

VIR-SLAM: Visual, Inertial, and Ranging SLAM for single and multi-robot systems

Yanjun Cao and Giovanni Beltrame

Abstract—Monocular cameras coupled with inertial measurements generally give high performance visual inertial odometry. However, drift can be significant with long trajectories, especially when the environment is visually challenging. In this paper, we propose a system that leverages ultra-wideband ranging with one static anchor placed in the environment to correct the accumulated error whenever the anchor is visible. We also use this setup for collaborative SLAM: different robots use mutual ranging (when available) and the common anchor to estimate the transformation between each other, facilitating map fusion. Our system consists of two modules: a double layer ranging, visual, and inertial odometry for single robots, and a transformation estimation module for collaborative SLAM. We test our system on public datasets by simulating an ultra-wideband sensor as well as on real robots. Experiments show our method can outperform state-of-the-art visual-inertial odometry by more than 20%. For visually challenging environments, our method works even the visual-inertial odometry has significant drift. Furthermore, we can compute the collaborative SLAM transformation matrix at almost no extra computation cost.

I. INTRODUCTION

Robot localization is a fundamental topic in any mobile robot application. Recent advances in robot hardware and software have boosted the opportunity and demand for multi-robot systems for their inherent benefits, such as high efficiency and robustness.

With a monocular camera and low-cost inertial measurement units, Visual Inertial Odometry (VIO) is accepted as the minimal sensor configuration for single robot state estimation and navigation considering size, weight, power and cost [1]. Recent technical advances [2]–[4] make VIO more and more robust and stable in many conditions. However, the drift caused by accumulated error is still hard to control without loop closures. Although loop closure is a natural part for SLAM system, the requirements to close the loop rely much on the trajectory and environment. For example, generating high-quality closures requires revisiting the same location with a similar viewpoint. Furthermore, perception outliers caused by illumination, self-similar environments, etc. are challenging for loop closure.

In this paper, we use single extra sensor, a static ultra-wideband (UWB) anchor, to improve the performance of robot localization. UWB technology has attracted a lot of attention recent years for its accurate ranging performance and long-distance support. For example, the latest Apple iPhone at the time of writing is equipped with a UWB

chip (actually, it includes all sensors needed in this paper). Most available UWB systems (e.g. [5]) use several (at least four for 3D and three for 2D) calibrated anchors as a Global Positioning System (GPS) for specific areas. This type of infrastructure is not applicable for the exploration in unknown environments, which is one of the primary objectives of SLAM. Therefore, we design our system to rely only on one anchor, which can be dropped off at any moment by a robot during its mission. Our experiments show this one anchor can improve the localization accuracy significantly.

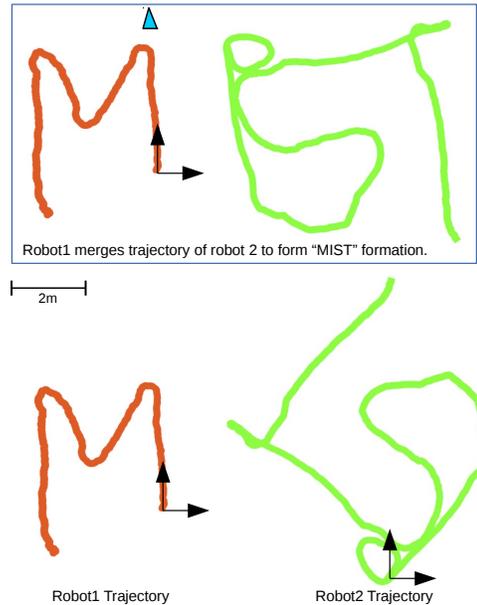


Fig. 1. Two robots are manually controlled independently. The first robot with a Realsense T265 and an UWB module draws an "M" in its own coordinate frame. The other robot, carrying a Realsense D435, a UWB module and a Pixracer flight controller, draws "IST" in its frame. One static UWB anchor is placed in the environment (the blue triangle shown in Robot1 frame). Their trajectories are shown at the bottom. With two inter-robot communications and measurements, robot 1 can estimate the transformation matrix from robot 2, and then map robot 2's trajectory in its own frame as the top picture shows, forming "MIST" (our lab).

This setup has another significant benefit for multi-robot SLAM. Multi-robot SLAM is a promising field, but much more challenging. One issue is the estimation of the transformation matrix between robots [6]. Most available works rely on common features to extract the relative pose of robots, which is a type of inter-robot loop closure, which leads to similar constraints as single robot loop closures. An additional challenge for inter-robot loop closure is the need

M. Cao and Prof. Beltrame are with the Department of Computer and Software Engineering, École Polytechnique de Montréal, 2900 Boul Édouard-Montpetit, Québec CA E-mail: (name.surname@polymtl.ca).

to exchange the information required for loop closure among all robots, which can be significant.

