# EDI: ESKF-based Disjoint Initialization for Visual-Inertial SLAM Systems

Weihan Wang[a], Jiani Li[b], Yuhang Ming[c], Philippos Mordohai[a]

*Abstract*— Visual-inertial initialization can be classified into joint and disjoint approaches. Joint approaches tackle both the visual and the inertial parameters together by aligning observations from feature-bearing points based on IMU integration then use a closed-form solution with visual and acceleration observations to find initial velocity and gravity. In contrast, disjoint approaches independently solve the Structure from Motion (SFM) problem and determine inertial parameters from up-to-scale camera poses obtained from pure monocular SLAM. However, previous disjoint methods have limitations, like assuming negligible acceleration bias impact or accurate rotation estimation by pure monocular SLAM. To address these issues, we propose EDI, a novel approach for fast, accurate, and robust visual-inertial initialization. Our method incorporates an Error-state Kalman Filter (ESKF) to estimate gyroscope bias and correct rotation estimates from monocular SLAM, overcoming dependence on pure monocular SLAM for rotation estimation. To estimate the scale factor without prior information, we offer a closed-form solution for initial velocity, scale, gravity, and acceleration bias estimation. To address gravity and acceleration bias coupling, we introduce weights in the linear least-squares equations, ensuring acceleration bias observability and handling outliers. Extensive evaluation on the EuRoC dataset shows that our method achieves an average scale error of 5.8% in less than 3 seconds, outperforming other state-of-the-art disjoint visual-inertial initialization approaches, even in challenging environments and with artificial noise corruption.

## I. INTRODUCTION

The combination of a single camera and an Inertial Measurement Unit (IMU) in Visual-Inertial Navigation Systems (VINS) is a cost-efficient and low-power solution for robot perception and AR/VR applications. The camera provides a rich representation of the environment while the IMU measures acceleration and angular velocity, making it robust to fast-motion and texture-less images. This combination makes them ideal for complementing each other. The initialization process is crucial for VINS, as it requires a good initial estimation for the scale, gravity, initial velocity, acceleration, and gyroscope biases, but remains a challenge as it requires fast and accurate recovery of observable parameters from visual and inertial measurements without prior knowledge. Poor initialization can lead to trajectory drift and hinder the convergence of subsequent optimization. A prolonged



Fig. 1: A diagram of the components of EDI. The green blocks represent steps 1 to 3 of EDI, and the different colored arrows represent different output flows from their corresponding steps.

initialization process is also impractical for both robotics and AR/VR applications.

Visual-inertial initialization is classified into two categories [11], [12]: joint [1]-[6] and disjoint approaches [7], [8], [9], [13]. Joint approaches tackle both the visual and the inertial parameters together by aligning observations from feature-bearing points based on IMU integration. Joint visual-inertial initialization methods begin by finding a closed-form solution to the visual-inertial problem. A closed-form solution to calculate the initial velocity and position of feature points using only the visual information from three consecutive image frames and a single feature point was first introduced by Kneip et al. [1]. Martinelli [2] later proposed a closed-form solution that also takes into account the scale, gravity, acceleration bias and feature points' depth, and analyzes the necessary conditions for the solution to be attainable. Kaiser et al. [3] apply Martinelli's solution and utilize non-linear optimization to account for the impact of gyroscope bias on the system. These methods track multiple points in all images and use a system of equations to minimize the 3D error of feature points in space. Dong-Si and Mourikis [4], [5] proposed a closed-form solution for estimating the attitude, velocity, feature positions, and camera-IMU extrinsic calibration. In their work, they also specifically discuss two methods for recovering the relative rotation between the camera and IMU under different scenarios with varying numbers of tracked features. Campos et al. [6] further improve joint methods by leveraging preintegration to reduce the computational cost of the closed-form initialization and conducting two rounds of visual-inertial Bundle Adjustment (VI-BA) to increase the precision of depth feature estimates, gyroscope bias,

[a]Stevens Institute of Technology, Hoboken, NJ, USA, 07030, {wwang103,pmordoha}@stevens.edu

[b]Meta, jiani.li@meta.com

[c]School of Computer Science, Hangzhou Dianzi University, Hangzhou 310018, China, yuhang.ming@ieee.org

TABLE I: Assumptions underlying joint and disjoint initialization methods. The fewer assumptions the methods make, the more versatile and practical they become.

| | Methods | Noiseless IMU and camera measurement | All features tracked in all frames | Acceleration bias is negligible | Gyroscope bias is negligible | Known Camera-IMU extrinsic calibration | Monocular SLAM accurately estimates camera pose |
|---|---|---|---|---|---|---|---|
| Joint | Kneip et al [1] | ✓ | ✓ | ✓ | ✓ | ✓ | × |
| | Martinelli [2] | ✓ | ✓ | × | ✓ | ✓ | × |
| | Kaiser et al [3] | ✓ | ✓ | ✓ | × | ✓ | × |
| | Dong-Si and Mourikis [4], [5] | ✓ | ✓ | × | × | × | × |
| | Campos et al. [6] | ✓ | × | × | × | ✓ | × |
| Disjoint | ORB-SLAM-VI [7] | × | × | × | × | ✓ | ✓ |
| | Huang and Liu [8] | × | × | × | × | × | ✓ |
| | VINS-Mono [9] | × | × | ✓ | × | ✓ | ✓ |
| | ORB-SLAM3 [10] | × | × | × | × | ✓ | ✓ |
| | EDI (Ours) | × | × | × | × | ✓ | × |

gravity, and initial velocity. However, these joint methods are based on the assumption that there is no noise in the IMU and camera measurements, and that all feature points are tracked correctly in all frames. Even though the work by Campos et al. [6] relaxes the requirement for tracking feature points in all frames, it still faces challenges with low recall rate, leading to extended initialization times.

In contrast, disjoint approaches aim to solve the SfM problem independently first, and then determine the inertial parameters based on up-to-scale camera poses obtained from a pure monocular simultaneous localization and mapping (SLAM) system. This approach is made possible by the use of monocular SLAM, which performs local bundle adjustment and takes into consideration both photometric and geometric consistency in feature point tracking, providing a more precise state estimation. One method in this category, used in ORB-SLAM-VI [7], is proposed by Mur-Artal and Tardós. It runs monocular SLAM for a few seconds, assuming the sensor undergoes a motion that makes all variables observable and divides the initialization process into four sub-problems.

