

# Stereo-Inertial Pose Estimation and Online Sensors Extrinsic Calibration

Fumin Pang and Tianmiao Wang, *Member, IEEE*

**Abstract**—The fusion of visual and inertial measurement has been popular in mobile robotics community for decades due to the complementary properties of two sensors. The combination of these two sensors offers rich texture of environment and accurate short-time motion prediction, making it particularly suitable for pose estimation, especially in GPS-denied unknown environment. In this paper, we propose a method which fuses stereo visual and inertial cues based on Multi-State Constraint Kalman Filter (MSCKF), to estimate 6DOF pose of mobile robot. Stereo vision offers real-metric perception of surroundings, giving a better initial scale estimation of the visual-inertial system. Compared with another class of methods based on batch nonlinear optimization, this filter-based method is more suitable for resource-constrained mobile platforms. In addition, the finite precision of sensors extrinsic calibration often makes estimator inconsistent. The proposed method includes extrinsic parameters in state vector to do online calibration. Experimental results on real-world datasets demonstrate that proposed method is consistent and substantially improves the accuracy of pose estimation, as well as the calibration between sensors.

## I. INTRODUCTION

Accurate 6DOF pose estimation in unknown environment is a critical component in autonomous navigation task. It is also one of the most popular topics in robotic community. Visual and inertial sensors have been more and more widely-applied to estimate the pose of mobile robots due to their complementary properties, low cost and safety. Visual measurement offers rich texture of surroundings. Each visual feature can always be tracked by a camera from a sequence of consecutive poses, which provide multiple constraint of camera motion. Inertial measurement can provide accurate motion prediction within a short time at a higher frequency. In Simultaneous localization and mapping (SLAM), fusion of visual and inertial measurement is used to estimate robot pose and build a map of environment [1] [2]. Another pose estimation paradigm is Visual-Inertial Odometry (VIO), which merely estimate robot motion typically.

In this paper, we are interested in visual-inertial odometry. The most common VIO methods can be classified into two categories: recursive filter-based and batch-based nonlinear optimization approaches. Filter-based approaches are the earlier ones used to solve SLAM and VIO problems. In this class of methods, inertial sensors are regarded as interoceptive sensors to provide linear accelerator and angular velocity, which are always integrated to form prior position and orientation estimation. Meanwhile, the visual sensors, as exteroceptive sensors at a lower frequency, provide visual

measurement to calculate the innovation of the filter. Thus, posterior estimation is achieved recursively. Batch based methods are inspired by Structure From Motion (SFM) [3] and come into sight later, but achieve some impressing results in recent years [4] [5] [6]. These methods promise results of higher accuracy compared with filtering approaches as they employ re-linearization at each iteration to get a better deal with the nonlinearity of the measurement models. They always implement bundle adjustment in a sliding window of states, using multiple iterations to minimize cost function, which results in increased computational cost, however. Thus, for a long-time running system, the lack of computational resources makes recursive algorithms a favorable choice for online estimation, especially in resource-constrained systems, such as micro aerial vehicles, mobile phones, and Augmented Reality (AR) devices [7]. Therefore, in this paper, we focus on the filter-based method.

Filter-based methods can also be divided into two main categories: loosely-coupled methods and tightly-coupled methods. Loosely-coupled ones are the most intuitive method for fusing visual and inertial measurement. This kind of methods deals with visual and inertial measurement separately. Visual pipeline, always treated as a black box, estimate poses independently. These poses are regard as the measurement of EKF and then fused with the inertial measurement recursively [8]. Weiss and Siegwart use a monocular SLAM framework PTAM to estimate 6DOF poses of a micro aerial vehicle mounted with a down-looking camera [9] [10]. Then, the poses are fused with the prior state vector driven by inertial reading. Lynen and Achtelik implement a EKF-based framework for dealing with loosely-coupled multiple sensor Fusion [11]. This kind of methods reduces computational complexity, but it will ignore some crucial information. For instance, processing IMU measurements separately does not allow for optimal estimation of sensor biases, and using feature measurements for motion estimation between pairs of images ignores the correlations between consecutive timesteps [12].

Tightly-coupled methods are ones directly fuse the visual and inertial data for achieving higher precision. The most common EKF estimator to tightly fuse visual and inertial data is EKF-based SLAM, in which current camera pose and features positions are jointly included in estimator state vector [13] [14]. EKF-based SLAM suffers from the cubic computational complexity of feature number in state vector. To overcome this limitation, Mourikis and Roumeliotis come up with multi-state constraint Kalman filter (MSCKF) [12]. Instead of feature positions, a sliding window of camera

poses is included in MSCKF estimator. By directly making use of the geometric constraints between multiple camera poses, it avoids the computational burden and loss of information associated with pairwise displacement estimation. To improve the consistency of the MSCKF, Li and Mourikis include Camera-IMU extrinsic calibration parameter in estimator state vector, giving a more accurate result [15].

In this paper, we propose a stereo visual inertial navigation system. As we know, almost all MSCKF-based implementations are based on monocular camera-IMU rig. Stereo visual measurement can give a real metric perception of environment, which results in a better initialization of the estimator. A key contribution of proposed method is we give the derivation of stereo visual measurement model, which is different from the existing monocular one. In addition, to improve the performance, the system includes online extrinsic calibration between IMU and two cameras. We test the proposed method on real-world dataset. The experimental results show that the proposed method is consistent, and that it attains substantially higher accuracy than monocular MSCKF. And it gives a effective estimation about extrinsic parameters.

