

# Visual-Inertial Direct SLAM

Alejo Concha<sup>1</sup>, Giuseppe Loianno<sup>2</sup>, Vijay Kumar<sup>2</sup>, and Javier Civera<sup>1</sup>

**Abstract**—The so-called direct visual SLAM methods have shown a great potential in estimating a semidense or fully dense reconstruction of the scene, in contrast to the sparse reconstructions of the traditional feature-based algorithms. In this paper, we propose for the first time a direct, tightly-coupled formulation for the combination of visual and inertial data.

Our algorithm runs in real-time on a standard CPU. The processing is split in three threads. The first thread runs at frame rate and estimates the camera motion by a joint non-linear optimization from visual and inertial data given a semidense map. The second one creates a semidense map of high-gradient areas only for camera tracking purposes. Finally, the third thread estimates a fully dense reconstruction of the scene at a lower frame rate. We have evaluated our algorithm in several real sequences with ground truth trajectory data, showing a state-of-the-art performance.

## I. INTRODUCTION

Simultaneous Localization and Mapping (usually referred by its acronym SLAM) has become a key technology for several potential applications like robotics, autonomous cars and augmented and virtual reality. Its goal is to estimate, from a stream of sensor data, a model of the surroundings of the sensor and its egomotion with respect to it.

In the latest decades, there has been an intense research on the use of visual sensors for SLAM, but its application has been limited by the sparsity of the scene estimation. The traditional techniques –denoted as feature-based– rely on the correspondences between image point features; that can only be reliably established for the image points that present a high degree of saliency. [1], [2] are two open-source examples of such feature-based monocular SLAM systems.

Recently, [3]–[5], have developed algorithms for real-time, online and dense scene reconstruction from monocular images. These algorithms are based, in contrast to the feature-based methods, on the direct minimization of the photometric pixel values and on a regularization function. These so-called direct methods are able to produce semidense and dense reconstructions and open the doors to a wider applicability of visual SLAM. On the other hand, their maturity is still low. For example, it was recently shown that their current accuracy is lower than the feature-based techniques [6].

Our contribution within direct SLAM is a novel algorithm for the fusion of inertial and photometric data. To our knowledge, this is the first algorithm that integrates the inertial information in direct visual SLAM in a jointly manner. The key advantage of the mono(cular)-inertial configuration

with respect to the pure monocular one is the recovery of the real scale of the scene and the camera motion. This might be essential in certain applications, for example in robotics and control. We also believe that the compactness, low cost and low consumption of the hardware and the dense reconstructions we achieve might be relevant for constrained systems like quadrotors. We validated the accuracy and real-time performance of our algorithm in real visual-inertial sequences with Vicon ground truth for the trajectory.

The rest of the paper is organized as follows. Section II describes the related work. Section III presents our visual-inertial fusion for camera tracking with the Jacobians of the optimization approach detailed in the appendix VII. Section IV describes the mapping algorithm we use in our paper. Finally, section V shows the experimental results and section VI presents the conclusions.

## II. RELATED WORK

SLAM from a monocular camera has a scale ambiguity that limits its use in certain applications. A usual configuration to overcome that is a stereo camera [7], [8]. However, the accuracy of the scale is limited by the ratio between the scene depth and the stereo baseline. Feature-based visual motion estimation has been also combined with wheel odometry (e.g., in [9]) and non-holonomic constraints [10], but such combination is constrained to terrestrial robots. Visual-GPS SLAM (e.g., [11]) also resolves the scale ambiguity, but it is constrained to outdoors.

The visual-inertial fusion also resolves the scale ambiguity in pure monocular SLAM. Such fusion is usually divided into two classes. Loosely-coupled approaches model the two inputs as independent (e.g., [12]). Tightly-coupled ones, that estimate jointly all the states, have proven a superior performance [13], [14] at an additional complexity in the optimization –but still within real time in modern computers. The recent works [15], [16] present two loosely-coupled direct stereo visual-inertial systems. In the first system, the inertial and stereo measurements are loosely-coupled fused with respect to the last keyframe, and in the latest one, they are fused in a –loosely coupled– pose-graph manner. Our algorithm is the first tightly-coupled one for direct SLAM methods. Also differently from [15], [16], we use monocular vision instead of stereo vision.

## III. VISUAL-INERTIAL JOINT OPTIMIZATION

The main idea is to take advantage of the fast rate of the Inertial Measurement Unit (IMU) that provides inertial data and propose a joint optimization where the integration of IMU measurements is able to infer additional constraints

<sup>1</sup>The authors are with the I3A, Universidad de Zaragoza, Spain {alejocb, jcivera}@unizar.es

<sup>2</sup>The authors are with the GRASP Lab, University of Pennsylvania, 3330 Walnut Street, 19104 Philadelphia, USA. {loiannog, kumar}@seas.upenn.edu

between camera poses. The IMU can contribute to speed up the estimation and obtain the velocity and the absolute scale information, which is not available using only the camera information.

### A. Integration of IMU Measurements

Let us consider, without loss of generality, that the IMU reference frame is coincident with the camera frame.

The rotation  $\mathbf{R}_j^w$ , translation  $\mathbf{t}_j^w$  and velocity  $\mathbf{v}_j^w$  of the current frame  $j$  in the world reference frame  $w$  are calculated from a previous frame  $i$  using an Euler forward integration.

$$\mathbf{R}_j^w = \mathbf{R}_i^w \mathbf{R}_j^i \quad (1)$$

$$\mathbf{v}_j^w = \mathbf{v}_i^w + \mathbf{v}_{ij}^w \quad (2)$$

$$\mathbf{t}_j^w = \mathbf{t}_i^w + \mathbf{t}_{ij}^w \quad (3)$$

Where  $\mathbf{R}_j^i$  is the relative rotation between the frames  $i$  and  $j$ ,  $\mathbf{v}_{ij}^w$  is the incremental velocity from  $i$  to  $j$  and  $\mathbf{t}_{ij}^w$  is the translation vector between  $i$  and  $j$ . These three integrate the IMU measurements as follows

$$\mathbf{R}_j^i = \prod_{p=k}^{k+N-1} \exp_{SO(3)}([\boldsymbol{\omega}(p) + \mathbf{b}_\omega(p)]^\wedge T) \quad (4)$$

$$\mathbf{v}_{ij}^w = \sum_{p=k}^{k+N-1} (\mathbf{R}_p^w (\mathbf{a}(p) + \mathbf{b}_a(p)) - \mathbf{g}) T \quad (5)$$

$$\begin{aligned} \mathbf{t}_{ij}^w &= N \mathbf{v}_i^w T + \\ &\frac{1}{2} \sum_{p=k}^{k+N-1} (2(k+N-1-p)+1) (\mathbf{R}_p^w (\mathbf{a}(p) + \mathbf{b}_a(p)) - \mathbf{g}) T^2 \end{aligned} \quad (6)$$

The IMU measurements at each step  $p$  are the angular velocity  $\boldsymbol{\omega} \in \mathbb{R}^3$  and the linear acceleration  $\mathbf{a} \in \mathbb{R}^3$  in the local frame. These measurements are affected by the biases  $\mathbf{b}_\omega \in \mathbb{R}^3$  –for the angular velocity– and  $\mathbf{b}_a \in \mathbb{R}^3$  –for the linear acceleration.  $T$  denotes the time step size,  $\mathbf{g} \in \mathbb{R}^3$  the gravity vector,  $k$  the time step of frame  $i$  and  $k+N$  the time step of frame  $j$ . The  $\exp_{SO(3)}$  operator maps an element of  $so(3)$  to an element of  $SO(3)$ . The inverse operation is denoted with  $\log_{SO(3)}$ .

