

Towards Consistent Visual-Inertial Navigation

Guoquan Huang, Michael Kaess, and John J. Leonard

Abstract—Visual-inertial navigation systems (VINS) have prevailed in various applications, in part because of the complementary sensing capabilities and decreasing costs as well as sizes. While many of the current VINS algorithms undergo inconsistent estimation, in this paper we introduce a new extended Kalman filter (EKF)-based approach towards consistent estimates. To this end, we impose both state-transition and observability constraints in computing EKF Jacobians so that the resulting linearized system can best approximate the underlying nonlinear system. Specifically, we enforce the propagation Jacobian to obey the semigroup property, thus being an appropriate state-transition matrix. This is achieved by parametrizing the orientation error state in the *global*, instead of *local*, frame of reference, and then evaluating the Jacobian at the *propagated*, instead of the *updated*, state estimates. Moreover, the EKF linearized system ensures correct observability by projecting the most-accurate measurement Jacobian onto the observable subspace so that no spurious information is gained. The proposed algorithm is validated by both Monte-Carlo simulation and real-world experimental tests.

I. INTRODUCTION

Over the past decades, inertial navigation systems (INS) [1] have been extensively used for estimating the 6 degrees-of-freedom (d.o.f.) poses of sensing platforms (a.k.a. robots) in GPS-denied environments, such as underwater, indoor, in the urban canyon, and on other planets. Most INS rely on an inertial measurement unit (IMU) that measures the 3 d.o.f. rotational velocity and 3 d.o.f. linear acceleration of the sensing platform on which it is rigidly attached. Unfortunately, simple integration of IMU measurements that are corrupted by noise and bias, often results in pose estimates unreliable for long-term navigation. On the other hand, a camera is small, light-weight, inexpensive, and energy efficient while providing rich information. We hence aid an INS with a monocular camera whose measurements are used to provide motion information of the sensor pair, i.e., visual-inertial navigation system (VINS). In this paper, we aim to develop a consistent estimation algorithm for this problem.

Various algorithms are available for VINS problems including visual-inertial simultaneous localization and mapping (SLAM) [2] and visual-inertial odometry (VIO) [3], such as the extended Kalman filter (EKF) [2], [4], [5], the unscented Kalman filter (UKF) [6], and the batch or incremental smoothers [7], [8], among which the EKF-based approach remains arguably the most popular because of

its efficiency. However, similar to 2D SLAM [9]–[11], the standard EKF produces inconsistent estimates when applied to VINS problems, primarily due to the mismatch of observability properties between the EKF linearized VINS and the underlying nonlinear system [3], [12]–[16]. This significantly limits a long-term deployment of VINS in critical scenarios. As defined in [17], a state estimator is *consistent* if the estimation errors are zero-mean, and the estimated covariance is equal to the true covariance. Consistency is one of the primary criteria for evaluating the performance of any estimator; if an estimator is inconsistent, then the accuracy of the computed state estimates is unknown, which in turn makes the estimator unreliable. In this paper, we also study the VINS problem within the EKF framework, while focusing on improving the filter consistency from the perspective of both state-transition *and* observability properties of the EKF linearized system.

In particular, as shown in [3], [12]–[16], the standard EKF-based VINS where the propagation and measurement Jacobians are evaluated at the latest state estimates, has different observability properties from the underlying nonlinear system (or the ideal linearized system where Jacobians are computed using the true states). This was shown to be one of main causes for the filter inconsistency. Furthermore, we analytically show for the first time that the propagation Jacobian in the standard EKF-based VINS violates the semigroup property of a state-transition matrix [18]. If such a Jacobian is used as the “state-transition” matrix to represent the underlying dynamical system, the produced state estimates conceivably may drift away from the solutions of the system, and thus become inconsistent or even diverge. To address the aforementioned two (observability and state-transition) issues, in the proposed algorithm, termed state-transition and observability constrained (STOC)-VINS, we first impose correct observability constraints as in [12]–[15]; and moreover, we *explicitly* enforce the propagation Jacobian to obey the semigroup property. This is achieved by parametrizing the orientation error state in the *global*, instead of *local*, frame of reference, and then directly evaluating the propagation Jacobian at the *propagated*, instead of the *updated*, state estimates. In addition, since in many practical cases the camera-IMU extrinsic calibration is not known perfectly, we include this 6 d.o.f. relative transformation as a part of the state vector and perform *online* calibration, which in effect contributes to improving consistency [16].

The remainder of the paper is organized as follows: After an overview of related work in the next section, the EKF-based VINS and its observability properties are described in Section III. In Section IV, we present the proposed STOC-VINS to improve filter consistency by enforcing appropriate

This work was partially supported by the ONR grants N00014-12-1-0093, N00014-10-1-0936, N00014-11-1-0688 and N00014-13-1-0588, and the NSF grant IIS-1318392, which we gratefully acknowledge.

G. Huang and J. Leonard are with the Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. Email: {gqhuang|jleonard}@mit.edu

M. Kaess is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA. Email: kaess@cmu.edu

state-transition and observability constraints. In Sections V and VI, the proposed approach is validated on both Monte-Carlo simulations and real-world experiments. Finally, Section VII outlines the main conclusions of this work, as well as possible future research directions.

II. RELATED WORK

Visual-inertial navigation has recently prevailed in robot localization in 3D (e.g., [2]–[8], [12]–[16], [19]–[26]), which can be broadly categorized into loosely-coupled and tightly-coupled approaches. The former processes the IMU measurements and/or images separately in a front end, and subsequently fuses them in a back end (e.g., [8], [23]). However, although this type of methods have advantage of computational efficiency, the decoupling results in information loss [16]. The latter seamlessly fuses the visual and inertial measurements by processing them in a single estimation thread (e.g., [3], [5], [12]–[16], [25], [26]). The approach proposed in this paper falls into the latter category, aiming at consistent VINS.

As system observability plays an important role in the proposed approach, we note that some work has recently studied the VINS observability properties under different scenarios. In particular, in [26], [27], nonlinear observability of IMU-camera extrinsic calibration was analyzed based on Lie derivatives and the conditions under which the IMU-camera transformation is observable were determined. In [25], the VINS observability was studied by examining the system’s indistinguishable trajectories [28] under different sensor configurations. Similarly, Martinelli [21] employed the concept of continuous symmetries [28] to show that in VINS, the IMU biases, 3D velocity, and absolute roll and pitch angles are observable.

Recently, similar to robot localization in 2D [9]–[11], consistency of EKF-based VINS has been investigated in [3], [12]–[16] from the perspective of observability. Specifically, Li and Mourikis [3], [16] studied the impact of filter inconsistency due to the VINS observability properties, and leveraged the first-estimates-Jacobian methodology [9] to mitigate the inconsistency. In [12]–[15], following the observability-based methodology proposed in [11], [29], the observability-constrained (OC)-VINS was introduced, which can employ any linearization method to ensure correct observability of the linearized system. While the same observability-based idea is used in the proposed STOC-VINS, we further explicitly enforce the propagation Jacobian to satisfy the semigroup property and thus to be a valid state transition matrix, which results in an alternative way of computing propagation Jacobians to that of the OC-VINS.

III. VISUAL-INERTIAL NAVIGATION

In this section, we first describe the IMU propagation and camera measurement models within the EKF framework, which govern the VINS. In the sequel, we briefly overview the observability properties of the linearized VINS, which will be useful for the design of our approach. For concise presentation of the analysis, we hereafter consider the case where only a single feature is included in the state vector,

while the results can be easily generalized to the case of multiple features.