## II. ESTIMATOR DESCRIPTION

As illustrated in Fig.1, we affix stereo camera-IMU rig body frame  $\{B\}$  to IMU, to track the 6D motion with respect to a global coordinate frame,  $\{G\}$ . Two camera coordinate frames, CAM0 and CAM1, are  $\{C_0\}$  and  $\{C_1\}$ .

### A. MSCKF State Parametrization

The full MSCKF state representation can be partitioned into three parts. The first is the evolving current body state. As we affix body frame  $\{B\}$  to IMU, we use  $\mathbf{x}_I$  to present body state. We parametrize the body state at time  $k$  as the 16-dimensional vector.

$$\mathbf{x}_{I,k} := \begin{bmatrix} {}^B_B\tilde{\mathbf{q}}_k^T & {}^G\mathbf{p}_{B,k}^T & {}^G\mathbf{v}_{B,k}^T & \mathbf{b}_{g,k}^T & \mathbf{b}_{a,k}^T \end{bmatrix}^T \quad (1)$$

Where  ${}^B_B\tilde{\mathbf{q}}_k$  is the unit quaternion representing the rotation which rotate vectors from the global frame  $\{G\}$  to the body frame  $\{B\}$ . In this paper, all quaternions follow JPL convention [16].  ${}^G\mathbf{p}_{B,k}$  is the vector from the origin of  $\{G\}$  to the origin of  $\{B\}$  expressed in  $\{G\}$  (i.e., the position of body in the global frame).  ${}^G\mathbf{v}_{B,k}$  is the vector representing original velocity of frame  $\{B\}$  expressed in  $\{G\}$ .  $\mathbf{b}_{g,k}$  is the bias on the gyro measurements  $\boldsymbol{\omega}_m$ ,  $\mathbf{b}_{a,k}$  is the bias on the velocity measurements  $\mathbf{a}_m$ .

The second state part is extrinsic calibration parameters (rotation and translation) between IMU and two cameras. We parametrize it as a 14-dimensional vector.

$$\mathbf{x}_{Calib,k} := \begin{bmatrix} {}^{C_0}_B\tilde{\mathbf{q}}_k^T & {}^B\mathbf{p}_{C_0,k}^T & {}^{C_1}_B\tilde{\mathbf{q}}_k^T & {}^B\mathbf{p}_{C_1,k}^T \end{bmatrix}^T \quad (2)$$

Where  ${}^{C_i}_B\tilde{\mathbf{q}}_k$  rotate vectors from body frame  $\{B\}$  to camera frame  $\{C_i\}$ .  ${}^B\mathbf{p}_{C_i,k}$  is the original point position

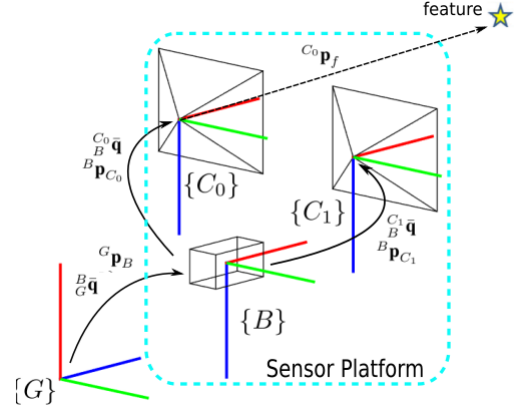


Fig. 1. Coordinate frames involved in sensor platform and visual feature.

of camera frame  $\{C_i\}$  expressed in body frame  $\{B\}$ ,  $i = 0, 1$ .

The third part of the full state is a sliding window of  $N$  past body states, in which active feature tracks were visible. The  $i$ -th body state,  $i = 0 \dots N - 1$ , including 6D pose and velocity, is a 10-dimensional vector.

$$\mathbf{x}_{B_i,k} := \begin{bmatrix} {}^{C_0}_B\tilde{\mathbf{q}}_k^T & {}^B\mathbf{p}_{B_i,k}^T & {}^G\mathbf{v}_{B_i,k}^T \end{bmatrix}^T \quad (3)$$

At time  $k$ , the full state of MSCKF is a  $(30 + 10 * N)$  vector consisting of current body state estimate, calibration parameters and  $N$  last body states.

$$\hat{\mathbf{x}}_k := \begin{bmatrix} \hat{\mathbf{x}}_{I,k}^T & \hat{\mathbf{x}}_{Calib,k}^T & \hat{\mathbf{x}}_{B_0,k}^T & \dots & \hat{\mathbf{x}}_{B_{N-1},k}^T \end{bmatrix}^T \quad (4)$$

We can define *error state* of the estimator at time  $k$ :

$$\tilde{\mathbf{x}}_k := \begin{bmatrix} \tilde{\mathbf{x}}_{I,k}^T & \tilde{\mathbf{x}}_{Calib,k}^T & \tilde{\mathbf{x}}_{B_0,k}^T & \dots & \tilde{\mathbf{x}}_{B_{N-1},k}^T \end{bmatrix}^T \quad (5)$$

where

$$\tilde{\mathbf{x}}_{I,k} := \begin{bmatrix} {}^G\delta\boldsymbol{\theta}_I^T & {}^G\tilde{\mathbf{p}}_{B,k}^T & {}^G\tilde{\mathbf{v}}_{B,k}^T & \tilde{\mathbf{b}}_{g,k}^T & \tilde{\mathbf{b}}_{a,k}^T \end{bmatrix}^T \quad (6)$$

$$\tilde{\mathbf{x}}_{B_i,k} := \begin{bmatrix} {}^G\delta\boldsymbol{\theta}_I^T & {}^G\tilde{\mathbf{p}}_{B_i,k}^T & {}^G\tilde{\mathbf{v}}_{B_i,k}^T \end{bmatrix}^T \quad (7)$$

