# Accurate and Interactive Visual-Inertial Sensor Calibration with Next-Best-View and Next-Best-Trajectory Suggestion

Christopher L. Choi[1], Binbin Xu[1] and Stefan Leutenegger[1,2]

*Abstract*— **Visual-Inertial (VI) sensors are popular in robotics, self-driving vehicles, and augmented and virtual reality applications. In order to use them for any computer vision or state-estimation task, a good calibration is essential. However, collecting *informative* calibration data in order to render the calibration parameters observable is not trivial for a non-expert. In this work, we introduce a novel VI calibration pipeline that guides a non-expert with the use of a graphical user interface and information theory in collecting *informative* calibration data with Next-Best-View and Next-Best-Trajectory suggestions to calibrate the intrinsics, extrinsics, and temporal misalignment of a VI sensor. We show through experiments that our method is faster, more accurate, and more consistent than state-of-the-art alternatives. Specifically, we show how calibrations with our proposed method achieve higher accuracy estimation results when used by state-of-the-art VI Odometry as well as VI-SLAM approaches. The source code of our software can be found on: `https://github.com/chutsu/yac`.**

## I. INTRODUCTION

In order to use Visual-Inertial (VI) sensors in computer vision or state-estimation tasks the calibration parameters must first be obtained. Conventionally, VI sensors are calibrated by an expert who would often collect calibration data by positioning and moving the sensors in front of a calibration target such as a checkerboard or grid of fiducial markers, then use an offline calibration tool such as Kalibr [1] to estimate the sensor calibration parameters. Good calibration results, however, may only be achieved, if the right kind and right amount of data is collected. More specifically, two potential practical issues arise during data capture: first, the choice of calibration views and the range of motions needed is not immediately clear to the non-expert. Secondly, the amount of data the user has to collect for calibration is also unclear, often collecting too much or too little data. A common practice to address these issues is to collect *multiple* calibration data sequences, however, this is impractical in the field and identifying which calibration is optimal becomes a tedious and time-consuming task.

A straight forward solution to this problem would be to mount the VI sensor on a robot arm and perform a rehearsed or optimal calibration "dance", such as in [2]. However, this requires extra hardware and is not a practical solution for many applications. As an alternative to classic offline calibration methods, one can estimate the calibration parameters within a state-estimation framework such as

[1]Smart Robotics Lab, Department of Computing, Imperial College London, United Kingdom. {`christopher.choi`, `b.xu17`, `s.leutenegger`}`@imperial.ac.uk`
[2]Smart Robotics Lab, Technical University of Munich, Germany. `stefan.leutenegger@tum.de`
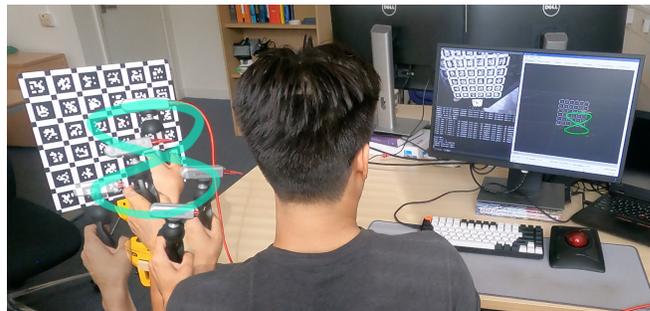
Fig. 1: Our system interactively suggests next-best-actions to collect calibration data.

OKVIS [3], VINS-MONO [4], and OpenVINS [5] in real-time. Note, however, that any of these frameworks require some form of sufficiently accurate initial calibration, as well as sufficient visual features and motion excitation, therefore suffering from similar issues as offline calibration. Furthermore, natural keypoints and the lack of precise knowledge of the corresponding 3D positions may not produce the best possible results.

In this work, we present an interactive VI sensor calibration pipeline that helps guide a non-expert in collecting *informative* calibration data for a VI sensor *once* through Next-Best-View (NBV) and Next-Best-Trajectory (NBT) suggestions (as shown in Fig. 1) in order to efficiently obtain sound calibrations. We show through extensive quantitative experiments on calibration sequences and several self-collected VICON real-world datasets that calibration parameters optimised through our system are more accurate and consistent than Kalibr by testing on state-of-the-art VI-SLAM ORBSLAM3 [6]. In summary our contributions are:

- A complete and open-sourced interactive VI-camera calibration tool that supports any number of cameras;
- An information-theoretic procedure to identify the most informative Next-Best-View (NBV) and Next-Best-Trajectory (NBT) among a pre-defined set of viewpoints and trajectory primitives;
- An interactive graphical user interface for guiding the user through the calibration data collection process;
- Through experiments we show that our proposed method is faster, more accurate and more reliable compared to state of the art traditional *non-guided* calibration methods, such as Kalibr [1], even when used by *novices*.

## II. RELATED WORK

**Offline Methods**. In the robotics community, early works in VI-sensor calibration methods such as [7], [8], [9], [10] showed that it is possible to calibrate the extrinsics between a camera and IMU, with Kalibr [1] regarded as the current state-of-the-art tool. It is an offline method capable of calibrating a multi-camera system, as well as a VI system. However, the use of this tool requires expert knowledge, as the result is highly dependent on the quality of the calibration data captured. Therefore, the calibration process may in practice have to be repeated until desired results are reached.

**Online Methods**. State of the art state-estimation framework such as OKVIS [3], VINS-MONO [4], and Open-VINS [5] can in practice estimate the calibration parameters in real-time. However, these frameworks require sufficiently accurate initial calibrations, as well as sufficient visual features and motion excitation in order to operate accurately.

**Reinforcement Learning Methods**: There has been a growing interest in using reinforcement learning for calibration such as [2], [11], [12] whereby the goal is to learn informative trajectories to render the VI-calibration parameters observable. However, the requirement of a robot arm to perform these motions is not always practical in the field. Further, these works do not provide quantitative results through a SLAM system to verify the optimality of the calibrated parameters.

**Information-Theoretic Methods**. The first calibration tool with an emphasis in guiding the user through capturing a good calibration sequence for a monocular camera is AprilCal [13]. The method used a quality metric to find and suggest the NBV in real-time during the camera calibration process. AprilCal, however, only supports calibrating the intrinsics of a single monocular camera.

A more recent work that uses an information-theoretic approach for VI sensor calibration is the work of [14], [15], where they proposed a segment-based method for calibrating a VI sensor system in a AR / VR headset and self-driving car setting. The idea is to extract informative data during online state-estimation using an information-theoretic metric, and then perform a full-batch optimisation to update the calibration parameters offline. This approach, however, relies on the fact that the VI sensors are calibrated well initially. Secondly, the available data does not guarantee informative segments for calibration.