Later, Huang and Liu [8] adopted this concept and expanded on it by incorporating an estimate of the camera-IMU extrinsic parameters. Qin et al. in VINS-Mono [9] align the camera's trajectory and orientation with the IMU preintegration measurements by solving linear least-squares equations to obtain values for the scaled velocity, scale, gravity, and gyroscope biases, resulting in high-performance output compared to other state-of-the-art visual-inertial odometry systems such as OKVIS[14], SVO [15], and ROVIO [16], as evaluated on various datasets [17]. Both methods determine inertial parameters by solving a set of linear equations using least-squares, but they differ in the steps involved. ORB-SLAM-VI does not take velocity into account, while VINS-Mono does not estimate the acceleration bias when solving a set of linear equations. However, both methods have limitations. For example, ORB-SLAM-VI requires 15 seconds of initialization to make the acceleration bias observable, and both methods assign equal weights to the residues without considering IMU measurement uncertainty when solving linear least-squares. To address these limitations in the above disjoint methods, Campos et al. [13] proposed a new disjoint initialization method in ORB-SLAM3 [10] by

formulating the visual-inertial initialization as a maximum a-posteriori (MAP) problem and demonstrate that their method outperforms the best joint and previous disjoint initialization methods. However, this method is sensitive to the scale factor. The second step in this approach, which involves the iterative inertial-only optimization to estimate the scale factor, requires a reliable initial estimate of scale. Furthermore, the method requires an empirical prior residual for the biases. Although the third step of the ORB-SLAM3 initialization phase, which involves joint visual-inertial bundle adjustment (VI-BA), aims to improve the previous estimate from the second step, the quality of the second step can greatly impact the convergence time of the third step and the ability of the VI-BA to produce an optimal solution. Additionally, the success of the disjoint method is highly dependent on the performance of the pure monocular SLAM system.

To overcome the limitations of the previous disjoint methods, we propose EDI, an innovative disjoint approach for initializing a visual-inertial system. The main contributions of the proposed initialization method are:

- Eliminating the need for computationally intensive visual-inertial bundle adjustment (VI-BA) while ensuring accuracy, improving the efficiency of the method, and increasing robustness in challenging conditions.
- Proposing a new method that incorporates an Error-state Kalman Filter (ESKF) [18] to estimate gyroscope bias and correct rotation estimates with monocular SLAM, considering the probabilistic model of IMU noise.
- Providing a closed-form solution for estimating initial velocity, scale, gravity, and acceleration bias, along with weights to handle outliers.

The summary of assumptions underlying joint and disjoint initialization methods, including ours, are shown in Table I.

## II. PRELIMINARIES

### A. Notation

In this paper, the world frame is represented by $(\cdot)^w$, the body frame is represented by $(\cdot)^b$, and the camera frame is represented by $(\cdot)^c$. $(\hat{\cdot})$ represents a posterior estimate after being corrected by ESKF, while an up-to-scale estimate is denoted as $(\bar{\cdot})$. We use $\mathbf{R}$ for rotation matrices and $\mathbf{g}$ for the gravity vector. The rotation, translation and velocity from the body frame to the world frame are represented by $\mathbf{R}_b^w$,

$\mathbf{p}_b^w$ and $\mathbf{v}_b^w$ respectively. $b_k$ is the body frame while taking the $k$-th image, and $c_k$ is the camera frame while taking the $k$-th image. Acceleration bias and gyroscope bias in the local body frame are represented by $\mathbf{b}_a$ and $\mathbf{b}_g$ respectively. The nominal state vectors are represented by $(\cdot)$, true state vectors are represented by $(\cdot)_t$, the error-state is represented by $\delta(\cdot)$ and $\mathcal{X}$ is a state variable.

### B. The Error-State Kalman Filter

The Error-state Kalman filter, or ESKF [18], is a type of filter used to estimate the true state of a system while taking into account measurement noise and model uncertainty. It is widely used in control systems, navigation and signal processing. It offers advantages such as minimal representation of state variables in rotation processing, and operating near the origin to avoid linearization approximation issues and gimbal lock problems. The state variables in ESKF are minimal, allowing for the omission of second-order variables, and the Jacobian matrices are also straightforward and can even be substituted with identity matrices.

In ESKF, there are three state variables: true, nominal, and error state. The nominal state integrates with noise and other potential model flaws, leading to the accumulation of errors and drift. The error state accounts for various noise sources and biases. The relationship among the true state ($\mathcal{X}_t$), the nominal state ($\mathcal{X}$) and the error state ($\delta\mathcal{X}$) is defined as

$$\mathcal{X}_t = \mathcal{X} \oplus \delta\mathcal{X}$$

where $\oplus$ indicates a generic composition.

The procedure of the ESKF outlined in this paper is as follows: upon receipt of an IMU measurement, it is integrated and incorporated into the nominal state variables. The error state takes into account the noise term and biases in the ESKF, providing a Gaussian distribution for the error state. The ESKF also incorporates a prediction and correction process, utilizing observations from sensors other than the IMU. Following correction, the ESKF yields a posterior error Gaussian distribution, and the error is incorporated into the nominal state variables, resetting the ESKF. This process is repeated iteratively.

## III. PROPOSED APPROACH

This section describes the proposed online initialization method, EDI, aiming to estimate precise initial values for the body velocities, gravity direction, scale factor, gyroscope and acceleration bias.

We begin by discussing the techniques used in two state-of-the-art disjoint initialization methods and comparing them to our own. ORB-SLAM3 [10] uses a three-step process: vision-only MAP estimation, inertial-only MAP estimation, and joint visual-inertial optimization for further refining the solution. While this method can provide accurate results, it is computationally demanding due to the extensive non-linear optimization procedure. On the other hand, VINS-Mono [9], a linear loosely-coupled initialization method, is quick to compute but relies on the assumption of zero acceleration bias.

Our method offers the best of both worlds - it is as efficient as VINS-Mono and produces an estimation comparable to ORB-SLAM3 without the need for multiple non-linear optimization steps or strong assumptions. Furthermore, it is more robust in challenging conditions. Our method, as shown in Fig. 1, is composed of four steps:

- **Step 0. Pure Monocular SLAM**: Obtain the initial keyframe poses with an unknown scale.
- **Step 1. ESKF-based Gyroscope Bias Estimation**: Utilize the ESKF to estimate the gyroscope bias and combine rotation estimates from IMU prediction and pure monocular SLAM.
- **Step 2. Linear Solver**: Align the IMU trajectory and pure monocular SLAM trajectory by determining the scale factor, keyframes' velocities, and gravity.
- **Step 3. Refinement**: Use the solution from the previous step as the initial estimate to obtain acceleration bias, refined scale, keyframes' velocities and gravity estimate.