The full error state has  $(27 + 9N)$  dimensions. In the above,  $\tilde{x}$  is the difference between the true value  $x$  and the estimated value  $\bar{x}$ . For position, velocity and bias, it is defined as  $\tilde{x} = x - \bar{x}$ . The quaternion error is defined according to

$$\delta\tilde{\mathbf{q}} := \hat{\mathbf{q}}^{-1} \otimes \tilde{\mathbf{q}} \approx \begin{bmatrix} \frac{1}{2}\delta\boldsymbol{\theta}^T & 1 \end{bmatrix}^T \quad (8)$$

Accordingly, the MSCKF error state covariance  $\mathbf{P}$  is a  $(27 + 9N) \times (27 + 9N)$  matrix. It can be partitioned into 4 blocks:

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{I,I} & \mathbf{P}_{I,Cali-B} \\ \mathbf{P}_{I,Cali-B}^T & \mathbf{P}_{Cali-B,Cali-B} \end{bmatrix} \quad (9)$$

where  $\mathbf{P}_{I,I}$  is the  $15 \times 15$  covariance matrix of current body state.  $\mathbf{P}_{Cali-B,Cali-B}$  is the combined  $(12 + 9N) \times (12 + 9N)$  covariance matrix of calibration state and past body states.  $\mathbf{P}_{I,Cali-B}$  is the cross-covariance between the current body state and combination of calibration and past body states.

### B. Filter Propagation

The IMU measurement is used to propagate the state estimates. As mentioned before, IMU measurement provide rotational velocity  $\omega_m$  and  $\mathbf{a}_m$ , described as below equations:

$$\omega_m = {}^G\omega + \mathbf{b}_g + \mathbf{n}_g \quad (10)$$

$$\mathbf{a}_m = {}^B\mathbf{R}({}^G\mathbf{a} - {}^G\mathbf{g}) + \mathbf{b}_a + \mathbf{n}_a \quad (11)$$

where  ${}^B\mathbf{R}$  is the rotation matrix cooresponding to  ${}^I_G\mathbf{q}$ .  ${}^G\mathbf{g}$  is the gravitational acceleration.  $\mathbf{n}_g$  and  $\mathbf{n}_a$  are zero-mean white Gaussian noise vectors. Using these measurements, we can write the dynamics of the filter state vector as:

$${}^I_G\dot{\mathbf{q}}(t) = \frac{1}{2}\Omega(\omega_m(t) - \mathbf{b}_g(t) - \mathbf{n}_g(t)) {}^I_G\bar{\mathbf{q}} \quad (12)$$

$${}^G\dot{\mathbf{v}}_B(t) = {}^B\mathbf{R}^T(t)(\mathbf{a}_m(t) - \mathbf{b}_a(t) - \mathbf{n}_a(t)) + {}^G\mathbf{g} \quad (13)$$

$${}^G\dot{\mathbf{p}}_B(t) = {}^G\mathbf{v}_B(t) \quad (14)$$

$$\dot{\mathbf{b}}_g(t) = \mathbf{n}_{wg}(t) \quad \dot{\mathbf{b}}_a(t) = \mathbf{n}_{wa}(t) \quad (15)$$

$${}^{C_i}_B\dot{\mathbf{q}}(t) = \mathbf{0} \quad {}^B\dot{\mathbf{p}}_{C_i}(t) = \mathbf{0} \quad (16)$$

where  $\Omega$  is the quaternion multiplication matrix corresponding to the angular velocity vector  $\omega$ . In the above, the first three equations describe the dynamics of the current body motion. The fourth line models the biases as random walk processes. The fifth line describes the fact that camera-IMU transformations do not change in time. As for the past body states, they remain constant once the corresponding image is captured.

Equation (12)-(16) describe the the continuous-time evolution of the true states. We follow the approach described in (High-precision, consistent EKF-based visual-inertial odometry) for propagating the state estimates in a discrete-time implementation. Further, a fifth-order Runge-Kutta procedure is used for propagate quaternion.

Besides the current body position, velocity, and orientation, all other state estimates remain unchanged during propagation. We can also examine the linearized continuous-time model of the current state error state,  $\tilde{\mathbf{x}}_I$ :

$$\dot{\tilde{\mathbf{x}}}_I = \mathbf{F}\tilde{\mathbf{x}}_I + \mathbf{G}\mathbf{n}_I \quad (17)$$

where  $\mathbf{F}$  is the continuous-time error-state transition matrix.  $\mathbf{G}$  is given by

$$\mathbf{G} = \begin{bmatrix} -\mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & -{}^B\mathbf{R}^T & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (18)$$

and  $\mathbf{n}_I = [\mathbf{n}_g^T \ \mathbf{n}_a^T \ \mathbf{n}_{wg}^T \ \mathbf{n}_{wa}^T]^T$  is the IMU process noise, which has covariance matrix  $\mathbf{Q}_I$ .

