

An iterative Kalman smoother for robust 3D localization and mapping

Dimitrios G. Kottas and Stergios I. Roumeliotis

Abstract In this paper, we present an iterative Kalman smoother (IKS) for robust 3D localization and mapping, using visual and inertial measurements. Contrary to extended Kalman filter (EKF) methods, smoothing increases the convergence rate of critical parameters (e.g., IMU’s velocity and camera’s clock drift), improves the positioning accuracy during challenging conditions (e.g., scarcity of visual features), and allows the immediate processing of visual observations. As opposed to existing smoothing approaches to VINS, based on the inverse filter (INVF), the proposed IKS exhibits superior numerical properties, allowing efficient implementations on mobile devices. Furthermore, we propose a classification of visual observations, for smoothing algorithms applied to VINS, based on their: (i) Track length, allowing their efficient processing as multi-state constraints, when possible and (ii) First observation, allowing their optional re-processing. Finally, we demonstrate the robustness of the proposed approach, over challenging indoor VINS scenarios, including, system (re)-initialization, and scarcity of visual observations.

1 Introduction and Related Work

Over the past decade, localization systems fusing inertial data, from an inertial measurement unit (IMU), with visual observations, from a camera [i.e., vision-aided inertial navigation systems (VINS)] have become a popular choice for GPS-denied navigation (e.g., indoors [1], or in space [2]). The dramatic upswing in manufacturing of low-cost, miniature IMUs and cameras, combined with the increasing capabilities of embedded computers, have made mobile devices (e.g., cell phones) excellent platforms for VINS, but raised also new challenges, when designing algorithms with improved robustness and efficiency characteristics.

The Maximum a Posteriori (MAP) estimator for VINS, corresponds to a batch least-squares (BLS) problem, over the platform’s trajectory, and the map of the environment. Unfortunately, BLS approaches cannot serve time-critical navigation ap-

The authors are with the Department of Computer Science and Engineering, University of Minnesota, USA (e-mail: dkottas|stergios@cs.umn.edu). This work was supported by the National Science Foundation (IIS-1328722) and UMN’s Doctoral Dissertation Fellowship.

plications (e.g., augmented reality), due to their unbounded processing and memory requirements, as the problem size increases with time [3].

Filtering approaches, on the other hand, maintain bounded processing time, by optimizing over a sliding window of the most recent visual and inertial measurements, while absorbing past measurements into a prior cost term. Depending on their representation of the prior information, filtering methods can be classified into extended Kalman filters (EKFs) and inverse filters (INVs) [4].

EKF-based algorithms exhibit excellent numerical stability, and have been demonstrated in real-time VINS implementations [5, 6]. Despite their efficiency, however, they do not allow re-processing the visual and/or inertial measurements within the optimization window, which may severely affect their performance under challenging conditions. Specifically, consider visual observations, that arrive as feature tracks, spanning a set of images within the estimator’s sliding window. EKF-based approaches (e.g., the MSC-KF [2]), postpone their processing till they have reached their maximum length, which causes a delay in the state correction equal to the time between when a feature track is available for processing (i.e., when it spans at least two images), and when it is actually processed. Although small (typically less than 2 sec for personal localization), such a delay can affect the accuracy of time-critical navigation systems. Furthermore, under scenarios with increased linearization errors [e.g., (re)-initialization (after failure) of the system’s operation, when the available estimates of critical quantities, such as the IMU’s velocity and biases, are of low accuracy], re-processing visual and inertial measurements, can significantly improve accuracy and robustness, by increasing their rate of convergence.

Such appealing capabilities of filtering algorithms, are supported by the INV (e.g., [7, 8, 9, 10]), which allows re-processing all inertial and visual measurements within the optimization window considered. Unfortunately, due to the high condition number of the Hessian, INVs typically require 64-bit precision, reducing their efficiency, especially when considering that most current mobile devices feature ARM processors and NEON co-processors that provide a 4-fold processing speedup when using 32-bit precision. Note also that, with the exception of [11], existing INVs for VINS do not classify visual observations based on their track length, not allowing their efficient processing (e.g., as in the MSC-KF) when possible.

To overcome the limitations of the EKF and INV when applied to visual-inertial odometry (VIO), in [12], we introduced an iterative Kalman smoother (IKS), that shares the advantages of both approaches. In this paper, we extend the methodology of [12], to the most general case for VINS, which allows (re)-processing visual data either using a VIO approach (as in the MSC-KF), or as SLAM landmarks when their track length exceeds the estimator’s sliding window. In particular, we introduce a sliding window IKS for VINS with the following key advantages:

- The IKS iteratively re-linearizes both inertial and camera measurements within the estimator’s window, and re-processes visual data over multiple overlapping sliding-window epochs, thus improving robustness and increasing accuracy.
- The IKS employs a covariance matrix, as well as a set of linearized constraints (instead of a Hessian) for representing prior information, thus inheriting the superior numerical properties of the EKF and leading to efficient implementations.

Besides the detailed analysis and description of the proposed IKS algorithm, we demonstrate its robustness and assess its accuracy in simulations and experiments over challenging indoor VINS scenarios, including, filter initialization and scarcity of feature tracks due to sudden turns or camera occlusions. Finally, we provide a timing comparison between the IKS and the EKF, using a mobile processor.

2 Vision-aided Inertial Navigation System (VINS)

The system state¹ at time t_k is given by $\mathbf{x}_k = [\mathbf{x}_{I_k}^T \ ^{I_k}\boldsymbol{\ell}_k^T]^T$ where \mathbf{x}_{I_k} contains all kinematic quantities, describing the motion of the IMU's frame of reference $\{I_k\}$, while $^{I_k}\boldsymbol{\ell}_k$ comprises landmarks of the environment. In particular, $\mathbf{x}_{I_k} = [{}^{I_k}\mathbf{q}_G^T \ ^G\mathbf{p}_{I_k}^T \ ^G\mathbf{v}_{I_k}^T \ \mathbf{b}_{a_k}^T \ \mathbf{b}_{g_k}^T]^T$, where $^{I_k}\mathbf{q}_G$ is the quaternion representation of the orientation of the global frame $\{G\}$ in $\{I_k\}$, $^G\mathbf{v}_{I_k}$ and $^G\mathbf{p}_{I_k}$ are the velocity and position of $\{I_k\}$ in $\{G\}$ respectively, while \mathbf{b}_{a_k} and \mathbf{b}_{g_k} correspond to the gyroscope and accelerometer biases. Finally, $^{I_k}\boldsymbol{\ell}_k$ comprises N_k landmarks, i.e., $^{I_k}\boldsymbol{\ell}_k = [{}^{I_k}\mathbf{f}_1^T \dots \ ^{I_k}\mathbf{f}_{N_k}^T]^T$ where $^{I_k}\mathbf{f}_j$ denotes the inverse-depth parameterization of feature \mathbf{f}_j in $\{I_k\}$.

2.1 Inertial Measurements

The IMU provides measurements of the platform's rotational velocity and linear acceleration, contaminated by white Gaussian noise and time-varying biases. Let $\mathbf{u}_{k,k+1}$ denote the inertial measurements within the time interval $[t_k, t_{k+1}]$, which through integration, define a constraint (discrete-time process model) of the form:

$$\mathbf{x}_{I_{k+1}} = \mathbf{f}(\mathbf{x}_{I_k}, \mathbf{u}_{k,k+1} - \mathbf{w}_{k,k+1}) \quad (1)$$

where $\mathbf{w}_{k,k+1}$ is a discrete-time zero-mean white Gaussian noise process with covariance \mathbf{Q}_k .² Linearizing (1), at the state estimates corresponding to the two consecutive states, $\mathbf{x}_{I_k}^*$, $\mathbf{x}_{I_{k+1}}^*$, results in the following IMU measurement model, relating the error states $\tilde{\mathbf{x}}_{I_k}^*$ and $\tilde{\mathbf{x}}_{I_{k+1}}^*$:

$$\tilde{\mathbf{x}}_{I_{k+1}}^* = \mathbf{r}_{u_{k,k+1}}^* + \boldsymbol{\Phi}_{k+1,k}^* \tilde{\mathbf{x}}_{I_k}^* + \mathbf{G}_{k+1,k}^* \mathbf{w}_{k,k+1} \quad (2)$$

where $\mathbf{r}_{u_{k,k+1}}^* := \mathbf{f}(\mathbf{x}_{I_k}^*, \mathbf{u}_{k,k+1}) - \mathbf{x}_{I_{k+1}}^*$, and we have defined the error state $\tilde{\mathbf{x}}_{I_k}^*$ as the difference between the true state \mathbf{x}_{I_k} and the linearization point $\mathbf{x}_{I_k}^*$ (i.e., $\tilde{\mathbf{x}}_{I_k}^* = \mathbf{x}_{I_k} - \mathbf{x}_{I_k}^*$).