As a disjoint method, EDI uses a pure monocular SLAM [19] with an increased keyframe insertion rate for a short period of time (1 or 2 seconds) to obtain the initial keyframe poses with an unknown scale, in order to ensure observability of all inertial variables.

### A. ESKF-based Gyroscope Bias Estimation

To estimate the gyroscope bias and combine rotation estimates from the IMU prediction and monocular SLAM during the initialization stage with a window of $N$ keyframes, we only consider the rotation and gyroscope bias in the nominal state and error state in this step.

Consider two consecutive keyframes $b_k$ and $b_{k+1}$ with interval $\Delta t$ and denote the nominal state at keyframe $b_{k+1}$ as $\mathcal{X}_{b_{k+1}} = \left[ \mathbf{R}_{b_{k+1}}^w, \mathbf{b}_{g,b_{k+1}} \right]$ and the error state as $\delta\mathcal{X}_{b_{k+1}} = \left[ \delta\boldsymbol{\theta}_{b_{k+1}}, \delta\mathbf{b}_{g,b_{k+1}} \right]$. The error state of rotation and gyroscope bias is integrated, allowing the IMU measurements to make predictions for the ESKF as follows:

$$\delta\boldsymbol{\theta}_{b_{k+1}} = \text{Exp}(-(\boldsymbol{\omega}_m - \mathbf{b}_g)\Delta t)\delta\boldsymbol{\theta}_{b_k} - \delta\mathbf{b}_{g,b_k}\Delta t + \boldsymbol{\eta}_\theta, \quad (1)$$

$$\delta\mathbf{b}_{g,b_{k+1}} = \delta\mathbf{b}_{g,b_k} + \boldsymbol{\eta}_g, \quad (2)$$

where $\boldsymbol{\omega}_m$ is the raw gyroscope measurement, $\delta\boldsymbol{\theta}_{b_k}$ and $\delta\mathbf{b}_{g,b_k}$ are error state of keyframe $b_k$, $\boldsymbol{\eta}_\theta$ and $\boldsymbol{\eta}_g$ are white Gaussian noise applied to rotation and gyroscope bias estimation respectively, $\boldsymbol{\eta}_\theta \sim \mathcal{N}(0, \sigma_{w_n}^2 \Delta t^2 \mathbf{I})$, $\boldsymbol{\eta}_g \sim \mathcal{N}(0, \sigma_{w_w}^2 \Delta t \mathbf{I})$. Meanwhile, the covariance matrix of the prediction at $b_{k+1}$ is updated as follows:

$$\boldsymbol{P}_{\text{pred}} = \boldsymbol{F}\hat{\boldsymbol{P}}_{b_k}\boldsymbol{F}^\top + \boldsymbol{Q},$$

where $\boldsymbol{P}_{\text{pred}}$ is the predicted covariance matrix of keyframe $b_{k+1}$, $\hat{\boldsymbol{P}}_{b_k}$ is the corrected covariance matrix of keyframe $b_k$, $\boldsymbol{Q}$ is the covariance matrix of the perturbation impulses ($\boldsymbol{Q} = \text{diag}(\text{Cov}(\boldsymbol{\eta}_\theta), \text{Cov}(\boldsymbol{\eta}_g))$), and $\boldsymbol{F}$ is the Jacobian matrix with respect to the error state of keyframe $b_k$ based on Eq. (1) and Eq. (2):

$$\boldsymbol{F} = \begin{bmatrix} \text{Exp}(-(\boldsymbol{\omega}_m - \mathbf{b}_g)\Delta t) & -I\Delta t \\ 0 & I \end{bmatrix}.$$

In order to avoid drift in the IMU prediction from integrating the IMU measurement directly, we need to correct the IMU prediction with other complementary sensors. In this paper, we treat monocular SLAM as a sensor and use its rotation estimates as observations to fuse with our IMU predictions. We imagine this abstract sensor as a typical sensor that gives information based on the current state:

$$\boldsymbol{r}_{\mathrm{b}_{k+1}}^{\mathrm{w}} = \boldsymbol{h}(\mathcal{X}_t) + \boldsymbol{v},$$

where $\boldsymbol{r}_{\mathrm{b}_{k+1}}^{\mathrm{w}}$ is an orientation observation at keyframe $\mathrm{b}_{k+1}$ from pure monocular SLAM, $\boldsymbol{h}()$ is the observation function of the system, $\boldsymbol{v}$ is white Gaussian noise with co-variance $\boldsymbol{V}$, $\boldsymbol{v} \sim \mathcal{N}(0, \boldsymbol{V})$ and the orientation difference between prediction and observation is denoted as $\mathbf{e}_r$ ($\mathbf{e}_r = \mathrm{Log}(\boldsymbol{h}(\mathcal{X}_t)^\top \boldsymbol{r}_{\mathrm{b}_{k+1}}^{\mathrm{w}})$). In ESKF, our goal is to update the error state, so we need to calculate the Jacobian matrix $\boldsymbol{H}$, which is the matrix of partial derivatives of the observation with respect to the error state $\delta \mathcal{X}$:

$$\boldsymbol{H} = \frac{\partial \boldsymbol{h}}{\partial \delta \mathcal{X}_{\mathrm{b}_{k+1}}} = \frac{\partial \boldsymbol{h}}{\partial \mathcal{X}_t} \frac{\partial \mathcal{X}_t}{\partial \delta \mathcal{X}_{\mathrm{b}_{k+1}}} = \boldsymbol{J}_r^{-1}(\mathbf{e}_r),$$

where $\boldsymbol{J}_r^{-1}$ is the inverse of the right Jacobian.

Then, we calculate the Kalman gain and the update of the error state as follows:

$$\boldsymbol{K} = \boldsymbol{P}_{\mathrm{pred}} \boldsymbol{H}^\top (\boldsymbol{H} \boldsymbol{P}_{\mathrm{pred}} \boldsymbol{H}^\top + \boldsymbol{V})^{-1},$$

$$\delta \hat{\mathcal{X}}_{\mathrm{b}_{k+1}} = \boldsymbol{K}(\mathrm{Log}(\mathbf{e}_r)),$$

$$\hat{\boldsymbol{P}}_{\mathrm{b}_{k+1}} = (\boldsymbol{I} - \boldsymbol{K} \boldsymbol{H}) \boldsymbol{P}_{\mathrm{pred}}.$$

where $\boldsymbol{K}$ is the Kalman gain.