In addition to the state estimate, the MSCKF propagates the state covariance matrix, as follows:

$$\mathbf{P}(t_{k+1}) = \begin{bmatrix} \mathbf{P}_{I,I}(t_{k+1}) & \Phi_{I_k} \mathbf{P}_{I,Cali-B}(t_{k+1}) \\ \mathbf{P}_{I,Cali-B}^T(t_k) \Phi_{I_k}^T & \mathbf{P}_{Cali-B,Cali-B}(t_k) \end{bmatrix} \quad (19)$$

where

$$\mathbf{P}_{I,I}(t_{k+1}) = \Phi_{I_k} \mathbf{P}_{I,I}(t_k) \Phi_{I_k}^T + \mathbf{G} \mathbf{Q}_I \mathbf{G}^T \Delta T \quad (20)$$

and  $\Delta T$  is the IMU sample period.  $\Phi_{I_k}$  is the discrete-time error-state transition matrix at  $k$ . Instead of using the most common approximation  $\Phi = \mathbf{I} + \mathbf{F} \Delta T$ , we use the closed-form matrix by [15].

### C. State Augmentation

Upon recording a new image, a copy of current body rotation, position and velocity is appended to the state vector, and the covariance matrix of the MSCKF is augmented accordingly:

$$\mathbf{P}(t_k) \leftarrow \begin{bmatrix} \mathbf{I}_{27+9N} \\ \mathbf{J} \end{bmatrix} \mathbf{P}(t_k) \begin{bmatrix} \mathbf{I}_{27+9N} \\ \mathbf{J} \end{bmatrix}^T \quad (21)$$

where  $\mathbf{J}$  is given by

$$\mathbf{J} = \begin{bmatrix} \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_{3 \times (18+9N)} \\ \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_{3 \times (18+9N)} \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_{3 \times (18+9N)} \end{bmatrix} \quad (22)$$

### D. Filter Update

When a tracked feature  $\mathbf{f}_j$  is lost in current frame, it is selected for state update, the MSCKF estimates its position  ${}^G\hat{\mathbf{p}}_f$  using an inverse depth [17] least-squares Gauss-Newton optimization [18]. The procedure takes as input  $N$  pair of camera pose and  $N$  sets of feature measurements. As for stereo visual system, each past body pose has two corresponding camera poses, and the two camera poses can be calculated by

$${}^{C_i}_G\hat{\mathbf{q}} = {}^{C_i}_B\hat{\mathbf{q}} \otimes {}^B\hat{\mathbf{q}} \quad {}^G\hat{\mathbf{p}}_{C_i} = {}^G\hat{\mathbf{p}}_B + {}^B\hat{\mathbf{R}}^T {}^B\hat{\mathbf{p}}_{C_i} \quad (23)$$

Now that we have estimated the positions of any features which can be used in the state update, we can apply the corresponding motion constraints to the window of poses from which each feature was tracked. To present the measurement model, we consider one feature,  $\mathbf{f}_j$ , which is observed by a set of pose pairs,  $\mathcal{M}_j$ .  $z_{l,0}$  and  $z_{l,1}$  are observations of feature  $\mathbf{f}_j$  from stereo camera poses  $\mathbf{C}_{l,0}$  and  $\mathbf{C}_{l,1}$ ,  $l \in \mathcal{M}_j$ . The  $4 \times 1$  stereo visual measurement residual can be computed by

$$\mathbf{r}_l^{(j)} = \begin{bmatrix} z_{l,0}^{(j)} \\ z_{l,1}^{(j)} \end{bmatrix} - \begin{bmatrix} \hat{z}_{l,0}^{(j)} \\ \hat{z}_{l,1}^{(j)} \end{bmatrix} \quad (24)$$

where

$$\hat{z}_{l,i}^{(j)} = \frac{1}{{}^{C_{l,i}}\hat{\mathbf{Z}}^{(j)}} \left[ {}^{C_{l,i}}\hat{\mathbf{X}}^{(j)} \quad {}^{C_{l,i}}\hat{\mathbf{Y}}^{(j)} \right]^T \quad (25)$$

with

$$\begin{aligned} {}^{C_{l,i}}\hat{\mathbf{p}}_{f_j} &= \left[ {}^{C_{l,i}}\hat{\mathbf{Z}}^{(j)} \quad {}^{C_{l,i}}\hat{\mathbf{X}}^{(j)} \quad {}^{C_{l,i}}\hat{\mathbf{Y}}^{(j)} \right]^T \\ &= {}^{C_{l,i}}\hat{\mathbf{R}} \left( {}^G\hat{\mathbf{p}}_{f_j} - {}^G\hat{\mathbf{p}}_{C_{l,i}} \right) \end{aligned} \quad (26)$$

Linearizing about the estimates for the body pose and for the feature position, the residual of Eq. (24) can be approximated as:

$$\mathbf{r}_l^{(j)} = \mathbf{H}_{x_{B_l}}^{(j)} \tilde{\mathbf{x}} + \mathbf{H}_{f_l}^{(j)G} \tilde{\mathbf{p}}_{f_j} + \mathbf{n}_l^{(j)} \quad (27)$$

where the matrices  $\mathbf{H}_{x_{B_l}}^{(j)}$  and  $\mathbf{H}_{f_l}^{(j)}$  are the corresponding Jacobians of the measurement  $\hat{\mathbf{z}}_l^{(j)}$  with respect to the state and the feature position, respectively, and  ${}^G\tilde{\mathbf{p}}_{f_j}$  is the error in the position estimate of  $\mathbf{f}_j$ .  $\mathbf{n}_l^{(j)}$  is the  $4 \times 1$  visual measurement noise with the corresponding covariance matrix  $\mathbf{R}_l^{(j)} = \sigma_{im}^2 \mathbf{I}_4$ .  $\mathbf{H}_{x_{B_l}}^{(j)}$  is given by