The *hat* operator  $(\wedge)$  is used to convert a  $3 \times 1$  vector in a an element of of  $so(3)$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \wedge = \begin{bmatrix} 0 & -x_3 & x_2 \\ x_3 & 0 & -x_1 \\ -x_2 & x_1 & 0 \end{bmatrix}, \quad (7)$$

while the opposite operation is performed using the *vee* operator  $(\vee)$

$$\begin{bmatrix} 0 & -x_3 & x_2 \\ x_3 & 0 & -x_1 \\ -x_2 & x_1 & 0 \end{bmatrix} \vee = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}. \quad (8)$$

The IMU biases are modeled as random walk processes with variances  $\boldsymbol{\eta}_a$  and  $\boldsymbol{\eta}_\omega$  for the acceleration  $\mathbf{a}$  and angular velocity  $\boldsymbol{\omega}$  respectively

$$\mathbf{b}_a(k+1) = \mathbf{b}_a(k) + T \boldsymbol{\eta}_a, \quad (9)$$

$$\mathbf{b}_\omega(k+1) = \mathbf{b}_\omega(k) + T \boldsymbol{\eta}_\omega. \quad (10)$$

They frame  $i$  is assumed to be fixed and only the current frame  $j$  is optimized. This assumption is made due to the computational cost of direct visual SLAM. Specifically, without this assumption, the Jacobians of the photometric error have to be computed in every iteration and we could not take advantage of the inverse compositional approach.

### B. State vector and residual

The state vector  $\mathbf{x}$  is composed of 15 variables:

$$\mathbf{x} = \left[ \left[ \log_{SO(3)}(\mathbf{R}_j^w)^\vee \right]^\top \mathbf{t}_j^w{}^\top \mathbf{v}_j^w{}^\top \mathbf{b}_\omega^\top \mathbf{b}_a^\top \right]^\top, \quad (11)$$

The residual  $\mathbf{r}$  to minimize is divided into two parts, referred to as the photometric  $\mathbf{r}_{ph}$  and IMU –inertial–  $\mathbf{r}_{imu}$  residual

$$\mathbf{r} = \left[ \mathbf{r}_{ph}^\top \mathbf{r}_{imu}^\top \right]^\top \quad (12)$$

The IMU residual is

$$\mathbf{r}_{imu} = \begin{bmatrix} \log_{SO(3)} \left( \mathbf{R}_j^i{}^\top (\mathbf{R}_i^w)^\top \mathbf{R}_j^w \right)^\vee \\ \left( \mathbf{t}_j^w - \mathbf{t}_{ij}^w - \mathbf{t}_i^w \right) \\ \left( \mathbf{v}_j^w - \mathbf{v}_{ij}^w - \mathbf{v}_i^w \right) \\ \mathbf{b}_\omega^i - \mathbf{b}_\omega \\ \mathbf{b}_a^i - \mathbf{b}_a \end{bmatrix} \quad (13)$$

The intensity residual is

$$\mathbf{r}_{ph} = \sum_{k=1}^n \mathbf{I}_i \left( \pi \left( \mathbf{T}_w^i \mathbf{p}_t^k \right) \right) - \mathbf{I}_j \left( \pi \left( \mathbf{T}_w^j \mathbf{p}_t^k \right) \right) \quad (14)$$

$$\mathbf{T}_w^j = \begin{bmatrix} \mathbf{R}_w^j & \mathbf{t}_w^j \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}, \quad (15)$$

where  $\pi$  is the pinhole camera model and  $\mathbf{p}_t$  is each one of the  $n$  3D points tracked from  $i$  to  $j$ .  $\mathbf{I}_i$  is the image taken at step  $k$  and  $\mathbf{I}_j$  the image taken at step  $k+N$ .

### C. Visual-Inertial Direct Tracking

The visual-inertial thread minimizes the residual of eq (12) using standard Gauss-Newton optimization

$$\mathbf{J}^\top \mathbf{A} \mathbf{J} \boldsymbol{\delta} = -\mathbf{J}^\top \mathbf{A} \cdot \mathbf{r}, \quad (16)$$

where  $\boldsymbol{\delta}$  are the updates for the pose and the IMU  $\boldsymbol{\delta} = [\boldsymbol{\delta}_P^\top, \boldsymbol{\delta}_I^\top]^\top$  and  $\mathbf{A}$  is the inverse covariance. The general Jacobian  $\mathbf{J}$  is composed of the Jacobian of the IMU residual with respect to the IMU parameters (biases and velocity)  $\mathbf{J}_I^{r_{imu}}$ , the Jacobian of the IMU error with respect to the current pose of the camera  $\mathbf{J}_P^{r_{imu}}$  and the Jacobian of the photometric error with respect to the current pose  $\mathbf{J}_P^{r_{ph}}$ . Note than  $\mathbf{J}_I^{r_{ph}} = 0$ .  $\mathbf{J}_P^{r_{ph}}$  will be defined in section III-D and  $\mathbf{J}_I^{r_{imu}}$  and  $\mathbf{J}_P^{r_{imu}}$  will be defined in the appendix (see section VII).

The full Jacobian is formed as follows

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_P^{r_{ph}} & 0 \\ \mathbf{J}_P^{r_{imu}} & \mathbf{J}_I^{r_{imu}} \end{bmatrix}, \quad (17)$$

and reparametrized for convenience as

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_P^r & \mathbf{J}_I^r \end{bmatrix}, \quad (18)$$

$$\text{where } \mathbf{J}_P^r = \begin{bmatrix} \mathbf{J}_P^{rph} \\ \mathbf{J}_P^{rimu} \end{bmatrix} \text{ and } \mathbf{J}_I^r = \begin{bmatrix} \mathbf{J}_I^{rph} \\ \mathbf{J}_I^{rimu} \end{bmatrix}.$$

To solve it, we perform Gauss-Newton optimization and use the Schur complement to accelerate it. The Jacobian  $\mathbf{J}$  yields the following Hessian approximation

$$\mathbf{H} = \begin{bmatrix} \mathbf{J}_P^r \top \mathbf{\Lambda} \mathbf{J}_P^r & \mathbf{J}_P^r \top \mathbf{\Lambda} \mathbf{J}_I^r \\ \mathbf{J}_P^r \mathbf{\Lambda} \mathbf{J}_I^r \top & \mathbf{J}_I^r \top \mathbf{\Lambda} \mathbf{J}_I^r \end{bmatrix} \quad (19)$$

The linear system is therefore

$$\begin{bmatrix} \mathbf{H}_{PP}^r & \mathbf{H}_{PI}^r \\ \mathbf{H}_{PI}^r \top & \mathbf{H}_{II}^r \end{bmatrix} \begin{bmatrix} \boldsymbol{\delta}_P \\ \boldsymbol{\delta}_I \end{bmatrix} = \begin{bmatrix} \mathbf{B}_P \\ \mathbf{B}_I \end{bmatrix}, \quad (20)$$

where  $\mathbf{B}_I = -\mathbf{J}_I^r \top \mathbf{\Lambda} \mathbf{r}$  and  $\mathbf{B}_P = -\mathbf{J}_P^r \top \mathbf{\Lambda} \mathbf{r}$ ,  $\boldsymbol{\delta}_I$  and  $\boldsymbol{\delta}_P$  are the updates for the IMU and pose parameters respectively.

