

Simultaneous State Initialization and Gyroscope Bias Calibration in Visual Inertial aided Navigation

Jacques Kaiser¹, Agostino Martinelli¹, Flavio Fontana² and Davide Scaramuzza²

Abstract—State of the art approaches for visual-inertial sensor fusion use filter-based or optimization-based algorithms. Due to the nonlinearity of the system, a poor initialization can have a dramatic impact on the performance of these estimation methods. Recently, a closed-form solution providing such an initialization was derived in [17]. That solution determines the velocity (angular and linear) of a monocular camera in metric units by only using inertial measurements and image features acquired in a short time interval. In this paper, we study the impact of noisy sensors on the performance of this closed-form solution. We show that the gyroscope bias, not accounted for in [17], significantly affects the performance of the method. Therefore, we introduce a new method to automatically estimate this bias. Compared to the original method, the new approach now models the gyroscope bias and is robust to it. The performance of the proposed approach is successfully demonstrated on both simulated and real data from a quadrotor MAV.

I. INTRODUCTION

Autonomous mobile robots navigating in unknown environments have an intrinsic need to perform localization and mapping using only on-board sensors. Concerning Micro Aerial Vehicles (MAV), a critical issue is to limit the number of on-board sensors to reduce weight and power consumption. Therefore, a common setup is to combine a monocular camera with an inertial measurements unit (IMU). On top of being cheap, these sensors have very interesting complementarities. Additionally, they can operate in indoor environments where Global Positioning System (GPS) signals are shadowed. An open question is how to optimally fuse the information provided by these sensors.

Currently, most sensor-fusion algorithms are either filter-based or iterative. That is, given a current state and measurements, they return an updated state. While working well in practice, these algorithms need to be provided by an initial state.

The initialization of a filter based method is critical. Due to nonlinearities of the system, a poor initialization can result into converging towards local minima and providing faulty states with high confidence.

In this paper, we demonstrate the efficiency of a recent closed-form solution introduced in [16][17] that fuses

visual and inertial data to obtain the structure of the environment at the global scale along with the attitude and the speed of the robot. By nature, a closed-form solution is deterministic and, thus, does not require any initialization.

We implement this method in order to test it with real world data. This allow us to identify its limitations and bring modifications to overcome them. Specifically, we investigate the impact of biased inertial measurements. Although the case of biased accelerometer was originally studied in [17], we show that a large bias does not significantly worsen the performance. One major limitation of this method is the impact of biased gyroscope measurements. In other words, the performance becomes very poor in presence of a bias on the gyroscope and, in practice, the overall method can only be successfully used with a very precise - and expensive - gyroscope. We then introduce a simple method that automatically estimates this bias. By adding this new method for the bias estimation to the original method, we obtain results that are equivalent to the ones in absence of bias. Compared to the original method, the new method is now robust to the gyroscope bias and automatically calibrates the gyroscope.

II. RELATED WORK

The problem of fusing vision and inertial data has been extensively investigated in the past. However, most of the proposed methods require a state initialization. Because of the system nonlinearities, lack of precise initialization can irreparably damage the entire estimation process. In literature, this initialization is often guessed or assumed to be known [1][13][11][2][7]. Recently, this sensor fusion problem has been addressed by using optimization-based approaches [21][3]. These methods outperform filter-based algorithms in terms of accuracy due to their capability of relinearizing past states. On the other hand, the optimization process can be affected by the presence of local minima. We are therefore interested in a deterministic solution that analytically expresses the state in terms of the measurements provided by the sensors during a short time-interval.

In computer vision, several deterministic solutions have been introduced. These techniques, known as *Structure from Motion*, can recover the relative rotation and translation up to scale between two camera poses [14][10][18][9][12]. Such methods are currently used in state-of-the-art visual navigation methods on MAVs in order to initialize maps [22][7]. However, the knowledge of the absolute scale and, at least, the absolute roll and pitch angles is essential for many applications ranging from autonomous navigation in GPS-denied environments to 3D reconstruction and augmented reality. It is required to take the inertial measurements into consideration to compute these values deterministically.

A procedure to quickly re-initialize a MAV after a failure was presented in [6]. However, this method requires an altimeter to initialize the scale.

A landmark of known size was used in [8] to recover the initial pose of the MAV. This method is therefore not suitable for unknown environment.