$$\mathbf{H}_{x_{B_l}}^{(j)} = [\mathbf{0}_{4 \times 15} \quad \mathbf{\Pi}_l \quad \mathbf{0}_{4 \times 9(l-1)} \quad \mathbf{H}_{B_l} \quad \mathbf{0}_{4 \times 9(N-l)}] \quad (28)$$

where

$$\mathbf{\Pi}_l = \begin{bmatrix} \mathbf{\Pi}_{\theta_{l,0}} & \mathbf{\Pi}_{p_{l,0}} & \mathbf{0}_{2 \times 3} & \mathbf{0}_{2 \times 3} \\ \mathbf{0}_{2 \times 3} & \mathbf{0}_{2 \times 3} & \mathbf{\Pi}_{\theta_{l,1}} & \mathbf{\Pi}_{p_{l,1}} \end{bmatrix} \quad (29)$$

$$\mathbf{H}_{B_l} = \begin{bmatrix} \mathbf{H}_{\theta_{l,0}} & \mathbf{H}_{p_{l,0}} & \mathbf{0}_{2 \times 3} \\ \mathbf{H}_{\theta_{l,1}} & \mathbf{H}_{p_{l,1}} & \mathbf{0}_{2 \times 3} \end{bmatrix} \quad (30)$$

and for more details, the nonzero blocks in matrixs are the Jacobians with respect to the two camera-IMU rotations, camera-IMU translations, body rotation and body position, respectively:

$$\mathbf{\Pi}_{\theta_{l,i}} = \mathbf{J}_l^{(j,i)C_i} \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G [\mathbf{R}^G \mathbf{p}_{f_j} - {}^G\tilde{\mathbf{p}}_B] \times \quad (31)$$

$$\mathbf{\Pi}_{p_{l,i}} = \mathbf{J}_l^{(j,i)} \quad (32)$$

$$\mathbf{H}_{\theta_{l,i}} = \mathbf{J}_l^{(j,i)C_i} \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G [\mathbf{R}^G \mathbf{p}_{f_j} - {}^G\tilde{\mathbf{p}}_B] \times \quad (33)$$

$$\mathbf{H}_{p_{l,i}} = -\mathbf{J}_l^{(j,i)C_i} \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G \quad (34)$$

where  $[c \times]$  is the skew symmetric matrix corresponding to vector  $c$ . And  $\mathbf{H}_{f_l}^{(j)}$  is given by

$$\mathbf{H}_{f_l}^{(j)} = \mathbf{J}_l^{(j,i)} \begin{bmatrix} C_0 \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G \\ C_1 \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G \\ C_2 \hat{\mathbf{R}}_B^B \hat{\mathbf{R}}_G^G \end{bmatrix} \quad (35)$$

In the above,  $\mathbf{J}_l^{(j,i)}$  is the Jacobian of the perspective model:

$$\mathbf{J}_l^{(j,i)} = \frac{1}{C_{l,i} \hat{\mathbf{Z}}^{(j)}} \begin{bmatrix} 1 & 0 & \frac{C_{l,i} \hat{\mathbf{X}}^{(j)}}{C_{l,i} \hat{\mathbf{Z}}^{(j)}} \\ 0 & 1 & \frac{C_{l,i} \hat{\mathbf{Y}}^{(j)}}{C_{l,i} \hat{\mathbf{Z}}^{(j)}} \end{bmatrix} \quad (36)$$

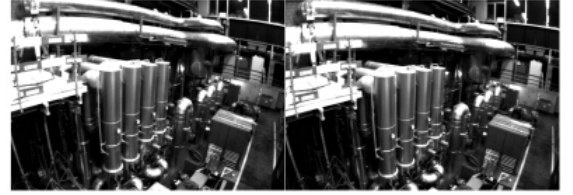
Stacking the residuals of all  $\mathcal{M}_j$  measurements of this feature, we can get:

$$\mathbf{r}^{(j)} = \mathbf{H}_{x_B}^{(j)} \tilde{\mathbf{x}} + \mathbf{H}_f^{(j)G} \tilde{\mathbf{p}}_{f_j} + \mathbf{n}^{(j)} \quad (37)$$

where  $\mathbf{r}^{(j)}$ ,  $\mathbf{H}_{x_B}^{(j)}$ ,  $\mathbf{H}_f^{(j)}$ , and  $\mathbf{n}^{(j)}$  are block vectors or matrices with elements  $\mathbf{r}_l^{(j)}$ ,  $\mathbf{H}_{x_{B_l}}^{(j)}$ ,  $\mathbf{H}_{f_l}^{(j)}$ , and  $\mathbf{n}_l^{(j)}$ , for  $l \in \mathcal{M}_j$ . And  $\mathbf{R}_l^{(j)} = \sigma_{im}^2 \mathbf{I}_{4\mathcal{M}_j}$ . Eq.(37) can not be directly used for measurement update in MSCKF. Because the term  $\mathbf{H}_f^{(j)G} \tilde{\mathbf{p}}_{f_j}$  is not part of state vector. In order to transform Eq.(37) into a standard form for update, we can compute a semi-unitary matrix  $\mathbf{A}$  whose columns form the basis of the



(a)



(b)

Fig. 2. Sample images from dataset used in pose estimation experiment. (a) is the first pair of stereo images. (b) is the last pair.

left nullspace of  $\mathbf{H}_f^{(j)}$ , and multiply both sides of Eq.(37) by  $\mathbf{A}$ .