In this paper, we place heavy emphasis on collecting *informative* calibration data by using an information-theoretic metric to find the NBV and NBT in real-time, and by *interactively* guiding the user in collecting them in order to calibrate the intrinsics, extrinsics, and time shift of a VI sensor. This is in contrast to current state-of-the-art calibration tools such as Kalibr [1] that assume the collected calibration data has sufficient views and range of motion.

## III. NOTATION

We employ the following notation throughout this work. Let $\mathcal{F}_W$ denote the world reference frame. A 3D point $P$ in the world frame $\mathcal{F}_W$ with respect to the origin is written as a position vector ${}_W\mathbf{r}_{WP}$. A rigid body transformation from the body frame, $\mathcal{F}_B$, to the world frame, $\mathcal{F}_W$, is represented by a homogeneous transformation matrix, $\mathbf{T}_{WB}$. Its rotation matrix component is written as, $\mathbf{C}_{WB}$, and the corresponding Hamiltonian quaternion is written as, $\mathbf{q}_{WB} = [\boldsymbol{\eta}^T, \epsilon]^T \in \mathcal{S}^3$, where $\epsilon$ and $\boldsymbol{\eta}$ are the real and imaginary parts.

In general, the state vector we will be estimating lives on a manifold and thus we define an operator $\boxplus$ that will be used to perturb the states in tangent space such that $\mathbf{x} = \bar{\mathbf{x}} \boxplus \delta\mathbf{x}$, where $\bar{\mathbf{x}}$ is the state estimate and $\delta\mathbf{x}$ is the local perturbation. Vector quantities such as positions, velocities, biases are updated via standard vector addition. Rotation components on the other hand such as a quaternion are updated via a combination of the group operator $\otimes$ (quaternion multiplication) and exponential map $\mathrm{Exp}(\cdot)$, such that $\mathbf{q} \boxplus \delta\boldsymbol{\alpha} = \mathrm{Exp}(\delta\boldsymbol{\alpha}) \otimes \mathbf{q}$. As a result we will be using a minimal coordinate representation approach similar to [3]. A comprehensive introduction to differential calculus is beyond of the scope of this paper, the reader is therefore encouraged to review [16], [17] for a more detailed treatment on the subject.

## IV. BACKGROUND

In robotics, the maximum a posteriori (MAP) estimator is commonly used to solve the camera and VI calibration problem,

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\mathrm{argmax}}\ p(\mathbf{x}|\mathbf{z}), \qquad (1)$$

where $\mathbf{x}$ is the state vector which may be comprised of poses, velocities, IMU biases and calibration parameters we are interested in jointly estimating, given the measurements, $\mathbf{z}$. Assuming Gaussian measurements resulting in independent error terms $\mathbf{e}_i$, maximising Eq. (1) is equivalent to solving the sum of nonlinear least squares using a nonlinear optimisation algorithm such as the Gauss Newton method,

$$\sum_i \mathbf{E}_i^T \mathbf{W}_i \mathbf{E}_i\ \Delta\mathbf{x} = \sum_i -\mathbf{E}_i^T \mathbf{W}_i \mathbf{e}_i(\mathbf{x}), \qquad (2)$$

where $\Delta\mathbf{x}$ is the update vector, $\mathbf{e}_i(\mathbf{x})$ is the $i^{\text{th}}$ error term evaluated at the current estimate $\mathbf{x}$, $\mathbf{E}_i$ is the Jacobian matrix of the error term and $\mathbf{W}_i$ the measurement information.

At convergence of the optimisation, we may approximate the posterior distribution as a Gaussian with mean $\mathbf{x}$ and find the covariance matrix $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}}$ by inverting the quantity $\sum_i \mathbf{E}_i^T \mathbf{W}_i \mathbf{E}_i$, also known as the Fisher Information matrix. However, recall that the state vector $\mathbf{x}$ not only contains calibration parameters $\boldsymbol{\theta}$, but also other state variables not related to the calibration parameters which we denote as $\boldsymbol{\gamma}$. In the context of calibration, we are only interested in the estimated calibration parameters, $\boldsymbol{\theta}$, and the covariance of the calibration parameters $\boldsymbol{\Sigma}_{\boldsymbol{\theta}\boldsymbol{\theta}}$. Expressing $\mathbf{x}$ and $\boldsymbol{\Sigma}_{\mathbf{x}}$ in partition form,

$$\mathbf{x} = \begin{bmatrix} \boldsymbol{\theta} \\ \boldsymbol{\gamma} \end{bmatrix},\ \ \boldsymbol{\Sigma}_{\mathbf{x}} = \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\theta}\boldsymbol{\theta}} & \boldsymbol{\Sigma}_{\boldsymbol{\theta}\boldsymbol{\gamma}} \\ \boldsymbol{\Sigma}_{\boldsymbol{\gamma}\boldsymbol{\theta}} & \boldsymbol{\Sigma}_{\boldsymbol{\gamma}\boldsymbol{\gamma}} \end{bmatrix}, \qquad (3)$$

we can employ marginalisation on Normal distributions to get $p(\boldsymbol{\theta}|\mathbf{z}) = \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}\boldsymbol{\theta}})$, by extracting the corresponding blocks in Eq. (3).
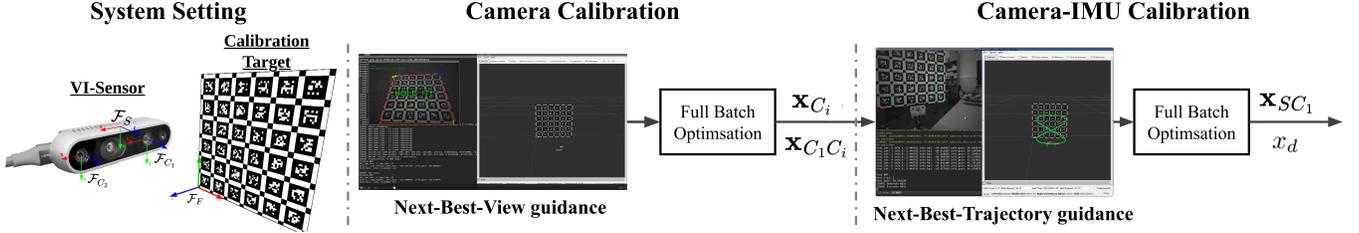
Fig. 2: An overview of our VI calibration pipeline.