After the ESKF update, the posterior error state $\delta \hat{\mathcal{X}}_{\mathrm{b}_{k+1}}$, is incorporated into the nominal state. Afterwards, $\delta \hat{\mathcal{X}}_{\mathrm{b}_{k+1}}$ is reset to zero, and its corresponding covariance matrix is updated to reflect this reset accordingly. The best true-state estimate at keyframe $\mathrm{b}_{k+1}$ is obtained using the appropriate compositions as follows:

$$\hat{\mathbf{R}}_{\mathrm{b}_{k+1}}^{\mathrm{w}} = \mathbf{R}_{\mathrm{b}_{k+1}}^{\mathrm{w}} \mathrm{Exp}(\delta \hat{\boldsymbol{\theta}}_{\mathrm{b}_{k+1}}),$$

$$\hat{\mathbf{b}}_{\mathrm{g},\mathrm{b}_{k+1}} = \mathbf{b}_{\mathrm{g},\mathrm{b}_{k+1}} + \delta \hat{\mathbf{b}}_{\mathrm{g},\mathrm{b}_{k+1}}.$$

Initially, the gyroscope bias is assumed to be zero. By using the ESKF with $N$ keyframes, the estimated gyroscope bias in the last keyframe in the window is considered the most accurate true-state estimate. We use this as the final estimated gyroscope bias. Additionally, our method is different from other initialization methods in that it fuses the rotation estimate from the IMU prediction with the rotation estimate from monocular SLAM for each keyframe within the window. This allows for a more accurate estimation of the rotation when monocular SLAM is not accurate and robust enough.

### B. Linear Solver

This step aims to obtain an optimal estimate of the keyframes' velocities, gravity and scale factor of the pure monocular SLAM. As the pure monocular SLAM system [19] only estimates keyframe poses without recovering the scale, the estimation of acceleration bias, correction of keyframe translations, velocities and gravity using ESKF is not possible. To overcome this limitation, we solve a set of linear equations to obtain the following estimates:

$$\mathcal{X}_1 = \left[ \mathbf{v}_{\mathrm{b}_0:\mathrm{b}_{N-1}}^{\mathrm{w}}, \mathbf{g}^{\mathrm{w}}, \mathrm{s} \right]^\top.$$

Considering two consecutive keyframes $\mathrm{b}_k$ and $\mathrm{b}_{k+1}$, we have the following relationships:

$$\Delta \mathbf{p}_{k,k+1} = \mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}(\mathrm{s}(\bar{\mathbf{p}}_{\mathrm{b}_{k+1}}^{\mathrm{w}} - \bar{\mathbf{p}}_{\mathrm{b}_k}^{\mathrm{w}}) - \mathbf{v}_{\mathrm{b}_k}^{\mathrm{w}} \Delta \mathbf{t}_{k,k+1} - \frac{1}{2}\mathbf{g}^{\mathrm{w}} \Delta \mathbf{t}_{k,k+1}^2), \tag{3}$$

$$\Delta \mathbf{v}_{k,k+1} = \mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}(\mathbf{v}_{\mathrm{b}_{k+1}}^{\mathrm{w}} - \mathbf{v}_{\mathrm{b}_k}^{\mathrm{w}} - \mathbf{g}^{\mathrm{w}} \Delta \mathbf{t}_{k,k+1}), \tag{4}$$

$$\mathrm{s}\bar{\mathbf{p}}_{\mathrm{b}_k}^{\mathrm{w}} = \mathrm{s}\bar{\mathbf{p}}_{\mathrm{c}_k}^{\mathrm{w}} - \mathbf{R}_{\mathrm{b}_k}^{\mathrm{w}} \mathbf{p}_{\mathrm{c}}^{\mathrm{b}}, \tag{5}$$

where $\Delta \mathbf{p}_{k,k+1}$ and $\Delta \mathbf{v}_{k,k+1}$ are the preintagration of translation and velocity respectively from $k$-th to $k+1$-th keyframes. We combine Eq. (3)$\sim$ Eq. (5) into the following linear equation:

$$\mathcal{A}_{k,k+1}\mathcal{X}_1 = \mathcal{B}_{k,k+1}, \tag{6}$$

$\mathcal{A}_{k,k+1} =$
$$\begin{bmatrix} 0_{3\times3k} & \boldsymbol{\alpha}_{k,k+1}^a & 0_{3\times3} & 0_{3\times3(N-k-2)} & \boldsymbol{\alpha}_{k,k+1}^b & \boldsymbol{\alpha}_{k,k+1}^c \\ 0_{3\times3k} & \boldsymbol{\beta}_{k,k+1}^a & \boldsymbol{\beta}_{k,k+1}^b & 0_{3\times3(N-k-2)} & \boldsymbol{\beta}_{k,k+1}^c & 0_{3\times1} \end{bmatrix},$$

$$\mathcal{B}_{k,k+1} = \begin{bmatrix} \Delta \mathbf{p}_{k,k+1} - \mathbf{p}_{\mathrm{c}}^{\mathrm{b}} + \mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k} \mathbf{R}_{\mathrm{b}_{k+1}}^{\mathrm{w}} \mathbf{p}_{\mathrm{c}}^{\mathrm{b}} \\ \Delta \mathbf{v}_{k,k+1} \end{bmatrix},$$

$$\begin{aligned} \boldsymbol{\alpha}_{k,k+1}^a &= -\mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k} \Delta \mathbf{t}_{k,k+1}, & \boldsymbol{\alpha}_{k,k+1}^b &= -\frac{1}{2}\mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k} \Delta \mathbf{t}_{k,k+1}^2, \\ \boldsymbol{\alpha}_{k,k+1}^c &= \mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}(\bar{\mathbf{p}}_{\mathrm{c}_{k+1}}^{\mathrm{w}} - \bar{\mathbf{p}}_{\mathrm{c}_k}^{\mathrm{w}}), & \boldsymbol{\beta}_{k,k+1}^a &= -\mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}, \\ \boldsymbol{\beta}_{k,k+1}^b &= \mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}, & \boldsymbol{\beta}_{k,k+1}^c &= -\mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k} \Delta \mathbf{t}_{k,k+1}, \end{aligned}$$