¹ Without loss of generality, we assume that the IMU-camera extrinsic calibration is the identity transformation and the clocks of the two sensors are perfectly synchronized. In practice, both are included in the system's state following the methodologies described in [13] and [6], respectively.

² The interested reader is referred to [14, 1, 5], and references there-in, for details on IMU integration.

$\mathbf{x}_{I_k}^*$).³ The Jacobians $\Phi_{k+1,k}^*$ and $\mathbf{G}_{k+1,k}^*$ are evaluated at $\mathbf{x}_{I_k}^*$, $\mathbf{x}_{I_{k+1}}^*$, and are available in numerical or analytical form [1]. Although the corresponding cost term,

$$c_{u_{k,k+1}}(\tilde{\mathbf{x}}_{k:k+1}^*) = \|\mathbf{r}_{u_{k,k+1}}^* - [-\Phi_{k+1,k}^* \mathbf{I}] \begin{bmatrix} \tilde{\mathbf{x}}_{I_k}^* \\ \tilde{\mathbf{x}}_{I_{k+1}}^* \end{bmatrix}\|_{\mathbf{Q}_k^*}^2$$

where $\tilde{\mathbf{x}}_{k:k+1}^* := [\tilde{\mathbf{x}}_{I_k}^{*T} \tilde{\mathbf{x}}_{I_{k+1}}^{*T}]^T$ and $\mathbf{Q}_k^* = \mathbf{G}_{k+1,k}^* \mathbf{Q}_k \mathbf{G}_{k+1,k}^{*T}$, can be re-linearized multiple times within the INVF framework, current EKF-based approaches linearize (1) only *once*, during state and covariance propagation (i.e., every time a new inertial measurement becomes available). This limitation may negatively affect performance when the linearization errors are large (e.g., during system initialization).

2.2 Visual Observations

In order for a camera to provide kinematic information, a feature-tracking pipeline is required. A common choice comprises a feature-extraction method (e.g., the Harris corner [15]) along with a tracking algorithm (e.g., the Kanade-Lucas-Tomasi (KLT) [16]). Once a new image arrives, point features from the previous image, are tracked to the new one, while new features are extracted from areas that just entered the camera's field of view [1]. An example of the feature tracks generated from such an image-processing pipeline is shown in Fig. 1.

Consider a point feature \mathbf{f}_j , observed in camera poses $\mathbf{x}_{I_{k+1:k+N_j}}$, where $N_j \leq M$ and M is the window's length. We represent \mathbf{f}_j , using its homogeneous coordinates and inverse distance in $\{I_{k+1}\}$, hereafter denoted by ${}^{I_{k+1}}\mathbf{f}_j$, or using its Cartesian coordinates, ${}^{I_{k+1}}\mathbf{p}_{f_j}$. For the m -th measurement, $m \in [1, \dots, N_j]$, the observation $\mathbf{z}_{k+m,j}$ acquired by a calibrated camera is:

$$\mathbf{z}_{k+m,j} = \pi(C({}^{I_{k+m}}\mathbf{q}_{I_{k+1}})^{I_{k+1}}\mathbf{p}_{f_j} + {}^{I_{k+m}}\mathbf{p}_{I_{k+1}}) + \mathbf{n}_{k+m,j} \quad (3)$$

where, $\pi(\begin{bmatrix} x & y & z \end{bmatrix}^T) = \begin{bmatrix} \frac{x}{z} & \frac{y}{z} \end{bmatrix}^T$, while $\mathbf{n}_{k+m,j} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_2)$, and \mathbf{I}_2 is the 2×2 identity matrix. Linearizing (3), yields:

$$\tilde{\mathbf{z}}_{k+m,j}^* = \mathbf{H}_{R,k+1,j}^* \tilde{\mathbf{x}}_{I_{k+1}}^* + \mathbf{H}_{R,k+m,j}^* \tilde{\mathbf{x}}_{I_{k+m}}^* + \mathbf{F}_{k+m,j}^* {}^{I_{k+1}}\tilde{\mathbf{f}}_j + \mathbf{n}_{k+m,j}$$

Collecting all N_j observations of feature \mathbf{f}_j in one vector, yields:

$$\tilde{\mathbf{z}}_j^* = \mathbf{H}_{R,j}^* \tilde{\mathbf{x}}_{I_{k+1:k+N_j}}^* + \mathbf{F}_j^{*I_{k+1}} \tilde{\mathbf{f}}_j + \mathbf{n}_j \quad (4)$$

which corresponds to the cost term:

$$c_{z_{f_j}}(\tilde{\mathbf{x}}_{I_{k+1:k+N_j}}^*, {}^{I_{k+1}}\tilde{\mathbf{f}}_j) = \|\tilde{\mathbf{z}}_j^* - \mathbf{H}_{R,j}^* \tilde{\mathbf{x}}_{I_{k+1:k+N_j}}^* - \mathbf{F}_j^{*I_{k+1}} \tilde{\mathbf{f}}_j\|_{\sigma^2 \mathbf{I}}^2 \quad (5)$$

Consider an orthonormal matrix Θ_j , partitioned as $\Theta_j = [\mathbf{S}_j \mathbf{U}_j]$, where the 3 columns of \mathbf{S}_j span the column space of \mathbf{F}_j^* , while the $2N_j - 3$ columns of \mathbf{U}_j , its left null space. Projecting (4) onto Θ_j , partitions $c_{z_{f_j}}$ into two parts:

³ For the quaternion \mathbf{q} we employ a multiplicative error model $\tilde{\mathbf{q}} = \mathbf{q} \otimes \mathbf{q}^{*-1} \simeq [\frac{1}{2} \delta \boldsymbol{\theta}_{3 \times 1}^T \ 1]^T$.

$$\begin{aligned}
c_{z_{f_j}}(\tilde{\mathbf{x}}_{k+1:k+N_j}^*, I_{k+1}\tilde{\mathbf{f}}_j) &= \|\Theta_j^T (\mathbf{z}_j^* - \mathbf{H}_{R,j}^* \tilde{\mathbf{x}}_{k+1:k+N_j}^* - \mathbf{F}_j^{*I_{k+1}} \tilde{\mathbf{f}}_j)\|_{\sigma^2 \mathbf{I}_{2N_j}}^2 \\
&= \|\mathbf{r}_j^{K*} - \mathbf{H}_j^{K*} \tilde{\mathbf{x}}_{k+1:k+N_j}^*\|_{\sigma^2 \mathbf{I}_{2N_j-3}}^2 + \|\mathbf{r}_j^{M*} - \mathbf{H}_j^{M*} \tilde{\mathbf{x}}_{k+1:k+N_j}^* - \mathbf{R}_j^{M*I_{k+1}} \tilde{\mathbf{f}}_j\|_{\sigma^2 \mathbf{I}_3}^2 \\
&= c_{z_{f_j}}^K(\tilde{\mathbf{x}}_{k+1:k+N_j}^*) + c_{z_{f_j}}^M(\tilde{\mathbf{x}}_{k+1:k+N_j}^*, I_{k+1}\tilde{\mathbf{f}}_j)
\end{aligned} \tag{6}$$

with $\mathbf{r}_j^{K*} = \mathbf{U}_j^T \tilde{\mathbf{z}}_j^*$, $\mathbf{H}_j^{K*} = \mathbf{U}_j^T \mathbf{H}_{R,j}^*$, and $\mathbf{r}_j^{M*} = \mathbf{S}_j^T \tilde{\mathbf{z}}_j^*$, $\mathbf{H}_j^{M*} = \mathbf{S}_j^T \mathbf{H}_{R,j}^*$, $\mathbf{R}_j^{M*} = \mathbf{S}_j^T \mathbf{F}_{R,j}^*$. The second term, $c_{z_{f_j}}^M$ contains *all* information regarding feature f_j , while $c_{z_{f_j}}^K$ defines a multi-state constraint *only* among the poses $\mathbf{x}_{I_{k+1:k+N_j}}$. For feature tracks that do not exceed the estimator’s window, as in the MSC-KF [2], we consider *only* the cost term $c_{z_{f_j}}^K$. Specifically, since \mathbf{R}_j^{M*} is an invertible square matrix and for any $\tilde{\mathbf{x}}_{k+1:k+N_j}^*$ there exists a $I_{k+1}\tilde{\mathbf{f}}_j$, such that $c_{z_{f_j}}^M$ is exactly zero, minimizing (6) is equivalent to minimizing $c_{z_{f_j}}^K(\tilde{\mathbf{x}}_{k+1:k+N_j}^*)$. As we describe later on, for feature tracks whose span exceeds the sliding window (i.e., SLAM landmarks), $c_{z_{f_j}}^M$ allows initializing them into the estimator’s map, and subsequently optimizing their measurements across *non-overlapping* epochs of the sliding window.