The system is solved using the Schur complement rearrangement. The lower part of the system of the eq. (20) results in the smaller equation system for the IMU update

$$(\mathbf{H}_{II} - \mathbf{H}_{PI}^r \top \mathbf{H}_{PP}^{-1} \mathbf{H}_{PI}^r) \boldsymbol{\delta}_I = \mathbf{B}_I - \mathbf{H}_{PI}^r \top \mathbf{H}_{PP}^{-1} \mathbf{B}_P, \quad (21)$$

This linear system is typically solved using Cholesky factorization or SVD decomposition. Once the update for the IMU parameters is found we can solve the upper system for the pose update:

$$\boldsymbol{\delta}_P = \mathbf{H}_{PP}^{-1} \mathbf{B}_P - \mathbf{H}_{PP}^{-1} \mathbf{H}_{PI}^r \boldsymbol{\delta}_I. \quad (22)$$

#### D. Visual Tracking Jacobian

The transformation from the current camera frame to the global reference frame  $\mathbf{T}_w^j$  is estimated based on the photometric reprojection error  $\mathbf{r}_{ph}$  using the inverse compositional approach [17]. The photometric error for the 3D points  $\mathbf{p}_t$  is

$$\mathbf{r}_{ph} = \sum_{k=1}^n \left( \mathbf{I}_i \left( \pi \left( \hat{\mathbf{T}} \mathbf{T}_w^i \mathbf{p}_t^k \right) \right) - \mathbf{I}_j \left( \pi \left( \mathbf{T}_w^j \mathbf{p}_t^k \right) \right) \right). \quad (23)$$

The tracking thread only uses a subset of image points, composed of high-gradient points –and superpixel contours if desired. Section IV details how the 3D position for those is obtained. We seek to estimate the transformation  $\hat{\mathbf{T}}$  from the closest keyframe  $\mathbf{I}_i$  to the current frame  $\mathbf{I}_j$ . The seed for  $\mathbf{T}_w^j$  comes from the raw inertial readings before the optimization.

For the optimization we use a minimal parametrization for the camera pose. The rotation  $\mathbf{R}$  is mapped to the tangent space  $so(3)$  of the rotation group  $SO(3)$  at the identity. The increments for the camera pose –the angular increment  $\delta\omega$  and the increment for the translation  $\delta\mathbf{t}$ – are defined as follows:

$$\hat{\mathbf{T}} = \begin{bmatrix} \exp_{SO(3)}([\delta\omega]^\wedge) & \delta\mathbf{t} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}, \quad (24)$$

The update for the current camera pose  $\mathbf{T}_w^j$  is as follows

$$\mathbf{T}_w^j = \hat{\mathbf{T}}^{-1} \mathbf{T}_w^i. \quad (25)$$

The Jacobian of the photometric residual  $\mathbf{J}_P^{rph}$  is obtained via the chain rule:

$$\mathbf{J}_P^{rph} = \frac{\partial \mathbf{r}_{ph}}{\partial \boldsymbol{\epsilon}} = \mathbf{J}_F^{rph} \mathbf{J}_{T_w^i}^F \mathbf{J}_{\boldsymbol{\epsilon}}^{T_w^i}, \quad (26)$$

where  $\mathbf{J}_F^{rph}$  are the image gradients of the reference keyframe,  $\mathbf{J}_{T_w^i}^F$  is the derivative of the projection model with respect to  $\mathbf{T}_w^i$  and  $\mathbf{J}_{\boldsymbol{\epsilon}}^{T_w^i}$  is the derivative of the transformation  $\mathbf{T}_w^i$  with respect to the motion  $\boldsymbol{\epsilon} = (\delta\omega, \delta\mathbf{t})$ .

Using the inverse compositional approach the Jacobians are always referred to the last keyframe, not being updated until a new frame becomes a keyframe  $-\mathbf{T}_w^i$ ,  $\mathbf{p}_t$  and  $\mathbf{J}_F^{rph}$  are constant during the optimization. This accelerates significantly the camera tracking.

To bootstrap our system we integrate the inertial measurements and we estimate the first map –obtaining the real scale– using the approach of section IV. After the initialization, we apply the joint optimization explained in the section III.

## IV. MAPPING

### A. Semidense Mapping

Rapid camera motions require high-frequency map updates for the camera to be tracked successfully. Similarly to [18]–[20], our system maintains a semidense map of high-gradient pixels that is quickly updated and serves for camera tracking. The semidense map is not only used for tracking, but also for the estimation of the fully dense map described in section IV-B.

The inverse depth  $\rho$  for every high-gradient pixel  $\mathbf{u}^*$  in a keyframe  $\mathbf{I}_i$  is estimated by minimizing its photometric error  $\mathbf{r}_{ph}^o$  with respect to several overlapping views  $\mathbf{I}_o$ . The specific optimization is

$$\hat{\rho} = \arg \min_{\rho} \mathbf{r}_{ph}^o, \quad (27)$$

with

$$\mathbf{r}_{ph}^o = \sum_o \|(\mathbf{I}_i(\mathbf{s}_{u^*}) - \mathbf{I}_o(G(\mathbf{s}_{u^*}, \mathbf{T}_i, \mathbf{T}_o, \rho)))\|_2^2. \quad (28)$$

$\mathbf{s}_{u^*}$  are the pixel coordinates of the template (7x7 pixels square) around the pixel  $\mathbf{u}^*$  and  $G$  is the function that backprojects the template from the new keyframe  $\mathbf{I}_i$  to the 3D world and projects it back to each overlapping image  $\mathbf{I}_o$ .

The reader is referred to [20] for more details on the semidense mapping algorithm used in this paper.

### B. Superpixels Mapping

We use the algorithm in [20] to extract the piecewise planar textureless regions to fill the semidense map of the previous section. First, we segment a keyframe  $\mathbf{I}_i$  into a set of superpixels  $\mathcal{S}_k = \{\mathbf{s}_1, \dots, \mathbf{s}_i, \dots, \mathbf{s}_m\}$  using the algorithm of [21].