$$\begin{aligned} \mathbf{r}_o^{(j)} &= \mathbf{A}^T \mathbf{r}^{(j)} = \mathbf{A}^T \mathbf{H}_{x_B}^{(j)} \tilde{\mathbf{x}} + \mathbf{0} + \mathbf{A}^T \mathbf{n}^{(j)} \\ &= \mathbf{H}_o^{(j)} \tilde{\mathbf{x}} + \mathbf{n}_o^{(j)} \end{aligned} \quad (38)$$

Now, we obtain the useful form for filter update.  $\mathbf{H}_{x_B}^{(j)}$  has full column rank, accordingly,  $\mathbf{A}$  has dimension  $4\mathcal{M}_j \times (4\mathcal{M}_j - 3)$  and  $\mathbf{r}_o^{(j)}$  has dimension  $(4\mathcal{M}_j - 3) \times 1$ . The covariance matrix of  $\mathbf{n}_o^{(j)}$  is  $\mathbf{R}_o^{(j)} = \mathbf{A}^T \mathbf{R}^{(j)} \mathbf{A}$ . We can now stack all the errors  $\mathbf{r}_o^{(j)}$  for all the features selected for update.

$$\mathbf{r}_o = \mathbf{H}_o \tilde{\mathbf{x}} + \mathbf{n}_o \quad (39)$$

To reduce the computational complexity of the MSCKF update, a QR-decomposition of  $\mathbf{H}_o$  is employed.

$$\mathbf{H}_o = [\mathbf{Q}_0 \quad \mathbf{Q}_1] \begin{bmatrix} \mathbf{H}_T \\ \mathbf{0} \end{bmatrix} \quad (40)$$

where  $\mathbf{Q}_0$ ,  $\mathbf{Q}_1$  are unitary matrices and  $\mathbf{H}_T$  is an uppertriangular matrix. Substituting this result into (39), we obtain:

$$\begin{bmatrix} \mathbf{Q}_0^T \mathbf{r}_o \\ \mathbf{Q}_1^T \mathbf{r}_o \end{bmatrix} = \begin{bmatrix} \mathbf{H}_T \\ \mathbf{0} \end{bmatrix} \tilde{\mathbf{x}} + \begin{bmatrix} \mathbf{Q}_0^T \mathbf{n}_o \\ \mathbf{Q}_1^T \mathbf{n}_o \end{bmatrix} \quad (41)$$

We discard term  $\mathbf{Q}_1^T \mathbf{r}_o$ , because it is only noise. And we re-define error term used for update:

$$\mathbf{r}_n := \mathbf{Q}_0^T \mathbf{r}_o = \mathbf{H}_T \tilde{\mathbf{x}} + \mathbf{Q}_0^T \mathbf{n}_o = \mathbf{H}_T \tilde{\mathbf{x}} + \mathbf{n}_n \quad (42)$$

Then, the covariance matrix of  $\mathbf{n}_n$  is  $\mathbf{R}_n = \mathbf{Q}_0^T \mathbf{R}_o \mathbf{Q}_0$ . Now, we arrive at the final form of measurement equation. We can finally formulate the Kalman gain and correction equations to obtain the updated estimates for the filter state

and covariance:

$$K = P_{k+1} H_T^T (H_T P_{k+1} H_T^T + R_n)^{-1} \quad (43)$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_{k+1} + K \mathbf{r}_n \quad (44)$$

$$P_{k+1} \leftarrow (I_{27+9N} - K H_T) P_{k+1} (I_{27+9N} - K H_T)^T + K R_n K^T \quad (45)$$

### III. EXPERIMENTS

We tested the proposed method and compared it with the monocular-based MSCKF method on the EuRoC MAV dataset recorded by ETH Autonomous Systems Lab. This dataset is recorded using VI-sensor and includes data streams from stereo camera and IMU [19]. The carefully calibrated extrinsic parameters and millimeter accurate position ground-truth is available in this dataset.

We implemented and tested both algorithms in MATLAB 2014b on a Lenovo laptop with a 2.4 GHz Intel Core i7 processor and 12 GB of DDR3L RAM. We extracted between 100 and 200 salient point features using Oriented FAST and Rotated BRIEF (ORB) detector [20] of OpenCV [21] from the stereo pairs and tracked them temporally using Kanade-Lucas-Tomasi (KLT) tracking [22]. Outliers are rejected using a 5-point Random Sample Consensus (RANSAC). We conducted two experiments on the dataset to compare the accuracy of the the proposed method with the monocular-based one and show the result of the calibration between sensors. We will discuss both of these in turn.

#### A. Pose Estimation

To compare the performance between monocular-based MSCKF (m-MSCKF) and the proposed method, we implemented a m-MSCKF method based [12] and carried out indoor experiment using real-world dataset. The total duration of the experiment is 130s and the trajectory length is approximately 59 m. The IMU sample rate is 200Hz and images are recorded at 20 Hz. We experiment based on *Machine Hall 01* data sequence. Fig.2 shows sample images from the dataset.

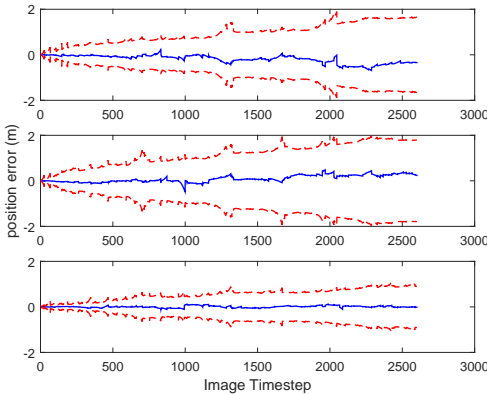


Fig. 3. The position error (blue) and  $\pm 3\sigma$  bounds (red) about axes x, y, and z.