A. IMU propagation model

The EKF uses the IMU (gyroscope and accelerometer) measurements for state propagation, and the state vector consists of the IMU states \mathbf{x}_I and the feature position ${}^G\mathbf{p}_f$:¹

$$\begin{aligned} \mathbf{x} &= [\mathbf{x}_I^T \quad {}^G\mathbf{p}_f^T]^T \\ &= [{}^I_G\bar{\mathbf{q}}^T \quad \mathbf{b}_g^T \quad {}^G\mathbf{v}^T \quad \mathbf{b}_a^T \quad {}^G\mathbf{p}^T \quad {}^G\mathbf{p}_f^T]^T \end{aligned} \quad (1)$$

where ${}^I_G\bar{\mathbf{q}}$ is the unit quaternion that represents the rotation from the global frame of reference $\{G\}$ to the IMU frame $\{I\}$ (i.e., different parametrization of the rotation matrix $\mathbf{C}({}^I_G\bar{\mathbf{q}}) =: {}^I_G\mathbf{C}$); ${}^G\mathbf{p}$ and ${}^G\mathbf{v}$ are the IMU position and velocity in the global frame; and \mathbf{b}_g and \mathbf{b}_a denote the gyroscope and accelerometer biases, respectively.

By noting that the feature is static (with trivial dynamics), as well as using the IMU motion dynamics [30], the continuous-time dynamics of the state (1) is given by:

$$\begin{aligned} {}^I_G\dot{\bar{\mathbf{q}}}(t) &= \frac{1}{2}\boldsymbol{\Omega}({}^I\boldsymbol{\omega}(t)) {}^I_G\bar{\mathbf{q}}(t), \quad {}^G\dot{\mathbf{p}}(t) = {}^G\mathbf{v}(t), \quad {}^G\dot{\mathbf{v}}(t) = {}^G\mathbf{a}(t) \\ \dot{\mathbf{b}}_g(t) &= \mathbf{n}_{wg}(t), \quad \dot{\mathbf{b}}_a(t) = \mathbf{n}_{wa}(t), \quad {}^G\dot{\mathbf{p}}_f(t) = \mathbf{0}_{3 \times 1} \end{aligned} \quad (2)$$

where ${}^I\boldsymbol{\omega} = [\omega_1 \quad \omega_2 \quad \omega_3]^T$ is the rotational velocity of the IMU, expressed in $\{I\}$, ${}^G\mathbf{a}$ is the IMU acceleration in $\{G\}$, \mathbf{n}_{wg} and \mathbf{n}_{wa} are the white Gaussian noise processes that drive the IMU biases, and $\boldsymbol{\Omega}(\boldsymbol{\omega})$ is defined by:

$$\boldsymbol{\Omega}(\boldsymbol{\omega}) = \begin{bmatrix} -[\boldsymbol{\omega} \times] & \boldsymbol{\omega} \\ -\boldsymbol{\omega}^T & 0 \end{bmatrix}, \quad [\boldsymbol{\omega} \times] = \begin{bmatrix} 0 & -\omega_3 & \omega_2 \\ \omega_3 & 0 & -\omega_1 \\ -\omega_2 & \omega_1 & 0 \end{bmatrix}$$

A typical IMU provides gyroscope and accelerometer measurements, $\boldsymbol{\omega}_m$ and \mathbf{a}_m , both of which are expressed in the IMU local frame $\{I\}$ and given by:

$$\boldsymbol{\omega}_m(t) = {}^I\boldsymbol{\omega}(t) + \mathbf{b}_g(t) + \mathbf{n}_g(t) \quad (3)$$

$$\mathbf{a}_m(t) = \mathbf{C}({}^I_G\bar{\mathbf{q}}(t)) ({}^G\mathbf{a}(t) - {}^G\mathbf{g}) + \mathbf{b}_a(t) + \mathbf{n}_a(t) \quad (4)$$

where ${}^G\mathbf{g}$ is the gravitational acceleration expressed in $\{G\}$, and \mathbf{n}_g and \mathbf{n}_a are zero-mean, white Gaussian noise.

Linearization of (2) at the current state estimate yields the continuous-time state-estimate propagation model [5]:

$$\begin{aligned} {}^I_G\hat{\bar{\mathbf{q}}}(t) &= \frac{1}{2}\boldsymbol{\Omega}({}^I\hat{\boldsymbol{\omega}}(t)) {}^I_G\hat{\bar{\mathbf{q}}}(t), \quad {}^G\hat{\mathbf{p}}(t) = {}^G\hat{\mathbf{v}}(t), \quad {}^G\hat{\mathbf{v}}(t) = {}^G\hat{\mathbf{a}}(t) \\ \hat{\mathbf{b}}_g(t) &= \mathbf{0}_{3 \times 1}, \quad \hat{\mathbf{b}}_a(t) = \mathbf{0}_{3 \times 1}, \quad {}^G\hat{\mathbf{p}}_f(t) = \mathbf{0}_{3 \times 1} \end{aligned} \quad (5)$$

where $\hat{\mathbf{a}} = \mathbf{a}_m - \hat{\mathbf{b}}_a$ and $\hat{\boldsymbol{\omega}} = \boldsymbol{\omega}_m - \hat{\mathbf{b}}_g$. The error state of dimension 18×1 is hence defined as follows [see (1)]:

$$\tilde{\mathbf{x}}(t) = \begin{bmatrix} {}^I\tilde{\boldsymbol{\theta}}^T(t) & \tilde{\mathbf{b}}_g^T(t) & {}^G\tilde{\mathbf{v}}^T(t) & \tilde{\mathbf{b}}_a^T(t) & {}^G\tilde{\mathbf{p}}^T(t) & {}^G\tilde{\mathbf{p}}_f^T(t) \end{bmatrix}^T \quad (6)$$

where we have employed the multiplicative error model for a quaternion [30]. That is, the error between the quaternion $\bar{\mathbf{q}}$ and its estimate $\hat{\bar{\mathbf{q}}}$ is the 3×1 angle-error vector, ${}^I\tilde{\boldsymbol{\theta}}$,

¹Throughout this paper the subscript $\ell|j$ refers to the estimate of a quantity at time-step ℓ , after all measurements up to time-step j have been processed. \hat{x} is used to denote the estimate of a random variable x , while $\tilde{x} = x - \hat{x}$ is the error in this estimate. \mathbf{I}_n and $\mathbf{0}_n$ are the $n \times n$ identity and zero matrices, respectively. Finally, the left superscript denotes the frame of reference with respect to which the vector is expressed.

implicitly defined by the error quaternion: $\delta\bar{\mathbf{q}} = \bar{\mathbf{q}} \otimes \hat{\mathbf{q}} \simeq \begin{bmatrix} 1 \\ \frac{1}{2} \mathbf{I} \tilde{\boldsymbol{\theta}} \\ 1 \end{bmatrix}$, where $\delta\bar{\mathbf{q}}$ describes the small rotation that causes the true and estimated attitude to coincide. The advantage of this parametrization permits a minimal representation, 3×3 covariance matrix $\mathbb{E} \left[\begin{smallmatrix} \tilde{\boldsymbol{\theta}} \\ \tilde{\boldsymbol{\theta}}^T \end{smallmatrix} \right]$, for the attitude uncertainty. It is important to note that the orientation error, $\tilde{\boldsymbol{\theta}}$, satisfies the following rotation-matrix relation [30]:

$$\mathbf{C}(^L_G \bar{\mathbf{q}}) \simeq \left(\mathbf{I}_3 - [^L \tilde{\boldsymbol{\theta}} \times] \right) \mathbf{C}(^L_G \hat{\mathbf{q}}) \quad (7)$$

Now the continuous-time error-state propagation is:

$$\dot{\tilde{\mathbf{x}}}(t) = \mathbf{F}_c(t) \tilde{\mathbf{x}}(t) + \mathbf{G}_c(t) \mathbf{n}(t) \quad (8)$$

where $\mathbf{n} = [\mathbf{n}_g^T \ \mathbf{n}_{wg}^T \ \mathbf{n}_a^T \ \mathbf{n}_{wa}^T]^T$ is the system noise, \mathbf{F}_c is the continuous-time error-state transition matrix, and \mathbf{G}_c is the input noise matrix, which are given by (see [30]):

$$\mathbf{F}_c = \begin{bmatrix} -[\hat{\boldsymbol{\omega}} \times] & -\mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ -\mathbf{C}^T(^L_G \hat{\mathbf{q}}) [\hat{\mathbf{a}} \times] & \mathbf{0}_3 & \mathbf{0}_3 & -\mathbf{C}^T(^L_G \hat{\mathbf{q}}) & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \end{bmatrix} \quad (9)$$

$$\mathbf{G}_c = \begin{bmatrix} -\mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ -\mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ -\mathbf{I}_3 & \mathbf{0}_3 & -\mathbf{C}^T(^L_G \hat{\mathbf{q}}) & \mathbf{0}_3 \\ -\mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \\ -\mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \end{bmatrix} \quad (10)$$

The system noise is modelled as zero-mean white Gaussian process with autocorrelation $\mathbb{E} [\mathbf{n}(t) \mathbf{n}(\tau)^T] = \mathbf{Q}_c \delta(t - \tau)$, which depends on the IMU noise characteristics.