In our system, we can solve this challenge using the UWB sensor. When any two robots move into their respective UWB ranging radius, they can exchange their mutual ranging and anchor information and they can estimate their respective coordinate systems transformations. Having the transformation matrix, a robot can correctly project all the information from its neighbors onto its own frame. This solution greatly simplifies multi-robot SLAM, minimizing the need for inter-robot loop closure, which requires the exchange of feature databases, the identification of loop closures and distributed pose graph optimization.

The last but not least reason for choosing UWB is its inherent advantage in data association. When a range measurement is received, the identity of the sender is recognized without extra effort, which can be of great assistance for multi-robot systems. Moreover, UWB can provide low bandwidth communication while ranging. For example, the state estimation of poses can be carried by UWB packets while performing ranging. In the future, we believe a monocular camera, inertial measurement units (IMU), and UWB will be a standard minimal configuration for multi-robot systems.

II. RELATED WORK

1) *Single Robot SLAM with Visual, Inertial and Ranging Measurements*: SLAM has been the subject of intense research for more than thirty years [7]. Monocular visual-inertial odometry is a popular choice as it provides good state estimation performance with a minimal sensor configuration. Although state-of-the-art VIO algorithms (e.g. SVO [3], VINS-Mono [4], DSO [8]) can reach very high accuracy in relative translation and orientation, the accumulated drift can still be an issue: any small orientation error can lead to large end-point error. Our system leverages UWB ranging measurements to correct the accumulated error. We developed our system based on VINS-Mono [4], which is a robust and versatile state estimator which uses a sliding window tightly-coupled nonlinear optimization for visual and IMU measurements.

UWB technology, as a localization solution on its own, has attracted a lot of attention in recent years both in research and industry for its decimeter localization accuracy. However, most results are based on a well-calibrated multi-anchor setup [9]–[12], which is not applicable for navigation in unexplored, unstructured environments. Wang et al. [13] propose a system using camera, IMU and UWB to bypass the complexity of loop closure. However, they still use multiple pre-configured UWB anchors. Their UWB module provides coarse drift-free global position and VIO identifies the local trajectory. On the contrary, in this paper we use only one anchor, placed in an unspecified location. Shi et al. [14] have a similar spirit to ours. They start with one UWB anchor and keep dropping anchors from a moving robot. Unfortunately, their experimental results are available only in simulation for one sequence of the EuRoC dataset [15], with five simulated anchors. In our work, we focus on the use of a single anchor

setup. We design a double layer sliding window algorithm, which effectively fuses accurate VIO and range constraints along the trajectory.

2) *Multi-Robot SLAM*: Multi-robot SLAM has gained recent attention for the increased viability and accessibility of multi-robot systems. Saeedi et al. [16] give a comprehensive review of multi-robot SLAM and points out one key issue: relative pose estimation. Most current multi-robot SLAM systems solve this issue by analyzing inter-robot loop closures, either in centralized [17] or distributed [18], [19] fashion. The distributed approach is more robust, but it is harder to implement in practice: robots need to exchange map data to get the feature database for future loop closures, and distributed optimization usually requires additional communication and computation.

Ranging measurements can assist in relative pose estimation. Trawny et al. [20] provide theoretical proofs and simulations that show how six range measurements can be used to get the transformation matrix between two robots. Martel et al. [21] extend the work and adapt it to 4DOF relative pose estimation with a UWB setup, and use it for merging maps for VR applications. Both methods use range measurements over long trajectories. In our solution, robots can estimate the transformation matrix as soon as they can get two measurements from their neighbors, which meets the requirement of real-time transformation estimation during robot rendezvous. These two methods [20], [21] represent a good solution when no common anchor is present and can be combined with inter-robot loop closures to improve the transformation results.

In practice, most works in literature have explored the combination of multiple UWB anchors with vision and/or IMU. In this paper, we focus on the use of a single anchor in an arbitrary location, which is trivial to deploy as a beacon in real exploration tasks. We propose a double layer sliding window technique to combine VIO with UWB ranging, which produces drift-free state estimation by leveraging VIO for its accurate short time relative pose estimation, and range constraints for longer trajectories. Moreover, recorded anchor ranges can help robots find the transformation matrix efficiently with only two range measurements in multi-robot settings. Our contributions can be summarized as:

- an UWB-aided SLAM system for single robots which outperforms state-of-the-art VIO;
- a double layer sliding window algorithm that combined relative pose estimation from VIO and UWB ranging constrains;
- an efficient method to estimate the transformation matrix between multiple robots.