Recently, a closed-form solution has been introduced in [16]. From integrating inertial and visual measurements over a short time-interval, this solution provides the absolute scale, roll and pitch angles, initial velocity, and distance to features. Specifically, all the physical quantities are obtained by simply inverting a linear system. The solution of the linear system can be refined with a quadratic equation assuming the knowledge of the gravity magnitude. This closed-form was improved in [13] to work with unknown camera-IMU calibration. Since the problem cannot be solved by simply inverting a linear system, a method to determine the six parameters that characterize the camera-IMU transformation was proposed. This method is therefore independent of external camera-IMU calibration, hence suitable for power-on-and-go systems.

A more intuitive expression of this closed-form solution was derived in [17]. While being mathematically sound, this closed-form solution is not robust to noisy sensor data. For this reason, to the best of our knowledge, it has never been used in an actual application. In this paper, we perform an analysis to find out its limitations. We start by reminding the reader the basic equations that characterize this solution (section III). In section IV, we show that this solution is resilient to the accelerometer bias but strongly affected by the gyroscope bias. We then introduce a simple method that automatically estimates the gyroscope bias (section V). By adding this new method for the bias estimation to the original method, we obtain results which are equivalent to the ones in absence of bias. Compared to the original method, the new method is now robust to

the gyroscope bias and also calibrates the gyroscope. We validate our new method against real world data to prove its robustness against noisy sensors in section VI. Finally, we provide our conclusion in section VII.

III. CLOSED-FORM SOLUTION

In this section, we provide the basic equations that characterize the closed-form solution proposed in [17]. We also provide the main features of this solution¹.

Let us refer to a short interval of time (of the order of 3 seconds). We assume that during this interval of time the camera observes simultaneously N point-features and we denote by t_1, t_2, \dots, t_{n_i} the times of this interval at which the camera provides an image of these points. Without loss of generality, we can assume that $t_1 = 0$. The following equation holds (see [17] for its derivation):

$$S_j = \lambda_1^i \mu_1^i - V t_j - G \frac{t_j^2}{2} - \lambda_j^i \mu_j^i \quad (6)$$

With:

- μ_j^i the normalized bearing of point feature i at time t_j in the local frame at time t_1 ;
- λ_j^i the distance to the point feature i at time t_j ;
- V the velocity in the local frame at time t_1 ;
- G the gravity in the local frame at time t_1 ;
- S_j the integration in the interval $[t_1, t_j]$ of the rotated linear acceleration data (i.e., the integration of the inertial measurements).

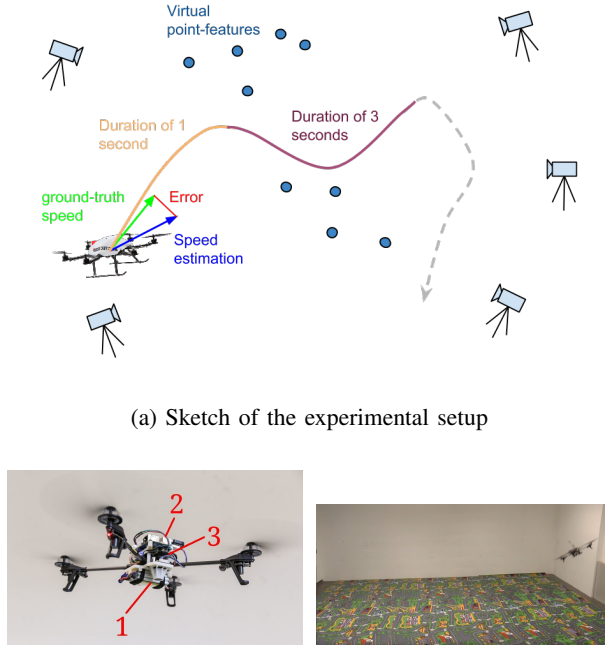
The local frame refers to a frame of reference common to the IMU and the camera. In a real application, we would work in the IMU frame and have some additional constant terms accounting for the camera-IMU transformation. We do not express these constant calibration terms explicitly here for clarity reasons.

The unknowns of Equation 6 are the scalars λ_j^i and the vectors V and G . Note that the knowledge of G is equivalent to the knowledge of the roll and pitch angles. The vectors μ_j^i are fully determined by visual and gyroscope measurements², and the vectors S_j are determined by accelerometer and gyroscope measurements.

Equation 6 provides three scalar equations for each point feature $i = 1, \dots, N$ and each frame starting from the second one $j = 2, \dots, n_i$. We therefore have a linear system consisting of $3(n_i - 1)N$ equations in $6 + Nn_i$ unknowns. Indeed, note that when the first frame occurs, at $t_1 = 0$, Equation 6 is always satisfied thus does not

¹Note that in this paper we do not provide a new derivation of this solution for which the reader is addressed to [17], section 3.

