{p}_{G,k}^{IG}))^\times \end{bmatrix}^T, \quad (5)$$

$$\mathbf{G}_{\mathbf{f},k} = \begin{bmatrix} \mathbf{G}_{\mathbf{f},k}^1 & & \\ & \mathbf{G}_{\mathbf{f},k}^2 & \\ & \vdots & \\ & & \mathbf{G}_{\mathbf{f},k}^M \end{bmatrix}, \quad (6)$$

j th block column position given by feature ID $j \in [1, M]$

$$\mathbf{G}_{\mathbf{f},k}^j = \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \bigg|_{\bar{\mathbf{p}}_{C_k}^{f_j C_k}} \mathbf{C}_{CI} \mathbf{C}_{IG,k}, \quad (7)$$

with

$$\bar{\mathbf{p}}_{C_k}^{f_j C_k} := \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \mathbf{C}_{CI} \left(\mathbf{C}_{IG,k} (\mathbf{p}_G^{f_j G} - \mathbf{p}_{G,k}^{IG}) - \mathbf{p}_I^{CI} \right), \quad (8)$$

$$\frac{\partial \mathbf{g}}{\partial \mathbf{p}} \bigg|_{\bar{\mathbf{p}}_{C_k}^{f_j C_k}} = \frac{1}{z^2} \begin{bmatrix} f_u z & 0 & -f_u x \\ 0 & f_v z & -f_v y \end{bmatrix}. \quad (9)$$

The weight matrix \mathbf{T} is given by

$$\mathbf{T} := \text{diag} \{ \mathbf{T}_1, \dots, \mathbf{T}_K \}, \quad (10)$$

where

$$\mathbf{T}_k := \begin{bmatrix} \mathbf{H}_{\mathbf{w},k} \mathbf{Q}_k \mathbf{H}_{\mathbf{w},k}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{G}_{\mathbf{n},k} \mathbf{R}_k \mathbf{G}_{\mathbf{n},k}^T \end{bmatrix}, \quad (11)$$

and \mathbf{R}_k , \mathbf{Q}_k are observation and motion model noise covariances, and $\mathbf{G}_{\mathbf{n},k}$, $\mathbf{H}_{\mathbf{w},k}$ are the observation and motion model Jacobians with respect to the noise. An important advantage of the SWF is that its computational cost depends on the number of features in the current window rather than the number of features in the entire map. By varying the spatial or temporal extent of the sliding window, the computational cost of the algorithm can be tailored to fit a given compute envelope, which makes the algorithm suitable for online operation over paths of arbitrary length.

However, the hard cut-off of the SWF may result in only a subset of measurements of each feature contributing to the optimization. As a result, the filter may not maximally constrain some poses, and hence localization may be less accurate than one would expect from the full batch solution.

IV. MULTI-STATE CONSTRAINT KALMAN FILTER

The Multi-State Constraint Kalman Filter (MSCKF) [1], [2] can be thought of as a hybrid of EKF-SLAM and the SWF. The key idea of the MSCKF is to maintain a variable window of vehicle poses and to simultaneously update each pose in the window using batch-optimized estimates of features that are visible across the entire window. This update step typically occurs when a tracked feature goes out of view of the camera, but it may also be triggered if the length of a feature track exceeds a preset threshold. Except where otherwise noted (i.e., state parametrization in IV-B and integration method in IV-C), we implemented the MSCKF as presented in [1], [2].

A. MSCKF State Parametrization

We evaluated the MSCKF on datasets in which the IMU ‘measures’ gravity-corrected linear velocities rather than raw linear accelerations (see Section V). To accommodate this alternative sensor configuration, we parametrize the IMU state at time k as the 13-dimensional vector

$$\mathbf{x}_{I,k} := [\mathbf{q}_{IG,k}^T \quad \mathbf{b}_{\omega,k}^T \quad \mathbf{b}_{\mathbf{v},k}^T \quad \mathbf{p}_{G,k}^{IG,T}]^T \quad (12)$$

where $\mathbf{q}_{IG,k}$ is the unit quaternion representing the rotation from the global frame \mathcal{F}_G to the IMU frame \mathcal{F}_I , $\mathbf{b}_{\omega,k}$ is the bias on the gyro measurements ω_m , $\mathbf{b}_{\mathbf{v},k}$ is the bias on the velocity measurements \mathbf{v}_m , and $\mathbf{p}_{G,k}^{IG}$ is the vector from the origin of \mathcal{F}_G to the origin of \mathcal{F}_I expressed in \mathcal{F}_G (i.e., the position of the IMU in the global frame).