III. SYSTEM DESIGN

In this section, we explain the preliminaries and symbols used in this paper. We then detail the formulation of our ranging-aided visual inertial SLAM, focusing on the cost factors of the optimization problem. Finally, we introduce our solution for the estimation of the transformation matrix between robots.

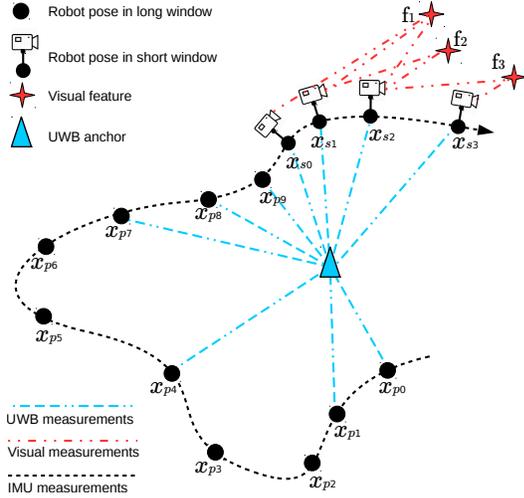


Fig. 2. Illustration of the scenario of VIR-SLAM, with a single anchor setup. Measurements from sensors are shown with different lines. Robot poses with the camera and IMU measurements belong to the short sliding window. All poses associated with UWB ranging to the anchor are kept for the long sliding window.

A. Single Robot Estimation Preliminaries

We assume a robot carries three kinds of sensors (a monocular camera, an IMU and UWB module), and moves in a 3D environment, as shown in Fig. 2. An UWB anchor is placed in the environment with an unknown initial position. We define the world frame of the robot i as $[\]^{iw}$, which is usually aligned with the first camera frame when the robot starts its mission. The position of the anchor is expressed in the robot world frame, denoted by \mathbf{P}_A^{iw} . We use $[\]^{ib}$ to indicate the body frame of robot i , and similarly $[\]^{ic}$ for the camera frame. Note that we do not define the UWB sensor frame because it is a scalar measurement. The UWB ranging measurement is transferred to body frame by considering the 3D offset of the UWB antenna in the body frame.

Classical VIO proposes an optimization formulation of states over a sliding window with size n as:

$$\mathcal{X} = [\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n, l_0, l_1, \dots, l_m] \quad (1)$$

$$\mathbf{x}_k = [\mathbf{p}_{b_k}^{iw}, \mathbf{v}_{b_k}^{iw}, \mathbf{q}_{b_k}^{iw}, \mathbf{b}_{a_k}^{ib}, \mathbf{b}_{w_k}^{ib}] \quad k \in [0, n]$$

which includes the state \mathbf{x} for all n frames and the visual feature inverse depth l . The k th frame state \mathbf{x}_k includes the position $\mathbf{p}_{b_k}^{iw}$, velocity $\mathbf{v}_{b_k}^{iw}$, and orientation in quaternions $\mathbf{q}_{b_k}^{iw}$ for robot i in its world frame, plus accelerometer bias $\mathbf{b}_{a_k}^{ib}$ and gyroscope bias $\mathbf{b}_{w_k}^{ib}$ in its body frame. l_i is the inverse depth of the i th feature among m features from the visual observations over the sliding window. If the sliding window includes all the camera frames since the start of the mission, the optimization becomes a full smoothing estimation. Although full smoothing offers the best accuracy, it is not scalable in reality, so we use a key frame approach that discards similar frames while not losing tracking. However, as the trajectory becomes longer, the size of state keeps growing. For this reason, we fix the size of key frames and marginalize older key frames into a prior factor

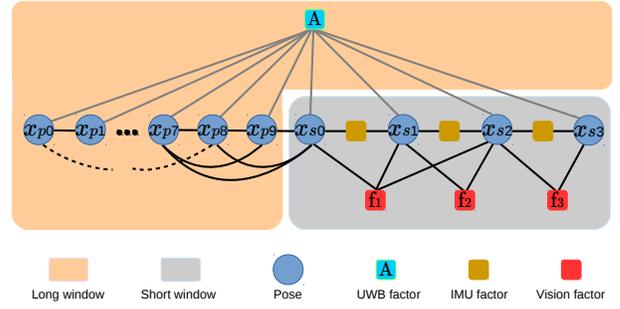


Fig. 3. Factor graph for the scenario in Fig. 2. The factor graph includes poses, UWB factor, IMU factors and visual factors. The edges are the error. We have a UWB factor for the poses with range measurements. Further, smooth edges (curves) are added between poses.

in the optimization. The drift from visual inertial odometry is still hard to avoid, and pose graph optimization with loop closure becomes the only chance to correct the accumulating error.