2.3 Visual Observations in Sliding-Window Estimators

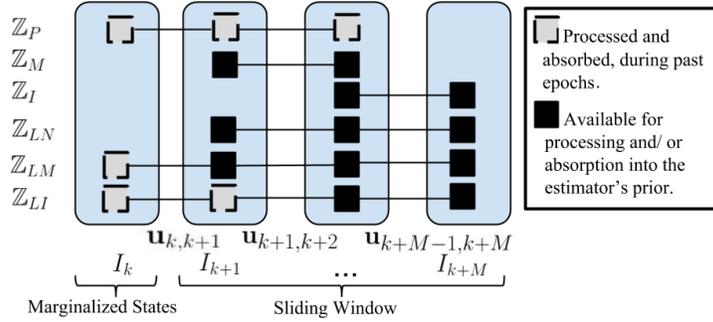


Fig. 1: At t_{k+M} , there exist 6 categories of feature tracks described in Section 2.3.

We classify visual observations, in sliding-window estimators, based on their:

- Track length** which distinguishes SLAM landmarks from MSC-KF features, since for the latter, their track length does not exceed the estimator’s window. Thus, optimizing over the MSC-KF feature’s coordinates is not required for minimizing their corresponding cost terms [see (6)].
- Earliest observation** which if it involves the sliding window’s “tail” (oldest) state, it does not allow postponing their absorption into the estimator’s prior.

In particular, we distinguish the following categories of visual observations:

- **Past features** (\mathbb{Z}_P) corresponding to visual observations that were absorbed in a past epoch of the sliding window and cannot be re-processed by any filter.
- **MSC-KF features** whose tracking length does not exceed the estimator’s window, i.e., $N_j \leq M$, hence they are not mapped but rather used only for providing multi-state constraints involving the camera poses observing them. Based on their earliest observation, we further classify MSC-KF features into 2 sets:
 - **Mature features** (\mathbb{Z}_M): These are MSC-KF features that have reached their maximum length. Both the EKF and the INVf linearize, process, and absorb them in a single step. Note also that the INVf, as well as the Iterative EKF (I-EKF) [4] can re-linearize these observations.
 - **Immature features** (\mathbb{Z}_I): This set represents feature tracks that have already entered the image-processing pipeline, but have not reached their maximum length, yet. Although the INVf can use (and re-linearize) these measurements multiple times, across overlapping epochs of the sliding window (from the second time they are observed till the track exits the optimization window), the EKF is not able to do so. This limitation introduces a delay, between the “birth” of a feature track and its impact on the filter’s state estimates; a potential drawback of EKF-based approaches for time-critical applications.
- **SLAM landmarks** corresponding to features, whose track length exceeds the estimator’s optimization window, i.e., $N_j > M$, and hence their optimal processing requires including them into the estimator’s map of the environment, ${}^{l_{k+1}}\ell_{k+M}$. Within a single epoch of the sliding window, observations to SLAM landmarks can be partitioned into 3 categories:
 - **Mature landmarks** (\mathbb{Z}_{LM}): These are observations to previously initialized SLAM landmarks, which include the estimator’s “tail” pose, hence their absorption into the estimator’s prior, cannot be postponed for future epochs of the sliding window.
 - **Immature landmarks** (\mathbb{Z}_{LI}): These correspond to measurements of previously initialized SLAM landmarks, which do not involve the estimator’s “tail” pose, thus their absorption can be postponed to later epochs, till one of them includes the estimator’s “tail” pose, i.e., they become mature landmarks.
 - **New SLAM landmarks** (\mathbb{Z}_{LN}): These are feature tracks that have not yet reached their maximum length, while their present observations span all cameras within the estimator’s optimization window. Hence, they will be processed, absorbed, and initialized as new landmarks in the estimator’s map.

3 Estimation Algorithm Description

When designing the proposed IKS our objective is two-fold: (i) Process all inertial and visual observations within the current epoch of the sliding window $\mathbf{x}_{k+1:k+M}$ (i.e., inertial measurements $\{\mathbf{u}_{\ell,\ell+1}\}$, for $k+1 \leq \ell \leq k+M-1$, and feature tracks \mathbb{Z}_{LM} , \mathbb{Z}_{LI} , \mathbb{Z}_{LN} , \mathbb{Z}_M , and \mathbb{Z}_I), and (ii) Allow future epochs to re-process all measurements

that are independent of the sliding window’s “tail” state $\mathbf{x}_{I_{k+1}}$ (i.e., the inertial measurements $\{\mathbf{u}_{\ell,\ell+1}\}$, for $k+2 \leq \ell \leq k+M-1$, and visual observations to immature MSC-KF features and SLAM landmarks (i.e., \mathbb{Z}_I and \mathbb{Z}_{LI} , respectively).

3.1 IKS Algorithm: Input

Before image $k+M$ arrives, the proposed IKS maintains:

- 1) **A set of linearization points**, over the estimator’s sliding-window of camera poses $\mathbf{x}_{I_{k+1:k+M-1}}^*$, and landmarks ${}^{I_{k+1}}\boldsymbol{\ell}_{k+M-1}^*$ that represent the estimator’s best estimates given all measurements up to t_{k+M-1} .
- 2) **A prior** comprising:
 - (a) The pdf of the oldest state, $\mathbf{x}_{I_{k+1}}$, within the sliding window approximated as a Gaussian $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1}}^\ominus, \mathbf{P}_{I_{k+1}}^\ominus)$.
 - (b) A set of N_L linearized constraints relating the oldest state, $\mathbf{x}_{I_{k+1}}^*$, with the rest of the poses within the sliding window, expressed as:

$$\mathbf{r}_L^{*\ominus} = \mathbf{H}_L^{*\ominus} \tilde{\mathbf{x}}_{I_{k+1:k+M-1}}^* + \mathbf{n}_L, \quad \mathbf{n}_L \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{N_L}). \quad (7)$$

- (c) A set of $3N_{k+M-1}$ linearized constraints relating the oldest state, $\mathbf{x}_{I_{k+1}}^*$, with the rest of the poses within the sliding window and the SLAM landmarks ${}^{I_{k+1}}\boldsymbol{\ell}_{k+M-1}$ expressed as:

$$\mathbf{r}_M^{*\ominus} = \mathbf{H}_M^{*\ominus} \tilde{\mathbf{x}}_{I_{k+1:k+M-1}}^* + \mathbf{F}_M^{*\ominus} {}^{I_{k+1}}\tilde{\boldsymbol{\ell}}_{k+M-1}^* + \mathbf{n}_M, \quad \mathbf{n}_M \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{3N_{k+M}}). \quad (8)$$

The ensemble of the pdf $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1}}^\ominus, \mathbf{P}_{I_{k+1}}^\ominus)$ and the linearized constraints $\{\mathbf{r}_L^{*\ominus}, \mathbf{H}_L^{*\ominus}\}$ and $\{\mathbf{r}_M^{*\ominus}, \mathbf{H}_M^{*\ominus}, \mathbf{F}_M^{*\ominus}\}$ in (7) and (8) represent all information for the poses $\mathbf{x}_{I_{k+1:k+M-1}}$ and the landmarks ${}^{I_{k+1}}\boldsymbol{\ell}_{k+M-1}$, accumulated through absorption of past visual observations (i.e., \mathbb{Z}_P in Fig. 1) and inertial measurements (i.e., $\{\mathbf{u}_{\ell,\ell+1}, \ell \leq k\}$).

