

Depth Enhanced Visual-Inertial Odometry Based on Multi-State Constraint Kalman Filter

Fumin Pang, Zichong Chen, Li Pu, and Tianmiao Wang, *Member, IEEE*

Abstract—There have been increasing demands for developing robotic system combining camera and inertial measurement unit in navigation task, due to their low-cost, lightweight and complementary properties. In this paper, we present a Visual Inertial Odometry (VIO) system which can utilize sparse depth to estimate 6D pose in GPS-denied and unstructured environments. The system is based on Multi-State Constraint Kalman Filter (MSCKF), which benefits from low computation load when compared to optimization-based method, especially on resource-constrained platform. In this paper, we enhance the features with depth information to form 3D landmark position measurements in space, which reduces uncertainty of position estimate. And we derive measurement model to access compatibility with both 2D and 3D measurements. In experiments, we evaluate the performance of the system in different in-flight scenarios, both cluttered room and industry environment. The results suggest that the estimator is consistent, substantially improves the accuracy compared with original monocular-based MSCKF and achieves competitive accuracy with other research.

I. INTRODUCTION

Accurate 6D pose estimate in unknown environment from a set of sensor measurements is one of the critical problems in robotics navigation task. Combination of complementary information from visual and inertial sensors is ubiquitous in mobile robot application with the requirement of high dynamic operation. Inertial sensors readings can be used to compute relatively accurate 6D pose in short period and give a real-metric estimate for absolute scale[1]. While visual sensor with the ability to provide rich texture of environment and bearing measurements for salient landmarks make it an efficient exteroceptive sensor to correct the motion and structure prediction. Each visual feature can always be tracked by a camera from a sequence of consecutive poses, which provide multiple constraint of camera motion and landmarks. VIO can be considered a subproblem of Simultaneous Localization and Mapping (SLAM), while it pays more focus on efficient 6D pose estimate. This method has achieved success in navigating Micro Aerial Vehicles(MAV) and cars.

To date, the majority of algorithms proposed for real-time VIO can be classified into two categories: extended Kalman filter-based methods and methods utilizing optimization approach. Filter-based approaches are the earlier ones used to solve SLAM and VIO problems. Davison et al. proposed one of the first real-time 3D monocular SLAM

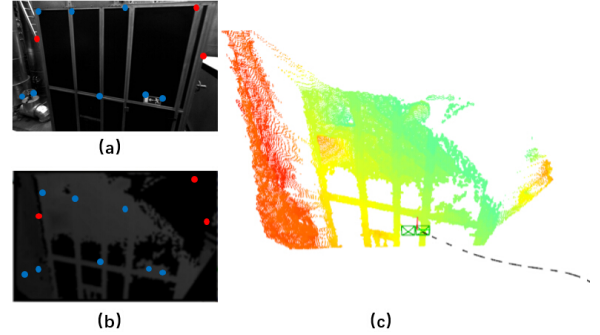


Fig. 1: (a) Features tracked at an image frame. The blue dots represent features whose depth comes from stereo matches, and the red dots represent features without depth. The proposed method uses both types of features in estimating 6D pose. (b) A disparity map from stereo a corresponding to (a), which shows depth information is only available in the vicinity of the camera. (c) A local point cloud built from (b), colors encode distance. This local map shows the sparse nature for space landmarks.

framework based on EKF in computer vision[2]. Based on this work, Roussillon et al. used IMU as a interoceptive to propagate position, orientation and velocity of platform in high dynamic environment[3]. Jones and Soatto[4] built a EKF using Lie derivatives and studied observability in visual-inertial system. Their results showed that the IMU biases, 3D velocity, absolute roll and pitch angles and IMU-camera transformation are observable. Some Similar conclusions were also drawn by Kelly et al.[5].

Meanwhile, optimization-based methods generally attain higher accuracy, as they re-linearize at each iteration to better deal with their nonlinear measurement models. Leutenegger et al. describe a tightly coupled approach in which the robot poses and sparse 3D landmarks are estimated through minimize a joint optimization problem using inertial error terms as well as the reprojection error of the landmarks in stereo image[6]. However, these methods always implement bundle adjustment in a sliding window of states, using multiple iterations to minimize cost function, which results in increased computational cost. Thus, filter-based methods might be a better choice for long-time running tasks, especially in resource-constrained systems, such as micro aerial vehicles and Augmented Reality (AR) devices[7].

Multi-State Constraint Kalman Filter improve efficiency of filter-based method further[8]. Instead of including the 3D feature position in the state vector of the EKF, MSCKF

Fumin Pang, Zichong Chen and Li Pu are with Segway Robotics Inc. Beijing, China. E-mail: {fumin.pang, zichong.chen, li.pu}@segwayrobotics.com

Tianmiao Wang is with the School of Mechanical Engineering and Automation, Beihang University, Beijing, China. E-mail: itm@buaa.edu.cn

maintains a sliding window of historical poses in state vector. The landmarks positions are triangulated from multiple observations from different camera poses when they slip away from field of view. Because the number of poses in sliding window is usually much less than the number of features, the computational complexity is reduced. Li et al. correct the observation properties of original MSCKF using First-Estimate Jacobian (FEJ)[9].

The sparsity of depth mainly lies on two aspects. First, with limitation of the sensors, scenarios with large depth variation leave large area in the images where depth is unavailable spatially. Second, as in MSCKF paradigm, there need to exist multiple measurements for one feature from a series of different robot poses before its 3D position is triangulated. However, only a subset of measurements for one feature has corresponding depth information temporally in general. The inevitable sparsity makes people unable to use ICP-like method in large scale environment to estimate poses.