To objectively quantify whether the next VI measurements are informative for the VI calibration problem, we used the Mutual Information (MI) defined in [18],

$$I(\boldsymbol{\theta}_1; \tilde{\mathbf{z}}_2) = \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_{\boldsymbol{\theta}_1 \boldsymbol{\theta}_1}|}{|\boldsymbol{\Sigma}_{\boldsymbol{\theta}_1 \boldsymbol{\theta}_1 | \mathbf{z}_2}|}, \quad (4)$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_1 \boldsymbol{\theta}_1}$ is the covariance estimate of $\boldsymbol{\theta}$ using measurements $\mathbf{z}_1$ alone, and $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_1 \boldsymbol{\theta}_1 | \mathbf{z}_2}$ is the covariance estimate of $\boldsymbol{\theta}$ using measurements $\mathbf{z}_1$ and $\mathbf{z}_2$, finally $|\cdot|$ is the matrix determinant. In summary, with Eq. (4) we can measure the amount of information $\mathbf{z}_2$ (next VI-sensor measurements) conveys to our current estimate $\boldsymbol{\theta}|\mathbf{z}_1$.

## V. SYSTEM OVERVIEW

An overview of our proposed calibration system is illustrated in Fig. 2[1]. It consists of two stages. The first stage aims to perform vision-only camera intrinsics and extrinsics calibration employing Next-Best-View (NBV) feedback. In the second stage the camera-IMU extrinsics are found by using Next-Best-Trajectory (NBT) feedback, with the camera intrinsics and extrinsics obtained in the previous stage fixed. Both stages of the calibration process require the use of a static fiducial marker grid of known size as a calibration target. Specifically, we use a planar calibration target grid of AprilTags [13] introduced by Kalibr [1]. Throughout this work, the VI sensor to be calibrated is assumed to capture images and inertial measurements with the same clock source.

## VI. CAMERA INTRINSICS AND EXTRINSICS CALIBRATION

In the following, we detail our approach of using Mutual Information (MI) and Next-Best-View (NBV) to calibrate intrinsics and extrinsics of all cameras.

### A. States

For the camera calibration problem, the states to be estimated consist of the camera poses relative to the fiducial target coordinate frame $\mathcal{F}_F$ as $\mathbf{x}_{FC_1}$, camera extrinsics relative to reference camera 1, $\mathbf{x}_{C_1 C_i}$, and camera intrinsics, $\mathbf{x}_{C_i}$, of the form:

$$\begin{aligned}
\mathbf{x}_{FC_1} &= \begin{bmatrix} F\mathbf{r}_{FC_1}^T & \mathbf{q}_{FC_1}^T \end{bmatrix}^T \in \mathbb{R}^3 \times \mathcal{S}^3, \\
\mathbf{x}_{C_1 C_i} &= \begin{bmatrix} C_1\mathbf{r}_{C_1 C_i}^T & \mathbf{q}_{C_1 C_i}^T \end{bmatrix}^T \in \mathbb{R}^3 \times \mathcal{S}^3, \quad (5) \\
\mathbf{x}_{C_i} &= \begin{bmatrix} f_x & f_y & c_x & c_y & k_1 & k_2 & p_1 & p_2 \end{bmatrix}^T \in \mathbb{R}^8,
\end{aligned}$$

[1]The pipeline is demonstrated in details in the supplementary video.

where $\mathcal{F}_{C_i}$ denotes the coordinate frame of the $i^{\text{th}}$ camera on the sensor assembly. We used the Radial-Tangential camera model consisting of focal lengths $f_x, f_y$, centre $c_x, c_y$, radial distortion parameters $k_1, k_2$, and tangential distortion parameters $p_1, p_2$ as the camera intrinsics. Note that any other projection model could be supported in principle. The full state vector for camera calibration thus becomes,

$$\mathbf{x} = \begin{bmatrix} \underbrace{\mathbf{x}_{FC_1}^{T,1} \cdots \mathbf{x}_{FC_1}^{T,k}}_{\text{Reference Camera 1 Poses}} & \underbrace{\mathbf{x}_{C_1 C_1}^T \cdots \mathbf{x}_{C_1 C_i}^T}_{\text{Camera Extrinsics}} & \underbrace{\mathbf{x}_{C_1}^T \cdots \mathbf{x}_{C_i}^T}_{\text{Camera Intrinsics}} \end{bmatrix}^T. \quad (6)$$

### B. Calibration Formulation

To estimate the camera calibration parameters we used a nonlinear least squares framework to minimise the cost function, $J_{\text{camera}}$, containing reprojection errors, $\mathbf{e}_r$, and the information matrix of the respective camera measurement, $\mathbf{W}_r$. The cost function has the form:

$$J_{\text{camera}}(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^I \sum_{k=1}^K \sum_{j \in \mathcal{J}(i,k)} \mathbf{e}_r^{i,j,k^T} \mathbf{W}_r^{i,j,k} \mathbf{e}_r^{i,j,k}, \quad (7)$$

where $i$ is the camera index, $k$ denotes the camera frame index, and $j$ denotes the fiducial target corner index. Finally, $\mathcal{J}(i,k)$ denotes the set of observable fiducial corner indices in the $i^{\text{th}}$ camera index and $k^{\text{th}}$ camera frame index.

Here, the standard reprojection error, $\mathbf{e}_r$, was used:

$$\mathbf{e}_r^{i,j,k} = \tilde{\mathbf{z}}^{i,j,k} - \mathbf{h}_i(\mathbf{T}_{C_i C_1} \; \mathbf{T}_{FC_1}^{-1} \; {}_F\mathbf{r}_{FF_j}, \; \mathbf{x}_{C_i}), \quad (8)$$

whereby $\mathbf{h}_i(\cdot)$ denotes the camera projection and distortion model. It needs as an input the fiducial corner, ${}_F\mathbf{r}_{FF_j}$, camera pose, $\mathbf{T}_{FC_1}$, camera extrinsics, $\mathbf{T}_{C_i C_1}$, and camera intrinsics $\mathbf{x}_{C_i}$. Lastly, $\tilde{\mathbf{z}}^{i,j,k}$ is the observed fiducial corner measurement.