We have described the continuous-time propagation model using IMU measurements. However, in any practical EKF implementation, the discrete-time state-transition matrix, $\Phi(t_{k+1}, t_k)$, is required in order to propagate the error covariance from time t_k to t_{k+1} . Typically it is found by solving the following matrix differential equation:

$$\dot{\Phi}(t_{k+1}, t_k) = \mathbf{F}_c(t_{k+1}) \Phi(t_{k+1}, t_k) \quad (11)$$

with the initial condition $\Phi(t_k, t_k) = \mathbf{I}_{18}$. Its solution has the following structure:

$$\Phi_k := \Phi(t_{k+1}, t_k) = \begin{bmatrix} \Phi_{k,11} & \Phi_{k,12} & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,31} & \Phi_{k,32} & \mathbf{I}_3 & \Phi_{k,34} & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,51} & \Phi_{k,52} & \delta t_k \mathbf{I}_3 & \Phi_{k,54} & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (12)$$

where $\delta t_k = t_{k+1} - t_k$. This matrix (12) can be found either numerically [5], [30] or analytically [3], [13], [14], [16]. Once it is computed, the EKF propagates the error covariance in a standard way [31]:

$$\mathbf{P}_{k+1|k} = \Phi_k \mathbf{P}_{k|k} \Phi_k^T + \mathbf{Q}_{d,k} \quad (13)$$

where $\mathbf{Q}_{d,k}$ is the discrete-time system noise covariance matrix computed as follows:

$$\mathbf{Q}_{d,k} = \int_{t_k}^{t_{k+1}} \Phi(t_{k+1}, \tau) \mathbf{G}_c(\tau) \mathbf{Q}_c \mathbf{G}_c^T(\tau) \Phi^T(t_{k+1}, \tau) d\tau$$

B. Camera measurement model

The camera observes visual corner features, which are used to concurrently estimate the ego-motion of the sensing platform. Assuming a calibrated perspective camera, the measurement of the feature at time-step k is the perspective projection of the 3D point, ${}^C_k \mathbf{p}_f$, expressed in the current camera frame $\{C_k\}$, onto the image plane, i.e.,

$$\mathbf{z}_k = \frac{1}{z_k} \begin{bmatrix} x_k \\ y_k \end{bmatrix} + \mathbf{v}_k \quad (14)$$

$$\begin{bmatrix} x_k \\ y_k \\ z_k \end{bmatrix} = {}^C_k \mathbf{p}_f = \mathbf{C}(^C_I \bar{\mathbf{q}}) \mathbf{C}(^L_G \bar{\mathbf{q}}_k) ({}^G \mathbf{p}_f - {}^G \mathbf{p}_k) + {}^C \mathbf{p}_I \quad (15)$$

where \mathbf{v}_k is the zero-mean, white Gaussian measurement noise with covariance \mathbf{R}_k . In (15), $\{^C_I \bar{\mathbf{q}}, {}^C \mathbf{p}_I\}$ is the rotation and translation between the camera and the IMU. This transformation can be obtained, for example, by performing camera-IMU extrinsic calibration *offline* [27]. However, in practice when the perfect calibration is unavailable, it is beneficial to VINS consistency to include these calibration parameters in the state vector and concurrently estimate them along with the IMU/camera poses [16]. For this reason, we perform *online* camera-IMU calibration in the proposed STOC-VINS (see Section IV).

For the use of EKF, linearization of (14) yields the following measurement residual [see (6)]:

$$\tilde{\mathbf{z}}_k = \mathbf{H}_k \tilde{\mathbf{x}}_{k|k-1} + \mathbf{v}_k = \mathbf{H}_{\mathbf{I}_k} \tilde{\mathbf{x}}_{\mathbf{I}_k|k-1} + \mathbf{H}_{\mathbf{f}_k} {}^G \tilde{\mathbf{p}}_{f_k|k-1} + \mathbf{v}_k \quad (16)$$

where the measurement Jacobian \mathbf{H}_k is computed as:

$$\mathbf{H}_k = [\mathbf{H}_{\mathbf{I}_k} \ \mathbf{H}_{\mathbf{f}_k}] \quad (17)$$

$$= \mathbf{H}_{\text{proj}} \mathbf{C}(^C_I \bar{\mathbf{q}}) [\mathbf{H}_{\theta_k} \ \mathbf{0}_{3 \times 9} \ \mathbf{H}_{\mathbf{p}_k} \ \mathbf{C}(^L_G \hat{\mathbf{q}}_k)]$$

$$\mathbf{H}_{\text{proj}} = \frac{1}{z_k^2} \begin{bmatrix} \hat{z}_k & 0 & -\hat{x}_k \\ 0 & \hat{z}_k & -\hat{y}_k \end{bmatrix} \quad (18)$$

$$\mathbf{H}_{\theta_k} = [\mathbf{C}(^L_G \hat{\mathbf{q}}_k) ({}^G \hat{\mathbf{p}}_f - {}^G \hat{\mathbf{p}}_k) \times], \ \mathbf{H}_{\mathbf{p}_k} = -\mathbf{C}(^L_G \hat{\mathbf{q}}_k) \quad (19)$$

Once the measurement Jacobian and residual are computed, we can apply the standard EKF update equations to update the state estimates and error covariance [31].

C. Observability properties

Observability analysis has recently been performed for both nonlinear and linearized VINS [3], [13]. In particular, the observability matrix for the EKF linearized system over the time interval $[k_o \ k]$ is defined by [31]:

$$\mathbf{M} = \begin{bmatrix} \mathbf{H}_{k_o} \\ \mathbf{H}_{k_o+1} \Phi_{k_o} \\ \vdots \\ \mathbf{H}_k \Phi_{k-1} \cdots \Phi_{k_o} \end{bmatrix} \quad (20)$$

It has been shown in [3], [13] that the nullspace of \mathbf{M} (i.e., unobservable subspace) for the VINS *ideally* spans the following *four* directions:

$$\mathbf{M}\mathbf{N} = \mathbf{0} \Rightarrow \mathbf{N} = \begin{bmatrix} \mathbf{0}_3 & \mathbf{C}(^L_G \bar{\mathbf{q}}_k) {}^G \mathbf{g} \\ \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & -[{}^G \mathbf{v}_k \times] {}^G \mathbf{g} \\ \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{I}_3 & -[{}^G \mathbf{p}_k \times] {}^G \mathbf{g} \\ \mathbf{I}_3 & -[{}^G \mathbf{p}_f \times] {}^G \mathbf{g} \end{bmatrix} \quad (21)$$

Note that the first block column of \mathbf{N} in (21) corresponds to the global translation while the second block column corresponds to the global rotation about the gravity vector ${}^G\mathbf{g}$. When designing a nonlinear estimator for VINS, we would like the system model employed by the estimator to have an unobservable subspace spanned by these directions. However, this is not the case for the standard EKF as shown in [3], [13]–[16]. In particular, the standard EKF linearized system, which linearizes system and measurement functions at the current state estimate, has an unobservable subspace of *three*, instead of four, d.o.f. This implies that the filter gains non-existent information from available measurements, which may lead to filter inconsistency.