where the $\mathcal{A}_{k,k+1}$ matrix has dimensions $6 \times (3N + 4)$ and $\mathcal{B}_{k,k+1}$ is a $6 \times 1$ vector. The camera's up-to-scale translations at two consecutive keyframes, $\bar{\mathbf{p}}_{\mathrm{c}_k}^{\mathrm{w}}$ and $\bar{\mathbf{p}}_{\mathrm{c}_{k+1}}^{\mathrm{w}}$, are obtained from the pure monocular SLAM, as well as the orientations of the IMU with respect to the world frame, $\mathbf{R}_{\mathrm{b}_k}^{\mathrm{w}}$ and $\mathbf{R}_{\mathrm{b}_{k+1}}^{\mathrm{w}}$. It is assumed that the extrinsic calibration matrix $[\mathbf{R}_{\mathrm{c}}^{\mathrm{b}}|\mathbf{p}_{\mathrm{c}}^{\mathrm{b}}]$ is known, which allows for the transformation of camera poses to the IMU frame of reference. $\mathbf{R}_{\mathrm{w}}^{\mathrm{b}_k}$ is the transpose of $\mathbf{R}_{\mathrm{b}_k}^{\mathrm{w}}$.

We then obtain $\mathcal{X}_1$ by considering all relationships among $N$ keyframes and solving the linear least squares problem:

$$\min_{\mathcal{X}_1} \sum_{k \in \mathcal{K}} \|\mathcal{A}_{k,k+1}\mathcal{X}_1 - \mathcal{B}_{k,k+1}\|^2$$

where $\mathcal{K}$ indexes all $N$ keyframes.

### C. Refinement

This step aims to obtain refined estimates of the keyframes' velocities, gravity, and the scale factor from the previous step and to estimate the acceleration bias. The parameters that we want to estimate in this step are:

$$\mathcal{X}_2 = \left[ \mathbf{v}_{\mathrm{b}_0:\mathrm{b}_{N-1}}^{\mathrm{w}}, \mathbf{b}_{\mathrm{a}}, w_1, w_2, \mathrm{s} \right]^\top.$$

As previously noted in [8], distinguishing between acceleration bias and gravity can be challenging as they tend to be coupled and difficult to separate. As a result, VINS-Mono

disregards acceleration bias during initialization and assumes it to be zero. Other methods, such as [7], rely on waiting for a prolonged period of time to observe these values. In our approach, we aim to decouple them by refining the initial gravity estimate from the previous step in its tangent space, and adding a weight matrix $\mathcal{W}$ to Eq. (8) to keep acceleration bias at zero and handle outliers when the motion performed does not provide enough information or is blurred.

We refine the gravity estimate using an approach similar to VINS-Mono, which maintains the magnitude of the gravity vector and adjusts it with two variables in its tangent space. We also decouple the acceleration bias during this process. This allows us to represent the gravity vector in a more accurate way:

$$\mathbf{g}^{\mathrm{w}} = g\mathbf{g}_{unit}^{\mathrm{w}} + \delta\mathbf{g}, \quad \delta\mathbf{g} = w_1\mathbf{b}_1 + w_2\mathbf{b}_2, \qquad (7)$$

where $g$ is the known magnitude of gravity, $\mathbf{g}_{unit}^{\mathrm{w}}$ is a unit vector denoting the gravity direction obtained by the previous step. $\mathbf{b}_1$ and $\mathbf{b}_2$ are two orthogonal basis vectors spanning the tangent plane. The initial values of $w_1$ and $w_2$ are set to zero.

By substituting the value of $\mathbf{g}^{\mathrm{w}}$ from Eq. (7) into Eq. (3) and Eq. (4), and introducing the acceleration bias term by approximating the first order of $\Delta\mathbf{p}_{k,k+1}$ and $\Delta\mathbf{v}_{k,k+1}$, we can rewrite Eq. (6) to obtain a new equation with a weight matrix $\mathcal{W}_{k,k+1}$ between two consecutive keyframes $\mathrm{b}_k$ and $\mathrm{b}_{k+1}$ as follows:

$$\mathcal{W}_{k,k+1}\mathcal{H}_{k,k+1}\mathcal{X}_2 = \mathcal{W}_{k,k+1}\mathcal{Z}_{k,k+1}. \qquad (8)$$

$$\mathcal{W}_{k,k+1} = \begin{bmatrix} w_{k,k+1}^{\alpha} & 0_{3\times3} \\ 0_{3\times3} & w_{k,k+1}^{\beta} \end{bmatrix},$$

$$\mathbf{e}_{\boldsymbol{\alpha}} = \mathcal{H}_{k,k+1}[0:2]\mathcal{X}_2 - \mathcal{Z}_{k,k+1}[0:2],$$
$$\mathbf{e}_{\boldsymbol{\beta}} = \mathcal{H}_{k,k+1}[3:5]\mathcal{X}_2 - \mathcal{Z}_{k,k+1}[3:5],$$

$$w_{k,k+1}^{\alpha} = \mathrm{diag}(\exp(-\|\mathbf{e}_{\boldsymbol{\alpha}}\|), \exp(-\|\mathbf{e}_{\boldsymbol{\alpha}}\|), \exp(-\|\mathbf{e}_{\boldsymbol{\alpha}}\|)),$$
$$w_{k,k+1}^{\beta} = \mathrm{diag}(\exp(-\|\mathbf{e}_{\boldsymbol{\beta}}\|), \exp(-\|\mathbf{e}_{\boldsymbol{\beta}}\|), \exp(-\|\mathbf{e}_{\boldsymbol{\beta}}\|)),$$

where $\mathcal{W}_{k,k+1}$ is a $6\times6$ matrix, $\mathbf{e}_{\boldsymbol{\alpha}}$ and $\mathbf{e}_{\boldsymbol{\beta}}$ are $3\times1$ vectors. $\mathcal{H}_{k,k+1}[i:j]$ denotes the submatrix of $\mathcal{H}_{k,k+1}$ from row $i$ to row $j$, and $\mathcal{Z}_{k,k+1}[i:j]$ denotes the submatrix of $\mathcal{Z}_{k,k+1}$ from row $i$ to row $j$. The detailed forms of matrices $\mathcal{H}_{k,k+1}$ and $\mathcal{Z}_{k,k+1}$ are given in the appendix.

By utilizing the solution of keyframes' velocities, scale, and gravity from the previous step as a seed, we can obtain $\mathcal{X}_2$ by solving the following linear least squares with preconditioned conjugate gradient (PCG):

$$\min_{\mathcal{X}_2} \sum_{k\in\mathcal{K}} \|\mathcal{W}_{k,k+1}\mathcal{H}_{k,k+1}\mathcal{X}_2 - \mathcal{W}_{k,k+1}\mathcal{Z}_{k,k+1}\|^2$$

where $\mathcal{K}$ indexes all $N$ keyframes.