In this paper, we propose a method that can utilize sparse depth information along with the imagery based on MSCKF, inspired by [10] and [11]. We utilize 2D visual feature and depth to form a 3D measurement of a landmark. Measurement model is adapted to be compatible with both 2D and 3D measurements. In the update step of this depth enhanced MSCKF, reprojection measurement errors for landmarks where depth is unavailable and 3D position measurement errors are jointed to correct the pose estimate. The proposed method has been tested on real-world data. The results demonstrate a lower drift than mono-MSCKF. The attained accuracy is competitive with other research in this field.

In this work, we use stereo camera to obtain depth information. But it can be extended and adapted to various sensors which can provide bearing and range information.

II. ESTIMATOR DESCRIPTION

Four different coordinate frames are used throughout the paper: we affix camera-IMU rig body frame $\{B\}$ to IMU, to track the 6D motion with respect to a global coordinate frame, $\{G\}$. The camera coordinate frame, $CAM0$, is $\{C^0\}$. In this paper, it is $CAM0$ where features are tracked. An additional camera, $CAM1$, is added to form a stereo rig to estimate depth by potentially matching with features in $CAM0$. Its coordinate frame is $\{C^1\}$.

A. Overall Filter Structure and State Parametrization

The full state representation can be partitioned into two parts according to MSCKF paradigm. 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. The body state at time k is a 17-dimensional vector as follow:

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

In this world-centric presentation, ${}^B\tilde{\mathbf{q}}_k$ 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 [12]. ${}^G\mathbf{p}_{B,k}$ 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 . t_d is a scalar modelling the unknown *time offset* between IMU and cameras, which will help to improve accuracy [13]. In this paper, we assume $CAM0$ and the depth sensor are synchronized temporally.

The second 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. Velocity state involves in t_d estimate.

$$\mathbf{x}_{B_i,k} := \begin{bmatrix} {}^{B_i}\tilde{\mathbf{q}}_k^T & {}^{B_i}\mathbf{p}_{B_i,k}^T & {}^{G}\mathbf{v}_{B_i,k}^T \end{bmatrix}^T \quad (2)$$

Combining this two parts, the full state of MSCKF is a $(17 + 10 \cdot N)$ vector consisting of current body state estimate and N last body states at time k .

In MSCKF, *Error State Kalman Filter* (ESKF)[12] is used to avoiding singularity brought by quaternions which use 4 dimensions to describe 3 degrees of freedom. While the error (\tilde{x}) between the true value (x) and the estimated value (\hat{x}) for position, velocity, bias and t_d can be defined as $\tilde{x} = x - \hat{x}$ trivially, the error for quaternion is defined as:

$$\delta\tilde{\mathbf{q}} := \hat{\mathbf{q}}^{-1} \otimes \tilde{\mathbf{q}} \approx \begin{bmatrix} \frac{1}{2}\delta\boldsymbol{\theta}^T & 1 \end{bmatrix}^T \quad (3)$$

Using $\delta\boldsymbol{\theta}$ to represent orientations in the Kalman filter reduces their dimensionality to 3. As a result, the *error state* of the estimator at time k :

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

where

$$\tilde{\mathbf{x}}_{I,k} := \begin{bmatrix} {}^G\delta\boldsymbol{\theta}_I^T & {}^G\tilde{\mathbf{p}}_{B,k}^T & {}^G\tilde{\mathbf{v}}_{B,k}^T & \tilde{\mathbf{b}}_{g,k}^T & \tilde{\mathbf{b}}_{a,k}^T & \tilde{t}_{d,k} \end{bmatrix}^T \quad (5)$$

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

The full error state has $(16 + 9 \cdot N)$ dimensions. Accordingly, the MSCKF error state covariance \mathbf{P} is a $(16 + 9 \cdot N) \times (16 + 9 \cdot N)$ matrix.

B. Filter Propagation and Augment

Every time an IMU measurement is received, it is used to propagate the IMU state. As mentioned before, IMU measurement provide rotational velocity $\boldsymbol{\omega}_m$ and \mathbf{a}_m , described as below equations:

$$\boldsymbol{\omega}_m = {}^G\boldsymbol{\omega} + \mathbf{b}_g + \mathbf{n}_g \quad (7)$$

$$\mathbf{a}_m = {}^B\mathbf{R}({}^G\mathbf{a} - {}^G\mathbf{g}) + \mathbf{b}_a + \mathbf{n}_a \quad (8)$$

where ${}^B_G\mathbf{R}$ is the rotation matrix corresponding 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. We follow the approach described in [9] and [14] to propagate filter state in discrete time.

Linearized continuous-time model of current state error state, $\tilde{\mathbf{x}}_I$, can be drawn from (9):

$$\dot{\tilde{\mathbf{x}}}_I = \mathbf{F}\tilde{\mathbf{x}}_I + \mathbf{G}\mathbf{n}_I \quad (9)$$

where \mathbf{F} is the continuous-time error-state transition matrix. \mathbf{G} is the Jacobian of current body error state with respect to the noise vector $\mathbf{n}_I = [\mathbf{n}_g^T \mathbf{n}_a^T \mathbf{n}_{wg}^T \mathbf{n}_{wa}^T]^T$.

State covariance matrix also propagates. And upon recording a new image, a copy of current body state is appended to the state vector, and the covariance matrix of the MSCKF is augmented accordingly [8].

C. Mean and Covariance of 3D measurement

In this paper, visual measures from stereo camera form a 3D measurement of a landmark when the depth is available. Stereo measurement can efficiently reduce the uncertainty (covariance) of 3D position estimate. We consider one feature, \mathbf{f} , which is observed by stereo rig at timestep i . We will derivate its mean and covariance of 3D position measurement.