IV. STATE-TRANSITION-AND-OBSERVABILITY CONSTRAINED (STOC)-VINS

As discussed in the preceding section, the standard EKF-based VINS where the propagation and measurement Jacobians are evaluated at the latest state estimates, has different observability properties from the ideal linearized system where the Jacobians are computed using the true states. This was shown to be *one* of the main causes for filter inconsistency [3], [12]–[16]. In this section, we revisit this inconsistency problem and further find that the propagation Jacobian of the standard EKF-based VINS violates the semigroup property of a state-transition matrix [18]. If such a Jacobian is used as the state-transition matrix to represent the underlying dynamical system, the produced state estimates conceivably may deviate from the solutions of the dynamical system, and thus become inconsistent or even diverge. Therefore, when designing consistent VINS algorithms, besides imposing correct observability constraints in computing Jacobians, we explicitly enforce the propagation Jacobian to obey the semigroup property and thus to be a valid state-transition matrix. The resulting method is thus termed as state-transition and observability constrained (STOC)-VINS.

A. Computing propagation Jacobians

We know from control theory that a state-transition matrix *must* have the following properties [18]:

$$\dot{\Phi}(t_1, t_0) = \mathbf{F}_c(t_1)\Phi(t_1, t_0) \quad (22)$$

$$\Phi(t_0, t_0) = \mathbf{I}_{\dim(\mathbf{x})} \quad (23)$$

$$\Phi(t_1, t_0) = \Phi^{-1}(t_0, t_1) \quad (24)$$

$$\Phi(t_2, t_0) = \Phi(t_2, t_1)\Phi(t_1, t_0) \quad (25)$$

which hold for any t_0, t_1 and t_2 . Note that in VINS we have derived the analytical state-transition matrix by solving the matrix differential equation (11) with the self-mapping initial condition, which is identical to (22) and (23). Note also that given (25) and (23), the identity of (24) immediately holds. Therefore, we hereafter focus on examining (25) which is the so-called *semigroup* property [18]. However, we show for the first time that the propagation Jacobian of the standard EKF-based VINS is *not* a valid state-transition matrix:

Lemma 4.1: The propagation Jacobian (12) of the standard EKF-based VINS, computed using the current state

estimates, violates the semigroup property (25) for being a state-transition matrix, i.e., for some t_{k-1}, t_k , and t_{k+1} ,

$$\Phi(t_{k+1}, t_{k-1}) \neq \Phi(t_{k+1}, t_k)\Phi(t_k, t_{k-1}) \quad (26)$$

Proof: See Appendix I. ■

As a state-transition matrix is used to construct the general solution of the corresponding linear dynamical systems, an invalid state-transition matrix can result in an erroneous solution. Therefore, using the propagation Jacobian as the incorrect “transition” matrix for the EKF linearized VINS system conceivably may cause the filter producing inaccurate, or even inconsistent, estimates. To address this issue, we aim to construct the propagation Jacobian in such a way that enforces this Jacobian to be a valid state-transition matrix for the EKF linearized system, and in particular, to possess the semigroup property (25). The key idea of our approach is that we parametrize the IMU orientation error in the *global*, instead of *local* (as commonly used in the regular VINS formulation [5], [12], [13]), frame of reference; and then *analytically* compute the propagation Jacobian using the *propagated*, instead of *updated*, state estimates.

In particular, we first notice that [in contrast to (7)]:

$$\mathbf{C}({}^L_G\bar{\mathbf{q}}) \simeq \mathbf{C}({}^L_G\hat{\mathbf{q}}) \left(\mathbf{I}_3 - [{}^G\tilde{\boldsymbol{\theta}} \times] \right) \quad (27)$$

which results in the global orientation error state, ${}^G\tilde{\boldsymbol{\theta}} = \mathbf{C}^T({}^L_G\hat{\mathbf{q}}) {}^L\tilde{\boldsymbol{\theta}}$. With this parametrization, the error states except the biases are all in the global frame, which will be useful for our ensuing derivations [see (6)]:

$$\tilde{\mathbf{x}}' := \begin{bmatrix} {}^G\tilde{\boldsymbol{\theta}} \\ \tilde{\mathbf{b}}_g \\ {}^G\tilde{\boldsymbol{\nu}} \\ \tilde{\mathbf{b}}_a \\ {}^G\tilde{\mathbf{p}} \\ {}^G\tilde{\mathbf{p}}_f \end{bmatrix} = \text{Diag} \left(\mathbf{C}^T({}^L_G\hat{\mathbf{q}}), \mathbf{I}_{15} \right) \begin{bmatrix} {}^L\tilde{\boldsymbol{\theta}} \\ \tilde{\mathbf{b}}_g \\ {}^G\tilde{\boldsymbol{\nu}} \\ \tilde{\mathbf{b}}_a \\ {}^G\tilde{\mathbf{p}} \\ {}^G\tilde{\mathbf{p}}_f \end{bmatrix} =: \boldsymbol{\Lambda}^T \tilde{\mathbf{x}} \quad (28)$$

Now the new error-state propagation can be written as:

$$\tilde{\mathbf{x}}'_{k+1|k} = \boldsymbol{\Lambda}_{k+1}^T \tilde{\mathbf{x}}_{k+1|k} = \boldsymbol{\Lambda}_{k+1}^T \Phi_k \boldsymbol{\Lambda}_k \tilde{\mathbf{x}}'_{k|k} =: \Phi'_k \tilde{\mathbf{x}}'_{k|k} \quad (29)$$

where we have used the fact that $\boldsymbol{\Lambda}^{-1} = \boldsymbol{\Lambda}^T$ [see (28)]. Note that $\Phi'_k := \boldsymbol{\Lambda}_{k+1}^T \Phi_k \boldsymbol{\Lambda}_k$ is the propagation Jacobian for the new parametrization, and can be computed analytically based on the analytical expression of Φ_k (see (12) and [13]):

$$\Phi'_k := \Phi'(t_{k+1}, t_k) = \begin{bmatrix} \mathbf{C}^T({}^L_G\hat{\mathbf{q}}_{k+1})\Phi_{k,11}\mathbf{C}({}^L_G\hat{\mathbf{q}}_k) & \mathbf{C}^T({}^L_G\hat{\mathbf{q}}_{k+1})\Phi_{k,12} & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,31}\mathbf{C}({}^L_G\hat{\mathbf{q}}_k) & \Phi_{k,32} & \mathbf{I}_3 & \Phi_{k,34} & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,51}\mathbf{C}({}^L_G\hat{\mathbf{q}}_k) & \Phi_{k,52} & \delta t_k \mathbf{I}_3 & \Phi_{k,54} & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (30)$$

In particular, as compared to $\Phi(t_{k+1}, t_k)$, the only blocks changed are $\Phi'_{k,11}$, $\Phi'_{k,12}$, $\Phi'_{k,31}$, and $\Phi'_{k,51}$, which are computed in closed form as follows (also see [3], [13]):

$$\Phi'_{k,11} = \mathbf{C}^T({}^L_G\hat{\mathbf{q}}_{k+1})\mathbf{C} \left(\frac{I^{(k+1)}\hat{\mathbf{q}}}{I^{(k)}} \right) \mathbf{C}({}^L_G\hat{\mathbf{q}}_k) = \mathbf{I}_3 \quad (31)$$

$$\Phi'_{k,12} = - \int_{t_k}^{t_{k+1}} \mathbf{C}^T({}^L_G\hat{\mathbf{q}}(t_\tau)) d\tau \quad (32)$$

$$\Phi'_{k,31} = - [({}^G\hat{\boldsymbol{\nu}}_{k+1} - {}^G\hat{\boldsymbol{\nu}}_k - {}^G\mathbf{g}\delta t_k) \times] \quad (33)$$

$$\Phi'_{k,51} = - \left[\left({}^G\hat{\mathbf{p}}_{k+1} - {}^G\hat{\mathbf{p}}_k - {}^G\hat{\boldsymbol{\nu}}_k\delta t_k - \frac{1}{2}{}^G\mathbf{g}\delta t_k^2 \right) \times \right] \quad (34)$$

We now show that the propagation Jacobian, $\Phi'(t_{k+1}, t_k)$, can be constructed analytically so as to satisfy the semigroup property (25) for being a valid state-transition matrix.