### C. Real-time Estimation

Since Eq. (7) will grow in complexity with every camera frame added, it cannot be solved in real-time as the problem size increases. We therefore adopted a fixed-lag sliding window scheme similar to [3], whereby the sliding window is bounded by marginalising out old camera poses $\mathbf{x}_{FC_1}$ with the Schur Complement, leading to a respective linear prior that enters the cost. Note that this is only needed for the real-time feedback to the user, and we still solve the full batch problem offline for the final calibration solution.
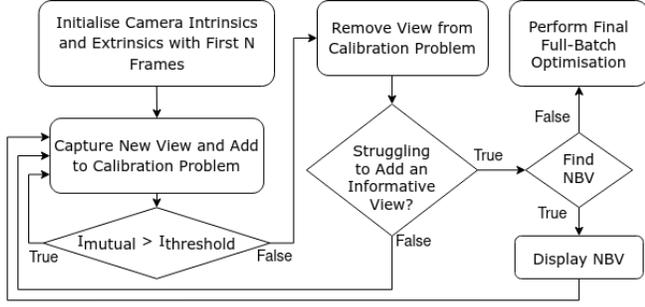
Fig. 3: Camera Calibration Pipeline

## D. Camera Calibration With Next-Best-View

In contrast to standard full-batch camera calibration, where the calibration data is first collected and then solved as a two step process, our method takes a more integrated approach, whereby data collection and solving the calibration problem are performed incrementally, until the addition of new data is no longer informative to the camera calibration problem (see Fig. 3).
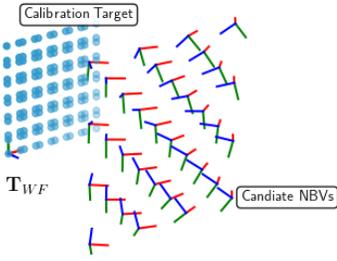


Fig. 4: NBV candidate poses in-front of the calibration target

First, the camera intrinsics and extrinsics are initialised with the first $N$ camera frames of a static fiducial marker of known size by minimising the cost function in Eq. (7). Once the camera parameters are initialised, the user is guided to maximise the calibration target measurement coverage over the image space. The information content of each camera view is evaluated using Eq. (4). Views which contain a MI score below the user-defined threshold, $I_{\text{threshold}}$ ($I_{\text{threshold}} = 0.2$, same as in Kalibr [1]), are removed from the calibration problem. If, however, the new candidate views are not informative enough (no new views added to the calibration problem in the last 3 frames), the calibration tool enters into "Find Next-Best-View" mode where it evaluates a set of possible NBVs. Similar to [13], NBVs are pre-determined by an expert ahead of time in order to reduce the search space and make the computation feasible in real-time (see Fig. 4). Using Eq. (4), the NBV is the one that has the highest mutual information. Once the NBV is determined, the calibration tool will guide the user to the NBV interactively through the graphical user-interface in capturing that view. If the mutual information of the NBV is found to be below $I_{\text{threshold}}$ the calibration tool stops capturing further measurements and proceeds to performing a final full batch optimisation to
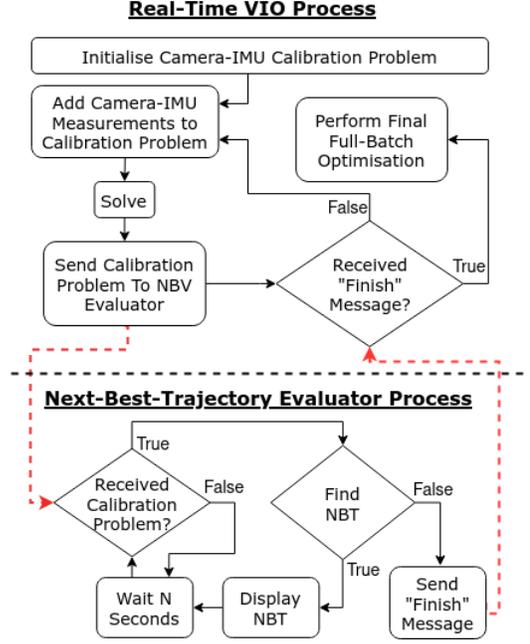


Fig. 5: Camera-IMU Calibration Pipeline

estimate the final calibration parameters.

## VII. CAMERA-IMU EXTRINSICS CALIBRATION

Once the camera intrinsics and extrinsics are known (from Sec. VI), we proceed to, without loss of generality, calibrate the extrinsics between the reference camera 1 and IMU, $\mathbf{T}_{SC_1}$, and camera-IMU delay, $t_d$, of a VI-sensor.

### A. States

The variables to be estimated are the VI sensor pose at discrete camera frame index $k$, $\mathbf{x}_{WS}^k$, fiducial target pose in the inertial frame $\mathbf{x}_{WF}$, extrinsics between reference camera 1 and IMU $\mathbf{x}_{WF}$, and camera-IMU time delay $x_d$:

$$
\begin{aligned}
\mathbf{x}_{WS} &= \begin{bmatrix} {_W}\mathbf{r}_{WS}^T & \mathbf{q}_{WS\ W}^T & \mathbf{v}_{WS}^T & \mathbf{b}_g^T & \mathbf{b}_a^T \end{bmatrix}^T \in \mathbb{R}^3 \times \mathcal{S}^3 \times \mathbb{R}^9, \\
\mathbf{x}_{WF} &= \begin{bmatrix} {_W}\mathbf{r}_{WF}^T & \mathbf{q}_{WF}^T \end{bmatrix}^T \in \mathbb{R}^3 \times \mathcal{S}^3, \\
\mathbf{x}_{SC_1} &= \begin{bmatrix} {_S}\mathbf{r}_{SC_1}^T & \mathbf{q}_{SC_1}^T \end{bmatrix}^T \in \mathbb{R}^3 \times \mathcal{S}^3, \\
x_d &= t_d \in \mathbb{R},
\end{aligned}
\tag{9}
$$

where the state vector $\mathbf{x}_{WS}$ holds the VI sensor position in the inertial frame $_W\mathbf{r}_{WS}$, the body orientation represented by a quaternion $\mathbf{q}_{WS}$, the velocity expressed in the sensor frame $\mathbf{v}_{WS}$, as well as the gyroscope and accelerometer biases $\mathbf{b}_g$ and $\mathbf{b}_a$. The state vectors $\mathbf{x}_{SC_1}$ and $\mathbf{x}_{WF}$ hold the sensor-camera relative pose and fiducial pose, respectively. The full state vector for camera-IMU calibration thus becomes,

$$
\mathbf{x} = \begin{bmatrix} \underbrace{\mathbf{x}_{WS_1}^{T,1} \dots \mathbf{x}_{WS_1}^{T,k}}_{\substack{\text{Sensor} \\ \text{Poses}}} & \underbrace{\mathbf{x}_{WF}^T}_{\substack{\text{Fiducial} \\ \text{Pose}}} & \underbrace{\mathbf{x}_{SC_1}^T}_{\substack{\text{Camera-IMU} \\ \text{Extrinsics}}} & \underbrace{x_d}_{\substack{\text{Camera-IMU} \\ \text{Time-Delay}}} \end{bmatrix}^T .
\tag{10}
$$