At time k , the full state of the MSCKF consists of the current IMU state estimate, and estimates of N 7-dimensional past camera poses in which active feature tracks were visible:

$$\hat{\mathbf{x}}_k := [\hat{\mathbf{x}}_{I,k}^T \quad \hat{\mathbf{q}}_{C_1G}^T \quad \hat{\mathbf{p}}_G^{C_1G,T} \quad \dots \quad \hat{\mathbf{q}}_{C_NG}^T \quad \hat{\mathbf{p}}_G^{C_NG,T}]^T.$$

We can also define the MSCKF error state at time k :

$$\tilde{\mathbf{x}}_k := [\tilde{\mathbf{x}}_{I,k}^T \quad \delta\boldsymbol{\theta}_{C_1}^T \quad \tilde{\mathbf{p}}_G^{C_1G,T} \quad \dots \quad \delta\boldsymbol{\theta}_{C_N}^T \quad \tilde{\mathbf{p}}_G^{C_NG,T}]^T$$

where

$$\tilde{\mathbf{x}}_{I,k} := [\delta\boldsymbol{\theta}_I^T \quad \tilde{\mathbf{b}}_{\omega,k}^T \quad \tilde{\mathbf{b}}_{\mathbf{v},k}^T \quad \tilde{\mathbf{p}}_{G,k}^{IG,T}]^T \quad (13)$$

is the 12-dimensional IMU error state. In the above, \tilde{x} denotes the difference between the true value and the estimated value of the quantity x . The rotational errors $\delta\boldsymbol{\theta}$ are defined according to

$$\delta\mathbf{q} := \hat{\mathbf{q}}^{-1} \otimes \mathbf{q} \simeq [\frac{1}{2}\delta\boldsymbol{\theta}^T \quad 1]^T \quad (14)$$

where \otimes denotes quaternion multiplication.

Accordingly, the MSCKF state covariance $\hat{\mathbf{P}}_k$ is a $(12 + 6N) \times (12 + 6N)$ matrix that may be partitioned as

$$\hat{\mathbf{P}}_k = \begin{bmatrix} \hat{\mathbf{P}}_{II,k} & \hat{\mathbf{P}}_{IC,k} \\ \hat{\mathbf{P}}_{IC,k}^T & \hat{\mathbf{P}}_{CC,k} \end{bmatrix} \quad (15)$$

where $\hat{\mathbf{P}}_{II,k}$ is the 12×12 covariance matrix of the current IMU state, $\hat{\mathbf{P}}_{CC,k}$ is the $6N \times 6N$ covariance matrix of the camera poses, and $\hat{\mathbf{P}}_{IC,k}$ is the $12 \times 6N$ cross-correlation between the current IMU state and the past camera poses.

B. MSCKF State Augmentation

When a new camera image becomes available, the MSCKF state must be augmented with the current camera pose. We obtain the camera pose by applying the known transformation $(\mathbf{q}_{CI}, \mathbf{p}_I^{CI})$ to a copy of the current IMU pose:

$$\hat{\mathbf{q}}_{C_{N+1}G} = \mathbf{q}_{CI} \otimes \hat{\mathbf{q}}_{IG,k} \quad (16)$$

$$\hat{\mathbf{p}}_G^{C_{N+1}G} = \hat{\mathbf{p}}_G^{IG} + \hat{\mathbf{C}}_{IG,k}^T \hat{\mathbf{p}}_I^{CI} \quad (17)$$

where $\hat{\mathbf{C}}_{IG,k}$ is the rotation matrix corresponding to $\hat{\mathbf{q}}_{IG,k}$.

Assuming the MSCKF state has already been augmented by N camera poses, we add the $(N+1)^{\text{th}}$ camera pose to the state according to

$$\hat{\mathbf{x}}_k \leftarrow \begin{bmatrix} \hat{\mathbf{x}}_k^T & \hat{\mathbf{q}}_{C_{N+1}G}^T & \hat{\mathbf{p}}_G^{C_{N+1}G,T} \end{bmatrix}^T \quad (18)$$

and augment the MSCKF state covariance according to

$$\hat{\mathbf{P}}_k \leftarrow \begin{bmatrix} \mathbf{1}_{12+6N} \\ \mathbf{J}_k \end{bmatrix} \hat{\mathbf{P}}_k \begin{bmatrix} \mathbf{1}_{12+6N} \\ \mathbf{J}_k \end{bmatrix}^T \quad (19)$$

where the Jacobian \mathbf{J}_k is given by

$$\mathbf{J}_k = \begin{bmatrix} \hat{\mathbf{C}}_{CI,k} & \mathbf{0}_{3 \times 6} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 6N} \\ \left(\hat{\mathbf{C}}_{IG,k}^T \mathbf{p}_{I,k}^{CI} \right)^\times & \mathbf{0}_{3 \times 6} & \mathbf{1}_3 & \mathbf{0}_{3 \times 6N} \end{bmatrix} \quad (20)$$

with $\mathbf{1}_m$ denoting the m -dimensional identity matrix and $\mathbf{0}_{n \times p}$ an $n \times p$ matrix of zeros.