²The gyroscope measurements in the interval $[t_1, t_j]$ are needed to express the bearing at time t_j in the frame at time t_1



(b) A closeup of our quadrotor: 1) down-looking camera, with an OptiTrack motion-capture system (for ground-truth recording)
(c) Our flying arena equipped with an OptiTrack motion-capture system (for ground-truth recording)

Fig. 2: Experimental setup for identifying the limitations of the performance. The drone is equipped with an IMU, and the visual measurements are simulated.

We represent this setup in Fig. 2. We define the relative error as the euclidean distance between the estimation and the ground truth, normalized by the ground truth. We measure our error on the absolute scale by computing the mean error over all estimated distances to point features λ_j^i .

We set the frame rate of the simulated camera at 10Hz. If the frame rate is too high (above 30Hz), the interval of time between two frames becomes too small for the integrations we perform in 6, thus yielding numerical instability. When using a real high-frequency camera, we can force the frame rate to 10Hz by discarding most of the provided frames.

Reducing the number of considered frames also reduces the size of the matrices, and thus, speeds up the computations. As an example, over a time interval of 3 seconds, we obtain 31 distinct frames. When observing 7 features, solving the closed-form solution is equivalent to inverting a linear system of $3 \times 30 \times 7 = 630$ equations and $6 + 7 \times 31 = 223$ unknowns (see section III).

The method we use to solve the overconstrained

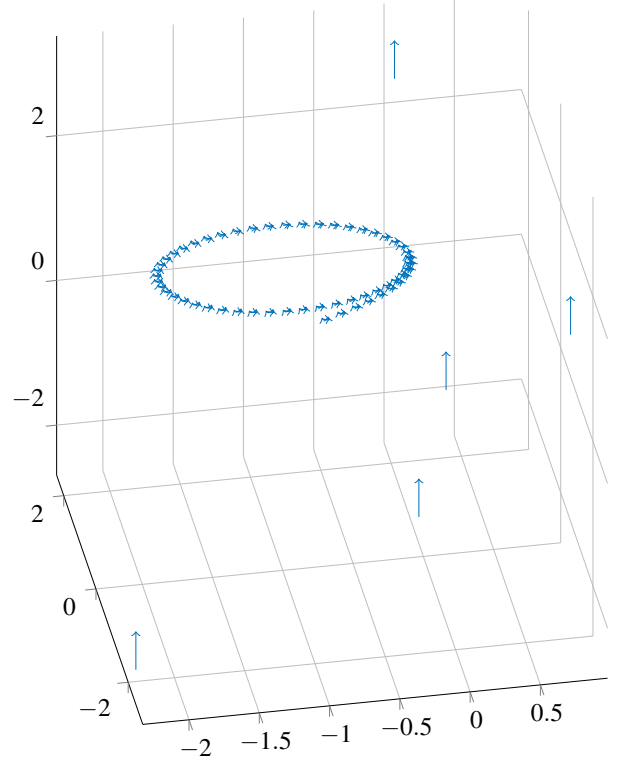


Fig. 3: Motion performed by the drone in 5 seconds.

linear system $\Xi X = S$ is a Singular Value Decomposition (SVD) since it yields numerically robust solutions.

In the next section, we will present the results obtained with the original closed-form solution on real IMU data, with simulated camera measurements. Our goal is to identify its performance limitations and introduce modifications to overcome them.

B. Original closed-form solution performance

The original closed-form solution described in Equation (9) will be used as a basis for our work. In the following, this closed-form solution will be denoted by *CF*. Moreover, we can also use the knowledge of the gravity magnitude to refine our results (Equation (10)). In this case, we are minimizing a linear objective function with a quadratic constraint. In Fig. 4, we display the performance of the original closed-form solution in estimating speed, gravity in the local frame and distances to the features with and without this additional constraint.