Lemma 4.2: If the propagation Jacobian $\Phi'(t_{\ell+1}, t_\ell)$ is evaluated at the propagated state estimates, $\hat{\mathbf{x}}_{\ell+1|\ell}$ and $\hat{\mathbf{x}}_{\ell|\ell-1}$, then it satisfies the semigroup property (25), i.e.,

$$\Phi'(t_{k+1}, t_{k-1}) = \Phi'(t_{k+1}, t_k) \Phi'(t_k, t_{k-1}) \quad (35)$$

Proof: See Appendix II. ■

B. Computing measurement Jacobians

The measurement Jacobian with respect to the new error state (28) is calculated as follows [see (14) and (17)-(19)]:

$$\mathbf{H}'_k = [\mathbf{H}'_{\mathbf{I}_k} \quad \mathbf{H}'_{\mathbf{f}_k}] \quad (36)$$

$$= \mathbf{H}_{\text{proj}} \mathbf{C}({}^G \hat{\mathbf{q}}) [\mathbf{H}'_{\theta_k} \quad \mathbf{0}_{3 \times 9} \quad \mathbf{H}'_{\mathbf{p}_k} \quad \mathbf{C}({}^I \hat{\mathbf{q}}_k)]$$

$$\mathbf{H}'_{\theta_k} = [({}^G \hat{\mathbf{p}}_f - {}^G \hat{\mathbf{p}}_k) \times], \quad \mathbf{H}'_{\mathbf{p}_k} = -\mathbf{I}_3 \quad (37)$$

Note that by performing observability analysis similar to [13], the *ideal* linearized error-state system (i.e., Jacobians are computed using the true states) with the global orientation-error parametrization has the following unobservable subspace of 4 d.o.f. [see (21)]:

$$\mathbf{N}' = \begin{bmatrix} \mathbf{0}_3 & {}^G \mathbf{g} \\ \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & -[{}^G \mathbf{v}_k \times] {}^G \mathbf{g} \\ \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{I}_3 & -[{}^G \mathbf{p}_k \times] {}^G \mathbf{g} \\ \mathbf{I}_3 & -[{}^G \mathbf{p}_f \times] {}^G \mathbf{g} \end{bmatrix} \quad (38)$$

However, in analogy to the case of standard EKF-based VINS (see Section III-C), it is not difficult to show that if we compute the measurement Jacobian using the current best state estimates as for the standard EKF, while computing the propagation Jacobian using the propagated state estimates as devised in the previous section, the resulting linearized error-state system has an unobservable subspace of only 3 (instead of 4) d.o.f. This may result in inconsistent estimates.

To address this issue, we impose appropriate observability constraints when computing measurement Jacobians, by following our prior observability-constrained methodology for designing consistent SLAM estimators [11], which was also exploited in [12]–[15]. Specifically, when computing the measurement Jacobian, we enforce that each block row of the observability matrix (20) has the same nullspace, i.e.,

$$\min_{\mathbf{H}'_k} \|\mathbf{H}'_k - \mathbf{H}_k\|_F^2 \quad (39)$$

$$\text{subject to } \mathbf{H}'_k \Phi'_{k-1} \cdots \Phi'_{k_0} \mathbf{N}' = \mathbf{0} \quad (40)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. Ideally, \mathbf{H}_k in (39) is the measurement Jacobian computed using the true states, which, however, is not realizable in practice. Hence, we employ the latest, and thus the best, state estimates to compute this Jacobian as for the standard EKF, i.e., $\mathbf{H}_k = \mathbf{H}_k(\hat{\mathbf{x}}_{k|k-1})$. On the other hand, \mathbf{N}' in (40) defines the desired nullspace. Although we would like to have the same one as in (38) computed using the true states (which is not realizable in practice), we select the nullspace that has the same structure as in (38) while computing it with the first available state estimates, i.e., $\mathbf{N}' = \mathbf{N}'(\hat{\mathbf{x}}_{k_0|k_0})$.

Once the choice of the nullspace \mathbf{N}' is made, we find the optimal solution to the above problem (39)-(40) in *closed form* based on the following lemma:

Lemma 4.3: The optimal solution to the constrained minimization problem (39)-(40) is given by:

$$\mathbf{H}'_k = \mathbf{H}_k (\mathbf{I}_{\dim(\mathbf{x})} - \mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T) \quad (41)$$

where $\mathbf{U} = \Phi'_{k-1} \cdots \Phi'_{k_0} \mathbf{N}'$.

Proof: See [32]. ■

It is interesting to note that \mathbf{U} in the above lemma is the propagated unobservable subspace (nullspace) at time-step k , and $(\mathbf{I}_{\dim(\mathbf{x})} - \mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T)$ is the subspace orthogonal to \mathbf{U} , i.e., the observable subspace. Hence, as seen from (41), the measurement Jacobian of the proposed STOC-VINS is the projection of the most accurate measurement Jacobian onto the observable subspace.

C. Application to MSCKF

The multi-state constraint Kalman filter (MSCKF) [5], [33] is a well-known VINS algorithm that performs tightly-coupled VIO over a sliding window of m poses, and has complexity only linear in the number of observed features. The MSCKF utilizes all feature observations available within the sliding window to impose probabilistic constraints between poses, without building a map. In what follows, we apply the proposed STOC-VINS to the MSCKF framework to address the VIO problem [16], while our methodology is applicable to other VINS problems including SLAM.

The MSCKF state vector at time-step k augments the current IMU state by the past m poses where the images were taken (i.e., stochastic cloning [34]):

$$\mathbf{x}_{A_k} = [\mathbf{x}_{I_k}^T \quad \mathbf{y}_{k-1}^T \quad \cdots \quad \mathbf{y}_{k-m}^T]^T \quad (42)$$

where $\mathbf{y}_\ell^T = [{}^I \hat{\mathbf{q}}_\ell^T \quad {}^G \hat{\mathbf{p}}_\ell^T]$ is the IMU pose (quaternion and position) where the image is recorded at time-step ℓ . Since the nullspace, \mathbf{N}' , is required for computing the STOC-VINS measurement Jacobian [see (41)], we accordingly augment the nullspace with the ones corresponding to the cloning states as follows:

$$\mathbf{N}'_A = \begin{bmatrix} \mathbf{N}' \\ \mathbf{N}'_{\text{clone},1} \\ \vdots \\ \mathbf{N}'_{\text{clone},m} \end{bmatrix} = \begin{bmatrix} \mathbf{N}' & & & \\ \mathbf{0}_3 & {}^G \mathbf{g} & & \\ \mathbf{I}_3 & -[{}^G \hat{\mathbf{p}}_{k-1|k-2} \times] {}^G \mathbf{g} & & \\ & & \ddots & \\ \mathbf{0}_3 & & & {}^G \mathbf{g} \\ \mathbf{I}_3 & -[{}^G \hat{\mathbf{p}}_{k-m|k-m-1} \times] {}^G \mathbf{g} & & \end{bmatrix} \quad (43)$$

During the MSCKF propagation, the current state estimates evolve forward in time by integrating (2), while the cloning-state estimates remain static. On the other hand, the augmented covariance is propagated as follows [see (13)]:

$$\mathbf{P}_{A_{k+1|k}} = \text{Diag}(\Phi'_k, \mathbf{I}_{6m}) \mathbf{P}_{A_{k|k}} \text{Diag}(\Phi_k^T, \mathbf{I}_{6m}) + \text{Diag}(\mathbf{Q}_{d,k}, \mathbf{0}_{6m}) \quad (44)$$

where the propagation Jacobian Φ'_k is computed using the propagated state estimates as devised in Section IV-A.