First, in frame $\{C^0\}$, landmark position ${}^{C^0}_i\mathbf{p}_f = [{}^{C^0}_iX \ {}^{C^0}_iY \ {}^{C^0}_iZ]^T$ is calculated by stereo triangulation. We present its 3D measurement using a pixel bearing vector $[\alpha \ \beta]^T$ and a inverse depth scalar ρ :

$${}^{C^0}_i\mathbf{z}_3 = \begin{bmatrix} {}^{C^0}_i\alpha \\ {}^{C^0}_i\beta \\ {}^{C^0}_i\rho \end{bmatrix} = \boldsymbol{\pi}_3 \begin{pmatrix} {}^{C^0}_iX \\ {}^{C^0}_iY \\ {}^{C^0}_iZ \end{pmatrix} = \begin{bmatrix} \frac{{}^{C^0}_iX}{{}^{C^0}_iZ} \\ \frac{{}^{C^0}_iY}{{}^{C^0}_iZ} \\ \frac{1}{{}^{C^0}_iZ} \end{bmatrix} \quad (10)$$

Similarly, in $\{C_1\}$, the 3D measurement ${}^{C_1}_i\mathbf{z}_3$ of ${}^{C^0}_i\mathbf{p}_f$ is

$${}^{C_1}_i\mathbf{z}_3 = \begin{bmatrix} {}^{C_1}_i\alpha \\ {}^{C_1}_i\beta \\ {}^{C_1}_i\rho \end{bmatrix} = \boldsymbol{\pi}_3 \left({}^{C^0}_i\mathbf{R} \begin{pmatrix} {}^{C^0}_iX \\ {}^{C^0}_iY \\ {}^{C^0}_iZ \end{pmatrix} + {}^{C^0}_i\mathbf{p}_{C^0} \right) \quad (11)$$

As above, α, β are given by pixel coordinates in image with measurement noise $\text{diag}(\sigma_\alpha^2, \sigma_\beta^2)$. ρ can be calculated by inverse depth. Following paper, we assume $\rho \sim \mathcal{N}(\rho, \sigma_\rho^2)$, where σ_ρ is set manually.

From (16)-(17), the covariance of Gaussian distribution in Euclidean space of ${}^{C^0}_i\mathbf{p}_f$ can be fused by the two measurements uncertainties:

$$\boldsymbol{\Sigma}_{p^j} = \mathbf{J}_{p^j} \text{diag}(\sigma_{\alpha^j}^2, \sigma_{\beta^j}^2, \sigma_{\rho^j}^2) \mathbf{J}_{p^j}^T \quad (12)$$

where $j = 0, 1$ and $\mathbf{J}_{p^j} = \frac{\partial {}^{C^0}_i\mathbf{p}_f}{\partial {}^{C^j}_i\mathbf{z}_3}$.

Given two 3D measurement covariances $\boldsymbol{\Sigma}_{p^0}$ and $\boldsymbol{\Sigma}_{p^1}$, we can fuse the position covariance using a Bayesian update step [15]:

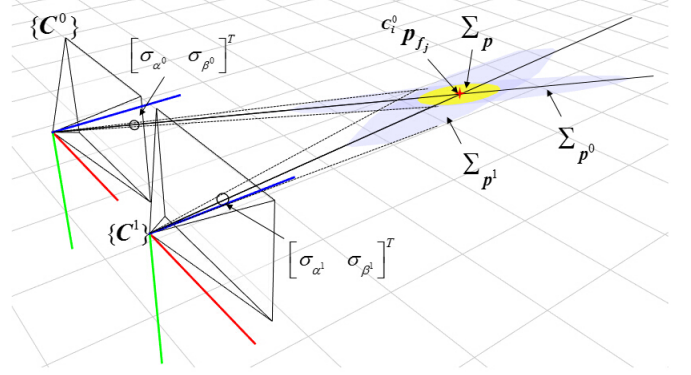


Fig. 2: 3D measurement can be achieved by fusing two visual measurements from stereo camera. Mean of 3D measurement can be calculate from triangulation, and the covariance $\boldsymbol{\Sigma}_p$ comes from merging two 3D Gaussian distributions, $\boldsymbol{\Sigma}_{p^0}$ and $\boldsymbol{\Sigma}_{p^1}$. This Bayesian fusion results in smaller 3D uncertainty than either $\boldsymbol{\Sigma}_{p^0}$ or $\boldsymbol{\Sigma}_{p^1}$, visualized as a smaller ellipsoid. The pixel bearing vector uncertainty involve visual feature uncertainty and scale.

$${}^{C_i}_i\boldsymbol{\Sigma}_p = \boldsymbol{\Sigma}_{p^0} (\boldsymbol{\Sigma}_{p^0} + \boldsymbol{\Sigma}_{p^1})^{-1} \boldsymbol{\Sigma}_{p^1} \quad (13)$$

Fusing two 3-dimensional Gaussian distributions reduces the uncertainty of the distribution (see Fig.2). Finally, the 3D measurement mean and covariance can be derived:

$${}^{C_i}_i\mathbf{z}_3 = \begin{bmatrix} {}^{C_i}_i\alpha \\ {}^{C_i}_i\beta \\ {}^{C_i}_i\rho \end{bmatrix} \quad {}^{C_i}_i\boldsymbol{\Sigma}_3 = \mathbf{J}_{\pi_3} {}^{C_i}_i\boldsymbol{\Sigma}_p \mathbf{J}_{\pi_3}^T \quad (14)$$

where $\mathbf{J}_{\pi_3} = \frac{\partial {}^{C^0}_i\mathbf{z}_3}{\partial {}^{C^0}_i\mathbf{p}_f}$. The 3D covariance in Euclidean space is transformed to bearing and inverse depth parameterization space. Here, for simplicity, C_i is used to present camera frame C^0_i at timestep i . Thus, fusing the visual measures in stereo provides a more precise estimate on landmark position and reduces uncertainty by fusing the two Gaussian distributions, which will potentially improve performance of visual odometry accuracy. In addition, these 3D measurements make filter update more efficient compared to two separate updates in left and right image due to a lower-dimensional measurement jacobain.