3.2 IKS Algorithm: Overview

A single recursion of the IKS, involves the following steps:

1. Propagation: The prior pdf, $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1}}^\ominus, \mathbf{P}_{I_{k+1}}^\ominus)$, of $\mathbf{x}_{I_{k+1}}$ and the inertial measurements $\{\mathbf{u}_{\ell,\ell+1}\}$, for $k+1 \leq \ell \leq k+M-1$, are used for creating a prior $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+M}}^\ominus, \mathbf{P}_{I_{k+1:k+M}}^\ominus)$ for all the poses within the sliding window.

2. State Update: All available feature tracks, either from SLAM landmarks, i.e., $\mathbb{Z}_{LM}, \mathbb{Z}_{LN}$ and \mathbb{Z}_{LI} , and MSC-KF features, i.e., \mathbb{Z}_M and \mathbb{Z}_I , as well as the prior constraints $\{\mathbf{r}_L^{*\ominus}, \mathbf{H}_L^{*\ominus}\}$ and $\{\mathbf{r}_M^{*\ominus}, \mathbf{H}_M^{*\ominus}, \mathbf{F}_M^{*\ominus}\}$ are processed for updating the current state estimates $\mathbf{x}_{I_{k+1:k+M}}^*$. This state-optimization can be performed iteratively.

3. Landmark Propagation: All landmarks are propagated to the next “tail” of the sliding window, $\mathbf{x}_{I_{k+2}}$, i.e., ${}^{I_{k+2}}\boldsymbol{\ell}_{k+M} = [{}^{I_{k+2}}\mathbf{f}_1^T \dots {}^{I_{k+2}}\mathbf{f}_{N_{k+M}}^T]^T$.

4. Covariance Update: The measurements, $\mathbb{Z}_{LM}, \mathbb{Z}_{LN}, \mathbb{Z}_M$, and $\mathbf{u}_{k+1,k+2}$, which are about to be absorbed, are used to compute the posterior covariance $\mathbf{P}_{I_{k+2}}^\oplus$ of $\mathbf{x}_{I_{k+2}}$,

which will become the new ‘‘tail’’ of the sliding window.

5. Construction of the next epoch’s prior constraints: The prior constraints $\{\mathbf{r}_L^{\star\ominus}, \mathbf{H}_L^{\star\ominus}\}$ and $\{\mathbf{r}_M^{\star\ominus}, \mathbf{H}_M^{\star\ominus}, \mathbf{F}_M^{\star\ominus}\}$ are updated so that they become independent of the state to be marginalized, $\mathbf{x}_{I_{k+1}}$, and instead reflect the new constraints between $\mathbf{x}_{I_{k+2}}$, $\mathbf{x}_{I_{k+3:k+M}}$, and ${}^{I_{k+2}}\boldsymbol{\ell}_{k+M}$.

3.3 IKS Algorithm: Detailed Description

In order to allow for a direct comparison with the INVF and the EKF, we follow a two-level presentation of the IKS: We first describe the effect that each step has on the cost function being minimized and then present the corresponding IKS equations. We start by stating that the IKS (iteratively) minimizes the cost function:

$$c_{k+M}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^{\star}, {}^{I_{k+1}}\tilde{\boldsymbol{\ell}}_{k+M}^{\star}) = c_{P_{I_{k+1}}^{\ominus}} + c_u + c_L + c_{z_M}^K + c_{z_I}^K + c_{z_{LM}}^K + c_{z_{LI}}^K + c_{z_{LN}}^K \quad (9) \\ + c_M + c_{z_{LM}}^M + c_{z_{LI}}^M + c_{z_{LN}}^M$$

where $c_{P_{I_{k+1}}^{\ominus}}$ corresponds to the prior pdf of the oldest state within the sliding window, \mathbf{x}_{k+1} , $c_u = \sum_{\ell=k+1}^{k+M-1} c_{u_{\ell,\ell+1}}$ to the inertial measurements $\mathbf{u}_{k+1:k+M}$ [see (2)], c_L to prior information about the poses $\mathbf{x}_{k+1:k+M-1}$ [see (7)], $c_{z_M}^K$, $c_{z_I}^K$, $c_{z_{LM}}^K$, $c_{z_{LI}}^K$, and $c_{z_{LN}}^K$ to geometric constraints between the poses from all available visual observations [see (6)], c_M to prior information about the SLAM landmarks ${}^{I_{k+1}}\boldsymbol{\ell}_{k+M-1}$, $c_{z_{LM}}^M$ [see (8)] and $c_{z_{LI}}^M$ to feature constraints between the poses *and* the SLAM landmarks [see (6)], and finally $c_{z_{LN}}^M$ to feature constraints for the *new* SLAM landmarks, ${}^{I_{k+1}}\boldsymbol{\ell}_{\mathcal{N}}$ [see (6)].

Hereafter, we employ the cost terms in (9) to describe the four main steps of the proposed IKS (see Section 3.2).

3.3.1 Prior Propagation

The prior pdf $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1}}^{\ominus}, \mathbf{P}_{I_{k+1}}^{\ominus})$ and the inertial measurements $\mathbf{u}_{k+1:k+M}$ are used to generate a prior pdf $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+M}}^{\ominus}, \mathbf{P}_{I_{k+1:k+M}}^{\ominus})$ over *all* the states, $\mathbf{x}_{I_{k+1:k+M}}$, within the sliding window. Through this process, the cost function in (9) takes the form:

$$c_{k+M}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^{\star}, {}^{I_{k+1}}\tilde{\boldsymbol{\ell}}_{k+M}^{\star}) = c_{P_{I_{k+1:k+M}}^{\ominus}} + c_L + c_{z_M}^K + c_{z_I}^K + c_{z_{LM}}^K + c_{z_{LI}}^K + c_{z_{LN}}^K \quad (10) \\ + c_M + c_{z_{LM}}^M + c_{z_{LI}}^M + c_{z_{LN}}^M$$

where, $c_{P_{I_{k+1:k+M}}^{\ominus}} = c_{P_{I_{k+1}}^{\ominus}} + \sum_{\ell=k+1}^{k+M-1} c_{u_{\ell,\ell+1}}$ corresponds to the prior $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+M}}^{\ominus}, \mathbf{P}_{I_{k+1:k+M}}^{\ominus})$. The mean $\hat{\mathbf{x}}_{I_{k+1:k+M}}^{\ominus}$ is computed as:

$$\hat{\mathbf{x}}_{I_{k+i}}^{\ominus} = \begin{cases} \hat{\mathbf{x}}_{I_{k+1}}^{\ominus} & , i = 1 \\ \mathbf{f}(\mathbf{x}_{I_{k+i-1}}^*, \mathbf{u}_{k+i-1, k+i}) + \Phi_{k+i, k+i-1}^* \delta \mathbf{x}_{I_{k+i-1}}^{\ominus} & , 2 \leq i \leq M \end{cases} \quad (11)$$

where $\delta \mathbf{x}_{I_{k+i-1}}^{\ominus} = \hat{\mathbf{x}}_{I_{k+i-1}}^{\ominus} - \mathbf{x}_{I_{k+i-1}}^*$. Note that for the EKF inertial measurements are linearized and processed once, soon as they become available; hence, $\delta \mathbf{x}_{I_{k+i-1}}^{\ominus} = \mathbf{0}$ and the mean of the prior pdf, $\hat{\mathbf{x}}_{I_{k+i}}^{\ominus}$, coincides with the linearization point $\mathbf{x}_{I_{k+i}}^*$. In contrast, the IKS re-processes inertial measurements by re-computing the prior over the sliding window $\mathbf{x}_{I_{k+1:k+M}}$ through the process described in (11).