## B. Calibration Formulation

Similar to Sec. VI, we seek to formulate the VI calibration problem as one joint nonlinear-optimisation of a cost function $J_{\text{imu-cam}}(\mathbf{x})$ containing both (weighted) reprojection errors $\mathbf{e}_r$ and (weighted) temporal error term from the IMU $\mathbf{e}_s$:

$$J_{\text{imu-cam}}(\mathbf{x}) = \underbrace{\frac{1}{2} \sum_{i=1}^{I} \sum_{k=1}^{K} \sum_{j \in \mathcal{J}(i,k)} \mathbf{e}_r^{i,j,k^T} \mathbf{W}_r^{i,j,k} \mathbf{e}_r^{i,j,k}}_{\text{visual}} \quad (11)$$

$$+ \underbrace{\frac{1}{2} \sum_{k=1}^{K-1} \mathbf{e}_s^{k^T} \mathbf{W}_s^k \mathbf{e}_s^k}_{\text{inertial}},$$

where $i$ is the camera index of the VI sensor, $k$ denotes the camera frame index, and $j$ denotes the fiducial target corner index. The set $\mathcal{J}(i,k)$ represents the indices of fiducial target corners observed in the $k^{\text{th}}$ frame and the $i^{\text{th}}$ camera.

The reprojection error was used to estimate the camera-IMU extrinsics $\mathbf{T}_{SC_1}$, sensor pose in the world frame $\mathbf{T}_{WS}$ and fiducial target in the world frame $\mathbf{T}_{WF}$:

$$\mathbf{e}_r = \tilde{\mathbf{z}}^{i,j,k} - \mathbf{h}_i(\mathbf{T}_{C_1 C_i}^{-1} \mathbf{T}_{SC_1}^{-1} \mathbf{T}_{SW}^k \mathbf{T}_{WF\,F} \mathbf{r}_{FF_j}, \ \mathbf{x}_{C_i}), \ (12)$$

where $\mathbf{h}_i(\cdot)$ denotes the $i^{\text{th}}$ camera projection model which includes distortion, $_F\mathbf{r}_{FF_j}$ denotes the $j^{\text{th}}$ fiducial target corner point and $\tilde{\mathbf{z}}^{i,j,k}$ denotes the corresponding measurement seen in camera $i$ and image frame $k$ in image coordinates. The camera-intrinsics $\mathbf{x}_{C_i}$ and camera-extrinsics $\mathbf{T}_{C_1 C_i}$ estimated in Sec. VI are fixed.

The fiducial target in the world frame $\mathbf{T}_{WF}$ is first initialised using initial measurements from the IMU and camera assuming low acceleration, where the measured acceleration vector corresponds to (inverse) acceleration due to gravity–yielding the camera pose $\mathbf{T}_{WC_i}$. Without loss of generality, we set the camera position and yaw around the world-z axis to zero. Next, the relative pose between the fiducial target and the $i^{\text{th}}$ camera, $\mathbf{T}_{FC_i}$, is computed with fiducial corner measurements using 3D-2D RANSAC and bundle adjustment, after which we can compose $\mathbf{T}_{WF} = \mathbf{T}_{WC_i} \mathbf{T}_{C_i F}$.

For the IMU error term, we adopted the pre-integration scheme in [19], where the error is the difference between the predicted relative state and the actual relative state, with the exception of orientation, where a simple multiplicative minimal error was used:

$$\mathbf{e}_S^k(\mathbf{x}_S^k, \mathbf{x}_S^{k+1}, \tilde{\mathbf{z}}_S^k) = \begin{bmatrix} {}_S\hat{\mathbf{r}}_{WS}^{k,k+1}(t_d) - {}_S\mathbf{r}_{WS}^{k,k+1} \\ 2\left[\mathbf{q}_S^{k,k+1} \otimes \hat{\mathbf{q}}_S^{k,k+1}(t_d)\right]_{1:3} \\ {}_S\hat{\mathbf{v}}^{k,k+1}(t_d) - {}_S\mathbf{v}^{k,k+1} \\ \hat{\mathbf{b}}_g^{k+1}(t_d) - \mathbf{b}_g^{k+1} \\ \hat{\mathbf{b}}_a^{k+1}(t_d) - \mathbf{b}_a^{k+1} \end{bmatrix} \in \mathbb{R}^{15}.$$

$$(13)$$

In addition to estimating the relative state, we further include the camera-IMU time delay scalar $t_d$. Since it is only a 1 dimensional parameter, the $15 \times 1$ Jacobian was obtained through the central finite difference by perturbing the IMU timestamps.

## C. Real-time Estimation

To keep the problem in Eq. (11) bounded for real-time operation, we used the same approach as in Sec. VI-C and adopted a fixed-lag sliding window scheme, marginalising out old sensor poses $\mathbf{T}_{WS}$, velocities $_W\mathbf{v}_{WS}$, accelerometer biases $\mathbf{b}_a$ and gyroscope biases $\mathbf{b}_g$. A full batch optimisation using all measurements will be performed to obtain the final calibration solution. The camera-IMU time delay parameter is fixed during online guidance, and estimated in the final full batch optimisation.

## D. Next-Best-Trajectories

Similar to [12], given our goal is to provide intuitive, easy and real-time feedback for a non-expert user to calibrate the VI-sensor, we discretized the continuous search space and used the results of [20] to design 6 non-degenerate NBTs that are computationally feasible in real-time, easy to display and followed by the user (see Fig. 6). Our NBTs are observable as the fisher-information matrix has to be invertible in order to evaluate the information gain [12].

Inspired by the Lissajous curve equations, each NBT is parameterised as:

$$x = w_{\text{traj}} \sin(at + \delta) + 0.5 w_{\text{calib}},$$
$$y = h_{\text{traj}} \cos(bt) + 0.5 h_{\text{calib}}, \quad (14)$$
$$z = \sqrt{d_{\text{nbt}} - x^2 - y^2},$$

where $x$, $y$ and $z$ are the trajectory positions relative to the fiducial target frame $\mathcal{F}_F$ to form $_F\mathbf{r}_{FS}$, $d_{\text{nbt}}$ is the distance away from the fiducial target center, $w_{\text{traj}}$ and $h_{\text{traj}}$ are the trajectory max width and height, $w_{\text{calib}}$ and $h_{\text{calib}}$ are the fiducial target width and height, $\delta$ represents the phase angle offset, and finally $a$ and $b$ are constants that determine the shape of the trajectory (e.g. a ratio of $\frac{a}{b} = 2$ forms a figure of 8). Finally, the sensor's orientation are parameterised as Euler angles and designed such that it is always pointing towards the center of the calibration target:

$$\phi = \phi_{\text{bound}} \sin(2\pi t) + \pi,$$
$$\theta = \theta_{\text{bound}} \sin(2\pi t), \quad (15)$$
$$\psi = 0.0,$$

where $\phi$, $\theta$ and $\psi$ are Euler angles around the x, y and z-axis to form $\mathbf{C}_{FS}$, respectively, and $\phi_{\text{bound}}$ and $\theta_{\text{bound}}$ are the maximum rotation around x and y-axis, respectively.