## IV. EXPERIMENTAL RESULTS

We evaluate the proposed initialization method using the EuRoC dataset [20] and compare it to two state-of-the-art disjoint visual-inertial initialization methods: 1) a linear loosely-coupled method in VINS-Mono [9], and 2) a non-linear optimization method in ORB-SLAM3 [13],[10]. The

evaluation metrics include computation speed, accuracy, and robustness. Note that since, according to Campos et al. [13], disjoint initialization generally outperforms joint approaches, we compare our method to disjoint methods only. We run each sequence five times, select the run that achieves median accuracy, and use that as the final outcome for all the metrics. Overall, the proposed method achieves the best performance in terms of accuracy and robustness, at competitive computation speed.

### A. Experimental Setup

The EuRoC dataset provides accurate rotation and translation data for 11 sequences recorded by a Micro Aerial Vehicle (MAV). These sequences vary from slow flights under favorable visual conditions to dynamic flights under challenging conditions such as motion blur, poor illumination, and occlusion. The dataset features visual-inertial sensor units that are hardware time synchronized, including: 1) two global shutter, monochrome cameras recording at 2x20 FPS and 2) a MEMS IMU providing angular rate and acceleration data at 200 Hz. Additionally, the dataset includes camera intrinsic and camera-IMU extrinsic parameters.

To guarantee a fair comparison among the various initialization methods, the initialization part of VINS-Mono and EDI are integrated into ORB-SLAM3, enabling the evaluation of all methods using the same ORB-SLAM3 framework. All experiments are conducted on an Intel i7-10700K desktop with 32GB of RAM.

Specifically, EDI and the initialization method of VINS-Mono are integrated into the Local Mapping thread in ORB-SLAM3, without affecting the real-time performance of the tracking thread. The following parameters have been defined for EDI: the standard deviation of angular velocity measurement noise is set to $1.7e^{-4}$ [rad/s], the standard deviation of angular velocity random walk noise is $2e^{-5}[rad/s\sqrt{s}]$, the magnitude of gravity is 9.81 [m/$s^2$] and the number of iterations in the PCG is 4.

### B. Computation Speed Evaluation

A fast initialization is important to allow the SLAM system to proceed to the tracking step in real-time. EDI

TABLE II: Computation time for estimating inertial parameters for ORB-SLAM3 [10], VINS-Mono [9], and our method (EDI) during initialization. Our approach, EDI, is an inertial only initialization method that estimates scale, keyframe velocities, gravity direction, and IMU biases using only inertial residuals, without considering visual residuals.

| Seq name | EDI Time (ms) Inert. Only | ORB-SLAM3 Time (ms) Inert. Only | ORB-SLAM3 Time (ms) Inert. Only + VI-BA | VINS-Mono Time (ms) Inert. Only |
|---|---|---|---|---|
| V1_01_easy | 0.44 | 1.23 | 60.30 | **0.30** |
| V1_02_medium | 0.46 | 1.55 | 28.33 | **0.22** |
| V1_03_difficult | **0.46** | 1.43 | 41.24 | 0.48 |
| V2_01_easy | 0.53 | 1.20 | 49.79 | **0.52** |
| V2_02_medium | 0.66 | 1.43 | 41.22 | **0.33** |
| V2_03_difficult | **0.54** | 1.16 | 46.05 | 0.66 |
| MH_01_easy | 0.60 | 1.56 | 22.50 | **0.52** |
| MH_02_easy | **0.38** | 1.22 | 28.27 | 0.64 |
| MH_03_medium | **0.48** | 1.13 | 32.37 | **0.48** |
| MH_04_difficult | 0.62 | 1.50 | 26.71 | **0.44** |
| MH_05_difficult | 0.75 | 1.51 | 22.25 | **0.31** |
| Avg | 0.54 | 1.36 | 36.28 | **0.45** |

TABLE III: Scale error comparison after initialization without VI-BA. The results of ORB-SLAM3 and VINS-Mono were obtained by executing the publicly available code with its default configuration. (Campos et al. [13] report the runs with the highest accuracy, while we report the median of five runs.)

| Seq name | EDI Scale Error(%) Inert. Only | ORB-SLAM3 Scale Error(%) Inert. Only | ORB-SLAM3 Scale Error(%) Inert. Only + VI-BA | VINS-Mono Scale Error(%) Inert. Only |
|---|---|---|---|---|
| V1_01_easy | **0.3** | 2.9 | 1.4 | 19.5 |
| V1_02_medium | **0.4** | 9.6 | 4.6 | 18.9 |
| V1_03_difficult | **9.8** | 53.6 | 13.3 | 14.9 |
| V2_01_easy | **0.6** | 2.8 | 2.2 | 5.6 |
| V2_02_medium | 4.4 | 3.8 | 3.8 | **3.4** |
| V2_03_difficult | 4.8 | 7.7 | **4.3** | 21.9 |
| MH_01_easy | 3.3 | 5.8 | **3.0** | 7.2 |
| MH_02_easy | 5.6 | **5.4** | 10.4 | 5.8 |
| MH_03_medium | **3.4** | 129.4 | 158.5 | 23.7 |
| MH_04_difficult | 19.9 | 17.6 | **1.7** | 13.3 |
| MH_05_difficult | **14.2** | 1017.4 | 97.6 | 31.2 |
| Avg | **5.8** | 114.2 | 27.3 | 15.0 |

TABLE IV: Scale and full trajectory error comparison using different initialization methods after VI-BA. ATE(m) is the absolute trajectory error of the entire trajectory in meters without a Sim(3) transformation.