The block elements of the covariance $\mathbf{P}_{I_{k+1:k+M}}^{\ominus}$ are computed through the EKF covariance propagation recursion, using, however, the most recent state estimates:

$$\begin{aligned} \mathbf{P}_{I_{k+i}}^{\ominus} &= \Phi_{k+i, k+i-1}^* \mathbf{P}_{I_{k+i-1}}^{\ominus} \Phi_{k+i, k+i-1}^{*T} + \mathbf{Q}_k^*, \quad i = 2, \dots, M \\ \mathbf{P}_{I_{k+i, k+j}}^{\ominus} &= \Phi_{k+i, k+j}^* \mathbf{P}_{I_{k+j}}^{\ominus}, \quad i = 2, \dots, M, \quad j = 1, \dots, i-1 \end{aligned} \quad (12)$$

3.3.2 State Update

All cost terms, which provide multi-state (i.e., $c_L, c_M^K, c_{z_i}^K, c_{z_{LM}}^K, c_{z_{LI}}^K, c_{z_{LN}}^K$), as well as mapping (i.e., $c_M, c_{z_{LM}}^M, c_{z_{LI}}^M, c_{z_{LN}}^M$) constraints, are used for updating the state estimates for $\mathbf{x}_{I_{k+1:k+M}}$ and $I_{k+1} \tilde{\ell}_{k+M}$. Although each of these terms could have been used independently, in successive updates, we choose to first merge them into two cost terms, c_S^K and c_S^M , comprising multi-state geometric and mapping constraints, respectively, and process them in a batch form.

Measurement Compression: Combining all mapping terms, $c_M, c_{z_{LM}}^M, c_{z_{LI}}^M, c_{z_{LN}}^M$, into a single cost term, $c_S^{M'}$, yields:

$$\begin{aligned} c_S^{M'}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+1} \tilde{\ell}_{k+M}) &= \left\| \begin{bmatrix} \mathbf{r}_{LM}^{* \ominus} \\ \mathbf{r}_{LI}^{* \ominus} \\ \mathbf{r}_{LN}^{* \ominus} \end{bmatrix} - \begin{bmatrix} \mathbf{H}_{LM}^{* \ominus} & \mathbf{F}_{LM}^{* \ominus} \\ \mathbf{H}_{LI}^{* \ominus} & \mathbf{F}_{LI}^{* \ominus} \\ \mathbf{H}_{LN}^{* \ominus} & \mathbf{F}_{LN}^{* \ominus} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^* \\ I_{k+1} \tilde{\ell}_{k+M-1} \\ I_{k+1} \tilde{\ell}_{\mathcal{N}} \end{bmatrix} \right\|_{\sigma^2 \mathbf{I}} \\ &= \left\| \mathbf{r}_S^{M'} - [\mathbf{H}_S^{M'} \mathbf{F}_S^{M'}] \begin{bmatrix} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^* \\ I_{k+1} \tilde{\ell}_{k+M} \end{bmatrix} \right\|_{\sigma^2 \mathbf{I}} \end{aligned} \quad (13)$$

As in Section 2.2, we project the linearized constraints of (13) onto the column space and left nullspace of $\mathbf{F}_S^{M'}$, which partitions $c_S^{M'}$ into, $3N_{k+M}$ constraints, denoted by c_S^M , providing information *only* over the landmarks ℓ_{k+M} and a cost term $c_S^{K'}$, providing geometric constraints, *only* over the poses $\mathbf{x}_{I_{k+1:k+M}}$, i.e.,

$$\begin{aligned} c_S^{M'}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+1} \tilde{\ell}_{k+M}) &= \left\| \mathbf{r}_S^{M'} - [\mathbf{H}_S^{M'} \mathbf{F}_S^{M'}] \begin{bmatrix} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^* \\ I_{k+1} \tilde{\ell}_{k+M} \end{bmatrix} \right\|_{\sigma^2 \mathbf{I}_{3N_{k+M}}} \\ &\quad + \left\| \mathbf{r}_S^{K'} - \mathbf{H}_S^{K'} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^* \right\|_{\sigma^2 \mathbf{I}} \\ &= c_S^M(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+1} \tilde{\ell}_{k+M}) + c_S^{K'}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*) \end{aligned}$$

which allows c_{k+M} in (10) to take the form:

$$c_{k+M}(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^*) = c_{P_{I_{k+1:k+M}}^\ominus} + c_S^M + c_S^K \quad (14)$$

where,⁴ c_S^K comprises all geometric constraints, i.e., c_L , $c_S^{K'}$, $c_{z_M}^K$, $c_{z_I}^K$, $c_{z_{LM}}^K$, $c_{z_{LI}}^K$, and $c_{z_{LN}}^K$:

$$c_S^K(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*) = \|\mathbf{r}_S^{K*} - \mathbf{H}_S^{K*} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^*\|_{\sigma^2 \mathbf{I}}^2. \quad (15)$$

Note that, since \mathbf{F}_S^{M*} is a square full-rank matrix and for every $\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*$ there exists an $I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^*$ that minimizes c_S^M to zero, for minimizing (14), it suffices to first minimize $c_{P_{I_{k+1:k+M}}^\ominus} + c_S^K$, over $\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*$, and then solve for $I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^*$, using c_S^M . Specifically, through an I-EKF update step [4] we first update the poses $\mathbf{x}_{I_{k+1:k+M}}$:

$$\mathbf{x}_{I_{k+1:k+M}}^{\star\oplus} = \mathbf{x}_{I_{k+1:k+M}}^* + \tilde{\mathbf{x}}_{I_{k+1:k+M}}^o, \quad \tilde{\mathbf{x}}_{I_{k+1:k+M}}^o = \delta \mathbf{x}_{I_{k+1:k+M}}^\ominus + \mathbf{P}_{I_{k+1:k+M}}^\ominus \mathbf{H}_S^{K*T} \mathbf{d}_S \quad (16)$$

where \mathbf{d}_S is the solution to the linear system $\mathbf{S}^* \mathbf{d}_S = \mathbf{r}_S^{K*} - \mathbf{H}_S^{K*} \delta \mathbf{x}_{I_{k+1:k+M}}^\ominus$, with $\mathbf{S}^* = \mathbf{H}_S^{K*} \mathbf{P}_{I_{k+1:k+M}}^\ominus \mathbf{H}_S^{K*T} + \sigma^2 \mathbf{I}$, and $\delta \mathbf{x}_{I_{k+1:k+M}}^\ominus = \hat{\mathbf{x}}_{I_{k+1:k+M}}^\ominus - \mathbf{x}_{I_{k+1:k+M}}^*$.

Second, we substitute $\tilde{\mathbf{x}}_{I_{k+1:k+M}}^o$ into c_S^M , and solve⁵ the linear system

$$\mathbf{F}_S^{M*} I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^o = \mathbf{r}_S^{M*} - \mathbf{H}_S^{M*} \tilde{\mathbf{x}}_{I_{k+1:k+M}}^o \quad (17)$$

for updating the SLAM landmarks: $I_{k+1} \boldsymbol{\ell}^{\star\oplus} = I_{k+1} \boldsymbol{\ell}^* + I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^o$.

It is important to note that processing immature observations (i.e., \mathbb{Z}_I and \mathbb{Z}_{LI}) is *optional* for the IKS, allowing us to adjust its computational cost, based on the availability of computational resources.

3.3.3 Landmark Propagation

Before proceeding with computing the next epoch's prior pdf $\mathcal{N}(\hat{\mathbf{x}}_{k+2}^\oplus, \mathbf{P}_{k+2}^\oplus)$ and linearized constraints, $\{\mathbf{r}_L^{\star\oplus}, \mathbf{H}_L^{\star\oplus}\}$ and $\{\mathbf{r}_M^{\star\oplus}, \mathbf{H}_M^{\star\oplus}, \mathbf{F}_M^{\star\oplus}\}$, we express all landmarks $I_{k+1} \boldsymbol{\ell}_{k+M}$ w.r.t. the new "tail" pose of the sliding window, $\mathbf{x}_{I_{k+2}}$, in (9). As we describe in detail in [17], the landmark parameters in the two frames I_{k+1} and I_{k+2} , as well as the poses $\mathbf{x}_{I_{k+1:k+2}}$, are related through a non-linear *deterministic* constraint, $\mathbf{g}(\mathbf{x}_{I_{k+1:k+2}}, I_{k+1} \boldsymbol{\ell}_{k+M}, I_{k+2} \boldsymbol{\ell}_{k+M}) = \mathbf{0}$, which after linearization, becomes:

$$\begin{aligned} & \mathbf{G}_{I_{k+1:k+2}}^* \tilde{\mathbf{x}}_{I_{k+1:k+2}}^* + \mathbf{G}_{I_{k+1}, \ell}^* I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M} + \mathbf{G}_{I_{k+2}, \ell}^* I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M} = \mathbf{0} \\ \Leftrightarrow I_{k+1} \tilde{\boldsymbol{\ell}}_{k+M}^* &= -\mathbf{G}_{I_{k+1}, \ell}^{\star-1} (\mathbf{G}_{I_{k+1:k+2}}^* \tilde{\mathbf{x}}_{I_{k+1:k+2}}^* + \mathbf{G}_{I_{k+2}, \ell}^* I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*) \end{aligned} \quad (18)$$

Substituting in (9), transforms c_{k+M} to a function of $\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*$ and $I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*$.