Each 3D point  $\mathbf{p}_t^k$  from the semidense map is then projected on the keyframe  $u^* = F(\mathbf{T}_w^i \mathbf{p}_t^k)$ . The 3D points  $\mathbf{p}_t^k$  are assigned to the superpixels if their projections  $\mathbf{u}^*$  lie within a threshold  $\xi$  of the superpixel contour.

The 3D points associated to the contour of every superpixel  $\mathbf{p}_t^k \in \mathcal{C}(\mathbf{s}_i)$  are used to robustly fit a plane  $\pi_i$  using singular value decomposition.

Again, we refer the reader to [20] for more details on the 3D superpixels estimation.

## V. EXPERIMENTS

### A. Experimental setup

The experimental tests have been performed in the GRASP Lab [22] at The University of Pennsylvania. The considered working area is a volume of  $5 \times 4 \times 5$  m<sup>3</sup>. The Vicon motion capture system<sup>1</sup>, composed of 20 T040 cameras, provides a ground truth for our rotation and translation estimates [22].

The camera-IMU system setup is composed by Matrix Vision mvBlueFOX-MLC<sup>2</sup> camera and a Microstrain 3DM-GX4-25 IMU<sup>3</sup> as shown in Figure 1. The image processing runs on  $752 \times 480$  pixels at 30 Hz, whereas the IMU rate has been set to 100 Hz. The entire solution has been developed using the ROS<sup>4</sup> framework. The visual-inertial direct tracking thread (section III-C) runs at 30 Hz while the 2 mapping threads (section IV) run at a lower rate. The computer used is a 3.5 GHz Intel Core i7-3770K processor and 8.0 GB of RAM memory.

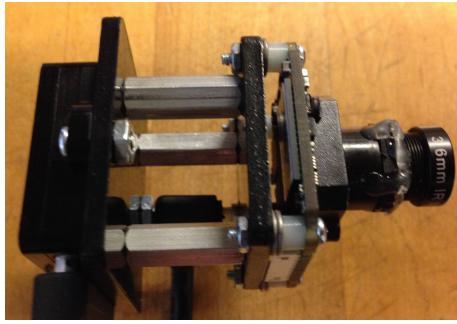


Fig. 1: The sensor setup used for the experiments.

### B. Results

We have evaluated our algorithm in 10 different sequences of visual and inertial data. The sequences were recorded with a hand-held camera while the first author walked around the room. Their duration goes from 20 to 55 seconds.

Table I shows the mean translational errors. Notice how the errors are of the same order than those reported for direct SLAM methods (for example, LSD-SLAM in [19] or in the results of [6]). But in our case we are able to recover the real scale of the scene and the trajectory thanks to the inertial fusion. For further details, figures 2 and 3 show respectively the angular and translational errors for each coordinate on the 10 experiments.

Experiment	1	2	3	4	5	6	7	8	9	10
Error [cm]	6.6	6.3	11.6	6.3	9.3	10.6	7.4	3.2	10.8	6.2

TABLE I: Mean trajectory error.

Figure 4 shows the computational cost of our algorithm –vertical axis– for every frame –horizontal axis– and for every sequence –in different colors. 99% of the frames have

<sup>1</sup><http://www.vicon.com>

<sup>2</sup><http://www.matrix-vision.com/>

<sup>3</sup><http://www.microstrain.com/>

<sup>4</sup><http://www.ros.org/>

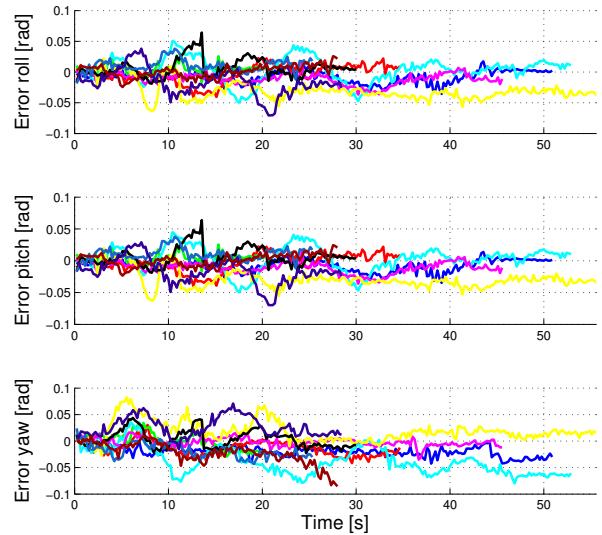


Fig. 2: Angular error for our 10 sequences. Each color is a different sequence.

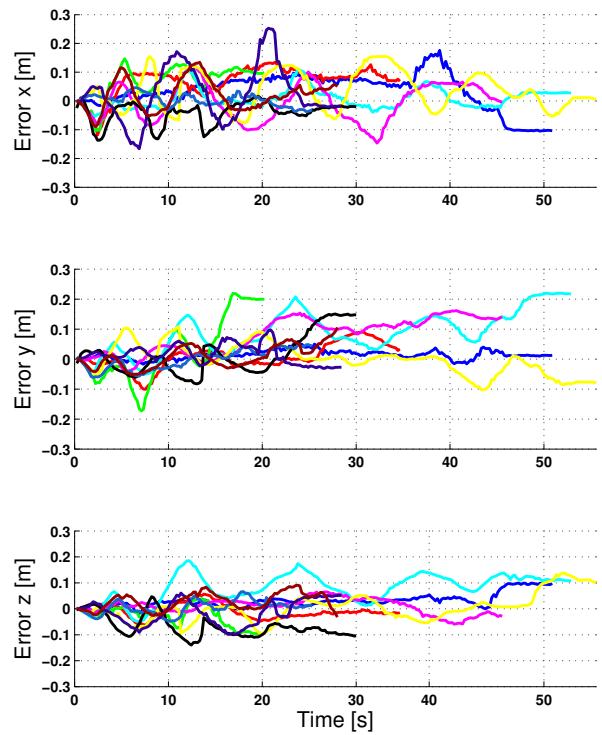


Fig. 3: Translational error for our 10 sequences. Each color is a different sequence.

a computational cost under 33 milliseconds, which is our frame period. We attribute the cost spikes in our results to other high priority tasks on the processor (we do not use a real-time OS). In any case, the sequences were played in real time, showing that our algorithm is resilient to a few lost frames.

Figure 5 shows some results of our algorithm. The top row

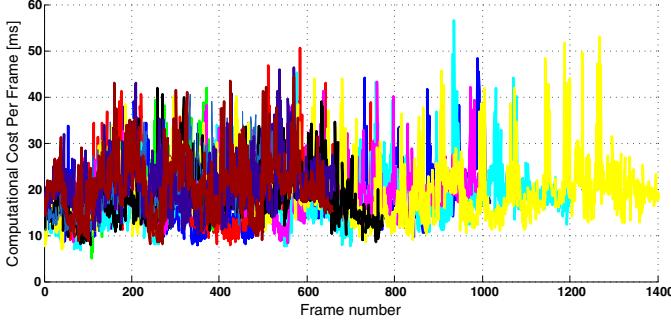


Fig. 4: Computational cost for our 10 sequences. Each color is a different sequence.

shows four sample images, and the bottom row several views of the estimated reconstruction and the camera trajectory in red.