D. Depth Enhanced Feature Position Estimation

In this work, once a tracked feature is lost in current frame, CAM_0 , the corresponding landmark position need to be estimated. Compared to original MSCKF which uses triangulation from a sliding window of 2D imagery measurements, we utilize both 2D and 3D measures.

To describe the position estimate process and update model, we consider one landmark, ${}^G\mathbf{p}_{f_j}$, which is observed at a set of poses, $\{T_i = \{{}^{C_i}_G\mathbf{R}, {}^G\mathbf{p}_{C_i}\} \in SE(3)\}$ with $i = 1, 2, 3, 4$ (see Fig.3). In C_1 and C_3 , depth is available. We initialize the optimization by estimating the position of feature f_j in camera frame C_1 using a linear least-squares

method [16]. We can express the feature position in camera frame C_i in terms of its position in camera frame C_1 as

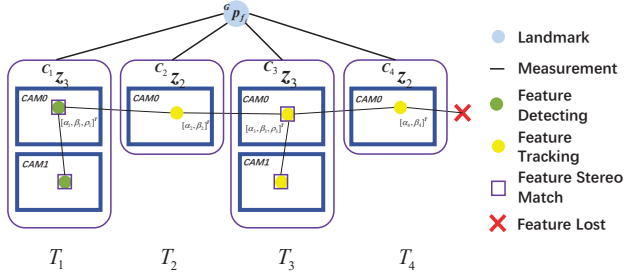


Fig. 3: An example indicates one feature f_j observed by camera at 4 pose. At T_1 and T_3 , depth is available which generates 3D measurements, while at T_2 and T_4 , no matches can be found in $CAM1$ resulting in two 2D measurements. Landmark position corresponding to f_j is estimated by minimizing either 2d or 3D measurement errors through least-square method.

$$C_i \hat{\mathbf{p}}_{f_j} = C_i \hat{\mathbf{R}}^{C_1} \hat{\mathbf{p}}_{f_j} + C_i \hat{\mathbf{p}}_{C_1} \quad (15)$$

we can write the 2D imagery measurement error as

$$e_{2,i}(C_i \hat{\mathbf{p}}_{f_j}) = C_i \mathbf{z}_{2,i} - \pi_2 \left(C_i \hat{\mathbf{R}}^{C_1} \hat{\mathbf{p}}_{f_j} + C_i \hat{\mathbf{p}}_{C_1} \right) \quad (16)$$

where $C_i \mathbf{z}_{2,i}$ is the pixel bearing vector given by feature coordinate and

$$\pi_2(\mathbf{h}) = \begin{bmatrix} \frac{h(1)}{h(3)} \\ \frac{h(2)}{h(3)} \end{bmatrix} \quad (17)$$

and 3D measurement error as

$$e_{3,i}(C_i \hat{\mathbf{p}}_{f_j}) = C_i \mathbf{z}_{3,i} - \pi_3 \left(C_i \hat{\mathbf{R}}^{C_1} \hat{\mathbf{p}}_{f_j} + C_i \hat{\mathbf{p}}_{C_1} \right) \quad (18)$$

By stacking the measurements together, the least-squares system then becomes

$$(\mathbf{J}^T \mathbf{W}^{-1} \mathbf{J}) \delta^{C_1} \mathbf{p}_{f_j}^* = -\mathbf{J}^T \mathbf{W}^{-1} \mathbf{e} \left(C_1 \hat{\mathbf{p}}_{f_j} \right) \quad (19)$$

where $\mathbf{e} \left(C_1 \hat{\mathbf{p}}_{f_j} \right)$ is stacked vector of four measurement errors and \mathbf{J} is stacked matrix of four Jacobians of measurement error with respect to $C_1 \hat{\mathbf{p}}_{f_j}$ and

$$\mathbf{W} = \text{diag} \left\{ C_1 \Sigma_3 \quad C_2 \Sigma_2 \quad C_3 \Sigma_3 \quad C_4 \Sigma_2 \right\} \quad (20)$$

with $C_i \Sigma_2 = \text{diag} \left\{ \sigma_{\alpha_i}^2, \sigma_{\beta_i}^2 \right\}$. $C_i \Sigma_3$ is calculated in previous section.

In practice, in order to improve numerical stability, inverse depth parameterization [17] is used. After iterative optimization, we get $C_1 \hat{\mathbf{p}}_{f_j}$ and further $^G \hat{\mathbf{p}}_{f_j}$. A simulation on how different quantities of 3D measurements in sliding window effects accuracy of position estimate is showed in Fig.4. The result shows depth information can observably benefits the estimation, especially along z-axis in local coordinate .