During the MSCKF update, we stack together all the feature measurements within the sliding window, and linearize

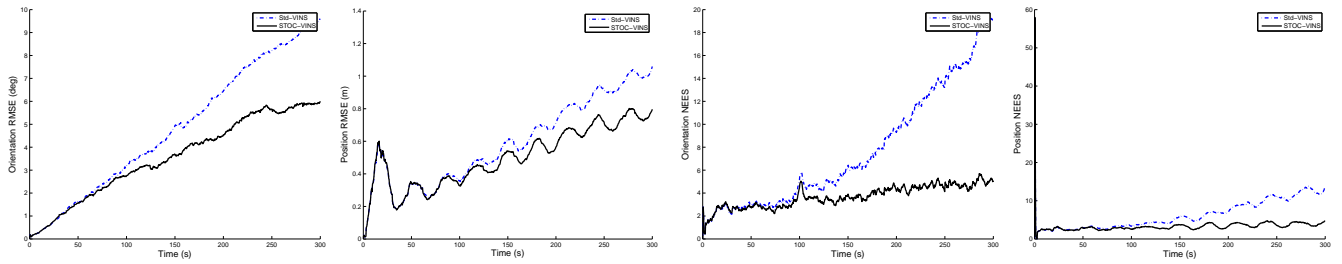


Fig. 1. Monte-Carlo simulation results for a VIO scenario. In these plots, the dash-dotted lines correspond to the standard VINS, the solid lines to the STOC-VINS. It is clear that the proposed STOC-VINS performs significantly better than the standard VINS, in terms of both NEES and RMSE.

them with respect to the augmented IMU states as well as the feature position [see (16)]:

$$\begin{aligned} \begin{bmatrix} \tilde{\mathbf{z}}_k \\ \vdots \\ \tilde{\mathbf{z}}_{k-m} \end{bmatrix} &= \begin{bmatrix} \mathbf{H}'_k \\ \vdots \\ \mathbf{H}'_{k-m} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}}_{A_{k|k-1}} \\ G \tilde{\mathbf{p}}_f \end{bmatrix} + \begin{bmatrix} \mathbf{v}_k \\ \vdots \\ \mathbf{v}_{k-m} \end{bmatrix} \\ &=: \mathbf{H}'_x \tilde{\mathbf{x}}_{A_{k|k-1}} + \mathbf{H}'_f G \tilde{\mathbf{p}}_f + \mathbf{v} \end{aligned} \quad (45)$$

where the measurement Jacobian \mathbf{H}'_k is computed as (41). Note that the feature position is not included in the MSCKF state vector (42), while we want to utilize the information contained in its measurements, we hence project (45) onto the left nullspace of \mathbf{H}'_f (i.e., $\mathbf{W}^T \mathbf{H}'_f = \mathbf{0}$) and have:

$$\begin{aligned} \mathbf{W}^T \tilde{\mathbf{z}} &= \mathbf{W}^T \mathbf{H}'_x \tilde{\mathbf{x}}_{A_{k|k-1}} + \mathbf{W}^T \mathbf{v} \\ \Leftrightarrow \tilde{\mathbf{z}}'' &= \mathbf{H}''_x \tilde{\mathbf{x}}_{A_{k|k-1}} + \mathbf{v}'' \end{aligned} \quad (46)$$

The EKF uses the above residual equation to update the state estimates and covariance [33].

As mentioned before, we include the camera-IMU extrinsic calibration parameters, $\{^C_I \bar{\mathbf{q}}, ^C \mathbf{p}_I\}$, in the state vector so as to perform this calibration *online*. As shown in [16], it is often unrealistic to assume the 6 d.o.f. camera-IMU transformation perfectly known, while using imperfect (known with finite precision) calibration as if it underestimates the uncertainty and thus harms the filter consistency. Interestingly, the inclusion of the calibration parameters in the state vector incurs minimal modifications to the MSCKF [16]. Since this transformation is static, it is easy to propagate over time (in a similar way as for the static feature). Linearization of the stacked measurements renders new Jacobian terms with respect to these parameters [see (45) and (36)], which can be easily used in the standard EKF update equations.

V. SIMULATION RESULTS

We conducted a series of Monte-Carlo simulations under realistic conditions to validate the the proposed STOC-VINS. The metrics used to evaluate filter performance are: (i) the root mean square error (RMSE), and (ii) the average normalized (state) estimation error squared (NEES) [17]. The RMSE provides us with a concise metric of the accuracy of a given filter, while the NEES is a standard criterion for evaluating the filter's consistency. By studying both metrics of the given filter, we obtain a comprehensive picture of the filter's performance.

In this simulation, we consider a VIO scenario [12], [16], where a robot equipped with a camera-IMU pair moves over

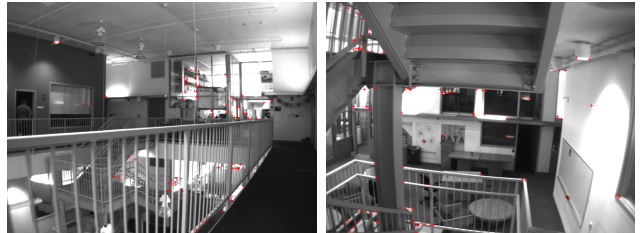


Fig. 2. Sample images with tracked features (red crosses) in the experiment.

a circular trajectory with radius 5 m at an average speed 0.6 m/sec. The camera with 45 deg field of view observes point features randomly distributed inside a circumscribing cylinder. The standard deviation of image noise was set to 1.5 pixels, while IMU measurements were modelled with MEMS sensor quality. We performed 50 Monte-Carlo simulations and compared our proposed STOC-VINS to the standard MSCK-based VINS [33]. Both filters use a sliding window with size of 10 camera poses, and during each run, process the same data to ensure a fair comparison.

The comparative Monte-Carlo results are presented in Fig. 1. As evident, the proposed STOC-VINS performs substantially better than that of the standard VINS, in terms of both RMSE (accuracy) and NEES (consistency). This is attributed to the fact that the proposed filter (i) by construction employs the linearized system models of correct observability properties, and (ii) explicitly enforces the semigroup property of the state transition matrix in the direct analytical computation of the propagation Jacobian.

VI. EXPERIMENTAL RESULTS

We further tested the proposed STOC-VINS in a real-world experiment, in which a hand-held camera/IMU platform travelled over two floors in the Stata Center at MIT. In this experiment, we were using a PointGrey Bumblebee2 stereo pair that records images of resolution 640×480 pixels at 30 Hz (only the right camera's images were used), and a MicroStrain IMU (3DM-GX3-25) which operates at 100 Hz. We employed the Shi-Tomasi corner detector [35] to extract point features from the first available image and track them over the subsequent images using the KLT tracking algorithm [36] (e.g., see Fig. 2). On average, approximately 100 features were tracked per image, while we initialize a new set of features when the number of successfully tracked features falls under a certain threshold. To remove outliers from the

$$\Phi(t_{k+1}, t_k)\Phi(t_k, t_{k-1}) = \begin{bmatrix} \Phi_{k,11}\Phi_{k-1,11} & \Phi_{k,11}\Phi_{k-1,12} + \Phi_{k,12} & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,31}\Phi_{k-1,11} + \Phi_{k-1,31} & \Phi_{k,31}\Phi_{k-1,12} + \Phi_{k,32} + \Phi_{k-1,32} & \mathbf{I}_3 & \Phi_{k,34} + \Phi_{k-1,34} & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi_{k,51}\Phi_{k-1,11} + \delta t_k \Phi_{k-1,31} + \Phi_{k-1,51} & \Phi_{k,51}\Phi_{k-1,12} + \Phi_{k,52} + \delta t_k \Phi_{k-1,32} + \Phi_{k-1,52} & (\delta t_k + \delta t_{k-1})\mathbf{I}_3 & \delta t_k \Phi_{k-1,34} + \Phi_{k,54} + \Phi_{k-1,54} & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (48)$$

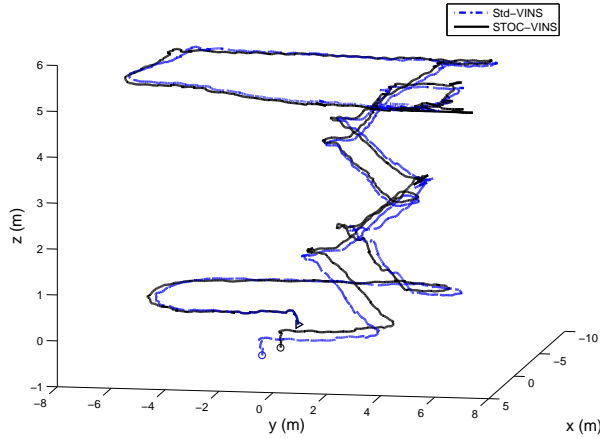


Fig. 3. The estimated trajectories of the two compared VINS algorithms in the real-world experiment conducted in the Stata Center at MIT. In this plot, \triangle denotes the starting position, while the estimated ending points are denoted by \circ . The dash-dotted lines correspond to the standard VINS, the solid lines to the STOC-VINS. Note that the two algorithms perform close to each other in some portions along the path, which makes the corresponding lines difficult to distinguish.

resulting tracks, which is a necessary step in practice, we used the RANSAC with five-point algorithm [37].