Note how the evaluations get better as we increase the integration time. Indeed, our equations come from an extended triangulation [16]. Therefore, it requires a significant difference in the measurements over time to

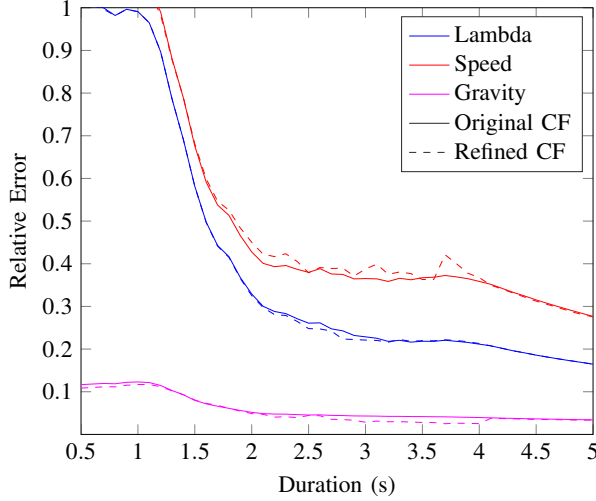


Fig. 4: Original closed-form solution estimations with and without gravity knowledge refinement. We are observing 12 features over a variable duration of integration.

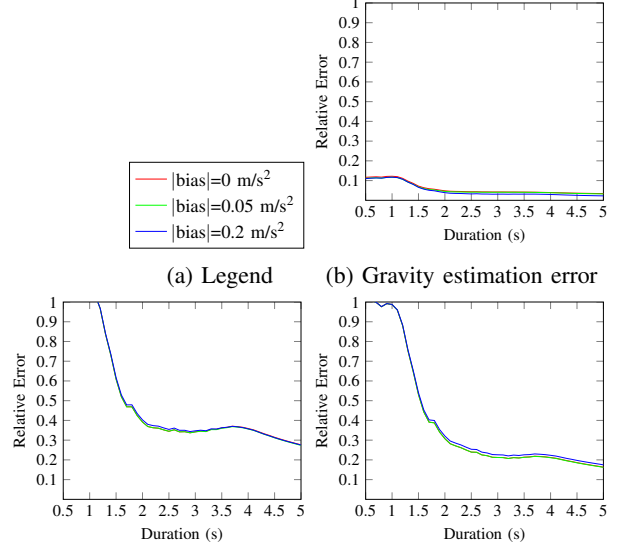
robustly estimate the state. Also note that the gravity is robustly estimated (around 5% error), whereas the speed and the distance to the features is more erroneous (above 10% error).

Since the gravity is well estimated with the original closed-form, it comes without surprise that constraining its magnitude does not improve the performance much. That is why the curves representing the performance with and without the gravity refinement are so close.

In the following sections, we will study the impact of biased inertial measurements on the performance of the closed-form solution without considering the gravity refinement.

C. Impact of accelerometer bias on the performance

In order to visualize the impact of the accelerometer bias on the performance, we corrupt the accelerometer measurements provided by the IMU by an artificial bias (Fig. 5). Despite a high accelerometer bias the closed-form solution still provides robust results. As seen on the Fig. 5, neither the estimation of the gravity, the velocity or the lambdas are impacted by the accelerometer bias. In [17], the author provides an alternative formulation of the closed-form solution including the accelerometer bias as an observable unknown of the system. However, the estimation of the accelerometer bias with this method is not robust since our system is only slightly affected by it. This is a counterintuitive result. Since our equations contain an integration of the acceleration, we also perform an integration of the accelerometer bias. We would



(c) Velocity estimation error (d) Lambda estimation error

Fig. 5: Impact of the accelerometer bias on the performance of the closed-form solution. We are observing 7 features over a variable duration of integration.

have expected the accelerometer bias to have a greater impact on the solutions yielded by the system.

D. Impact of gyroscope bias on the performance

To visualize the impact of the gyroscope bias on the performance, we corrupt the gyroscope measurements provided by the IMU by an artificial bias (Fig. 6). As seen in Fig. 6, the performance becomes very poor in presence of a bias on the gyroscope and, in practice, the overall method could only be successfully used with a very precise—and expensive—gyroscope.

V. ESTIMATING THE GYROSCOPE BIAS

Previous work has shown that the gyroscope bias is an observable mode when using an IMU and a camera, which means that it can be estimated [16].

Optimally, we would add the gyroscope bias in our unknown vector X and determine X by simply inverting the system $\Xi X = S$ as in the standard closed-form solution. However, we cannot express the gyroscope bias linearly with this system.

In this section, we propose a different approach to estimate the gyroscope bias using the closed-form solution.

A. Nonlinear minimization of the residual

Since our system of equations (6) is overconstrained, inverting it is equivalent to finding the vector X that minimizes the residual $\|\Xi X - S\|^2$.