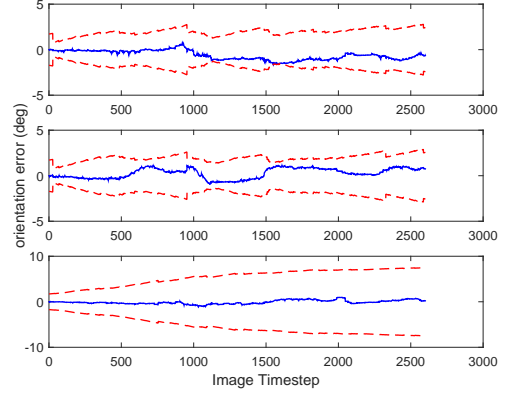


Fig. 4. The orientation error (blue) and  $\pm 3\sigma$  bounds (red) about axes x (roll), y (pitch), and z (yaw).

The results shown in Fig.3 and 4 show the error and  $\pm 3\sigma$  bounds corresponding to position and orientation estimate. The plots reveal that the error remains bounded by  $\pm 3\sigma$ , meaning that the filter is consistent. And the proposed estimator has correct observability properties for camera/IMU system. The roll and pitch are observable, the error bounds do not increase indefinitely. In contrast, those for the position and yaw do continuously increase because they are not observable [14].

In Fig.5, the estimated poses are displayed in X, Y and Z axis separately. It also gives a comparison between the m-MSCKF method. According to the ground truth, the final position is  $[3.8050 \ 1.3780 \ 0.7572]^T$  m. The final position estimation given by the m-MSCKF method is  $[4.6609 \ 0.9028 \ 0.5055]^T$  m, while the proposed method gives  $[3.6275 \ 1.5840 \ 0.7833]^T$  m. The average error between the trajectory estimated by proposed method and the ground truth is  $[0.2810 \ 0.1804 \ 0.0659]^T$  m in three axis. The results show that the proposed method owns better precision than m-MSCKF method.

#### B. Sensors Calibration

In order to validate the proposed filter algorithm is effective for estimating the IMU-camera transformation while pose estimation, we added initial alignment error for translation and rotation to the known extrinsic calibration between sensors.

For simplicity and due to limited space, in this section, we just discuss the calibration of transformation between IMU and CAM0. Because the calibration of IMU and CAM1 is in a similarity situation. The initial alignment error for translation is set to  ${}^B \tilde{\mathbf{p}}_{C_0} = [-0.03 \ 0.10 \ -0.09]^T$  m with a standard deviation of  $[0.0548 \ 0.0447 \ 0.0458]^T$  m in each axis. The initial alignment error for rotation is set to  ${}^B \tilde{\boldsymbol{\theta}} = [5^\circ \ 6^\circ \ -4^\circ]^T$  with a standard deviation of  $[3.3914^\circ \ 3.8455^\circ \ 3.1917^\circ]^T$  in each axis of rotation. Consequently, the filter state vector and error-state covariance matrix are initialized according to the process described in Section II.

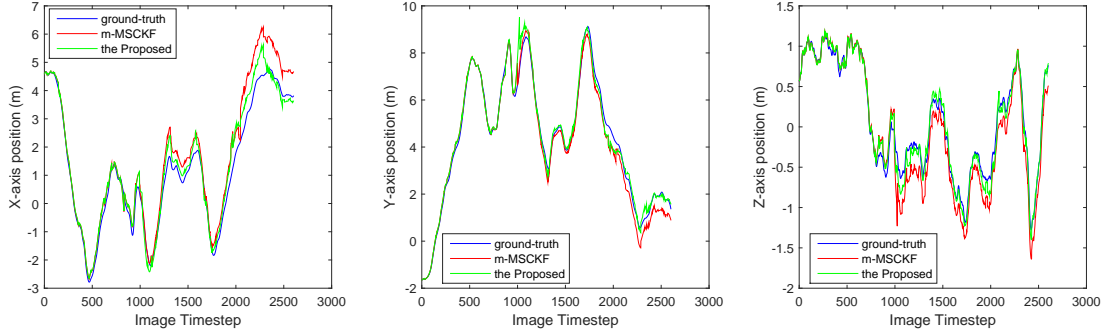


Fig. 5. Experimental results: comparison of ground-truth (blue), trajectories estimated by m-MSCKF (red) and the proposed method (green)

In Figs.6 and 7, the IMU-camera calibration errors and their  $3\sigma$  bounds for the 6-DOF transformation between the IMU and the camera are shown.

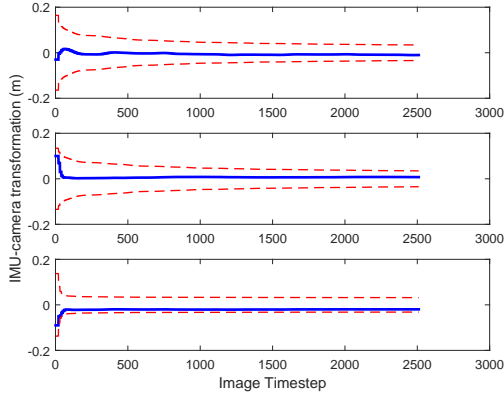


Fig. 6. IMU-camera translation error (blue) and  $3\sigma$  bounds (red). Translation along axes x, y, and z.