C. IMU State Estimate Propagation

The evolution of the mean estimated IMU state $\hat{\mathbf{x}}_I$ over time is described by a continuous-time motion model:

$$\begin{aligned} \dot{\hat{\mathbf{q}}}_{IG} &= \frac{1}{2} \boldsymbol{\Omega}(\hat{\boldsymbol{\omega}}) \hat{\mathbf{q}}_{IG} & \dot{\hat{\mathbf{b}}}_{\omega} &= \mathbf{0}_{3 \times 1} \\ \dot{\hat{\mathbf{p}}}_G^{IG} &= \hat{\mathbf{C}}_{IG}^T \hat{\mathbf{v}} & \dot{\hat{\mathbf{b}}}_{\mathbf{v}} &= \mathbf{0}_{3 \times 1} \end{aligned} \quad (21)$$

where $\hat{\mathbf{C}}_{IG}$ is the rotation matrix corresponding to $\hat{\mathbf{q}}_{IG}$,

$$\hat{\mathbf{v}} = \mathbf{v}_m - \hat{\mathbf{b}}_{\mathbf{v}}, \quad \hat{\boldsymbol{\omega}} = \boldsymbol{\omega}_m - \hat{\mathbf{b}}_{\omega},$$

$$\boldsymbol{\Omega}(\hat{\boldsymbol{\omega}}) = \begin{bmatrix} -\hat{\boldsymbol{\omega}}^\times & \hat{\boldsymbol{\omega}} \\ -\hat{\boldsymbol{\omega}}^T & 0 \end{bmatrix}, \quad \text{and} \quad \hat{\boldsymbol{\omega}}^\times = \begin{bmatrix} 0 & -\hat{\omega}_3 & \hat{\omega}_2 \\ \hat{\omega}_3 & 0 & -\hat{\omega}_1 \\ -\hat{\omega}_2 & \hat{\omega}_1 & 0 \end{bmatrix}.$$

In our implementation we propagate the motion model using a simple forward-Euler integration rather than the fifth-order Runge-Kutta procedure used in [2].

We can also examine the linearized continuous-time model of the IMU error state:

$$\dot{\tilde{\mathbf{x}}}_I = \mathbf{F} \tilde{\mathbf{x}}_I + \mathbf{G} \mathbf{n}_I \quad (22)$$

where the Jacobians \mathbf{F} , \mathbf{G} are given by

$$\mathbf{F} = \begin{bmatrix} -\hat{\boldsymbol{\omega}}^\times & -\mathbf{1}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ -\hat{\mathbf{C}}_{IG}^T \hat{\mathbf{v}}^\times & \mathbf{0}_{3 \times 3} & -\hat{\mathbf{C}}_{IG}^T & \mathbf{0}_{3 \times 3} \end{bmatrix} \quad (23)$$

$$\mathbf{G} = \begin{bmatrix} -\mathbf{1}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{1}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{1}_3 \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & -\hat{\mathbf{C}}_{IG}^T & \mathbf{0}_{3 \times 3} \end{bmatrix}, \quad (24)$$

and $\mathbf{n}_I = [\mathbf{n}_\omega^T \ \mathbf{n}_{\mathbf{b}_\omega}^T \ \mathbf{n}_v^T \ \mathbf{n}_{\mathbf{b}_v}^T]^T$ is the IMU process noise, which has covariance matrix \mathbf{Q}_I .

D. MSCKF State Covariance Propagation

With reference to the partitions defined in (15), we compute the predicted camera state covariances and IMU-camera cross-correlations as follows:

$$\hat{\mathbf{P}}_{CC,k+1}^- = \hat{\mathbf{P}}_{CC,k} \quad (25)$$

$$\hat{\mathbf{P}}_{IC,k+1}^- = \Phi(t_k + \Delta T, t_k) \hat{\mathbf{P}}_{IC,k} \quad (26)$$

where ΔT is the IMU sampling period.

The state transition matrix $\Phi(t_k + \Delta T, t_k)$ and the predicted IMU state covariance $\hat{\mathbf{P}}_{II,k+1}^-$ are computed according to [18]:

$$\Phi(t_k + \Delta T, t_k) = \mathbf{1}_{12} + \mathbf{F}\Delta T \quad (27)$$

$$\begin{aligned} \hat{\mathbf{P}}_{II,k+1}^- &= \Phi(t_k + \Delta T, t_k) \hat{\mathbf{P}}_{II,k} \Phi^T(t_k + \Delta T, t_k) \\ &+ \mathbf{G}\mathbf{Q}_I\mathbf{G}^T\Delta T. \end{aligned} \quad (28)$$

To improve numerical stability, we found it useful to enforce the positive semi-definiteness of $\hat{\mathbf{P}}_{II,k+1}^-$ by making the replacements

$$P_{ii} \leftarrow |P_{ii}| \quad \text{and} \quad P_{ij}, P_{ji} \leftarrow \frac{1}{2}(P_{ij} + P_{ji}), i \neq j,$$

where P_{ij} denotes the i^{th} row and j^{th} column of $\hat{\mathbf{P}}_{II,k+1}^-$.