⁴ For reducing the computational cost (linear in the number of MSC-KF features within the sliding window), the residual \mathbf{r}_S^{K*} and Jacobian matrix \mathbf{H}_S^{K*} are compressed using QR factorization [2].

⁵ Not surprisingly, the computational cost of this step is cubic in the number of SLAM landmarks, which are observed at the present epoch, as in the corresponding EKF and INVF state-update steps.

3.3.4 Covariance Update

Our objective is to compute the posterior $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+2}}^{\oplus}, \mathbf{P}_{I_{k+2}}^{\oplus})$, which will be used as the prior pdf during the next epoch. To do so, we will operate on those terms of the cost function in (9) that contain the state $\mathbf{x}_{I_{k+1}}$, which is about to be marginalized; that is the cost function:

$$c_{k+M}^M = c_{P_{I_{k+1}}^{\ominus}} + c_{u_{k+1:k+2}} + c_L + c_{z_M}^K + c_{z_{LM}}^K + c_{z_{LN}}^K + c_M + c_{z_{LM}}^M + c_{z_{LN}}^M \quad (19)$$

In particular, we follow a 4-step process:

Prior propagation and measurement compression: Following the same process, as in Sections 3.3.1 and 3.3.2, we use the prior $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1}}^{\ominus}, \mathbf{P}_{I_{k+1}}^{\ominus})$ and the inertial measurements $\mathbf{u}_{k+1:k+2}$ to compute the prior $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+2}}^{\ominus}, \mathbf{P}_{I_{k+1:k+2}}^{\ominus})$, and merge the linearized constraints into two terms, c_C^K and c_C^M , comprising multi-state and mapping constraints, respectively. Thus, (19) becomes:

$$c_{k+M}^M = c_{P_{I_{k+1:k+2}}^{\ominus}} + c_C^K + c_C^M \quad (20)$$

Partitioning of the linearized constraints: Following the same process as in (6), we partition c_C^K into $c_{C_1}^K$ and $c_{C_2}^K$, where the first term depends only on $\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*$, i.e., $c_C^K(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*) = c_{C_1}^K(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*) + c_{C_2}^K(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*)$, which after substitution in (20), yields:

$$c_{k+M}^M = c_{P_{I_{k+1:k+2}}^{\ominus}} + c_{C_1}^K(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*) + c_{C_2}^K(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*) + c_C^M. \quad (21)$$

Covariance Update: At this point, we combine the first two terms in (21), thus updating the *prior pdf* $\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+2}}^{\ominus}, \mathbf{P}_{I_{k+1:k+2}}^{\ominus})$:

$$c_{k+M}^M = c_{P_{I_{k+1:k+2}}^{\ominus}} + c_{C_2}^K + c_C^M. \quad (22)$$

$\mathcal{N}(\hat{\mathbf{x}}_{I_{k+1:k+2}}^{\oplus}, \mathbf{P}_{I_{k+1:k+2}}^{\oplus})$ has mean $\hat{\mathbf{x}}_{I_{k+1:k+2}}^{\oplus} = \hat{\mathbf{x}}_{I_{k+1:k+2}}^{\ominus} + \mathbf{P}_{I_{k+1:k+2}}^{\ominus} \mathbf{H}_{C_1}^{K*T} \mathbf{d}_C$, where \mathbf{d}_C is the solution to the linear system $\mathbf{S}_C \mathbf{d}_C = \mathbf{r}_{C_1}^{K*} - \mathbf{H}_{C_1}^{K*T} \delta \mathbf{x}_{k+1:k+2}^{\ominus}$, with $\mathbf{S}_C = \mathbf{H}_{C_1}^{K*} \mathbf{P}_{I_{k+1:k+2}}^{\ominus} \mathbf{H}_{C_1}^{K*T} + \sigma^2 \mathbf{I}_{N_1}$, and covariance $\mathbf{P}_{I_{k+1:k+2}}^{\oplus} = \mathbf{P}_{I_{k+1:k+2}}^{\ominus} \left(\mathbf{I} - \mathbf{H}_{C_1}^{K*T} \mathbf{S}_C^{-1} \mathbf{H}_{C_1}^{K*} \mathbf{P}_{I_{k+1:k+2}}^{\ominus} \right)$.

3.3.5 Construction of Next Epoch's Prior

During this last step of the IKS, we will bring c_{k+M}^M into a form whose minimization is independent of $\mathbf{x}_{I_{k+1}}$. To achieve this, we follow a 2-step process.

Partitioning of $c_{P_{I_{k+1:k+2}}^{\oplus}}$: By employing the Schur complement, the prior term $c_{P_{I_{k+1:k+2}}^{\oplus}}$ in (22) is partitioned into a prior over $\mathbf{x}_{I_{k+2}}$, $c_{P_{I_{k+2}}^{\oplus}}$, and a conditional term, $c_{I_{k+1|k+2}}$, representing linearized constraints between $\mathbf{x}_{I_{k+1}}$ and $\mathbf{x}_{I_{k+2}}$, i.e.,

$$\begin{aligned}
c_{P_{k+1:k+2}^\oplus}(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*) &= \|\tilde{\mathbf{x}}_{I_{k+2}}^* - \delta \mathbf{x}_{I_{k+2}}^\oplus\|_{\mathbf{P}_{I_{k+2}}^\oplus} + \|\delta \mathbf{x}_{I_{k+1:k+2}}^\oplus - [\mathbf{I} - \mathbf{P}_{I_{k+1:k+2}}^\oplus \mathbf{P}_{I_{k+2}}^{\oplus-1}] \begin{bmatrix} \tilde{\mathbf{x}}_{I_{k+1}}^* \\ \tilde{\mathbf{x}}_{I_{k+2}}^* \end{bmatrix}\|_{\mathbf{P}_{I_{k+1:k+2}}^\oplus}^2 \\
&= c_{P_{I_{k+2}}^\oplus}(\tilde{\mathbf{x}}_{I_{k+2}}^*) + c_{I_{k+1:k+2}}(\tilde{\mathbf{x}}_{I_{k+1:k+2}}^*)
\end{aligned} \tag{23}$$

where $\mathbf{P}_{I_{k+1:k+2}}^\oplus = \mathbf{P}_{I_{k+1}}^\oplus - \mathbf{P}_{I_{k+1,k+2}}^\oplus \mathbf{P}_{I_{k+2}}^{\oplus-1} \mathbf{P}_{I_{k+2,k+1}}^\oplus$. Substituting in (22), yields:

$$c_{k+M}^M = c_{P_{I_{k+2}}^\oplus} + c_{I_{k+1:k+2}} + c_{C_2}^K + c_C^M. \tag{24}$$

Marginalization of $\mathbf{x}_{I_{k+1}}$: Firstly, we combine all terms involving $\mathbf{x}_{I_{k+1}}$, i.e., $c_{I_{k+1:k+2}}$, $c_{C_2}^K$, and c_C^M into a single quadratic cost, corresponding to $15 + N_{C_2}^K + 3N_{k+M}$ linearized constraints:

$$c_J(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*) = \|\mathbf{b} - \mathbf{J}_1 \tilde{\mathbf{x}}_{I_{k+1}}^* - \mathbf{J}_2 \begin{bmatrix} \tilde{\mathbf{x}}_{I_{k+2:k+M}}^* \\ I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^* \end{bmatrix}\|_{\mathbf{I}_{15+N_{C_2}^K+3N_{k+M}}}^2 \tag{25}$$

Following the same process as in (6), we partition c_J into $c_{I_{k+1:k+2:k+M}}$, that contains all information regarding $\mathbf{x}_{I_{k+1}}$, and c_{L^\oplus} and c_{M^\oplus} , which are independent of $\mathbf{x}_{I_{k+1}}$:

$$c_J(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*) = c_{I_{k+1:k+2:k+M}} + c_{L^\oplus} + c_{M^\oplus}. \tag{26}$$

where the detailed analytical expressions for c_{L^\oplus} and c_{M^\oplus} are given in [17].