Finally, figures 6 and 7 illustrates the difference between the semidense and dense maps estimated by our algorithm. The top row in both images displays again sample frames from the sequence, and the bottom row a close view of the semidense and dense maps estimated by our algorithm. Notice in the figures 6(e) and 7(e) that the semidense map only contains the high gradient regions of the images, being of little use for some applications like robotic navigation. Notice the higher point density of the dense map in the figures 6(f) and 7(f), where the textureless areas have been reconstructed using a piecewise planar prior.

As shown in [23], a TV-regularization of a photometric map produces low quality results in low texture scenes. We also run a TV-regularization on the experiments of this paper and the qualitative accuracy of the maps was lower than the piecewise-planar approach we use.

## VI. CONCLUSIONS

We have presented in this paper a tightly-coupled visual-inertial SLAM algorithm for real-time, online, dense scene reconstruction and camera tracking. Our main contribution is the joint optimization of the IMU measurements and visual data using direct methods. Our full SLAM system is divided into 3 threads; a low-cost mapping thread produces semi-dense reconstructions for real-time camera tracking, a camera tracking thread responsible for the visual-inertial fusion, and a high-quality dense mapping thread producing dense maps at low frame rate. Our dense mapping algorithm uses a piecewise planar approximation based on 3D superpixels.

We have validated the accuracy and real-time performance of our approach with several real sequences, comparing the estimated pose with respect to a Vicon motion capture system. To our knowledge this is the first real-time, tightly-coupled, direct visual SLAM approach. Given the low cost and compactness of the mono-inertial hardware setting, our approach can be relevant for low-cost, low-consumption, lightweight robots like some quadrotors, or for augmented and virtual reality applications.

## VII. APPENDIX

### A. Background

The following result relates infinitesimal increments in  $\text{so}(3)$  with right hand-multiplications [14]:

$$\exp_{SO(3)}([\boldsymbol{\theta} + \delta\boldsymbol{\theta}]^\wedge) = \exp_{SO(3)}([\boldsymbol{\theta}]^\wedge) \exp_{SO(3)}([\mathbf{J}_r(\boldsymbol{\theta})\delta\boldsymbol{\theta}]^\wedge) \quad (29)$$

where  $\mathbf{J}_r(\boldsymbol{\theta})$  is the  $SO(3)$  Jacobian.

$$\mathbf{J}_r(\boldsymbol{\theta}) = \mathbf{I} - \frac{1 - \cos(\|\boldsymbol{\theta}\|)}{\|\boldsymbol{\theta}\|^2} \boldsymbol{\theta}^\wedge + \frac{\|\boldsymbol{\theta}\| - \sin(\|\boldsymbol{\theta}\|)}{\|\boldsymbol{\theta}\|^3} (\boldsymbol{\theta}^\wedge)^2 \quad (30)$$

A similar first-order approximation holds for the logarithm as shown in [14]

$$\log_{SO(3)}\left(\exp_{SO(3)}([\boldsymbol{\theta}]^\wedge) \exp_{SO(3)}([\delta\boldsymbol{\theta}]^\wedge)\right)^\vee = \boldsymbol{\theta} + \mathbf{J}_r(\boldsymbol{\theta})^{-1} \delta\boldsymbol{\theta}. \quad (31)$$

The following relation holds for  $SO(3)$

$$\exp_{SO(3)}([\boldsymbol{\theta}]^\wedge) \mathbf{R} = \mathbf{R} \exp_{SO(3)}\left(\left[\mathbf{R}^\top \boldsymbol{\theta}\right]^\wedge\right) \quad (32)$$

### B. Jacobians

#### a) Jacobians of the IMU residual with respect to the camera pose $J_P^{r_{imu}}$ :

Definition of the rotation error with respect to the increment

$$r_{imu}[1]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right) = \log_{SO(3)}\left(\mathbf{R}_j^{i^\top} \mathbf{R}_i^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)^\vee = \\ \log_{SO(3)}\left(\left(\prod_{p=k}^{k+N-1} \exp_{SO(3)}([\omega(p) + b\omega(p)]^\wedge T)\right)^\top \mathbf{R}_i^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)^\vee \quad (33)$$

By  $r_{imu}[1]$  we denote the first term of the IMU error (the rotation error in this case). Rearranging terms and using eq. (31)

$$r_{imu}[1]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right) = r_{imu}[1](\mathbf{R}_j^w) \\ + \mathbf{J}_r(r_{imu}[1](\mathbf{R}_j^w))^{-1} \delta\boldsymbol{\theta}_j \quad (34)$$

$$\frac{\partial r_{imu}[1]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)}{\partial \delta\boldsymbol{\theta}_j} = \mathbf{J}_r(r_{imu}[1](\mathbf{R}_j^w))^{-1}. \quad (35)$$

$$\begin{aligned} \text{Note} &\quad \text{that} & \frac{\partial r_{imu}[2]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)}{\partial \delta\boldsymbol{\theta}_j} &= \\ \frac{\partial r_{imu}[3]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)}{\partial \delta\boldsymbol{\theta}_j} &= \frac{\partial r_{imu}[4]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)}{\partial \delta\boldsymbol{\theta}_j} &= \\ \frac{\partial r_{imu}[5]\left(\mathbf{R}_j^w \exp_{SO(3)}([\delta\boldsymbol{\theta}_j]^\wedge)\right)}{\partial \delta\boldsymbol{\theta}_j} &= \mathbf{0}_{3 \times 3}. \end{aligned}$$

For the translation, we define the translation error  $r_{imu}[2]$  with respect to the  $SE(3)$  increment

$$r_{imu}[2]\left(\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w\right) = \left(\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w - \mathbf{t}_{ij} - \mathbf{t}_i^w\right) \quad (36)$$

$$\frac{\partial r_{imu}[2](\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w)}{\partial \delta\mathbf{t}_j^w} = \mathbf{R}_j^w \quad (37)$$

$$\begin{aligned} \text{Note} &\quad \text{that} & \frac{\partial r_{imu}[1](\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w)}{\partial \delta\mathbf{t}_j^w} &= \frac{\partial r_{imu}[3](\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w)}{\partial \delta\mathbf{t}_j^w} &= \\ \frac{\partial r_{imu}[4](\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w)}{\partial \delta\mathbf{t}_j^w} &= \frac{\partial r_{imu}[5](\mathbf{t}_j^w + \mathbf{R}_j^w \delta\mathbf{t}_j^w)}{\partial \delta\mathbf{t}_j^w} &= \mathbf{0}_{3 \times 3}. \end{aligned}$$