E. Feature Position Estimation

When a tracked feature f_j is selected for a state update, the MSCKF estimates its position $\hat{\mathbf{p}}_G^{f_j G}$ using an inverse-depth least-squares Gauss-Newton optimization. The procedure takes as input N camera poses and N sets of “ideal” pixel measurements, where “ideal” means that the pixel measurements have been corrected for the camera intrinsics:

$$\hat{\mathbf{z}}_i^{(j)} = [u'_i \ v'_i]^T = [(u_i - c_u)/f_u \ (v_i - c_v)/f_v]^T. \quad (29)$$

We initialize the optimization by estimating the position of feature f_j in camera frame C_1 using a linear least-squares method with measurements from the first two camera frames, C_1 and C_2 :

$$\hat{\mathbf{p}}_{C_1}^{f_j C_1} := [\hat{X}_{C_1}^{(j)} \ \hat{Y}_{C_1}^{(j)} \ \hat{Z}_{C_1}^{(j)}]^T = \lambda \rho_{C_1}^{f_j} \quad (30)$$

where

$$\rho_{C_i}^{f_j} := \frac{1}{\sqrt{u_i'^2 + v_i'^2 + 1}} [u'_i \ v'_i \ 1]^T \quad (31)$$

is the direction of the ray emanating from camera C_i along which feature f_j must lie, and

$$\lambda = [(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \hat{\mathbf{p}}_{C_1}^{C_2 C_1}]_1 \quad (32)$$

with

$$\mathbf{A} := [\rho_{C_1}^{f_j} \ -\rho_{C_2}^{f_j}]. \quad (33)$$

We can express the feature position in camera frame C_i in terms of its position in camera frame C_1 as

$$\hat{\mathbf{p}}_{C_i}^{f_j C_i} = \hat{\mathbf{C}}_{i1} \hat{\mathbf{p}}_{C_1}^{f_j C_1} + \hat{\mathbf{p}}_{C_i}^{C_1 C_i}. \quad (34)$$

By forming the vector

$$\hat{\mathbf{y}} := [\alpha \ \beta \ \gamma]^T := \frac{1}{\hat{Z}_{C_1}^{(j)}} [\hat{X}_{C_1}^{(j)} \ \hat{Y}_{C_1}^{(j)} \ 1]^T, \quad (35)$$

we can define three functions

$$\begin{bmatrix} h_1(\hat{\mathbf{y}}) \\ h_2(\hat{\mathbf{y}}) \\ h_3(\hat{\mathbf{y}}) \end{bmatrix} = \hat{\mathbf{C}}_{i1} \begin{bmatrix} \alpha \\ \beta \\ 1 \end{bmatrix} + \gamma \hat{\mathbf{p}}_{C_i}^{C_1 C_i} \quad (36)$$

and rewrite the camera measurement error as

$$\mathbf{e}(\hat{\mathbf{y}}) = \hat{\mathbf{z}}_i^{(j)} - \frac{1}{h_3(\hat{\mathbf{y}})} \begin{bmatrix} h_1(\hat{\mathbf{y}}) \\ h_2(\hat{\mathbf{y}}) \end{bmatrix}. \quad (37)$$

The least-squares system then becomes

$$(\mathbf{E}^T \mathbf{W}^{-1} \mathbf{E}) \delta \mathbf{y}^* = -\mathbf{E}^T \mathbf{W}^{-1} \mathbf{e}(\hat{\mathbf{y}}) \quad (38)$$

where

$$\mathbf{E} = \frac{\partial \mathbf{e}}{\partial \hat{\mathbf{y}}} \quad \text{and} \quad \mathbf{W} = \text{diag} \{ \mathbf{R}_1^{(j)}, \dots, \mathbf{R}_N^{(j)} \} \quad (39)$$

with $\mathbf{R}_i^{(j)} = \text{diag} \{ \sigma_{u'}^2, \sigma_{v'}^2 \}$.

F. MSCKF State Correction

Now that we have estimated the positions of any features to be used in the state update, we can apply the corresponding motion constraints to the window of poses from which each feature was observed. We begin by forming the exteroceptive measurement error corresponding to an observation $\mathbf{z}_i^{(j)}$ of feature f_j from camera pose C_i :

$$\mathbf{r}_i^{(j)} := \mathbf{z}_i^{(j)} - \hat{\mathbf{z}}_i^{(j)} \quad (40)$$

where

$$\hat{\mathbf{z}}_i^{(j)} = \frac{1}{\hat{Z}_{C_i}^{(j)}} [\hat{X}_{C_i}^{(j)} \ \hat{Y}_{C_i}^{(j)}]^T \quad (41)$$

with

$$\begin{aligned} \hat{\mathbf{p}}_{C_i}^{f_j C_i} &= [\hat{X}_{C_i}^{(j)} \ \hat{Y}_{C_i}^{(j)} \ \hat{Z}_{C_i}^{(j)}]^T \\ &= \hat{\mathbf{C}}_{C_i G} (\hat{\mathbf{p}}_G^{f_j G} - \hat{\mathbf{p}}_G^{C_i G}). \end{aligned} \quad (42)$$