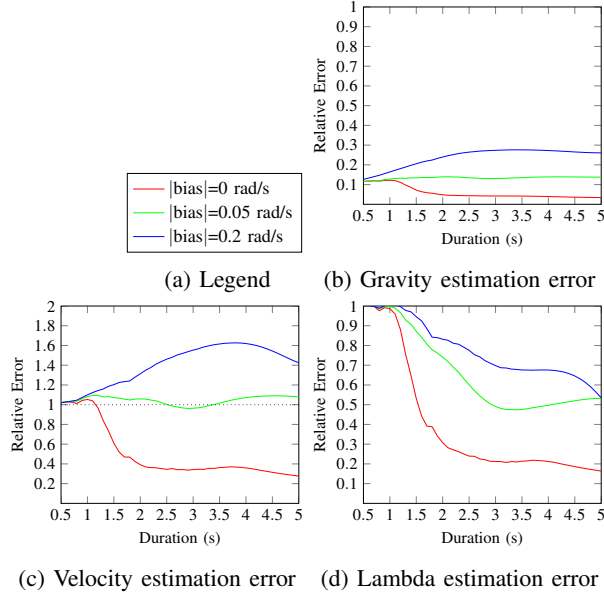


Fig. 6: Impact of the gyroscope bias on the performance of the closed-form solution. We are observing 7 features over a variable duration of integration.

Since we cannot express the gyroscope bias linearly, we define the following cost function:

$$cost(B) = \|\Xi X - S\|^2 \quad (1)$$

With:

- B the gyroscope bias;
- Ξ and S computed with respect to B .

By minimizing this cost function, we recover the gyroscope bias B and the unknown vector X which compensated for the gyroscope bias B . We can initialize the optimization process with $B = 0_3$ since the bias is usually a rather small quantity. Figure 7 displays the three components of the bias estimated by this minimization.

The optimized closed-form solution provides better results than the standard closed-form solution. Fig. 8 shows an improvement in precision of around 13% for the distance to the features, and around 30% for the speed after 4 seconds of integration.

Moreover, this method is robust even for high values of the gyroscope bias. Fig. 9 displays the performance of the proposed method in estimating speed, gravity in the local frame and distances to the features in presence of the same artificial gyroscope bias from Fig. 6.

As seen in Fig. 9, after a certain integration duration, the estimations agree no matter how high the bias is. In

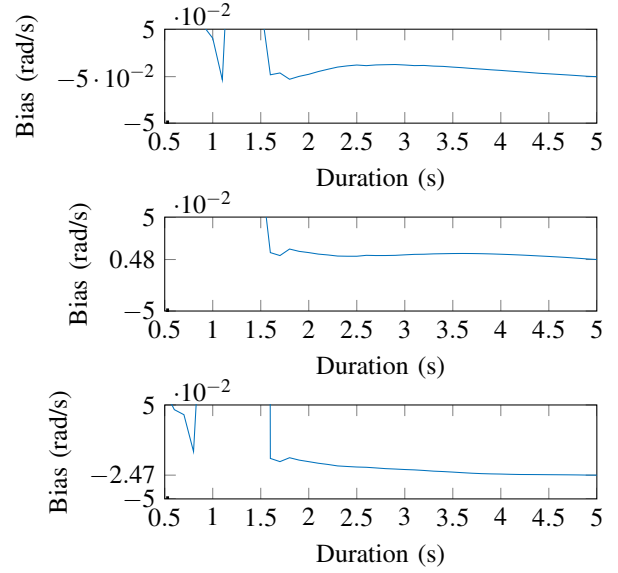


Fig. 7: Gyroscope bias estimation from nonlinear minimization of the residual. We are observing 12 features over a variable duration of integration.

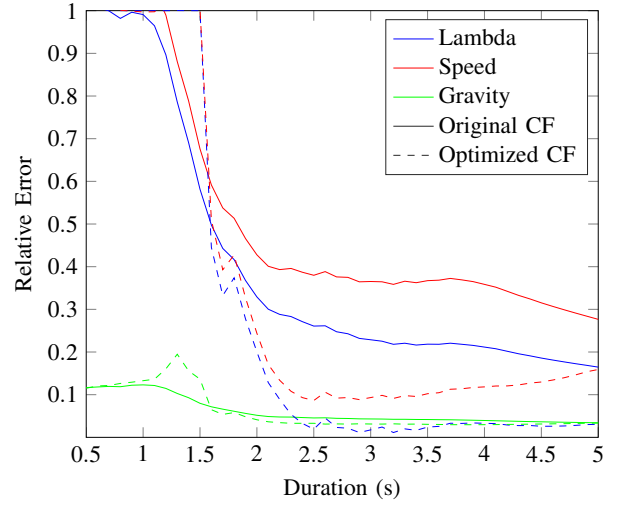


Fig. 8: Estimation error of the original closed-form solution against the optimized closed-form solution. We are observing 12 features over a variable duration of integration.