Substituting (26) in (24), yields:

$$\begin{aligned}
c_{k+M}^M(\tilde{\mathbf{x}}_{I_{k+1:k+M}}^*, I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*) &= c_{P_{I_{k+2}}^\oplus}(\tilde{\mathbf{x}}_{I_{k+2}}^*) + c_{L^\oplus}(\tilde{\mathbf{x}}_{I_{k+2:k+M}}^*) + c_{M^\oplus}(\tilde{\mathbf{x}}_{I_{k+2:k+M}}^*, I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*) \\
&\quad + c_{I_{k+1:k+2:k+M}}(\tilde{\mathbf{x}}_{I_{k+2:k+M}}^*, I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*).
\end{aligned} \tag{27}$$

The last term, $c_{I_{k+1:k+2:k+M}}$ in (27), is irrelevant for the minimization of c_{k+M}^M over $\tilde{\mathbf{x}}_{I_{k+2:k+M}}^*$ and $I_{k+2} \tilde{\boldsymbol{\ell}}_{k+M}^*$ since, for any of their values, there exists a $\tilde{\mathbf{x}}_{I_{k+1}}^o$ that minimizes $c_{I_{k+1:k+2:k+M}}$ to *exactly* zero. Hence, all prior information from the current to the next IKS recursion, is represented completely through the terms $c_{P_{k+2}^\oplus}$, c_{L^\oplus} , and c_{M^\oplus} all of which do *not* involve $\tilde{\mathbf{x}}_{I_{k+1}}^*$.

4 Simulations

Our simulations involved a MEMS-quality commercial grade IMU, similar to those present on current mobile devices, running at 100 Hz, and a wide (175° degrees) field of view camera with resolution 640 × 480. Visual observations were contaminated by zero-mean white Gaussian noise with $\sigma = 1.0$ pixel. We compared the Root Mean Square Error (RMSE) for the real-time estimates⁶ of the MSC-KF VINS (denoted as EKF), with that of the proposed iterative Kalman smoother (denoted as IKS). Both estimators maintained a sliding window of 10 camera poses. Feature tracks that spanned beyond the sliding window were initialized as SLAM landmarks

⁶ By real-time, we refer to the estimate for \mathbf{x}_{I_k} , right before, processing image k .

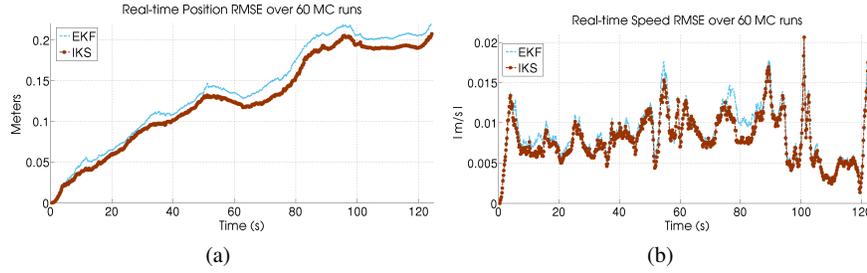


Fig. 2: Monte Carlo Simulation I: Comparison of the proposed IKS versus the EKF

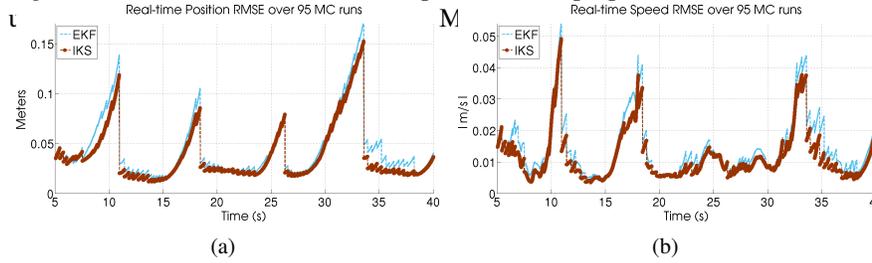


Fig. 3: Monte Carlo Simulation II: Comparison of the proposed IKS versus the EKF during camera occlusions: (a) Position RMSE; (b) Speed RMSE

(and marginalized when lost from the camera’s field of view), while shorter feature tracks were processed as MSC-KF features.

- **VINS under nominal conditions:** In Simulation I (see Fig. 2), the platform’s trajectory and dynamics resembled those of a person traversing 120 m of an indoor environment, while new camera poses were generated every 25 cm, and the rate that new features enter the camera’s field of view followed that of real-world experimental trials. As seen in Fig. 2(a), the performance difference between the EKF-based VINS and the proposed IKS is rather small, since in the presence of many visual measurements, both estimators are able to accurately track the system’s state. Note however, that even under these nominal conditions the IKS always maintained a more accurate estimate of the platform’s speed (see Fig. 2(b)), while for certain parts of the trajectory its estimate improved by $\sim 20\%$, over the EKF, due to the inability of the latter to process feature tracks immediately as they become available.

- **Camera Occlusions:** In Simulation II (see Fig. 3), we simulated the motion of a handheld device, “hovering” over the same scene, for 40 s, emulating a common scenario for augmented-reality applications. We then introduced 3 periods, of approximately 5 s each, during which the camera was occluded and no feature tracks were available. As evident from Figs. 3(a) and 3(b), by re-processing visual and inertial measurements, the IKS, converges faster to the correct position and velocity estimates, right after re-gaining access to camera measurements.

5 Experiments

We further validated the proposed IKS on real-world data, using a Project Tango developer tablet, and as ground truth the result of a batch least-squares (BLS) over the entire set of visual and inertial measurements. As in Section 4, we compared the real-time estimates of the proposed IKS, to those of the MSC-KF VINS, both of which processed measurements to SLAM landmarks, as well as MSC-KF feature tracks, and maintained a sliding window of length 14.

Initialization: In Experiment I (Fig. 4), we compared the initialization phase, of both estimators, when they start from inaccurate estimates of their initial velocity and IMU biases, as often happens in practice. As it is seen in Fig. 4(b), the IKS converged faster to its correct velocity estimates, which lead to a reduced position error, as compared to the EKF, for the rest of their trajectories [Fig. 4(a)]. Note that a BLS over a small set of initial poses, could have been used for system initialization, making both estimators equally accurate. Such a choice, however, would inevitably introduce a significant time delay, a key drawback for real-time applications, while the complexity of the corresponding implementation would significantly increase. For the IKS, however, an iterative optimization during its initial phase, seamlessly takes place, without the need to transition from BLS to filtering.

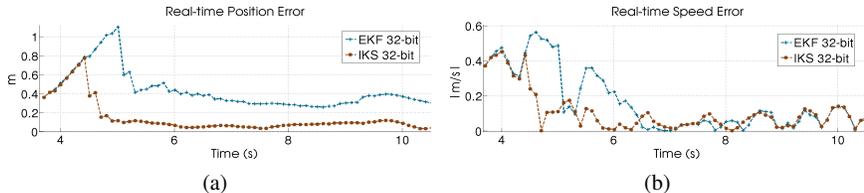


Fig. 4: Experiment I: Comparison of the proposed IKS versus the EKF during filter initialization: (a) Real-time position error; (b) Real-time speed error

Fast camera turns: In Experiment II (Fig. 5), we collected a “stress” dataset, during which the camera performed quick turns, inside an office area, causing abrupt reductions in the number, quality, and length of feature tracks for short periods of time [Fig. 5(a)]. As it is evident from Fig. 5(d), the inability of the EKF to re-process visual observations, caused sudden accumulations of yaw error, e.g., at 40 s. Furthermore, on par with our simulation results, the IKS maintained an improved real-time estimate, of the platform’s velocity, throughout the experiment, while at certain points, its velocity estimate was even 2 times better than the EKF’s [Fig. 5(c)].