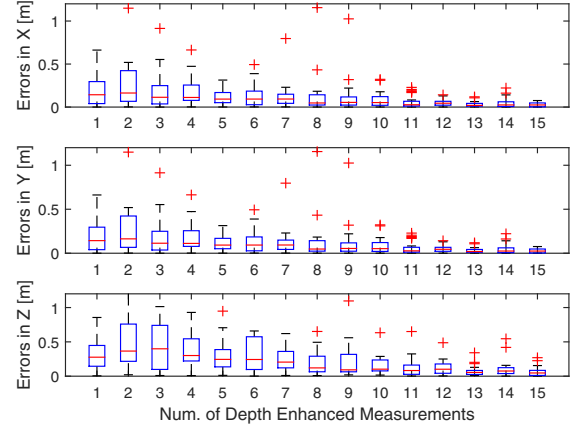


Fig. 4: Boxplot on simulation result about of 3-axis errors of \mathbf{p}_{f_j} presented in C_1 , in which different amount of 3D measurement is applied to estimate landmark position . We generate 50 features within a distance of 1-5 meters in front of cameras, and the sliding window size is 15. All feature positions are estimated with 3D measurement amounts varying from 1 to 15, respectively. The results indicate more 3D measurements reduce the measurements error, especially along z-direction in local coordinate.

E. Depth Enhanced Filter Update

Once the landmark position $^G \hat{\mathbf{p}}_{f_j}$ is estimated, we can calculated the measurement residuals for either 2D measurements or 3D measurements to update the system.

$$\mathbf{r}_i^{(j)} = C_i \mathbf{z}_i^{(j)} - C_i \hat{\mathbf{z}}_i^{(j)} \quad (21)$$

where for 2D measurement, $\mathbf{r}_i^{(j)}$ is 2×1 and

$$C_i \hat{\mathbf{z}}_{2,i}^{(j)} = \pi_2 \left(C_i \hat{\mathbf{R}} \left(^G \hat{\mathbf{p}}_{f_j} - ^G \hat{\mathbf{p}}_{C_i} \right) \right) \quad (22)$$

and for 3D measurement, $\mathbf{r}_i^{(j)}$ is 3×1 and

$$C_i \hat{\mathbf{z}}_{3,i}^{(j)} = \pi_3 \left(C_i \hat{\mathbf{R}} \left(^G \hat{\mathbf{p}}_{f_j} - ^G \hat{\mathbf{p}}_{C_i} \right) \right) \quad (23)$$

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

$$\mathbf{r}_i^{(j)} = \mathbf{H}_{x_{B_i}}^{(j)} \tilde{\mathbf{x}} + \mathbf{H}_{f_i}^{(j)G} \tilde{\mathbf{p}}_{f_j} + \mathbf{n}_i^{(j)} \quad (24)$$

where the matrices $\mathbf{H}_{x_{B_i}}^{(j)}$ and $\mathbf{H}_{f_i}^{(j)}$ are the corresponding Jacobians of the measurement $\hat{\mathbf{z}}_i^{(j)}$ 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 . $\mathbf{n}_i^{(j)}$ is the measurement noise with covariance either $C_i \Sigma_2$ or $C_i \Sigma_3$. $\mathbf{H}_{x_{B_i}}^{(j)}$ and $\mathbf{H}_{f_i}^{(j)}$ are given by

$$\mathbf{H}_{x_{B_i}}^{(j)} = \mathbf{J}_{\pi,i}^{(j)} \begin{bmatrix} \mathbf{0}_{3 \times 15} & \mathbf{M}_{t_d,i} & \mathbf{0}_{3 \times 9(i-1)} & \mathbf{M}_{B_i} & \mathbf{0}_{3 \times 9(N-l)} \end{bmatrix} \quad (25)$$

$$\mathbf{H}_{f_i}^{(j)} = \mathbf{J}_{\pi,iB}^{(j)C_0} \mathbf{R}_G^B \hat{\mathbf{R}}(t + \hat{t}_d) \quad (26)$$

where

$$\mathbf{M}_{t_d,i} = {}^{C_0}_B \mathbf{R} \left(\lfloor {}^B_G \hat{\mathbf{R}} \left({}^G \hat{\mathbf{p}}_{f_j} - {}^G \hat{\mathbf{p}}_{B_i} \right) \times \rfloor^B \boldsymbol{\omega} - {}^B_G \hat{\mathbf{R}}^G \hat{\mathbf{v}}_{B_i} \right) \quad (27)$$

$$\mathbf{M}_{B_i} = {}^{C_0}_B \mathbf{R}_G^B \hat{\mathbf{R}}(t + \hat{t}_d) \left[\left[({}^G \mathbf{p}_{f_j} - {}^G \tilde{\mathbf{p}}_B) \times \right] - \mathbf{I}_3 \quad \mathbf{0}_3 \right] \quad (28)$$

$\lfloor \mathbf{c} \times \rfloor$ is the skew symmetric matrix corresponding to vector \mathbf{c} . And for more details, the nonzero blocks in $\mathbf{H}_{x_{B_i}}^{(j)}$ are the Jacobians with respect to body rotation and body position, respectively. $\mathbf{J}_{\pi,i}^{(j)}$ is the Jacobian of the perspective model depending on 2D or 3D measurement type:

$$\mathbf{J}_{\pi_2,i}^{(j)} = \frac{1}{C_i \hat{\mathbf{Z}}^{(j)}} \begin{bmatrix} C_i \hat{\mathbf{X}}^{(j)} \\ 1 & 0 & -\frac{C_i \hat{\mathbf{Z}}^{(j)}}{C_i \hat{\mathbf{Y}}^{(j)}} \\ 0 & 1 & -\frac{C_i \hat{\mathbf{Z}}^{(j)}}{C_i \hat{\mathbf{X}}^{(j)}} \end{bmatrix} \quad (29)$$

$$\mathbf{J}_{\pi_3,i}^{(j)} = \frac{1}{C_i \hat{\mathbf{Z}}^{(j)}} \begin{bmatrix} C_i \hat{\mathbf{X}}^{(j)} \\ 1 & 0 & -\frac{C_i \hat{\mathbf{Z}}^{(j)}}{C_i \hat{\mathbf{Y}}^{(j)}} \\ 0 & 1 & -\frac{C_i \hat{\mathbf{Z}}^{(j)}}{C_i \hat{\mathbf{X}}^{(j)}} \\ 0 & 0 & -\frac{C_i \hat{\mathbf{Z}}^{(j)}}{C_i \hat{\mathbf{Z}}^{(j)}} \end{bmatrix} \quad (30)$$