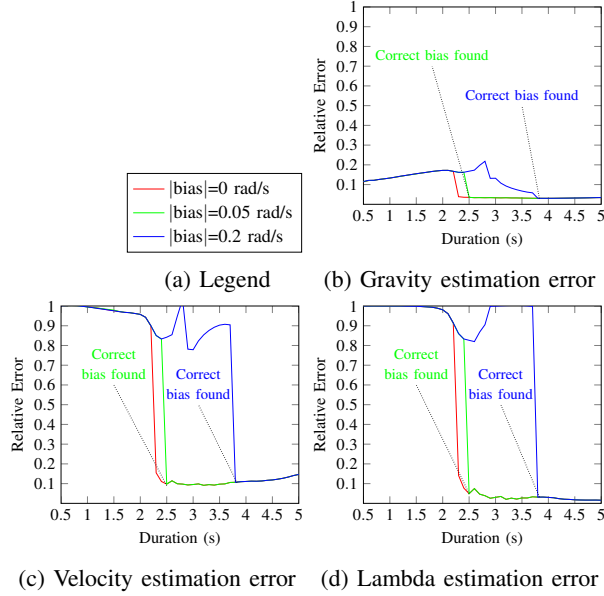


Fig. 9: Impact of the gyroscope bias on the performance of the optimized closed-form solution. We are observing 7 features over a variable duration of integration.

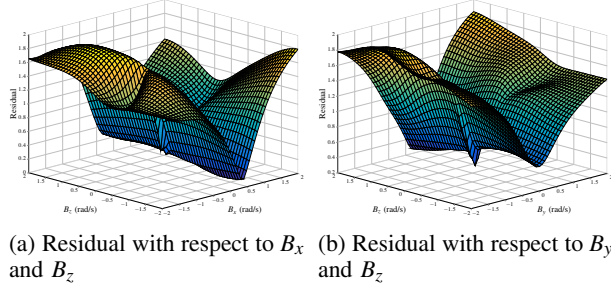


Fig. 10: Cost function (residual) with respect to the gyroscope bias for a small amount of available measurements (integration of 2 seconds while observing 7 features)

other words, given that the integration duration is long enough, this method is unaffected by the gyroscope bias.

However, for very short time of integration (< 1.7 seconds), the proposed method can fail by providing a gyroscope bias much larger than the correct one. To understand this misestimation, in Fig. 10 we plot the residual with respect to the bias, which is the cost function we are minimizing. We highlight a misestimation of the gyroscope bias by setting the duration of integration to 2 seconds while observing 7 features. We refer to the components of the gyroscope bias by $B = [B_x, B_y, B_z]$.

As we can see in Fig. 10, the cost function admits a symmetry with respect to B_z .

B. Removing the symmetry in the cost function

The symmetry in the cost function is induced by the strong weight of the gravity in the equation 6. In general, the residual is almost constant with respect to the component of the gyroscope bias along the direction \vec{u} when this direction \vec{u} is collinear with the gravity throughout the motion. In the real world data we had, the motion satisfies this constraint. Specifically, the gyroscope was strapped on the MAV such that the vector $[0, 0, 1]$ in the gyroscope frame was pointing upwards when the MAV was hovering. That is why the residual varies only slightly for certain time sequences with respect to this vector. Indeed, in normal operations, a MAV will often have a pose close to its hovering stance in order to stay stable.

If the MAV rotates such that the vector \vec{u} becomes noncollinear with the gravity, the cost function does not exhibit this symmetry anymore. In this case the gyroscope bias is well estimated.

A simple solution to avoid having that symmetry in our system would be to constrain the motion of our MAV while it is operating. Another way to artificially get rid of this symmetry is to tweak the cost function. Specifically, we can add a regularization term that penalizes high estimations of the gyroscope bias:

$$\text{cost}(B) = \|\Xi X - S\|^2 + \lambda \|B\| \quad (2)$$

The coefficient λ is the weight given to how much we want the bias to be small. For small values of λ , our cost function is similar to the previous one and the bias can grow arbitrarily high. For high values of λ , the estimations provided by the optimized closed-form solution are similar to the ones provided by the standard closed-form solution. Indeed, high values of λ force the estimation of the bias to 0_3 .