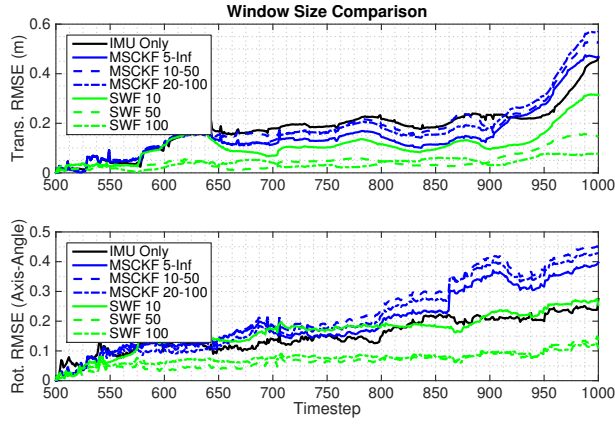
If we linearize (40) about the estimates for the camera pose and feature position, we obtain an estimate of the exteroceptive measurement error

$$\mathbf{r}_i^{(j)} \simeq \mathbf{H}_{\mathbf{x},i}^{(j)} \tilde{\mathbf{x}}_i + \mathbf{H}_{f,i}^{(j)} \tilde{\mathbf{p}}_G^{f_j G} + \mathbf{n}_i^{(j)} \quad (43)$$

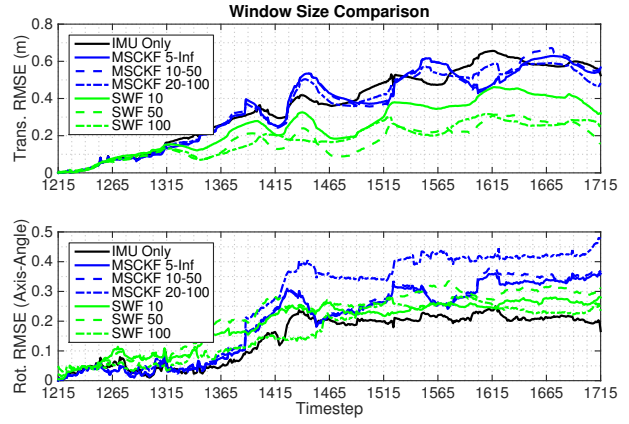
where $\mathbf{H}_{\mathbf{x},i}^{(j)}$ and $\mathbf{H}_{f,i}^{(j)}$ are the Jacobians of the measurement of feature f_j from camera pose C_i with respect to the filter state and the position of the feature, respectively:

$$\mathbf{H}_{\mathbf{x},i}^{(j)} = [\mathbf{0} \ \mathbf{J}_i^{(j)} (\hat{\mathbf{p}}_{C_i}^{f_j C_i})^\times \ -\mathbf{J}_i^{(j)} \hat{\mathbf{C}}_{C_i G} \ \mathbf{0}] \quad (44)$$

$$\mathbf{H}_{f,i}^{(j)} = \mathbf{J}_i^{(j)} \hat{\mathbf{C}}_{C_i G} \quad (45)$$



(a) At least three features are visible in 64% of the interval (500, 1000). The number of visible features is frequently exceeds ten.



(b) At least three features are visible in only 56% of the interval (1215, 1715). The number of visible features rarely exceeds ten.

Fig. 3. Comparison of MSCKF, SWF, and IMU integration for multiple parameter settings on two intervals of the “Starry Night” dataset. The numbers next to “MSCKF” in the legend refer to the minimum and maximum feature track lengths before an EKF update is triggered, while the numbers next to “SWF” refer to the number of states in the sliding window.

for both the sensor head motion and the feature positions. We conducted three experiments on this dataset to compare the effect of feature visibility and window size on each algorithm. We discuss each of these in turn.

1) *Several visible features*: In the first experiment, we compared the SWF and MSCKF for various parameter settings on an interval with many visible features. At least three features are visible in 64% of this interval, and the number of visible features is frequently in excess of ten.

Figure 3(a) shows translational and rotational root mean squared errors (RMSE) for pure IMU integration and various parameter settings for the MSCKF and SWF. For the MSCKF, we varied the maximum number of observations before triggering an update and the minimum number of observations for a feature track to be used in an update. For the SWF, we varied the number of states in the window. These parameters did not significantly affect the accuracy of either filter, however Figure 3(a) shows some small gains in accuracy for larger feature track lengths and window sizes.

On this interval, the SWF outperforms the MSCKF in terms of both translational and rotational RMSE. The MSCKF does not perform much better than pure IMU integration on this interval, likely due to the overall low number of feature tracks in this dataset.

2) *Fewer visible features*: In the second experiment, we compared the SWF and MSCKF for various parameter settings on an interval with few visible features. At least three features are visible in only 56% of this interval, and the number of visible features rarely exceeds ten. Moreover, the features that are visible do not remain in view of the camera for as long a time as in the first interval.