| Seq name | EDI Scale Error(%) | ATE(m) | ORB-SLAM3 Scale Error(%) | ATE(m) | VINS-Mono Scale Error(%) | ATE(m) |
|---|---|---|---|---|---|---|
| V1_01_easy | **0.9** | **0.031** | 1.1 | 0.032 | 1.9 | 0.044 |
| V1_02_medium | **0.0** | **0.061** | 0.6 | 0.064 | 0.6 | 0.075 |
| V1_03_difficult | 3.3 | **0.076** | 3.7 | 0.076 | **1.0** | 0.103 |
| V2_01_easy | **1.1** | 0.059 | 1.5 | 0.060 | 1.8 | 0.063 |
| V2_02_medium | **0.1** | 0.063 | 0.8 | **0.059** | 0.7 | 0.065 |
| V2_03_difficult | **0.2** | 0.060 | 0.7 | 0.063 | 0.3 | 0.075 |
| MH_01_easy | **1.5** | **0.085** | 2.2 | 0.093 | 2.2 | 0.104 |
| MH_02_easy | **0.3** | **0.077** | 1.6 | 0.081 | 1.0 | 0.081 |
| MH_03_medium | **0.4** | **0.066** | 0.5 | 0.070 | 1.3 | 0.198 |
| MH_04_difficult | **0.0** | 0.132 | 0.7 | **0.107** | 1.3 | 0.183 |
| MH_05_difficult | 1.0 | 0.126 | 1.3 | **0.110** | **0.7** | 0.382 |
| Avg | **0.8** | 0.076 | 1.3 | **0.074** | 1.3 | 0.125 |

TABLE V: Relative rotation error comparison after the initialization in four challenging sequences with added noise. The rotation error is calculated as the root mean square error of the relative rotation (RMSE).

| Seq name | EDI RMSE (rad) | ORB-SLAM3 RMSE (rad) | VINS-Mono RMSE (rad) |
|---|---|---|---|
| V1_03_difficult | **0.116** | 0.302 | 0.298 |
| V2_03_difficult | **0.160** | 0.374 | 0.291 |
| MH_04_difficult | **0.125** | 0.304 | 0.304 |
| MH_05_difficult | **0.127** | 0.297 | 0.297 |
| Avg | **0.132** | 0.319 | 0.298 |

eliminates the need for the computationally intensive step of visual-inertial bundle adjustment (VI-BA) during the initialization, similar to VINS-Mono [9], making it more efficient and faster. Table II presents the runtime comparison for estimating inertial parameters during the initialization stage involving 10 keyframes, with the best results for each sequence highlighted in bold. The results show that our method is 2-3 times faster than the ORB-SLAM3 [10] inertial-only method on average and has a comparable computational efficiency to VINS-Mono.

*C. Accuracy Evaluation*

To measure accuracy, we evaluate both scale and trajectory error. Scale error measures how closely the estimated scale factor aligns with the true scale, calculated as $|s^* - \hat{s}|/|s^*| \times 100\%$ , where $\hat{s}$ is the scale factor determined by aligning the estimated trajectory to the ground truth, and $s^*$ is 1.

Table III compares the scale error using 10 keyframe trajectories for initialization. Our method consistently achieves the highest accuracy, with a scale error of 5% or less for most tasks. In contrast, the other two methods struggle to achieve similar levels of accuracy. Although ORB-SLAM3 with VI-BA performs better on the MH_04_difficult task, it takes more than 67 times longer to run on average compared to our method. Furthermore, the effectiveness of using VI-BA in ORB-SLAM3 heavily depends on the quality of the inertial-only estimation used as a seed for finding the optimal solution, and the initialization method in ORB-SLAM3 is highly sensitive to the scale factor. As shown in Table IV, our method's trajectory accuracy is comparable to ORB-SLAM3, with an average scale error of less than 1% for entire trajectories after undergoing two rounds of VI-BA refinement, similar to what is used in ORB-SLAM3.

*D. Robustness Evaluation*

A robust initialization step is crucial for ensuring a reliable starting point for the SLAM system to build upon. To evaluate the robustness of our method, we test it in challenging conditions such as motion blur and illumination change using sequences MH_04_difficult, MH_05_difficult, V1_03_difficult, and V2_03_difficult. Additionally, to make the task more challenging, we introduce noise to the rotation estimates of keyframes used in the initialization, obtained from a monocular SLAM system. The added noise has a standard deviation of 0.1 radians for roll, pitch, and yaw, and simulates a scenario in which pose estimation from pure monocular SLAM is of poor quality and the initialization methods are challenged to maintain accuracy and stability.

As shown in Table V and Fig. 2, our method outperforms ORB-SLAM3 [10] and VINS-Mono [9] in rotation estimation, with a median error of 0.099-0.128 radians compared to ORB-SLAM3's 0.160-0.350 radians and VINS-Mono's 0.231-0.247 radians. It also has a lower root mean square error (rmse) of 0.116-0.127 radians compared to ORB-SLAM3's 0.297-0.374 radians and VINS-Mono's 0.291-0.304 radians in four challenging sequences with added noise. Our method, which uses ESKF during the initialization phase, is able to improve the accuracy of the rotation estimates, particularly when the visual estimates from a pure monocular SLAM system are not accurate. In terms of the full trajectory, in one of the four tasks (V1_03_difficult sequence), our method can run the entire trajectory with an Absolute Trajectory Error (ATE) of 0.253 meters, while the other two methods fail to run all of the four challenging sequences with added noise, which demonstrates the robustness of our proposed method in comparison to the others.

## V. CONCLUSIONS

Our proposed approach, EDI addresses the limitations of previous disjoint methods by utilizing an Error-state Kalman Filter (ESKF) to estimate gyroscope bias and correct rotation estimates, providing adaptability for application to other SLAM systems that use sensors such as GNSS and GPS. In addition, EDI offers a closed-form solution

**Fig. 2:** Box plots of the relative rotation error: comparison between our initialization method and ORB-SLAM3 for pure monocular SLAM's rotation estimation with added noise.

and introduces weights to handle outliers when estimating initial velocity, scale, gravity, and acceleration bias. EDI outperforms previous disjoint methods in terms of accuracy and robustness, at competitive computation speed, even in challenging environments with artificial noise. This new approach has promising potential for the development of efficient and reliable navigation systems in the future.