Stacking the residuals of all measurements either 2D or 3D, we can get:

$$\mathbf{r}^{(j)} = \mathbf{H}_{x_B}^{(j)} \tilde{\mathbf{x}} + \mathbf{H}_f^{(j)G} \tilde{\mathbf{p}}_{f_j} + \mathbf{n}^{(j)} \quad (31)$$

where $\mathbf{r}^{(j)}$, $\mathbf{H}_{x_B}^{(j)}$, $\mathbf{H}_f^{(j)}$, and $\mathbf{n}^{(j)}$ are block vectors or matrices with elements $r_i^{(j)}$, $\mathbf{H}_{x_{B_i}}^{(j)}$, $\mathbf{H}_{f_i}^{(j)}$, and $\mathbf{n}_i^{(j)}$. Eq.(37) can not be directly used for measurement update in MSCKF. Because the term ${}^G \tilde{\mathbf{p}}_{f_j}$ is not part of state vector. In order to transform (37) into a standard form for update, we can compute a semi-unitary matrix \mathbf{A} whose columns form the basis of the left nullspace of $\mathbf{H}_f^{(j)}$, and multiply both sides of (37) by \mathbf{A} [8].

$$\begin{aligned} \mathbf{r}_o^{(j)} &= \mathbf{A}^T \mathbf{r}^{(j)} = \mathbf{A}^T \mathbf{H}_{x_B}^{(j)} \tilde{\mathbf{x}} + \mathbf{0} + \mathbf{A}^T \mathbf{n}^{(j)} \\ &= \mathbf{H}_o^{(j)} \tilde{\mathbf{x}} + \mathbf{n}_o^{(j)} \end{aligned} \quad (32)$$

Now, we obtain the useful form for filter update. We assume the quantity of 2D measurements is $\mathcal{M}_2^{(j)}$, while $\mathcal{M}_3^{(j)}$ for 3D measurements. Thus \mathbf{A} has dimension $(2\mathcal{M}_2^{(j)} + 3\mathcal{M}_3^{(j)}) \times (2\mathcal{M}_2^{(j)} + 3\mathcal{M}_3^{(j)} - 3)$ and $\mathbf{r}_o^{(j)}$ has dimension $(2\mathcal{M}_2^{(j)} + 3\mathcal{M}_3^{(j)} - 3) \times 1$. The covariance matrix of $\mathbf{n}_o^{(j)}$ is $\boldsymbol{\Sigma}_o^{(j)} = \mathbf{A}^T \boldsymbol{\Sigma}^{(j)} \mathbf{A}$. We can now stack all the errors $\mathbf{r}_o^{(j)}$ for all the features selected for update.

$$\mathbf{r}_o = \mathbf{H}_o \tilde{\mathbf{x}} + \mathbf{n}_o \quad (33)$$

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

Finally, we calculate the Kalman gain. And the state vector and covariance are updated following the description in [8].

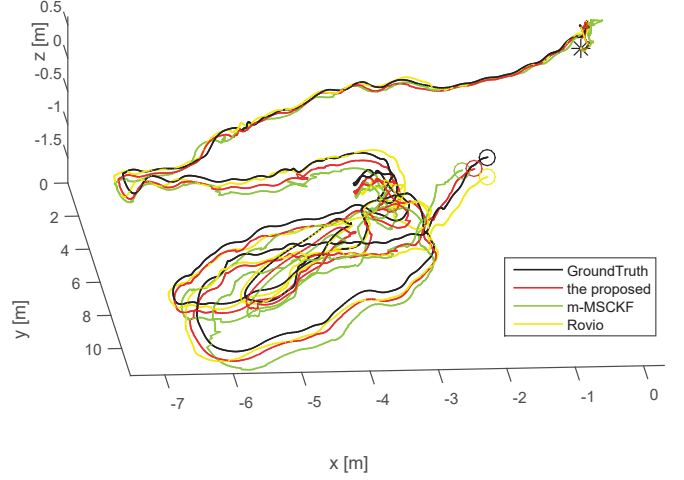


Fig. 5: Comparison on trajectory estimations of mono-MSCKF, Rovio and our method on *MH_01_easy* dataset.

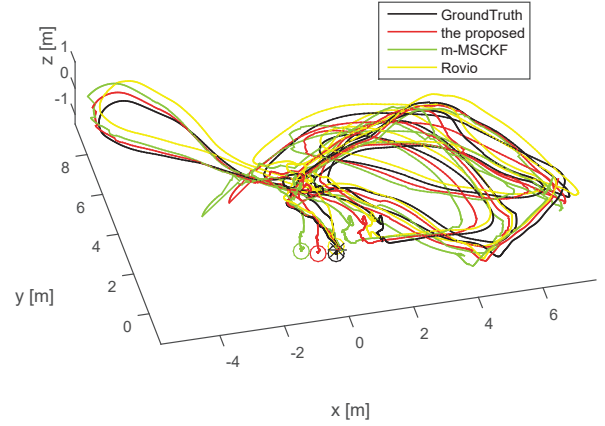


Fig. 6: Comparison on trajectory estimations of mono-MSCKF, Rovio and our method on *MH_03_medium* dataset.