Figure 3(b) shows translational and rotational RMSE for pure IMU integration and the same MSCKF and SWF parameter settings as in Experiment 1. Compared to Experiment 1, the difference in performance between the MSCKF and SWF is less pronounced, with the exception of the SWF’s clearly superior translational accuracy. Neither filter performs substantially better than pure IMU integration on

this interval, again likely due to the overall low number of feature tracks in this dataset.

3) *Effect of feature density in synthetic maps*: In order to investigate the effect of feature density on the performance of the MSCKF and SWF, we modified the “Starry Night” dataset by creating synthetic feature distributions with larger spatial extents and more features than the original dataset. We constructed each dataset so that the larger maps contained the same features as the smaller maps, plus additional features to make up the difference. In each dataset, we retained the IMU data from the original dataset and corrupted the synthetic camera measurements with zero-mean Gaussian noise. By generating longer feature tracks and increasing the number of visible features, these synthetic datasets allowed for a clearer comparison between the MSCKF and SWF.

Figure 4 compares translational and rotational RMSE of the MSCKF, SWF, and pure IMU integration on three synthetic datasets with 40, 60, and 100 features. With more features, both algorithms consistently outperform pure IMU integration, and the SWF outperforms the MSCKF by a wide margin on all three datasets.

Note that the MSCKF’s performance improves as the number of features increases, while the SWF’s performance is not significantly affected by increasing feature count. This result indicates that the MSCKF is much more sensitive to feature density than the SWF. This may be due to the fact that the MSCKF updates its state whenever a feature goes out of view and does not associate tracks corresponding to the same feature. If observations of a given feature are frequently interrupted, each track will not constrain the full set of poses from which the feature is visible. As feature density increases, the number of long feature tracks is likely to increase, and so the MSCKF benefits from more high-quality motion constraints. Since the SWF always associates all observations of each feature in a given window, it is not sensitive to the contiguity of the feature observations.

Table I summarizes each algorithm’s performance on the interval shown in Figure 4. As noted above, the SWF

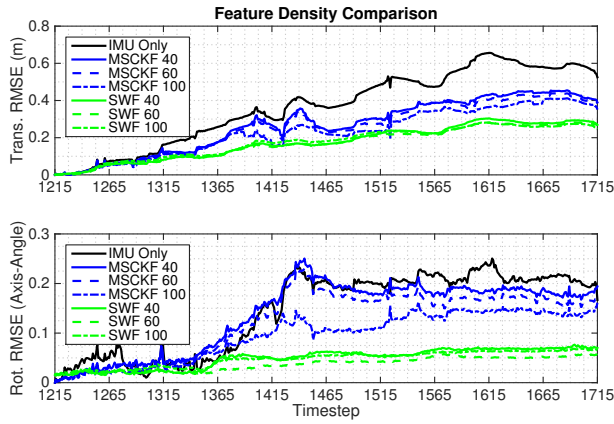


Fig. 4. Comparison of Root Mean Squared Error (RMSE) of the MSCKF, SWF, and pure IMU integration for three synthetic maps generated using the IMU data in the interval (1215, 1715). The numbers next to “MSCKF” and “SWF” refer to the number of features in the dataset (40, 60 or 100).

TABLE I

COMPARISON OF AVERAGE ROOT MEAN SQUARED ERROR (ARMSE), AVERAGE NORMALIZED ESTIMATION ERROR SQUARED (ANEES), AND COMPUTE TIME OF IMU INTEGRATION, MSCKF, AND SWF ON SYNTHETIC DATASETS FOR THE INTERVAL (1215, 1715).

		Feature Count		
		40	60	100
IMU Only	Trans. ARMSE	0.3679	0.3679	0.3679
	Rot. ARMSE	0.1452	0.1452	0.1452
	ANEES	0.2850	0.2850	0.2850
	Compute Time [†]	8.90 s	8.90 s	8.90 s
MSCKF (20-100)	Trans. ARMSE	0.2672	0.2550	0.2304
	Rot. ARMSE	0.1378	0.1247	0.0952
	ANEES	10.18	12.03	16.76
	Compute Time [†]	12.19 s	14.64 s	20.58 s
SWF (25)	Trans. ARMSE	0.1750	0.1687	0.1755
	Rot. ARMSE	0.0495	0.0377	0.0481
	ANEES	2280	2093	2013
	Compute Time [†]	114.3 s	175.9 s	245.3 s

[†] Running MATLAB 2014b on a MacBook Pro Retina (11,3) with a 2.3 GHz Intel Core i7 processor and 16 GB of DDR3L RAM.

outperforms the MSCKF in terms of average RMSE on both rotation and translation. It is worth noting, however, that, on average, the MSCKF achieves low Normalized Estimation Error Squared (NEES) values, which indicates that the filter scores well on consistency if not on accuracy. The average NEES values for the SWF are substantially higher because the SWF treats each window independently and has no built-in mechanism for propagating uncertainty from window to window. The authors of [5] describe how this can be remedied by marginalizing old poses, but we did not implement this technique because our intent was to directly compare the two algorithms in their simplest forms.