To ensure the velocity and angular velocity are realistic, we parameterise $t$ in Eq. (14) and Eq. (15) as a function of $k$ between $[0, t_{\text{nbt}}]$ such that the first derivative of both equations, velocity and angular velocity, start and end at 0, $t(k) = \sin^2(\pi k/2 \, t_{\text{nbt}})$, where $t_{\text{nbt}}$ is the time to complete a NBT. Differentiating both Eq. (14) and Eq. (15) enables us to simulate the camera and IMU measurements for evaluating NBTs using Eq. (4).

## E. Camera-IMU Calibration With Next-Best-Trajectory

The Camera-IMU calibration begins with two separate processes running in parallel, a real-time VI estimator solving Eq. (11) and a NBT evaluator (see Fig. 5). The real-time

(a) Vertical Figure-8    (b) Horizontal Figure-8

(c) Horizontal Pan    (d) Vertical Pan
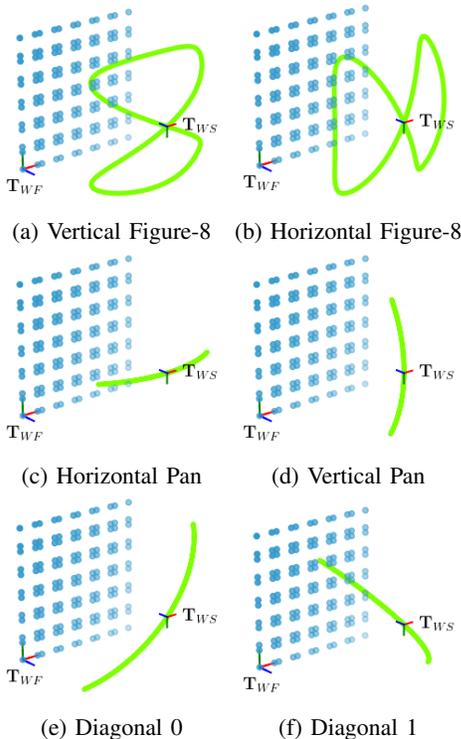
(e) Diagonal 0    (f) Diagonal 1

Fig. 6: Next-Best-Trajectory (NBT) candidates in-front of a calibration target

VI camera parameters are initialised using the parameters optimised in Section. VI and are fixed throughout. The fiducial pose, $\mathbf{T}_{WF}$ and camera-IMU extrinsics, $\mathbf{T}_{SC_1}$ on the other hand are initialised by solving Eq. (11) with the first $N$ camera frames, and IMU measurements between the first and last camera frame timestamps.

As the real-time VI estimator is solving the camera-IMU calibration problem, it periodically sends the calibration problem data to the NBT evaluator process. The NBT evaluator in turn would use the data to evaluate the MI of a set of pre-defined NBTs (Section. VII-D) using Eq. (4), find the NBT with the highest MI and guide the user in executing the NBT in order to render the calibration parameters optimally observable, i.e. reducing the expected uncertainty on the estimated camera-IMU extrinsics. If none of the candidate NBTs satisfies $I_{\text{mutual}} > I_{\text{threshold}}$ ($I_{\text{threshold}} = 0.2$) then the NBT evaluator sends a "finish" message to communicate to the real-time VI process that it should proceed to perform a final full-batch calibration.

## VIII. EXPERIMENTS

To evaluate our method, we conducted two sets of experiments. First, we evaluated our calibration pipeline in offline mode with the EuRoC [21] dataset to verify our calibration *accuracy* is competitive against that of Kalibr's *without the interactive component* of our system, and despite different approaches to solving the camera-IMU calibration problem, where Kalibr uses a continuous time full-batch optimisation

in contrast to our method which uses a discrete time full-batch optimisation. This is *independent* of our contributions regarding interactivity. With this we wanted to highlight our calibration tool without interactivity is at least as good as Kalibr.

Since the main motivation in this work is to provide non-experts with good calibration results for VIO/VI-SLAM systems, we further conducted experiments involving a small batch of graduate students to prove that our system can *efficiently* and *reliably* calibrate VI sensors, achieving superior performance for existing VIO and VI-SLAM systems. To compare the calibrations we used them in ORBSLAM3 [6] and evaluated the accuracy using the evaluation scheme of [22] with RMSE Absolute Trajectory Error (ATE) by aligning the estimated trajectory with the ground-truth.

All experiments were conducted on a Lenovo P52 Thinkpad laptop containing an Intel Core i7-8750H CPU at 2.2 Ghz with 16GB of memory running Ubuntu 20.04 and ROS Melodic. The experiments with the graduate students were conducted with the aim of calibrating an Intel RealSense D435i which contains a stereo IR global shutter depth sensor, a monocular RGB rolling shutter sensor, and additionally an IMU sensor running at 15Hz, 15Hz and 400Hz respectively. For our purposes, we *do not* use the RGB rolling shutter sensor. We have instead disabled the IR projector and used the stereo IR global shutter depth sensors as a standard gray-scale stereo camera.

During the camera and VI calibrations the default settings for Kalibr [1] were used to generate their results, whereas in our method we used Cauchy loss ($s = 1.5$) on the reprojection errors and a fixed-lag smoothing window size of 10 and 3 for the camera calibration and camera-IMU calibration stages respectively. The IMU parameters used for the camera-IMU calibration are: $\sigma_a = 2.52 \times 10^{-2} \frac{\text{m}}{\text{s}^2} \frac{1}{\sqrt{\text{Hz}}}$ for the accelerometer noise density, $\sigma_{ba} = 4.41 \times 10^{-3} \frac{\text{m}}{\text{s}^3} \frac{1}{\sqrt{\text{Hz}}}$ for the accelerometer drift noise density, $\sigma_g = 2.78 \times 10^{-3} \frac{\text{rad}}{\text{s}} \frac{1}{\sqrt{\text{Hz}}}$ for gyroscope noise density, and $\sigma_{bg} = 1.65 \times 10^{-5} \frac{\text{rad}}{\text{s}^2} \frac{1}{\sqrt{\text{Hz}}}$ gyroscope drift noise density.