Timing analysis: We used the Samsung S4 cell phone, as a testbed for comparing the processing time of the proposed IKS, *with and without processing of immature visual observations*, denoted by IKS w/, and IKS w/o, respectively, as well as, a reference EKF implementation. Albeit the 32-bit arithmetic precision, of the NEON co-processor, present on the 1.6 GHz Cortex-A15 ARM CPU of the Samsung S4, no numerical inaccuracies were introduced, when compared to 64-bit arithmetic precision. All estimators maintained a sliding window of 14 poses and on average 5 SLAM landmarks in their state vector. As it is evident from Table 1, the proposed IKS achieves real-time performance, even under re-linearization, while it

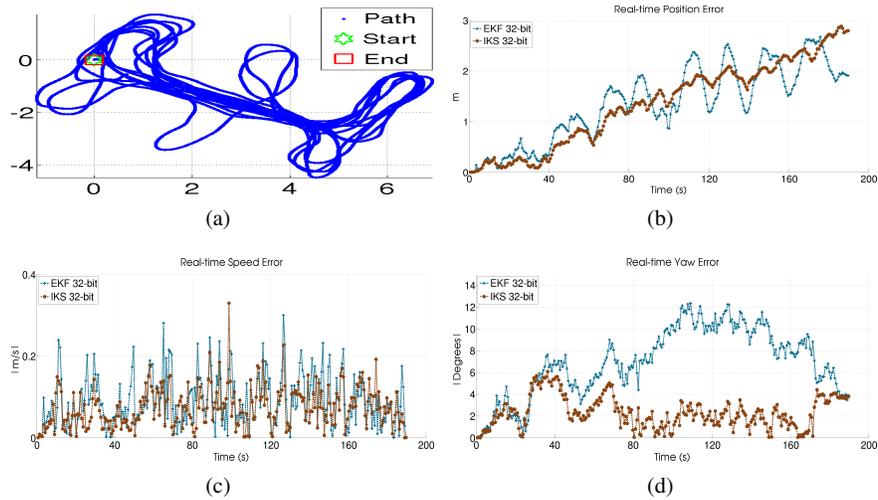


Fig. 5: Experiment II: Comparison of the proposed IKS versus the EKF during sudden turns: (a) 3D Trajectory; (b) Position error; (c) Speed error; (d) Yaw error

is able to bring its cost down to levels comparable to the EKF by temporary disabling the re-processing of visual observations.

| Step \ Algorithm | IKS w/ | IKS w/o | EKF |
|--------------------------|------------|-----------|-----------|
| Propagation | 14 | 14 | 1 |
| Jacobians Calculation | 101 | 12 | 9 |
| Measurement Compressions | 14 | 6 | 2 |
| State Update | 12 | 12 | 12 |
| Covariance Update | 0.5 | 0.5 | 18 |
| Prior Constraints Update | 1.5 | 1.5 | N/A |
| Total | 166 | 46 | 42 |

Table 1: Timing analysis on the Samsung S4: Numbers denote average time in ms

6 Conclusion

In this paper, we have presented an iterative Kalman smoother (IKS) for vision-aided inertial navigation that incorporates advantages of competing approaches. Through smoothing, the proposed IKS iteratively re-linearizes both inertial and visual measurements over a single, or multiple overlapping, sliding windows, thus improving robustness. At the same time, the IKS inherits the excellent numerical properties of the Kalman filter, making it amenable to very efficient implementations (4-fold speed up on ARM NEON co-processor) using single-precision (32 bit) arithmetic. As part of our validation process, we demonstrated the resilience and efficiency of the proposed approach, under adverse navigation conditions.

References

1. J. A. Hesch, D. G. Kottas, S. L. Bowman, and S. I. Roumeliotis, "Consistency analysis and improvement of vision-aided inertial navigation," *IEEE Trans. on Robotics*, vol. 30, pp. 158–176, Feb. 2014.
2. A. I. Mourikis, N. Trawny, S. I. Roumeliotis, A. E. Johson, A. Ansar, and L. Matthies, "Vision-aided navigation for spacecraft entry, descent, and landing," *IEEE Trans. on Robotics*, vol. 25, pp. 264–280, Apr. 2009.
3. E. D. Nerurkar, K. J. Wu, and S. I. Roumeliotis, "C-KLAM: Constrained Keyframe Localization and Mapping for long-term navigation," in *Proc. of the IEEE International Conference on Robotics and Automation*, (Hong Kong, China), pp. 3638–3643, May 31 – June 6 2013.
4. A. Jazwinski, *Stochastic processes and filtering theory*. No. 64 in Mathematics in science and engineering, New York, Academic Press, 1970.
5. J. A. Hesch, D. G. Kottas, S. L. Bowman, and S. I. Roumeliotis, "Towards consistent vision-aided inertial navigation," in *Proc. of the 10th International Workshop on the Algorithmic Foundations of Robotics (WAFR'12)*, (Cambridge, Massachusetts), pp. 559–574, June 13–15 2012.
6. C. Guo, D. Kottas, R. DuToit, A. Ahmed, R. Li, and S. Roumeliotis, "Efficient visual-inertial navigation using a rolling-shutter camera with inaccurate timestamps," in *Proceedings of Robotics: Science and Systems*, (Berkeley, USA), July 2014.
7. G. Sibley, L. Matthies, and G. Sukhatme, "Sliding window filter with application to planetary landing," *Journal of Field Robotics*, vol. 27, no. 5, pp. 587–608, 2010.
8. G. P. Huang, A. I. Mourikis, and S. I. Roumeliotis, "An Observability-Constrained Sliding Window Filter for SLAM," in *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, (San Francisco, CA), pp. 65–72, Sept. 25 – 30 2011.
9. H.-P. Chiu, S. Williams, F. Dellaert, S. Samarasekera, and R. Kumar, "Robust vision-aided navigation using Sliding-Window Factor Graphs," in *Proc. of the IEEE International Conference on Robotics and Automation*, (Karlsruhe, Germany), pp. 46–53, May 6 – 10 2013.
10. S. Leutenegger, P. T. Furgale, V. Rabaud, M. Chli, K. Konolige, and R. Siegwart, "Keyframe-Based Visual-Inertial SLAM using Nonlinear Optimization," in *Proceedings of Robotics: Science and Systems*, (Berlin, Germany), June 2013.
11. K. J. Wu, A. Medhat, G. Georgiou, and S. I. Roumeliotis, "A square root inverse filter for efficient vision-aided inertial navigation on mobile devices," in *Robotics: Science and Systems*, (Rome, Italy), July 13 - 17 2015.
12. D. G. Kottas and S. I. Roumeliotis, "An iterative Kalman smoother for robust 3D localization on mobile and wearable devices," in *Proc. of the IEEE International Conference on Robotics and Automation*, (Seattle, WA), pp. 6336–6343, May 26 – 30 2015.
13. F. M. Mirzaei and S. I. Roumeliotis, "A Kalman Filter-based algorithm for IMU-camera calibration: Observability analysis and performance evaluation," *IEEE Trans. on Robotics*, vol. 24, pp. 1143–1156, Oct. 2008.
14. N. Trawny and S. I. Roumeliotis, "Indirect Kalman Filter for 3D Attitude Estimation," tech. rep., University of Minnesota, Dept. of Comp. Sci. & Eng., March 2005.
15. C. Harris and M. Stephens, "A combined corner and edge detector," in *Proc. of the Alvey Vision Conference*, (Manchester, UK), pp. 147–151, Aug. 31 – Sept. 2 1988.
16. B. D. 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, British Columbia), pp. 674–679, Aug. 24 – 28 1981.
17. D. G. Kottas and S. I. Roumeliotis, "An iterative Kalman smoother for robust 3D localization and mapping, http://www-users.cs.umn.edu/~dkottas/pdfs/iks_slam_2014_tr.pdf," tech. rep., University of Minnesota, Dept. of Comp. Sci. & Eng., MARS Lab, Apr. 2015.