Another point in favour of the MSCKF is that the computational effort required was an order of magnitude smaller than for the SWF. The MSCKF may therefore be better suited to **vehicles** with limited computational resources, particularly if they are operating in feature-rich environments.

TABLE II

COMPARISON OF AVERAGE ROOT MEAN SQUARED ERROR (ARMSE) AND AVERAGE NORMALIZED ESTIMATION ERROR SQUARED (ANEES) OF IMU INTEGRATION, MSCKF, AND SWF ON KITTI TRAVERSES.

		KITTI Traverse			
		0001	0036	0051	0095
IMU Only	Trans. ARMSE	0.7197	0.5131	0.7834	1.039
	ANEES	0.1630	0.0092	0.1170	0.6254
MSCKF (5-Inf)	Trans. ARMSE	0.3492	0.4401	0.7530	0.8170
	ANEES	5.103	1.826	2.031	14.98
SWF (10)	Trans. ARMSE	0.3372	0.3778	0.5832	0.7196
	ANEES	358.3	703.2	1124	3767

B. KITTI Dataset

In addition to the “Starry Night” datasets, we also compared the performance of the MSCKF and SWF over four traverses from the KITTI dataset [3] totaling 534 m of urban driving. We selected these particular traverses because they contained few moving objects, which we found to be common failure cases for both filters. To accommodate our alternative state parametrization, we used the pre-integrated linear velocity measurements provided in the dataset instead of the raw linear accelerations from the IMU. We extracted between 50 and 100 salient point features from the left camera in the stereo pair using Speeded Up Robust Features (SURF) [19] and tracked them temporally using Kanade-Lucas-Tomasi (KLT) tracking [20]. We rejected outliers using an M-estimator Sample Consensus (MSAC) [21] procedure with a 2-point similarity transform.

Figure 5 shows the RMSE of pure IMU integration, the MSCKF, and the SWF over each of the four traverses we tested, representative images from which are shown in Figure 6. Table II summarizes these results and reports NEES values for each traverse. Similarly to the “Starry Night” results, both algorithms outperformed pure IMU integration, and the SWF slightly outperformed the MSCKF in terms of translational RMSE. We did not consider rotational RMSE for the KITTI traverses because ground truth was obtained from GPS and did not provide reliable estimates of the entire 6DOF vehicle pose. As expected, the reported ANEES values show that the MSCKF produced estimates that were more consistent than those of the SWF.

VI. CONCLUSIONS

On the datasets we tested, the SWF slightly outperformed the MSCKF, but the MSCKF improved in accuracy with additional features while the SWF was less sensitive to feature quantity. In relatively featureless environments, neither filter performed substantially better than pure IMU integration.

We stress that our inertial data was obtained from high-quality IMUs and has been sanitized to account for gravity, biases, and integration error. Using a consumer-grade IMU, we expect that pure IMU integration would have performed much worse, and that the benefit of the SWF and MSCKF would have been more apparent in these cases.

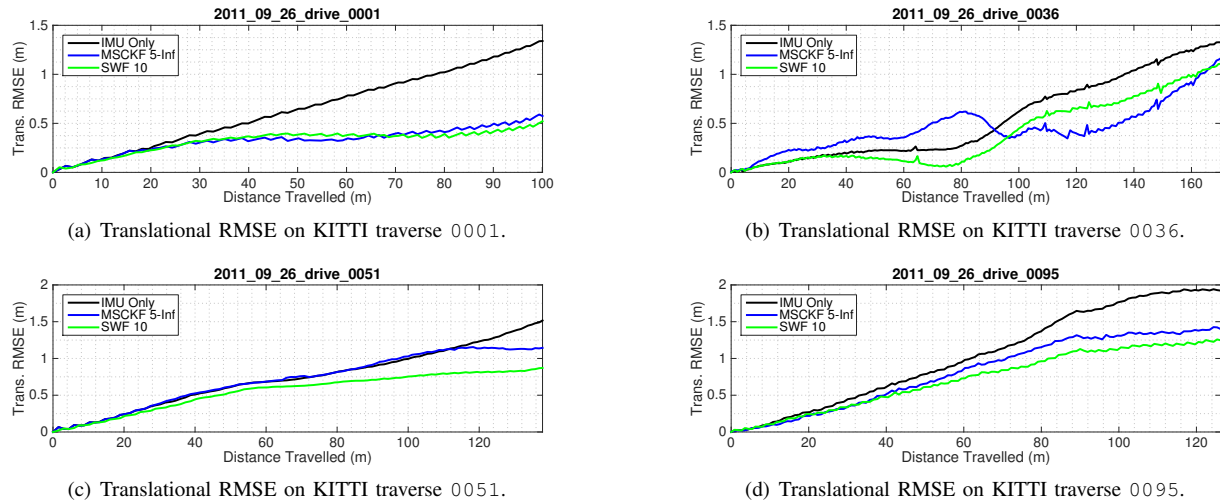


Fig. 5. Translational RMSE over four traverses from the KITTI dataset [3] totaling 534 m of urban driving.