Note that, instead of forcing the gyroscope bias to be close to 0_3 , we can easily force it to be close to any value. Therefore, if we have the knowledge of an approximately known gyroscope bias, we can use it to provide a better estimation of the gyroscope bias.

$$\text{cost}(B) = \|\Xi X - S\|^2 + \lambda \|B - B^{\text{approx}}\|$$

With B^{approx} the known approximate gyroscope bias. This methods allow us to reuse previously computed gyroscope bias since it is known to slowly vary over time.

Selecting a reliable and safe value for the regularization parameter λ is complicated. In this paper, we picked a value of λ by experimentation.

VI. EXPERIMENTS ON REAL DATA

We validate our method against a different dataset than the one we used in the previous sections to draw our conclusions. Specifically, it contains IMU and camera measurements along with ground truth. Therefore, we are now only relying on fully real world data. However, since we are no longer simulating the point features, we do not have the ground truth for the distance to the point features anymore. We can therefore only compare the performance of the evaluation of the speed and the gravity.

The drone is flying indoor at low altitude. The feature extraction and matching is done with FAST corner algorithm [19][20].

We compare the performance on the estimations of the gravity and the initial velocity obtained with three different methods:

- The original closed-form solution (Equation 9);
- Our modified closed-form solution (Equation 1);
- The loosely-coupled visual-inertial algorithm (MSF) in [15] using pose estimates from the Semi-direct Visual Odometry (SVO) package [7] (how to combine the sensor fusion in [15] with SVO can be found in [5]).

The reason we included the SVO in the validation is for having a reference to state-of-the-art pose estimation method. However, this method requires to be initialized with the knowledge of the absolute scale, whereas our method works without initialization.

We set the integration duration for the closed-form solution to 2.8 seconds. The camera provides 60fps, but we discard most of the frames and consider only one frame every 0.1 seconds. Indeed, considering frames that are separated by a time interval shorter than 0.1 seconds does not add significant information to our system.

As seen in Fig. 11, the original closed-form solution and the optimized closed-form solution have similar performance. Indeed, for this dataset the gyroscope bias was estimated to $B = [0.0003, 0.009, 0.001]$, which is very low ($\|B\| = 0.0091$).

Moreover, the performance is also similar to the one obtained with a well-initialized SVO. We remind the reader that unlike SVO, the closed-form solution does not require the knowledge of the absolute scale to be provided.

VII. CONCLUSION

In this paper, we studied the recent closed-form solution proposed by [17] that performs visual-inertial sensor fusion without requiring an initialization. We implemented this method in order to test it with real

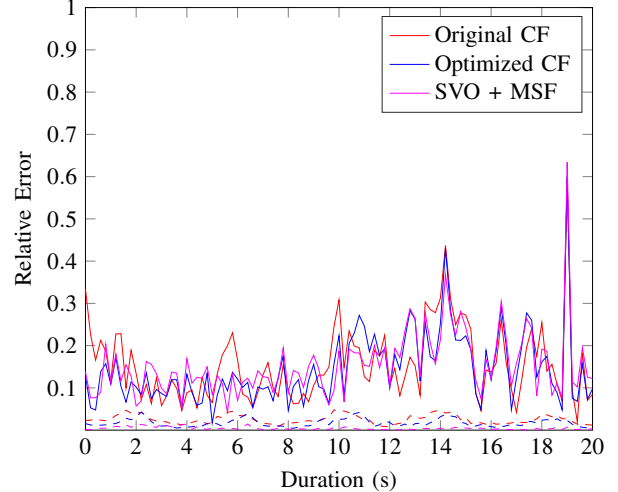


Fig. 11: Estimation error of the optimized closed-form solution against the original closed-form solution [17] and SVO [5]. The duration of integration is set to 2.8 seconds, and 10 point features are observed throughout the whole operation.

world data. This allowed us to identify its performance limitations and bring modifications to overcome them.

We investigated the impact of biased inertial measurements. Although the case of biased accelerometer was originally studied in [17], we showed that its low impact on the system makes it hard to estimate.

One major performance limitation of this method was due to the impact of biased gyroscope measurements. In other words, the performance becomes very poor in presence of a bias on the gyroscope and, in practice, the overall method could only be successfully used with a very precise (and expensive) gyroscope. We then introduced a simple method that automatically estimates this bias.