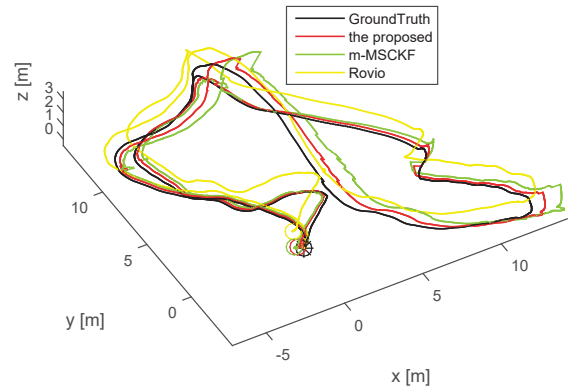


Fig. 7: Comparison on trajectory estimations of mono-MSCKF, Rovio and our method on *MH_05_difficult* dataset.

TABLE I: ATE on EuRoC / ASL Dataset

Dataset	mono-MSCKF			Rovio			The Proposed		
	Mean	Median	RMSE[m]	Mean	Median	RMSE[m]	Mean	Median	RMSE[m]
<i>MH_01.easy</i>	0.3536	0.3433	0.3875	0.2517	0.2560	0.2761	0.1716	0.1768	0.1936
<i>MH_03.medium</i>	0.7778	0.8026	0.8511	0.4053	0.3839	0.4524	0.4002	0.4134	0.4439
<i>MH_05.difficult</i>	0.8943	0.6438	1.1091	1.0490	1.0891	1.1250	0.4292	0.2984	0.5262
<i>V2_01.easy</i>	0.2500	0.1963	0.2998	0.2383	0.2437	0.2564	0.0959	0.0633	0.1370
<i>V2_02.medium</i>	0.4825	0.4576	0.5329	0.3660	0.2948	0.4296	0.1643	0.4525	0.4576

III. IMPLEMENTATION DETAILS

A. Feature Detection and Match

In this paper, We extracted 60-to-200 salient point features using Oriented FAST and Rotated BRIEF (ORB) detector in each frame [18]. ORB feature maintains a image pyramid to find the best corners independent of the scale. For compromise of precision and efficiency, we use a 3-level pyramid while the scale factor between adjacent levels is 1.2. The relationship between feature pixel bearing vector uncertainty and level in which this feature is detected is

$$\sigma_\alpha = \frac{\sigma_u}{f_u} (1.2)^{level} \quad \sigma_\beta = \frac{\sigma_v}{f_v} (1.2)^{level} \quad (34)$$

σ_u and σ_v are feature imagery coordinate isotropic uncertainties at 0-level image in pyramid. f_u and f_v are camera focal length.

Image is divide in cells of fixed size (e.g. 32×30 pixels). This results in evenly distributed features in image by limiting each cell only maintaining the strongest feature ,if it exists.

Each feature track always tries to estimate 3D position of the corresponding landmark in the frame where first detected. To make 2D-to-3D matching be available to accelerate feature search, first stereo measurements (if available) or first two *CAM0* measurement of one feature are used to triangulate this feature's 3D position . Reprojecting the this 3D local point ,which gives a good searching guess for feature matches in succeeding images, keeps the algorithm efficient.

B. Outlier Detection

After the local feature position is achieved when it is first detected, the position uncertainty is reprojected into succeeding images

$${}^{C_{cur}}S = J_h {}^{C_{first}}\Sigma_p J_h^T \quad (35)$$

where J_h is the 2×3 Jacobian of $h(\bullet)$ which projects 3D points in the frame where the feature is first detected to current image with respect to the feature position. Than the 2σ uncertainty ellipse derivated from 2×2 matrix ${}^{C_{cur}}S$ in image frame is used to predict potential matching and reject outliers.

Further, once $\mathbf{r}_o^{(j)}$ and $\mathbf{H}_o^{(j)}$ are calculated in (38), we proceed to carry out a Mahalanobis gating test for the residual $\mathbf{r}_o^{(j)}$. Here a 95th percentile of the χ^2 test with $(2\mathcal{M}_2^{(j)} + 3\mathcal{M}_3^{(j)} - 3)$ degrees of freedom is made.

IV. EXPERIMENTS

We implemented and tested both algorithms as a single thread program on a Lenovo laptop with a 2.4 GHz Intel Core i7-5500 CPU and 12 GB of DDR3L RAM. The experiments are mainly 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 with accurate state groundtruth [19].

A. Pose Estimation

We compare the the proposed method with monocular-based MSCKF (mono-MSCKF) and Rovio, another popular visual-inertial odometry based on EKF [20]. We implemented a mono-MSCKF method based [8] and carried out indoor experiment in different scenarios. In

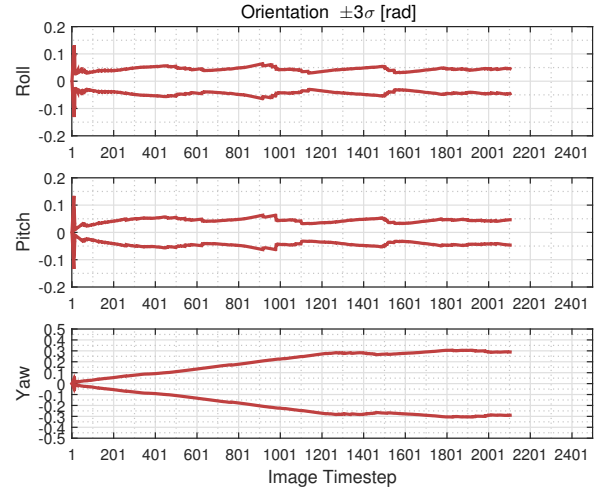


Fig. 8: Estimated orientation $\pm 3\sigma$ bounds of Roll, Pitch and Yaw

Fig.5-7, groundtruth and trajectories estimated from all three methods on *MH_01.easy*, *MH_03.medium* and *MH_05.difficult* dataset. TABLE.1 contains the absolute trajectory error (ATE) between the estimated and the reference trajectory on five EuRoC datasets which contain either fast or slow motion in either light or dark scenes from different methods. Means, medians and RMSEs are listed. In all cases both trajectories were aligned by a rigid transform that minimizes their distance. These comparison results demonstrate the proposed method is robust to different motion and scenarios. Compared to mono-MSCKF and

Rovio, the proposed method has a higher accuracy on 6D pose estimation.