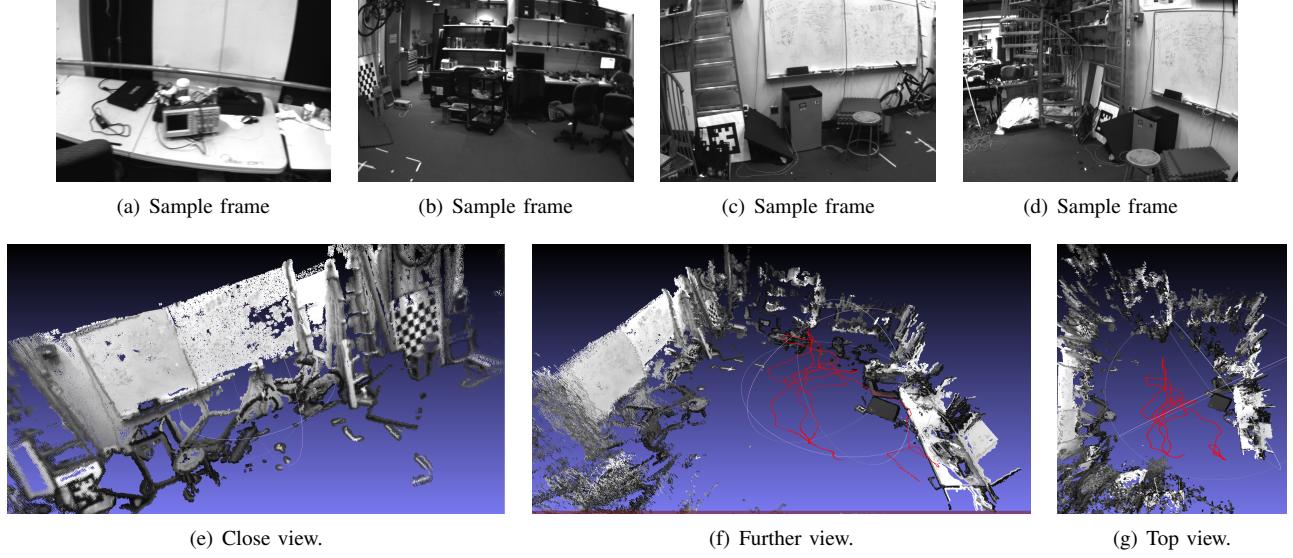


Fig. 5: Sample images and views of the reconstructed scene for the largest of our experiments, covering almost the whole GRASP Lab Vicon room. The estimated trajectory of the camera is shown in red. Notice, qualitatively, the accuracy and density of the estimated map.

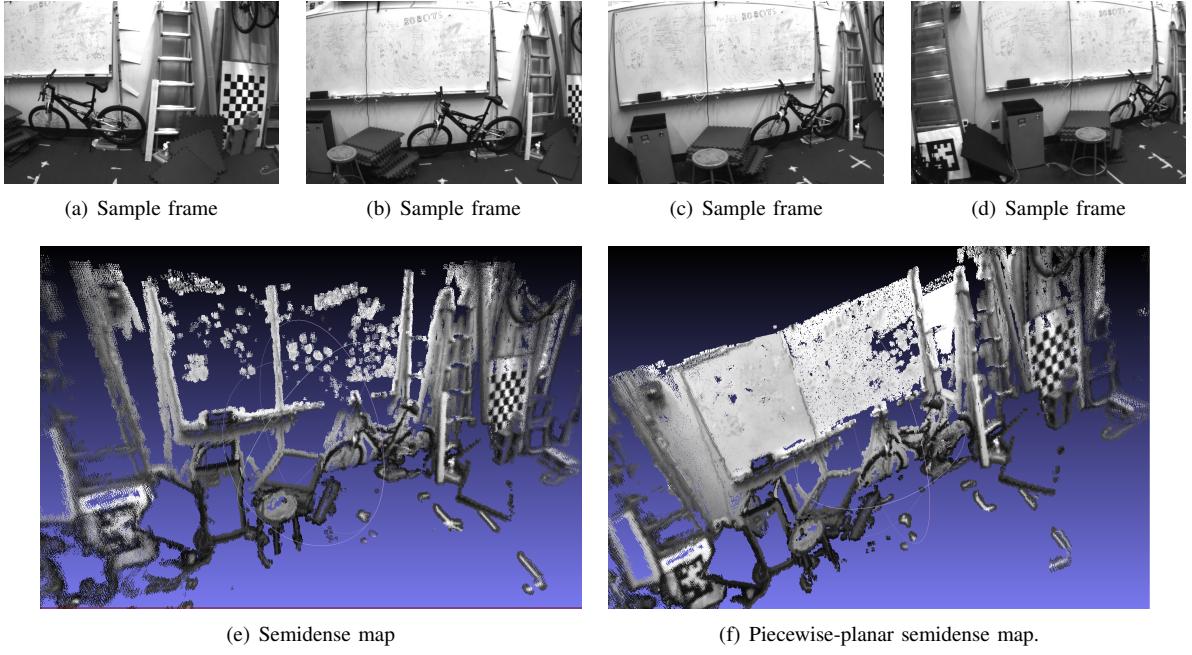


Fig. 6: Sample frames for one of our experiments, and close-ups of the semidense map used for camera tracking and the piecewise planar dense map. Notice how, in the latest, the low-textured areas have been correctly reconstructed.

**b) Jacobian of the IMU residual with respect to IMU parameters  $J_I^{imu}$ :**

- Jacobian with respect to the velocity in the frame j

$$\frac{\partial \mathbf{r}_{imu}[3](\mathbf{v}_j^w + \delta \mathbf{v}_j^w)}{\partial \delta \mathbf{v}_j^w} = \mathbf{I}_{3 \times 3}. \quad (38)$$

Note that  $\frac{\partial \mathbf{r}_{imu}[1](\mathbf{v}_j^w + \delta \mathbf{v}_j^w)}{\partial \delta \mathbf{v}_j^w} = \frac{\partial \mathbf{r}_{imu}[2](\mathbf{v}_j^w + \delta \mathbf{v}_j^w)}{\partial \delta \mathbf{v}_j^w} =$

$$\frac{\partial \mathbf{r}_{imu}[4](\mathbf{v}_j^w + \delta \mathbf{v}_j^w)}{\partial \delta \mathbf{v}_j^w} = \frac{\partial \mathbf{r}_{imu}[5](\mathbf{v}_j^w + \delta \mathbf{v}_j^w)}{\partial \delta \mathbf{v}_j^w} = \mathbf{0}_{3 \times 3}.$$

- Jacobian of the rotation error  $\mathbf{r}_{imu}[1]$  with respect to the angular velocity bias  $\mathbf{b}_\omega$

$$\begin{aligned} \mathbf{r}_{imu}[1](\mathbf{b}_\omega + \delta \mathbf{b}_\omega) &= \\ \log_{SO(3)} \left( \left( \prod_{p=k}^{k+N-1} \exp_{SO(3)}([\mathbf{\omega}(p) + \mathbf{b}_\omega(p) + \delta \mathbf{b}_\omega] \wedge T) \right)^T \mathbf{R}_i^w \top \mathbf{R}_j^w \right)^V. \end{aligned} \quad (39)$$

Applying eq. (29) in the increments and using eq. (32) to move the all increment terms to the