We validated this method by comparing its performance against the original method and the SVO described in [5] which is the state-of-the-art approach for pose estimation on MAV.

For future work, we see this optimized closed-form solution being implemented on a MAV to provide accurate state initialization. This would allow aggressive take-off maneuvers, such as hand throwing the MAV in the air, as already demonstrated in [6] with a range sensor. With our technique however, we could get rid of the altimeter sensor. The drone could therefore perform any motion right after the throw instead of having a required hovering stage to compute the absolute scale.

REFERENCES

- [1] L. Armesto, J. Tornero, and M. Vincze. Fast Ego-motion Estimation with Multi-rate Fusion of Inertial and Vision. *The International Journal of Robotics Research*, 26(6):577–589, 2007.
- [2] Marco Bibuli, Massimo Caccia, and Lionel Lapierre. Sliding Window Filter with Application to Planetary Landing. *IFAC Proceedings Volumes (IFAC-PapersOnline)*, 7(PART 1):81–86, 2007.
- [3] Forster C., Carlone L., Dellaert F., and Scaramuzza D. IMU Preintegration on Manifold for Efficient Visual-Inertial Maximum-a-Posteriori Estimation. In *Robotics: Science and Systems (RSS)*, 2015.
- [4] Jakob Engel and Daniel Cremers. Scale-Aware Navigation of a Low-Cost Quadcopter with a Monocular Camera. 2014.
- [5] Matthias Faessler, Flavio Fontana, Christian Forster, Elias Mueggler, Matia Pizzoli, and Davide Scaramuzza. Autonomous, vision-based flight and live dense 3d mapping with a quadrotor micro aerial vehicle. *Journal of Field Robotics*, 2015.
- [6] Matthias Faessler, Flavio Fontana, Christian Forster, and Davide Scaramuzza. Automatic Re-Initialization and Failure Recovery for Aggressive Flight with a Monocular Vision-Based Quadrotor. In *International Conference on Robotics & Automation*, 2015.
- [7] Christian Forster, Matia Pizzoli, and Davide Scaramuzza. SVO: Fast semi-direct monocular visual odometry. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [8] P. Gemeiner, P. Einramhof, and M. Vincze. Simultaneous Motion and Structure Estimation by Fusion of Inertial and Vision Data. *The International Journal of Robotics Research*, 26(6):591–605, 2007.
- [9] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004.
- [10] Richard I. Hartley. In defense of the eight-point algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(6):580–593, 1997.
- [11] Guoquan P. Huang, Anastasios I. Mourikis, and Stergios I. Roumeliotis. On the complexity and consistency of UKF-based SLAM. *Proceedings - IEEE International Conference on Robotics and Automation*, pages 4401–4408, 2009.
- [12] Hongdong Li and Richard Hartley. Five-point motion estimation made easy. *Proceedings - International Conference on Pattern Recognition*, 1:630–633, 2006.
- [13] M. Li and a. I. Mourikis. High-precision, consistent EKF-based visual-inertial odometry. *The International Journal of Robotics Research*, 32(6):690–711, 2013.
- [14] H. C. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections, 1981.
- [15] Simon Lynen, Markus W Ahtelik, Steven Weiss, Maria Chli, and Roland Siegwart. A robust and modular multi-sensor fusion approach applied to mav navigation. pages 3923–3929, 2013.
- [16] Agostino Martinelli. Vision and imu data fusion: Closed-form solutions for attitude, speed, absolute scale, and bias determination. *Transactions on Robotics*, 28(1):44–60, 2012.
- [17] Agostino Martinelli. Closed-form solution of visual-inertial structure from motion. *International Journal of Computer Vision*, 106(2):138–152, 2014.
- [18] D Nister. An efficient solution to the five point relative pose problem. pages 195–202, 2003.
- [19] Edward Rosten and Tom Drummond. Fusing points and lines for high performance tracking. *Proceedings of the IEEE International Conference on Computer Vision*, II:1508–1515, 2005.
- [20] Edward Rosten and Tom Drummond. Machine learning for high-speed corner detection. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 3951 LNCS:430–443, 2006.
- [21] Leutenegger Stefan, Lynen Simon, Bosse Michael, Siegwart Roland, and Furgale Paul. Keyframe-based visual-inertial odometry using nonlinear optimization. *The International Journal of Robotics Research*, 34(3):314–334, 2014.
- [22] Stephan M Weiss. Vision Based Navigation for Micro Helicopters (PhD Thesis - Weiss 2012). (20305), 2012.