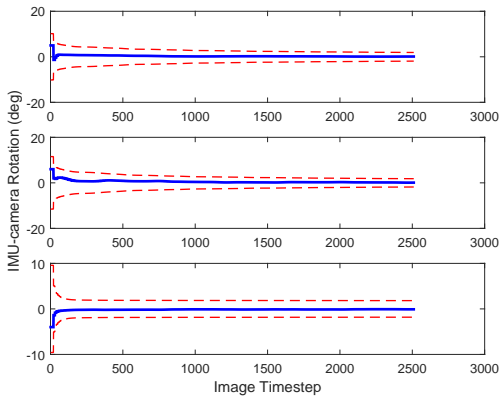


Fig. 7. IMU-camera rotation error (blue) and  $3\sigma$  bounds (red). Rotation about axes x (roll), y (pitch), and z (yaw).

Fianlly, we list the initial and final uncertainty of the IMU-camera parameters (translation and rotation), including CAM0 and CAM1 in Table I. The results shown here demonstrate that the proposed calibration is capable of operating in a real-world environment and this method is effective.

TABLE I  
FINAL UNCERTAINTY ( $3\sigma$ ) OF THE IMU-CAMERAS  
PARAMETERS AFTER 2500 IMAGE TIMESTEP

$3\sigma$		x(m)	y(m)	z(m)	r( $^\circ$ )	p( $^\circ$ )	y( $^\circ$ )
CAM0	Initial	0.16	0.13	0.134	10.17	11.53	9.57
	Final	0.03	0.03	0.03	1.90	1.82	1.81
CAM1	Initial	0.28	0.23	0.19	9.41	12.63	12.06
	Final	0.03	0.03	0.03	1.90	1.83	1.84

#### IV. CONCLUSIONS

In this paper, we presented a filter-based visual inertial sensor fusion algorithm for pose estimation. The main contribution is the derivation of a measurement model of stereo visual observation. And no other method of attempting to stereo visual information with inertial sensor in a MSCKF framework has been found on literature yet. In addition, the proposed method estimates estrinsic calibration parameters between sensors to improve the performamce. The experimental results indicate that our method has a better accuracy than a self-implemented m-MSCKF ant it is effective to estimate sensor extrinsic calibration online.

#### REFERENCES

- [1] E. Eade and T. Drummond, "Scalable monocular slam," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 1. IEEE, 2006, pp. 469–476.
- [2] D. Strelow and S. Singh, "Motion estimation from image and inertial measurements," *The International Journal of Robotics Research*, vol. 23, no. 12, pp. 1157–1195, 2004.
- [3] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment: a modern synthesis," in *International workshop on vision algorithms*. Springer, 1999, pp. 298–372.
- [4] 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.
- [5] H. Strasdat, J. Montiel, and A. J. Davison, "Scale drift-aware large scale monocular slam," *Robotics: Science and Systems VI*, 2010.
- [6] R. Mur-Artal, J. 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.
- [7] M. Li, "Visual-inertial odometry on resource-constrained systems," 2014.
- [8] M. Kleinert and S. Schleith, "Inertial aided monocular slam for gps-denied navigation," in *Multisensor Fusion and Integration for Intelligent Systems (MFI), 2010 IEEE Conference on*. IEEE, 2010, pp. 20–25.

- [9] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *Mixed and Augmented Reality, 2007. ISMAR 2007. 6th IEEE and ACM International Symposium on*. IEEE, 2007, pp. 225–234.
- [10] S. Weiss and R. Siegwart, "Real-time metric state estimation for modular vision-inertial systems," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 4531–4537.
- [11] S. Lynen, M. Achtelik, S. Weiss, M. Chli, and R. Siegwart, "A robust and modular multi-sensor fusion approach applied to mav navigation," in *Proc. of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2013.
- [12] A. I. Mourikis and S. I. Roumeliotis, "A multi-state constraint kalman filter for vision-aided inertial navigation," in *Proceedings 2007 IEEE International Conference on Robotics and Automation*. IEEE, 2007, pp. 3565–3572.
- [13] E. S. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *The International Journal of Robotics Research*, vol. 30, no. 4, pp. 407–430, 2011.
- [14] J. Kelly and G. S. Sukhatme, "Visual-inertial sensor fusion: Localization, mapping and sensor-to-sensor self-calibration," *The International Journal of Robotics Research*, vol. 30, no. 1, pp. 56–79, 2011.
- [15] M. Li and A. I. Mourikis, "High-precision, consistent ekf-based visual-inertial odometry," *The International Journal of Robotics Research*, vol. 32, no. 6, pp. 690–711, 2013.
- [16] J. Sola, "Quaternion kinematics for the error-state kf," *Laboratoire d'Analyse et d'Architecture des Systemes-Centre national de la recherche scientifique (LAAS-CNRS), Toulouse, France, Tech. Rep*, 2012.
- [17] J. Montiel, J. Civera, and A. J. Davison, "Unified inverse depth parametrization for monocular slam," *analysis*, vol. 9, p. 1, 2006.
- [18] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [19] 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. [Online]. Available: <http://ijr.sagepub.com/content/early/2016/01/21/0278364915620033.abstract>
- [20] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *2011 International conference on computer vision*. IEEE, 2011, pp. 2564–2571.
- [21] G. Bradski, *Dr. Dobb's Journal of Software Tools*.
- [22] S. Birchfield, "Klt: An implementation of the kanade-lucas-tomasi feature tracker," 2007.