Fig. 6. Sample frames from the four KITTI traverses shown in Figure 5.

Although the SWF appears to produce more accurate pose estimates than the MSCKF in many cases, the MSCKF is less computationally intensive than the SWF and has better consistency properties in its most basic form. However, in our experiments we found the MSCKF to be much more sensitive to tuning parameters than the SWF, sometimes diverging wildly with small changes in parameters. We conclude that the MSCKF may be a better choice of algorithm when computational resources are limited and the operational environment is feature-rich, but that the SWF may be preferable in situations where robustness is paramount.

In future work we would like to investigate the sensitivity of these algorithms to variations in other parameters such as scene geometry, frame rate, field of view, and camera type (e.g., stereo and omnidirectional cameras).

ACKNOWLEDGMENT

We would like to thank Prof. Timothy D. Barfoot for the use of the “Starry Night” dataset prepared for the graduate state estimation course, AER1513, taught at the University of Toronto.

REFERENCES

- [1] A. I. Mourikis, “A Multi-State Constraint Kalman Filter for Vision-aided Inertial Navigation (Tech Note),” pp. 1–20, 2006.
- [2] A. I. Mourikis and S. I. Roumeliotis, “A Multi-State Constraint Kalman Filter for Vision-aided Inertial Navigation,” in *Proc. ICRA*, 2007, pp. 3565–3572.
- [3] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The KITTI dataset,” *Int. J. Rob. Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [4] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. T. Furgale, “Keyframe-based visual-inertial odometry using nonlinear optimization,” *Int. J. Rob. Research*, vol. 34, no. 3, pp. 314–334, 2015.
- [5] G. Sibley, L. Matthies, and G. Sukhatme, “Sliding window filter with application to planetary landing,” *J. Field Rob.*, vol. 27, no. 5, pp. 587–608, 2010.
- [6] A. Huster, E. Frew, and S. Rock, “Relative position estimation for auvs by fusing bearing and inertial rate sensor measurements,” in *Proc. OCEANS*, vol. 3, 2002.
- [7] D. G. Kottas, K. J. Wu, and S. I. Roumeliotis, “Detecting and dealing with hovering maneuvers in vision-aided inertial navigation sys.”
- [8] J. Kim and S. Sukkariehb, “Real-time implementation of airborne inertial-slam,” in *Proc. RSS*, 2007.
- [9] M. Li and A. I. Mourikis, “High-precision, consistent EKF-based visual-inertial odometry,” *Int. J. Rob. Research*, vol. 32, no. 6, pp. 690–711, 2013.
- [10] S. You and U. Neumann, “Fusion of vision and gyro tracking for robust augmented reality registration.”
- [11] D. G. Kottas and S. I. Roumeliotis, “Exploiting urban scenes for vision-aided inertial navigation,” in *Proc. RSS*, 2013.
- [12] S. Ebcin and M. Veth, “Tightly-coupled image-aided inertial navigation using the unscented kalman filter,” in *Proc. ION GNSS*, 2007, pp. 1851–1860.
- [13] D. Fox, W. Burgard, and S. Thrun, “Markov localization for mobile robots in dynamic environments,” *J. Artif. Intelligence Research*, vol. 11, pp. 391–427, 1999.
- [14] M. Pupilli and A. Calway, “Real-time camera tracking using a particle filter,” in *Proc. BMVC*, 2005, pp. 519–528.
- [15] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon, “Bundle adjustment – a modern synthesis,” in *Vision Algorithms: Theory and Practice*. Springer Verlag, 2000, pp. 298–375.
- [16] D. Strelow and S. Singh, “Motion estimation from image and inertial measurements,” *Int. J. Rob. Research*, vol. 23, no. 12, pp. 1157 – 1195, December 2004.
- [17] A. I. Mourikis, M. Li, and B. H. Kim, “Real-time motion tracking on a cellphone using inertial sensing and a rolling-shutter camera,” in *Proc. ICRA*, 2013, pp. 4712–4719.
- [18] S. Leutenegger, P. T. Furgale, V. Rabaud, M. Chli, K. Konolige, and R. Siegwart, “Keyframe-Based Visual-Inertial SLAM using Nonlinear Optimization,” in *Proc. RSS*, 2013.
- [19] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, “Speeded-Up Robust Features (SURF),” *Comp. Vis. and Image Understanding*, vol. 110, no. 3, pp. 346–359, Jun. 2008.
- [20] C. Tomasi and T. Kanade, “Detection and tracking of point features,” Carnegie Mellon University, Tech. Rep. CMU-CS-91-132, 1991.
- [21] P. Torr and A. Zisserman, “Robust computation and parametrization of multiple view relations,” in *Proc. ICCV*, Jan 1998, pp. 727–732.