## VI. APPENDIX

Details of Eq. (8)

$$\mathcal{H}_{k,k+1} =$$
$$\begin{bmatrix} \boldsymbol{\alpha}^a & \boldsymbol{\alpha}^b & \boldsymbol{\alpha}^c & \boldsymbol{\alpha}^d & \boldsymbol{\alpha}^e & \boldsymbol{\alpha}^f & \boldsymbol{\alpha}^g \\ \boldsymbol{\beta}^a & \boldsymbol{\beta}^b & \boldsymbol{\beta}^c & \boldsymbol{\beta}^d & \boldsymbol{\beta}^e & \boldsymbol{\beta}^f & \boldsymbol{\beta}^g \end{bmatrix},$$

$$\mathcal{Z}_{k,k+1} = \begin{bmatrix} \Delta\mathbf{p}_{k,k+1} - \mathbf{p}_c^b + \mathbf{R}_w^{b_k} \mathbf{R}_{b_{k+1}}^w \mathbf{p}_c^b + \frac{1}{2}\mathbf{R}_w^{b_k} \Delta\mathbf{t}_{k,k+1}^2 \mathbf{g_0} \\ \Delta\mathbf{v}_{k,k+1} + \mathbf{R}_w^{b_k} \Delta\mathbf{t}_{k,k+1} \mathbf{g_0} \end{bmatrix},$$

$$\begin{aligned}
\boldsymbol{\alpha}^a &= 0_{3\times 3k}, & \boldsymbol{\beta}^a &= 0_{3\times 3k} \\
\boldsymbol{\alpha}^b &= -\mathbf{R}_w^{b_k}\Delta\mathbf{t}_{k,k+1} & \boldsymbol{\beta}^b &= -\mathbf{R}_w^{b_k} \\
\boldsymbol{\alpha}^c &= 0_{3\times 3} & \boldsymbol{\beta}^c &= \mathbf{R}_w^{b_k} \\
\boldsymbol{\alpha}^d &= 0_{3\times 3(N-k-2)} & \boldsymbol{\beta}^d &= 0_{3\times 3(N-k-2)} \\
\boldsymbol{\alpha}^e &= -\boldsymbol{J}_{\mathbf{b_a}}^{\Delta\mathbf{p}} & \boldsymbol{\beta}^e &= -\boldsymbol{J}_{\mathbf{b_a}}^{\Delta\mathbf{v}} \\
\boldsymbol{\alpha}^f &= -\frac{1}{2}\mathbf{R}_w^{b_k}\mathbf{b}\Delta\mathbf{t}_{k,k+1}^2 & \boldsymbol{\beta}^f &= -\mathbf{R}_w^{b_k}\mathbf{b}\Delta\mathbf{t}_{k,k+1} \\
\boldsymbol{\alpha}^g &= \mathbf{R}_w^{b_k}(\bar{\mathbf{p}}_{c_{k+1}}^w - \bar{\mathbf{p}}_{c_k}^w) & \boldsymbol{\beta}^g &= 0_{3\times 1}
\end{aligned}$$

where $\mathcal{H}_{k,k+1}$ has dimensions $6 \times (3N + 6)$ and $\mathcal{Z}_{k,k+1}$ is a $6 \times 1$ vector. $\mathbf{g_0} = g\mathbf{g}_{unit}^w$, where $\mathbf{g}_{unit}^w$ is the unit vector of the gravity in the world frame. The Jacobians $\boldsymbol{J}_{\mathbf{b_a}}^{\Delta\mathbf{p}}$ and $\boldsymbol{J}_{\mathbf{b_a}}^{\Delta\mathbf{v}}$ represent how the preintegration changes due to a small difference in bias estimation, and the vector of biases $\mathbf{b} = [\mathbf{b}_1, \mathbf{b}_2]^\top$ includes two bias terms, $\mathbf{b}_1$ and $\mathbf{b}_2$. These bias terms are used to perturb the gravity vector, and are chosen to be two orthogonal basis vectors on the tangent plane.

## REFERENCES

[1] L. Kneip, A. Martinelli, S. Weiss, D. Scaramuzza, and R. Siegwart, "Closed-form solution for absolute scale velocity determination combining inertial measurements and a single feature correspondence," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2011, pp. 4546–4553.

[2] A. Martinelli, "Closed-form solution of visual-inertial structure from motion," *International Journal of Computer Vision (IJCV)*, vol. 106, no. 2, pp. 138–152, 2014.

[3] J. Kaiser, A. Martinelli, F. Fontana, and D. Scaramuzza, "Simultaneous state initialization and gyroscope bias calibration in visual inertial aided navigation," *IEEE Robotics and Automation Letters*, vol. 2, no. 1, pp. 18–25, 2017.

[4] T. Dong-Si and A. I. Mourikis, "Estimator initialization in vision-aided inertial navigation with unknown camera-imu calibration," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2012, pp. 1064–1071.

[5] T. Dong-Si and A. I. Mourikis, "Closed-form solutions for vision-aided inertial navigation." in *Tech. rep. Dept. of Electrical Engineering, University of California, Riverside*, 2011.

[6] C. Campos, J. M. M. Montiel, and J. D. Tardós, "Fast and robust initialization for visual-inertial slam," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019, pp. 1288–1294.

[7] R. Mur-Artal and J. D. Tardós, "Visual-Inertial monocular SLAM with map reuse," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 796–803, 2017.

[8] W. Huang and H. Liu, "Online initialization and automatic camera-imu extrinsic calibration for monocular visual-inertial slam," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 5182–5189.

[9] T. Qin, P. Li, and S. Shen, "VINS-Mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.

[10] C. Campos, R. Elvira, J. J. Gómez, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM," *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.

[11] D. Zuñiga-Noël, F.-A. Moreno, and J. Gonzalez-Jimenez, "An analytical solution to the imu initialization problem for visual-inertial systems," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 6116–6122, 2021.

[12] J. Cheng, L. Zhang, and Q. Chen, "An improved initialization method for monocular visual-inertial SLAM," *Electronics*, vol. 10, no. 24, p. 3063, 2021.

[13] C. Campos, J. M. M. Montiel, and J. D. Tardós, "Inertial-only optimization for visual-inertial initialization," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 51–57.

[14] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, "Keyframe-based visual–inertial odometry using nonlinear optimization," *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 314–334, 2015.

[15] C. Forster, Z. Zhang, M. Gassner, M. Werlberger, and D. Scaramuzza, "SVO: Semidirect visual odometry for monocular and multicamera systems," *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 249–265, 2017.

[16] M. Bloesch, M. Burri, S. Omari, M. Hutter, and R. Siegwart, "Iterated extended kalman filter based visual-inertial odometry using direct photometric feedback," *The International Journal of Robotics Research*, vol. 36, no. 10, pp. 1053–1072, 2017.

[17] B. Joshi, S. Rahman, M. Kalaitzakis, *et al.*, "Experimental comparison of open source visual-inertial-based state estimation algorithms in the underwater domain," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 7227–7233.

[18] J. Sola, "Quaternion kinematics for the error-state kalman filter," *arXiv preprint arXiv:1711.02508*, 2017.

[19] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM: A versatile and accurate monocular slam system," *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.

[20] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart, "The euroc micro aerial vehicle datasets," *The International Journal of Robotics Research*, 2016.