In this test, the same two MSCKF-based VINS algorithms as in the preceding simulation were compared, and Fig. 3 shows the estimated trajectories. Note that in this test, the camera-IMU platform traversed about 140 m and returned to its starting position. At the end of the trajectory, the standard VINS has a position error of 1.6111 m, while the position error of the proposed STOC-VINS is only 1.0975 m. These errors respectively account for approximately 1.13% (standard VINS), and 0.79% (STOC-VINS) of the total distance travelled. It becomes clear from these results that the proposed STOC-VINS performs better than the standard VINS which agrees with the previous simulation results.

VII. CONCLUSIONS AND FUTURE WORK

In this paper, we have introduced a new EKF-based VINS algorithm, termed STOC-VINS, which ensures appropriate state-transition and observability properties of the linearized system so as to improve consistency and accuracy. In particular, we use the global, instead of local, parametrization for the orientation error state, which enables the direct analytical computation of the propagation Jacobian that fulfils the semigroup property of an appropriate state-transition matrix. Moreover, by adopting the observability-constrained methodology, we project the most accurate (canonical) measurement Jacobian – computed using the latest, and thus best, state

estimates – onto the observable subspace so that no spurious information is gained by the filter. As a result, the proposed STOC-VINS was shown to outperform the standard VINS algorithms, in terms of both consistency and accuracy. In the future, we will focusing on further improving VINS performance (including accuracy, consistency, and efficiency), e.g., how to efficiently integrate loop closure to enable long-term navigation while attaining bounded errors.

APPENDIX I PROOF OF LEMMA 4.1

Note first that the standard EKF computes the propagation Jacobian, $\Phi(t_k, t_\ell)$, using the current state estimates, $\hat{\mathbf{x}}_{k|k-1}$ and $\hat{\mathbf{x}}_{\ell|\ell}$. In order to verify the semigroup property (25), we substitute the pertinent state estimates to the analytical expressions of the IMU propagation Jacobians found in [13]. In particular, the multiplication of the two propagation Jacobians, $\Phi(t_{k+1}, t_k)$ and $\Phi(t_k, t_{k-1})$, does not alter the matrix structure as shown in (48).

Consider (1,1) entry of the above multiplication (48). Substitution of the current state estimates into the analytical expressions of $\Phi(t_{k+1}, t_k)$ and $\Phi(t_k, t_{k-1})$ [13] yields:

$$\begin{aligned} \Phi_{k,11}\Phi_{k-1,11} &= \mathbf{C} \begin{pmatrix} I^{(k+1|k)} \hat{\mathbf{q}} \\ I^{(k|k)} \hat{\mathbf{q}} \end{pmatrix} \mathbf{C}^T \begin{pmatrix} I^{(k-1|k-1)} \hat{\mathbf{q}} \\ I^{(k|k-1)} \hat{\mathbf{q}} \end{pmatrix} \\ &\neq \mathbf{C} \begin{pmatrix} I^{(k+1|k)} \hat{\mathbf{q}} \\ I^{(k-1|k-1)} \hat{\mathbf{q}} \end{pmatrix} = \Phi_{11}(t_{k+1}, t_{k-1}) \end{aligned} \quad (49)$$

where we have employed the fact that the propagated estimate of the IMU orientation generally is different from its updated estimate, i.e., $\mathbf{C} \begin{pmatrix} I \hat{\mathbf{q}}_{k|k-1} \\ I \hat{\mathbf{q}}_{k|k} \end{pmatrix} \neq \mathbf{C} \begin{pmatrix} I \hat{\mathbf{q}}_{k|k-1} \\ I \hat{\mathbf{q}}_{k|k} \end{pmatrix}$. The above inequality (49) immediately completes the proof.

APPENDIX II PROOF OF LEMMA 4.2

Due to the space constraint, we here prove only for the case *without* biases, while the proof for the general case with biases can be found in [32]. In this case, by removing the entries corresponding to the biases from (30), we have the propagation Jacobian $\Phi'(t_{k+1}, t_k)$ as follows:

$$\Phi'(t_{k+1}, t_k) = \begin{bmatrix} \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi'_{k,31} & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi'_{k,51} & \delta t_k \mathbf{I}_3 & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (50)$$

The product of the two consecutive propagation Jacobians assumes the following form [see (31)-(34)]:

$$\begin{aligned} \Gamma &:= \Phi'(t_{k+1}, t_k)\Phi'(t_k, t_{k-1}) = \\ &\begin{bmatrix} \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi'_{k,31} + \Phi'_{k-1,31} & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \Phi'_{k,51} + \delta t_k \Phi'_{k-1,31} + \Phi'_{k-1,51} & (\delta t_{k-1} + \delta t_k)\mathbf{I}_3 & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \end{aligned} \quad (51)$$

where

$$\begin{aligned}\Gamma_{21} &= \Phi'_{k,31} + \Phi'_{k-1,31} \\ &= - \left[\left({}^G \hat{\mathbf{v}}_{k+1|k} - {}^G \hat{\mathbf{v}}_{k-1|k-2} - {}^G \mathbf{g}(\delta t_{k-1} + \delta t_k) \right) \times \right] \\ &= \Phi'_{31}(t_{k+1}, t_{k-1})\end{aligned}\quad (52)$$

$$\begin{aligned}\Gamma_{31} &= \Phi'_{k,51} + \delta t_k \Phi'_{k-1,31} + \Phi'_{k-1,51} \\ &= - \left[\left({}^G \hat{\mathbf{p}}_{k+1|k} - {}^G \hat{\mathbf{p}}_{k|k-1} - {}^G \hat{\mathbf{v}}_{k|k-1} \delta t_k - \frac{1}{2} {}^G \mathbf{g} \delta t_k^2 \right) \times \right] \\ &\quad - \delta t_k \left[\left({}^G \hat{\mathbf{v}}_{k|k-1} - {}^G \hat{\mathbf{v}}_{k-1|k-2} - {}^G \mathbf{g} \delta t_{k-1} \right) \times \right] \\ &\quad - \left[\left({}^G \hat{\mathbf{p}}_{k|k-1} - {}^G \hat{\mathbf{p}}_{k-1|k-2} - {}^G \hat{\mathbf{v}}_{k-1|k-2} \delta t_{k-1} - \frac{1}{2} {}^G \mathbf{g} \delta t_{k-1}^2 \right) \times \right] \\ &= - \left[\left({}^G \hat{\mathbf{p}}_{k+1|k} - {}^G \hat{\mathbf{p}}_{k-1|k-2} - {}^G \hat{\mathbf{v}}_{k-1|k-2} (\delta t_k + \delta t_{k-1}) \right. \right. \\ &\quad \left. \left. - \frac{1}{2} {}^G \mathbf{g} (\delta t_k + \delta t_{k-1})^2 \right) \times \right] \\ &= \Phi'_{51}(t_{k+1}, t_{k-1})\end{aligned}\quad (53)$$

Note that the other (trivial) entries are easy to verify. Thus, this completes the proof for the case of no bias.