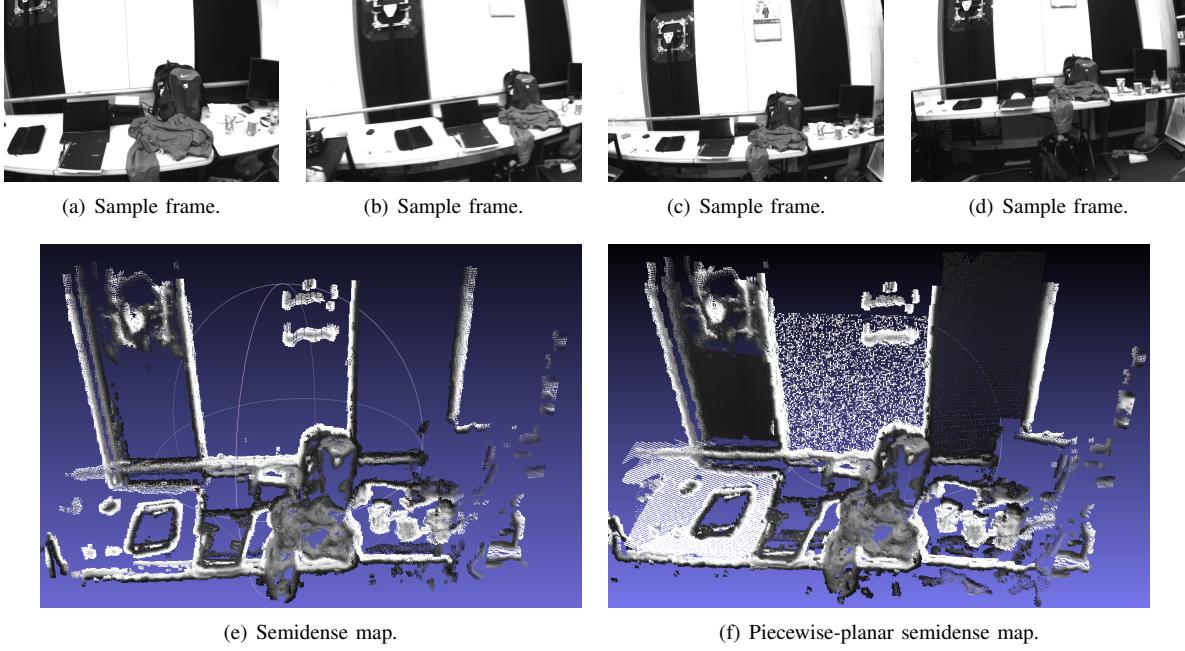


Fig. 7: Sample frames for one of our experiments, and close-ups of the semidense map used for camera tracking and the piecewise planar dense map. Notice how, in the latest, the low-textured areas have been correctly reconstructed.

right side, assuming that the increments are small  $\exp_{SO(3)}([\mathbf{a}]^\wedge) \exp_{SO(3)}([\mathbf{b}]^\wedge) = \exp_{SO(3)}([\mathbf{a} + \mathbf{b}]^\wedge)$ , rearranging terms and applying eq. (31) we obtain

$$\begin{aligned} \mathbf{r}_{imu}[1](\mathbf{b}_w + \delta\mathbf{b}_\omega) &= \\ \mathbf{r}_{imu}[1] - \mathbf{J}_r(\mathbf{r}_{imu}[1])^{-1} * \exp_{SO(3)}([\mathbf{r}_{imu}[1]]^\wedge)^\top \mathbf{A} \delta\mathbf{b}_w &\end{aligned} \quad (40)$$

where

$$\mathbf{A} = \sum_{p=k}^{k+N-1} \left[ \left[ \prod_{m=p+1}^{k+N-1} \exp_{SO(3)}([\mathbf{\omega}(m) + \mathbf{b}_\omega(m)]^\wedge T) \right]^\top \mathbf{J}_r^p * T \right] \quad (41)$$

$$\frac{\partial \mathbf{r}_{imu}[1](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = -\mathbf{J}_r(\mathbf{r}_{imu}[1])^{-1} * \exp_{SO(3)}([\mathbf{r}_{imu}[1]]^\wedge)^\top \mathbf{A} \quad (42)$$

Where

$$\mathbf{J}_r^p = \mathbf{J}_r((\mathbf{\omega}(p) + \mathbf{b}_\omega(p))T).$$

- Jacobian of the position error  $\mathbf{r}_{imu}[2]$  with respect to the angular velocity bias  $\mathbf{b}_w$ . Applying eq. (29) in the increments and using eq. (32) to move all increment terms to the right side, assuming that the increments are small  $\exp_{SO(3)}([\mathbf{a}]^\wedge) \exp_{SO(3)}([\mathbf{b}]^\wedge) = \exp_{SO(3)}([\mathbf{a} + \mathbf{b}]^\wedge)$ , using a first order approximation for the exponential  $\exp_{SO(3)}([\mathbf{a}]^\wedge) = \mathbf{I}_{3 \times 3} + [\mathbf{a}]^\wedge$  and rearranging terms we obtain

$$\frac{\partial \mathbf{r}_{imu}[2](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = \frac{1}{2} \sum_{p=k}^{k+N-1} \mathbf{C}(\mathbf{a}(p) + \mathbf{b}_a(p))^\wedge \mathbf{B} * T^2 \quad (43)$$

$$\begin{aligned} \mathbf{B} &= \sum_{l=k}^{p-1} \left[ \left[ \prod_{m=l+1}^{p-1} \exp_{SO(3)}([\mathbf{\omega}(m) + \mathbf{b}_\omega(m)]^\wedge T) \right]^\top \mathbf{J}_r^l * T \right] \\ \mathbf{C} &= \prod_{q=k}^{p-1} (\mathbf{R}_i^w \exp_{SO(3)}([\mathbf{\omega}(q) + \mathbf{b}_\omega(q)]^\wedge T)) \end{aligned} \quad (44)$$

- Jacobian of the velocity error  $\mathbf{r}_{imu}[3]$  with respect to the angular velocity bias  $\mathbf{b}_w$ . Similarly to the previous Jacobian we obtain

$$\frac{\partial \mathbf{r}_{imu}[3](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = \sum_{p=k}^{k+N-1} \mathbf{D}(\mathbf{a}(p) + \mathbf{b}_a(p))^\wedge \mathbf{B} * T \quad (45)$$

$$\mathbf{D} = \prod_{q=k}^{p-1} \mathbf{R}_i^w \exp_{SO(3)}([\mathbf{\omega}(q) + \mathbf{b}_\omega(q)]^\wedge T)$$

- Jacobian of the angular velocity bias error  $\mathbf{r}_{imu}[4]$  with respect to the angular velocity bias  $\mathbf{b}_w$

$$\frac{\partial \mathbf{r}_{imu}[4](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = -\mathbf{I}_{3 \times 3} \quad (46)$$

- Jacobian of the acceleration bias error  $\mathbf{r}_{imu}[5]$  with respect to the angular velocity bias  $\mathbf{b}_\omega$

$$\frac{\partial \mathbf{r}_{imu}[5](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = \mathbf{0}_{3 \times 3} \quad (47)$$