B. Orientation Observability

The results shown in Fig.8 show $\pm 3\sigma$ bounds corresponding to orientation estimate. The plots reveal the proposed estimator has correct observability properties for VINS. 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 [21].

C. Bias Estimation

In Fig.9, IMU biases are estimated. We can see the gyroscope biases exhibit a better convergence than the accelerometer biases. This result may come from the more direct link of rotational rates to visual errors [22].

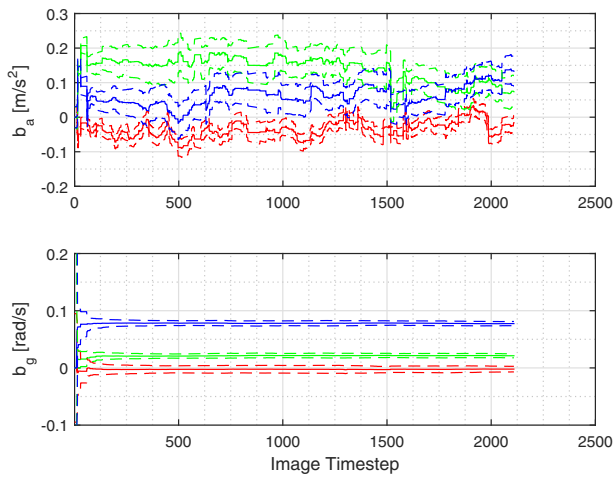


Fig. 9: Estimated IMU biases with their $\pm 3\sigma$ bounds. Top: accelerometer biases (red:x, green:y, blue:z). Bottom: gyroscope biases (red:x, green:y, blue:z).

V. CONCLUSIONS

In this paper, we have presented a depth enhanced visual-inertial odometry system based on MSCKF. By integrating visual feature and depth information, if available, a 3D measurement model of landmark position is derived, which reduces the estimation uncertainty. Combining 2D and 3D measurements results in more accurate estimation of landmark position. MSCKF update model is modified to gain compatibility of 2D and 3D measurement.

Experiments under datasets recorded in different scenarios underline the robustness of the proposed method on pose estimate. Compared to mono-MSCKF and other popular open-source method, the proposed gets competitive accuracy.

REFERENCES

- [1] 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.
- [2] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "Monoslam: Real-time single camera slam," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 6, 2007.

- [3] C. Roussillon and S. Lacroix, "High rate-localization for high-speed all-terrain robots," in *Communications, Computing and Control Applications (CCCA), 2012 2nd International Conference on*. IEEE, 2012, pp. 1–8.
- [4] 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.
- [5] J. Kelly, S. Saripalli, and G. S. Sukhatme, "Combined visual and inertial navigation for an unmanned aerial vehicle," in *Field and Service Robotics*. Springer, 2008, pp. 255–264.
- [6] 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.
- [7] M. Li, "Visual-inertial odometry on resource-constrained systems," 2014.
- [8] 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.
- [9] 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.
- [10] J. Zhang, M. Kaess, and S. Singh, "Real-time depth enhanced monocular odometry," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*. IEEE, 2014, pp. 4973–4980.
- [11] M. N. Galfond, "Visual-inertial odometry with depth sensing using a multi-state constraint kalman filter," Ph.D. dissertation, Massachusetts Institute of Technology, 2014.
- [12] N. Trawny and S. I. Roumeliotis, "Indirect kalman filter for 3d attitude estimation," *University of Minnesota, Dept. of Comp. Sci. & Eng., Tech. Rep.*, vol. 2, p. 2005, 2005.
- [13] M. Li and A. I. Mourikis, "3-d motion estimation and online temporal calibration for camera-imu systems," in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 5709–5716.
- [14] F. Pang and T. Wang, "Stereo-inertial pose estimation and online sensors extrinsic calibration," in *Robotics and Biomimetics (ROBIO), 2016 IEEE International Conference on*. IEEE, 2016, pp. 1377–1382.
- [15] A. W. Stroupe, M. C. Martin, and T. Balch, "Merging gaussian distributions for object localization in multi-robot systems," in *Experimental Robotics VII*. Springer, 2001, pp. 343–352.
- [16] L. E. Clement, V. Peretroukhin, J. Lambert, and J. Kelly, "The battle for filter supremacy: A comparative study of the multi-state constraint kalman filter and the sliding window filter," in *Computer and Robot Vision (CRV), 2015 12th Conference on*. IEEE, 2015, pp. 23–30.
- [17] J. Montiel, J. Civera, and A. J. Davison, "Unified inverse depth parametrization for monocular slam," *analysis*, vol. 9, p. 1, 2006.
- [18] 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.
- [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*, vol. 35, no. 10, pp. 1157–1163, 2016.
- [20] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct ekf-based approach," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 298–304.
- [21] 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.
- [22] S. M. Weiss, "Vision based navigation for micro helicopters," Ph.D. dissertation, 2012.