### A. Calibration Results on EuRoC Dataset

To assess our approach in offline mode, we used the calibration sequences from the EuRoC dataset [21] to calibrate the VI-sensor. The calibration process is split into two stages. First, the camera intrinsics and camera extrinsics are estimated. Then in the second stage, only the camera-IMU extrinsics and time-delay are estimated, with the camera intrinsics and camera extrinsics estimated in the first phase fixed.

The results show comparable calibration reprojection errors, where in the camera calibration stage our method obtained an RMSE reprojection error of $0.6042$ pixels compared to Kalibr's $0.6087$ pixels, and in the camera-IMU calibration stage the RMSE reprojection errors are $0.5569$ pixels and $0.5775$ pixels for our method and Kalibr respectively. Fig. 7 and Fig. 8 report RMSE ATE after running ORBSLAM3 [6] on the EuRoC dataset sequences 10 times

in Stereo-VO mode and Stereo-VIO mode, respectively. We did not change the ORB-SLAM3 EuRoC configuration that was orginally tuned for Kalibr calibration parameters. Both figures show that the calibrations produced by our method yielded better results on most sequences in Stereo-VO mode and all sequences in VIO mode, compared to Kalibr.

To verify our camera-IMU time delay estimation, we assumed the EuRoC dataset [21] has a camera-IMU time delay of $\approx 7\mu s$, as reported in [23], and perturbed the `imu_april` IMU timestamps with 100ms, 10ms and 1ms time offsets. With our offline camera-IMU calibration *without interactivity* we were able to recover the time offsets $100ms$, $9.97ms$ and $0.987ms$ respectively, thus showing that our offline camera-IMU calibrator is capable of accurately estimating the camera-IMU time delay.



Fig. 7: Comparison of ORBSLAM3 using calibrations from Kalibr and ours in Stereo-VO mode on EuRoC Dataset

### B. Trials with Graduate Students

To evaluate our calibration method, we conducted a series of tests involving 16 graduate students to measure the effectiveness of our approach compared to the state-of-the-art calibrator, Kalibr [1]. Our test-subjects were postgraduate students at Imperial College London. Of the 16 students, 4 reported some previous experience with camera calibration, and only 2 reported some previous experience with camera-IMU calibration. Each participant was asked to calibrate the same Intel RealSense D435i sensor by first collecting two calibration sequences for Kalibr (one for camera calibration and the second for camera-IMU calibration), and then another two with our calibration method.

Because we do not have ground truth for the calibration parameters, we evaluated the estimated calibration parameters by applying them in ORBSLAM3 [6] running in odometry mode (with loop-closure switched off) on 10 custom-collected Vicon room sequences where ground-truth poses were recorded with various motions.
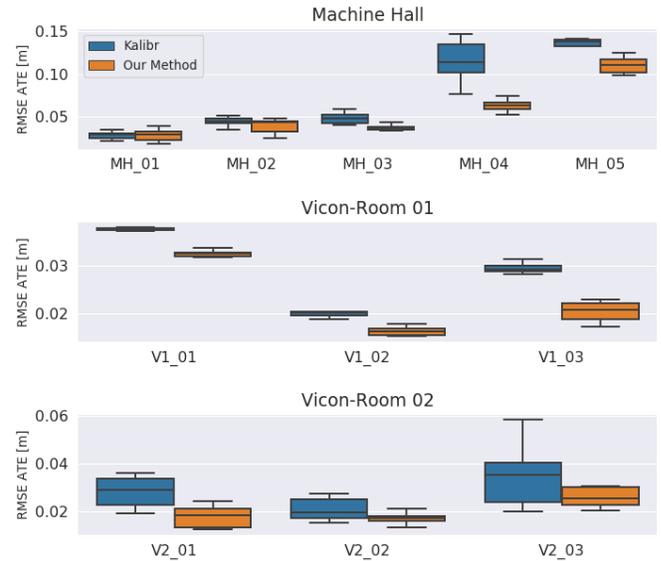


Fig. 8: Comparison of ORBSLAM3 using calibrations from Kalibr and ours in Stereo-VIO mode on EuRoC Dataset
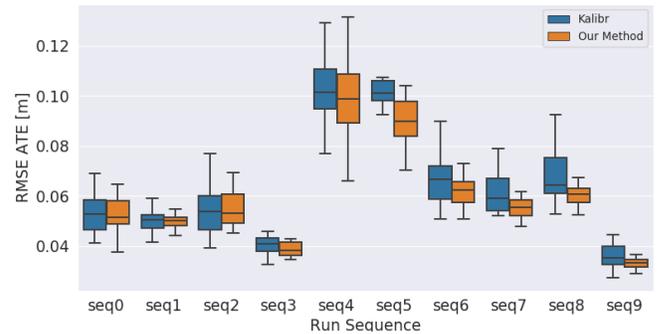


Fig. 9: Comparing calibrations by graduate students across 10 different evaluation VICON room sequences by running ORBSLAM3 in odometry mode

Our study shows that novices who have little to no experience in calibrating a VI sensor can obtain better calibrations using our approach compared to Kalibr. Out of 10 Vicon room sequences, the RMSE ATE error is lowest across all sequences using calibration parameters obtained through our method (see Fig. 9). Our calibration parameters also yielded overall smaller RMSE ATE variances, showing more consistent and reliable odometry accuracy, regardless of the experience of calibration users. The estimated camera-IMU time delay with our method is $3.07 \pm 0.932ms$, and Kalibr's estimate is $4.81 \pm 0.981ms$. Since ground-truth is not available we can only conclude our method is more consistent compared to Kalibr's result. The break-down of the median total time taken to calibrate the VI sensor between Kalibr and our method is shown in Fig. 10, where our method's median is 381.11 seconds compared to Kalibr's 455.44 seconds.