REFERENCES

- [1] D. Titterton and J. Weston, *Strapdown Inertial Navigation Technology*, 2nd ed. The Institution of Engineering and Technology, 2005.
- [2] J. Kim and S. Sukkariéh, "Real-time implementation of airborne inertial-SLAM," *Robotics and Autonomous Systems*, vol. 55, no. 1, pp. 62–71, Jan. 2007.
- [3] M. Li and A. I. Mourikis, "Improving the accuracy of EKF-based visual-inertial odometry," in *Proc. of the IEEE International Conference on Robotics and Automation*, Minneapolis, MN, May 2012, pp. 828–835.
- [4] M. Bryson and S. Sukkariéh, "Observability analysis and active control for airborne SLAM," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 44, no. 1, pp. 261–280, Jan. 2008.
- [5] A. Mourikis, N. Trawny, S. Roumeliotis, A. Johnson, A. Ansar, and L. Matthies, "Vision-aided inertial navigation for spacecraft entry, descent, and landing," *IEEE Transactions on Robotics*, vol. 25, no. 2, pp. 264–280, Apr. 2009.
- [6] S. Ebcin and M. Veth, "Tightly-coupled image-aided inertial navigation using the unscented Kalman filter," Air Force Institute of Technology, Dayton, OH, Tech. Rep., 2007.
- [7] D. Strelow, "Motion estimation from image and inertial measurements," Ph.D. dissertation, Carnegie Mellon University, 2004.
- [8] V. Indelman, S. Williams, M. Kaess, and F. Dellaert, "Information fusion in navigation systems via factor graph based incremental smoothing," *Robotics and Autonomous Systems*, vol. 61, no. 8, pp. 721–738, Aug. 2013.
- [9] G. P. Huang, A. I. Mourikis, and S. I. Roumeliotis, "Analysis and improvement of the consistency of extended Kalman filter-based SLAM," in *Proc. of the IEEE International Conference on Robotics and Automation*, Pasadena, CA, May 19–23, 2008, pp. 473–479.
- [10] —, "Observability-based rules for designing consistent EKF SLAM estimators," *International Journal of Robotics Research*, vol. 29, no. 5, pp. 502–528, Apr. 2010.
- [11] G. Huang, "Improving the consistency of nonlinear estimators: Analysis, algorithms, and applications," Ph.D. dissertation, University of Minnesota, Dept. of Computer Science and Engineering, 2012.
- [12] D. G. Kottas, J. A. Hesch, S. L. Bowman, and S. I. Roumeliotis, "On the consistency of vision-aided inertial navigation," in *Proc. of the 13th International Symposium on Experimental Robotics*, Quebec City, Canada, Jun. 17–20, 2012.
- [13] J. Hesch, D. Kottas, S. Bowman, and S. Roumeliotis, "Towards consistent vision-aided inertial navigation," in *Algorithmic Foundations of Robotics X*, ser. Springer Tracts in Advanced Robotics, E. Frazzoli, T. Lozano-Perez, N. Roy, and D. Rus, Eds. Springer Berlin Heidelberg, 2013, vol. 86, pp. 559–574.
- [14] —, "Consistency analysis and improvement of vision-aided inertial navigation," *IEEE Transactions on Robotics*, vol. 30, no. 1, pp. 158–176, Feb. 2013.
- [15] —, "Camera-IMU-based localization: Observability analysis and consistency improvement," *International Journal of Robotics Research*, vol. 33, pp. 182–201, 2014.
- [16] M. Li and A. I. Mourikis, "High-precision, consistent EKF-based visual-inertial odometry," *International Journal of Robotics Research*, vol. 32, no. 6, pp. 690–711, Jun. 2013.
- [17] Y. Bar-Shalom, X. R. Li, and T. Kirubarajan, *Estimation with applications to tracking and navigation*. New York: Wiley, 2001.
- [18] W. L. Brogan, *Modern Control Theory*. Upper Saddle River, NJ: Prentice Hall, 1991.
- [19] P. Corke, J. Lobo, and J. Dias, "An introduction to inertial and visual sensing," *International Journal of Robotics Research*, vol. 26, no. 6, pp. 519–535, Jun. 2007.
- [20] T. Lupton and S. Sukkariéh, "Visual-inertial-aided navigation for high-dynamic motion in built environments without initial conditions," *IEEE Transactions on Robotics*, vol. 28, no. 1, pp. 61–76, Feb. 2012.
- [21] A. Martinelli, "Vision and IMU data fusion: Closed-form solutions for attitude, speed, absolute scale, and bias determination," *IEEE Transactions on Robotics*, vol. 28, no. 1, pp. 44–60, 2012.
- [22] S. Shen, Y. Mulgaonkar, N. Michael, and V. Kumar, "Vision-based state estimation for autonomous rotorcraft mavs in complex environments," in *Proc. of the IEEE International Conference on Robotics and Automation*, Karlsruhe, Germany, May 6–10, 2013, pp. 1750–1756.
- [23] S. Weiss and R. Siegwart, "Real-time metric state estimation for modular vision-inertial systems," in *Proc. of the IEEE International Conference on Robotics and Automation*, 2011, pp. 4531–4537.
- [24] D. Kottas and S. Roumeliotis, "Exploiting urban scenes for vision-aided inertial navigation," in *Proc. of the Robotics: Science and Systems Conference*, Berlin, Germany, Jun. 24–28, 2013.
- [25] E. S. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *International Journal of Robotics Research*, vol. 30, no. 4, pp. 407–430, Apr. 2011.
- [26] J. Kelly and G. S. Sukhatme, "Visual-inertial sensor fusion: Localization, mapping and sensor-to-sensor self-calibration," *International Journal of Robotics Research*, vol. 30, no. 1, pp. 56–79, Jan. 2011.
- [27] F. M. Mirzaei and S. I. Roumeliotis, "A Kalman filter-based algorithm for IMU-camera calibration: Observability analysis and performance evaluation," *IEEE Transactions on Robotics*, vol. 24, no. 5, pp. 1143–1156, Oct. 2008.
- [28] A. Isidori, *Nonlinear Control Systems*. Springer, 1995.
- [29] G. P. Huang, A. I. Mourikis, and S. I. Roumeliotis, "A quadratic-complexity observability-constrained unscented Kalman filter for SLAM," *IEEE Transactions on Robotics*, vol. 29, no. 5, pp. 1226–1243, Oct. 2013.
- [30] N. Trawny and S. I. Roumeliotis, "Indirect Kalman filter for 3D attitude estimation," University of Minnesota, Dept. of Comp. Sci. & Eng., Tech. Rep., Mar. 2005.
- [31] P. S. Maybeck, *Stochastic Models, Estimation, and Control*, ser. Mathematics in Science and Engineering. London: Academic Press, 1979, vol. 141-1.
- [32] G. Huang, "Towards consistent visual-inertial navigation," Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, Tech. Rep., 2013. [Online]. Available: <http://people.csail.mit.edu/guang/paper/tr/stocvins.pdf>
- [33] A. I. Mourikis and S. I. Roumeliotis, "A multi-state constraint Kalman filter for vision-aided inertial navigation," in *Proceedings of the IEEE International Conference on Robotics and Automation*, Rome, Italy, Apr. 10–14, 2007, pp. 3565–3572.
- [34] S. I. Roumeliotis and J. W. Burdick, "Stochastic cloning: A generalized framework for processing relative state measurements," in *Proceedings of the IEEE International Conference on Robotics and Automation*, Washington, DC, May 11–15 2002, pp. 1788–1795.
- [35] J. Shi and C. Tomasi, "Good features to track," in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Jun. 21–23, 1994, pp. 593–600.
- [36] B. Lucas and T. Kanade, "An iterative image registration technique with an application to stereo vision," in *Proc. of the International Joint Conference on Artificial Intelligence*, Vancouver, BC, Aug. 1981, pp. 674–679.
- [37] D. Nister, "An efficient solution to the five-point relative pose problem," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 6, pp. 756–770, Jun. 2004.