- Jacobian of the rotation error  $\mathbf{r}_{imu}[1]$  with respect to the acceleration bias  $\mathbf{b}_a$

$$\frac{\partial \mathbf{r}_{imu}[1](\mathbf{b}_w + \delta\mathbf{b}_\omega)}{\partial \delta\mathbf{b}_\omega} = \mathbf{0}_{3 \times 3} \quad (48)$$

- Jacobian of the position error  $\mathbf{r}_{imu}[2]$  with respect to the acceleration bias  $\mathbf{b}_a$

$$\frac{\partial \mathbf{r}_{imu}[2](\mathbf{b}_a + \delta \mathbf{b}_a)}{\partial \delta \mathbf{b}_a} = -\frac{1}{2} \sum_{p=k}^{k+N-1} (2(k+N-1-p)+1) (\mathbf{R}_p^w) T^2 \quad (49)$$

- Jacobian of the velocity error  $\mathbf{r}_{imu}[3]$  with respect to the acceleration bias  $\mathbf{b}_a$

$$\frac{\partial \mathbf{r}_{imu}[3](\mathbf{b}_a + \delta \mathbf{b}_a)}{\partial \delta \mathbf{b}_a} = - \sum_{p=k}^{k+N-1} (\mathbf{R}_p^w) T, \quad (50)$$

- Jacobian of the angular velocity bias error  $\mathbf{r}_{imu}[4]$  with respect to the acceleration bias  $\mathbf{b}_a$

$$\frac{\partial \mathbf{r}_{imu}[4](\mathbf{b}_a + \delta \mathbf{b}_a)}{\partial \delta \mathbf{b}_a} = \mathbf{0}_{3 \times 3} \quad (51)$$

- Jacobian of the acceleration bias error  $\mathbf{r}_{imu}[5]$  with respect to the acceleration bias  $\mathbf{b}_a$

$$\frac{\partial \mathbf{r}_{imu}[5](\mathbf{b}_a + \delta \mathbf{b}_a)}{\partial \delta \mathbf{b}_a} = -\mathbf{I}_{3 \times 3} \quad (52)$$

## ACKNOWLEDGMENTS

This research has been partially funded by the Spanish government (projects DPI2012-32168 and DPI2015-67275), the Aragón regional government (Grupo DGA T04-FSE) and the University of Zaragoza (JIUZ-2015-TEC-03). In addition, this work was supported by the ARL grant W911NF-08-2-0004, ONR grants N00014-07-1-0829, N00014-14-1-0510, N00014-09-1-1051, N00014-09-1-103, NSF grant IIS-1426840, IIS-1138847, DARPA grants HR001151626, HR0011516850, and TerraSwarm, one of six centers of STARnet, a Semiconductor Research Corporation program sponsored by MARCO and DARPA

## REFERENCES

- [1] G. Klein and D. Murray, "Parallel tracking and mapping for small AR workspaces," in *Sixth IEEE and ACM International Symposium on Mixed and Augmented Reality*, 2007.
- [2] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós., "ORB-SLAM: A versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147-1163, October 2015, 2015.
- [3] J. Stühmer, S. Gumhold, and D. Cremers, "Real-time dense geometry from a handheld camera," in *Pattern Recognition*, 2010, pp. 11–20.
- [4] G. Graber, T. Pock, and H. Bischof, "Online 3d reconstruction using convex optimization," in *2011 IEEE International Conference on Computer Vision Workshops*, 2011, pp. 708–711.
- [5] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "DTAM: Dense tracking and mapping in real-time," in *2011 IEEE International Conference on Computer Vision (ICCV)*, 2011, pp. 2320–2327.
- [6] R. Mur-Artal and J. D. Tardós., "Probabilistic semi-dense mapping from highly accurate feature-based monocular slam," in *Robotics: Science and Systems*, 2015.
- [7] C. Mei, G. Sibley, M. Cummins, P. M. Newman, and I. D. Reid, "A constant-time efficient stereo slam system," in *BMVC*, 2009, pp. 1–11.
- [8] T. Pire, T. Fischer, J. Civera, P. D. Cristóforis, and J. J. Berles, "Stereo parallel tracking and mapping for robot localization," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015.
- [9] J. Civera, O. G. Grasa, A. J. Davison, and J. M. M. Montiel, "1-point ransac for EKF filtering: Application to real-time structure from motion and visual odometry," *Journal of Field Robotics*, vol. 27, no. 5, pp. 609–631, October 2010.
- [10] D. Scaramuzza, "1-point-ransac structure from motion for vehicle-mounted cameras by exploiting non-holonomic constraints," *International journal of computer vision*, vol. 95, no. 1, pp. 74–85, 2011.
- [11] D. Schleicher, L. M. Bergasa, M. Ocaña, R. Barea, and E. López, "Real-time hierarchical gps aided visual slam on urban environments," in *ICRA'09. IEEE International Conference on Robotics and Automation*, 2009. IEEE, 2009, pp. 4381–4386.
- [12] S. Weiss, M. W. Achtelik, S. Lynen, M. Chli, and R. Siegwart, "Real-time onboard visual-inertial state estimation and self-calibration of mavs in unknown environments," in *2012 IEEE International Conference on Robotics and Automation (ICRA)*, 2012, pp. 957–964.
- [13] 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.
- [14] C. Forster, L. Carbone, F. Dellaert, and D. Scaramuzza, "IMU preintegration on manifold for efficient visual-inertial maximum-a-posteriori estimation," in *IEEE Robotics: Science and Systems*, 2015.
- [15] S. Omari, M. Bloesch, P. Gohl, and R. Siegwart, "Dense visual-inertial navigation system for mobile robots," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 2634–2640.
- [16] L. Ma, J. M. Falquez, S. McGuire, and G. Sibley, "Large scale dense visual inertial slam," in *International Symposium on Experimental Robotics*, 2015.
- [17] S. Baker and I. Matthews, "Lucas-kanade 20 years on: A unifying framework," *Int. J. Comput. Vision*, vol. 56, no. 3, pp. 221–255, 2004.
- [18] C. Forster, M. Pizzoli, and D. Scaramuzza, "SVO: Fast semi-direct monocular visual odometry," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [19] J. Engel, T. Schöps, and D. Cremers, "LSD-SLAM: Large-scale direct monocular slam," in *ECCV 2014*, 2014, pp. 834–849.
- [20] A. Concha and J. Civera, "DPPTAM: Dense piecewise-planar tracking and mapping from a monocular sequence," in *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems* , 2015.
- [21] P. F. Felzenswalb and D. P. Huttenlocher, "Efficient graph-based image segmentation," *International Journal of Computer Vision*, vol. 59, no. 2, pp. 167–181, 2004.
- [22] N. Michael, D. Mellinger, Q. Lindsey, and V. Kumar, "The Grasp Multiple Micro-UAV Test Bed," *IEEE Robotics and Automation Magazine*, vol. 17, no. 3, pp. 56–65, 2010.
- [23] A. Concha, W. Hussain, L. Montano, and J. Civera, "Manhattan and piecewise-planar constraints for dense monocular mapping," in *Robotics:Science and Systems*, 2014.