In addition to showing that our method yields better SLAM results and calibrations faster, by inspecting the Shannon Entropy of the calibration parameters, a metric
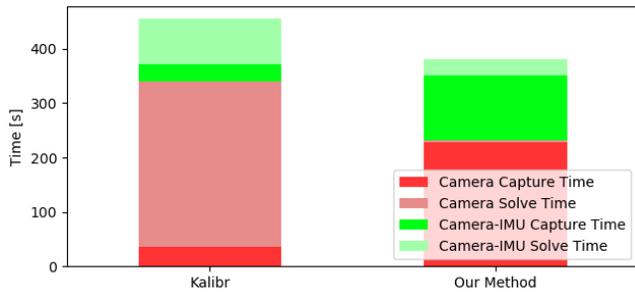
Fig. 10: Median total calibration time per person in seconds



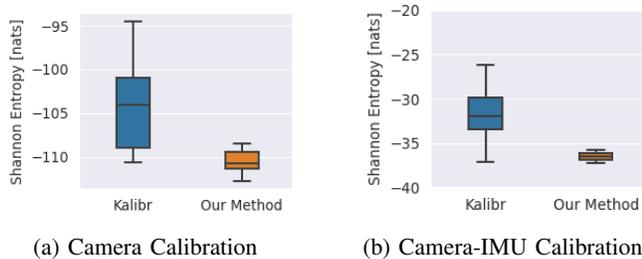(a) Camera Calibration  (b) Camera-IMU Calibration

Fig. 11: Shannon Entropy of camera and camera-IMU calibration across all calibrations by graduate students

used to measure the uncertainty of information content [24], we also observe a lower entropy (more certainty) with our method compared to Kalibr (see Fig. 11a and Fig. 11b). This means that our method can successfully guide a novice to collect a more informative calibration dataset for a good calibration.

## IX. CONCLUSIONS

The success of SoTA computer-vision and state-estimation algorithms often hinges on good VI calibrations. However, collecting high-quality VI calibration data is not trivial, especially since most existing calibration tools do not provide an interactive live feedback to the user which ultimately increases the risk of poor calibrations. In this work, we have introduced a novel visual-inertial calibration guidance system to provide real-time NBV and NBT suggestions to guide users in collecting informative calibration data. It achieves competitive calibration results against the SoTA offline calibrator, Kalibr [1], and produces faster, more accurate and more reliable calibrations for existing SoTA visual and VI SLAM systems, even when used by novices.

## ACKNOWLEDGMENT

## REFERENCES

[1] Paul Furgale, Joern Rehder, and Roland Siegwart. Unified temporal and spatial calibration for multi-sensor systems. In *IROS*, 2013.

[2] Yunke Ao, Le Chen, Florian Tschopp, Michel Breyer, Roland Siegwart, and Andrei Cramariuc. Unified data collection for visual-inertial calibration via deep reinforcement learning. In *ICRA*, 2022.

[3] Stefan Leutenegger, Simon Lynen, Michael Bosse, Roland Siegwart, and Paul Furgale. Keyframe-based visual-inertial odometry using non-linear optimization. *The International Journal of Robotics Research*, 34(3):314–334, 2014.

[4] Tong Qin, Peiliang Li, and Shaojie Shen. VINS-Mono: A robust and versatile monocular visual-inertial state estimator. *T-RO*, 34(4):1004–1020, 2018.

[5] Patrick Geneva, Kevin Eckenhoff, Woosik Lee, Yulin Yang, and Guoquan Huang. OpenVINS: A research platform for visual-inertial estimation. In *ICRA*, 2020.

[6] Carlos Campos, Richard Elvira, Juan J. Gómez, José M. M. Montiel, and Juan D. Tardós. ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM. *T-RO*, 37(6):1874–1890, 2021.

[7] Joao Alves, Jorge Lobo, and Jorge Dias. Camera-inertial sensor modelling and alignment for visual navigation. *Machine Intelligence and Robotic Control*, 5(3):103–112, 2003.

[8] Jorge Lobo and Jorge Dias. Relative pose calibration between visual and inertial sensors. *IJRR*, 26(6):561–575, 2007.

[9] F.M. Mirzaei and S.I. Roumeliotis. A Kalman Filter-Based Algorithm for IMU-Camera Calibration: Observability Analysis and Performance Evaluation. *T-RO*, 24(5):1143–1156, 2008.

[10] Tue-Cuong Dong-Si and Anastasios I Mourikis. Estimator initialization in vision-aided inertial navigation with unknown camera-IMU calibration. In *IROS*, 2012.

[11] Le Chen, Yunke Ao, Florian Tschopp, Andrei Cramariuc, Michel Breyer, Jen Jen Chung, Roland Siegwart, and Cesar Cadena. Learning trajectories for visual-inertial system calibration via model-based heuristic deep reinforcement learning. In *CoRL*, 2021.

[12] Fernando Nobre and Christoffer Heckman. Learning to calibrate: Reinforcement learning for guided calibration of visual–inertial rigs. *IJRR*, 38(12-13):1388–1402, 2019.

[13] Andrew Richardson, Johannes Strom, and Edwin Olson. AprilCal: Assisted and repeatable camera calibration. In *IROS*, 2013.

[14] Thomas Schneider, Mingyang Li, Michael Burri, Juan Nieto, Roland Siegwart, and Igor Gilitschenski. Visual-inertial self-calibration on informative motion segments. In *ICRA*, 2017.

[15] Thomas Schneider, Mingyang Li, Cesar Cadena, Juan Nieto, and Roland Siegwart. Observability-aware self-calibration of visual and inertial sensors for ego-motion estimation. *IEEE Sens. J*, 19(10):3846–3860, 2019.

[16] Michael Bloesch, Hannes Sommer, Tristan Laidlow, Michael Burri, Gabriel Nützi, Peter Fankhauser, Dario Bellicoso, Christian Gehring, Stefan Leutenegger, Marco Hutter, and Roland Siegwart. A Primer on the Differential Calculus of 3D Orientations. *CoRR*, abs/1606.0, 2016.

[17] J. Solà, J. Deray, and D. Atchuthan. A micro Lie theory for state estimation in robotics. *arXiv:1812.01537*, 2018.

[18] Jérôme Maye, Paul Furgale, and Roland Siegwart. Self-supervised calibration for robotic systems. In *IV*, 2013.

[19] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza. IMU preintegration on manifold for efficient visual-inertial maximum-a-posteriori estimation. In *RSS*, 2015.

[20] Yulin Yang, Patrick Geneva, Kevin Eckenhoff, and Guoquan Huang. Degenerate motion analysis for aided VINS with online spatial and temporal sensor calibration. *RAL*, 4(2):2070–2077, 2019.

[21] Michael Burri, Janosch Nikolic, Pascal Gohl, Thomas Schneider, Joern Rehder, Sammy Omari, Markus W Achtelik, and Roland Siegwart. The EuRoC Micro Aerial Vehicle Datasets. *IJRR*, 35(10):1157–1163, September 2016.

[22] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A Benchmark for the Evaluation of RGB-D SLAM Systems. In *IROS*, 2012.

[23] Janosch Nikolic, Joern Rehder, Michael Burri, Pascal Gohl, Stefan Leutenegger, Paul T Furgale, and Roland Siegwart. A synchronized visual-inertial sensor system with FPGA pre-processing for accurate real-time SLAM. In *ICRA*, pages 431–437, 2014.

[24] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley-Interscience, second edition, 2006.