Aiming at an accurate localization system and avoiding all rely on loop closure, we design a SLAM system that uses a novel double layer sliding window structure, with an implementation based on VINS-Mono [4].

B. Double Layer Tightly Coupled VIR Optimization

We propose a double layer sliding window tightly coupled SLAM optimization by considering following aspects: high accuracy relative pose estimation from VIO, less accurate but absolute UWB measurements, and the computation cost for these factors.

We design two sliding windows involving three kinds of variables, shown in Equ. 2. A short window Ψ , the same as the classical visual SLAM sliding window with size n from Equ.1. The variables in this window include \mathbf{x}_k (Equ. 3) and $[l_0, l_1, \dots, l_m]$, which have same meaning as the variables in Equ. 1. The novelty in our system lies in the addition of another long sliding window Ω , which carries state \mathbf{w}_t , as shown in Equ. 4. The size of the long slide window is s , which is much larger than n . The state in this window only contains the robot position \mathbf{p}_b^{iw} in robot world frame.

$$\mathcal{X} = [\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_s, \mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n, l_0, l_1, \dots, l_m] \quad (2)$$

$$\mathbf{x}_k = [\mathbf{p}_{b_k}^{iw}, \mathbf{v}_{b_k}^{iw}, \mathbf{q}_{b_k}^{iw}, \mathbf{b}_{a_k}^{ib}, \mathbf{b}_{w_k}^{ib}] \quad k \in \Psi[0, n] \quad (3)$$

$$\mathbf{w}_t = [\mathbf{p}_{b_t}^{iw}] \quad t \in \Omega[0, s] \quad (4)$$

Fig.3 shows the factor graph of our system, corresponding to the scenario of Fig. 2. The short window area in gray includes the state from the camera and IMU measurements. The long window area in orange contains the state from the UWB measurements. Following the state definition, we formulate a full nonlinear optimization problem as:

$$\min_{\mathcal{X}} \left\{ \underbrace{\sum_{t \in \Omega} \rho(\|\mathbf{r}_U(\hat{\mathbf{z}}_t, \mathcal{X})\|_{\mathbf{P}_{b_t}^{iw}})}_{\text{UWB factor}} + \underbrace{\sum_{k \in \Psi} \|\mathbf{r}_B(\hat{\mathbf{z}}_{b_{k+1}}^{b_k}, \mathcal{X})\|_{\mathbf{P}_{b_{k+1}}^{b_k}}^2}_{\text{IMU factor}} \right. \\ \left. + \underbrace{\sum_{(l,j) \in \mathcal{C}} \rho(\|\mathbf{r}_C(\hat{\mathbf{z}}_l^{c_j}, \mathcal{X})\|_{\mathbf{P}_l^{c_j}})}_{\text{Vision factor}} \right\} \quad (5)$$

This nonlinear optimization problem considers three factors, corresponding to the UWB factor, the IMU factor, and the vision factor. The UWB residuals are calculated for the long window Ω , while IMU and visual residuals for the short window Ψ .

1) *UWB Factor*: UWB localization uses different protocols: time of arrival (TOA), time difference of arrival (TDoA), and two-way ranging (TWR). In our system, we use TWR, which measures the distances between two transceivers by sending a packet back and forth. Although the number of nodes supported is limited because of the shared UWB communication medium, TWR can be used without device synchronization, which makes the protocol widely used. Since we are not considering hundreds of robot communicating simultaneously, we choose TWR to avoid synchronization, which can be difficult to implement in distributed multi-robot applications. We model the ranging measurement of UWB modules as:

$$\hat{d} = d + b_d + e$$

where \hat{d} is the UWB measurement, d is true distance, b_d is distance based bias, and e is the error following a Gaussian distribution $N(0, \sigma^2)$. Note that the b_d can be calibrated with a Gaussian process model beforehand by collecting data with ground truth. Then the model can be simplified as:

$$\hat{d} = d + e$$

With this simplified UWB model, we define the UWB factor in Equ. 5 as:

$$\mathbf{r}_U(\hat{\mathbf{z}}_t, \mathcal{X}) = \underbrace{\gamma_r \cdot \rho(\|\mathbf{p}_{b_t}^{iw} - \mathbf{P}_A^{iw}\| - \hat{d}_t)}_{\text{UWB ranging residual}} \\ + \underbrace{\gamma_s \cdot \left(\sum_{j \in (t, t+s]} \{\mathbf{p}_{b_j}^{iw} - \mathbf{p}_{b_t}^{iw}\} - \hat{\mathbf{z}}_{b_t}^{b_j} \right)}_{\text{Relative transformation residual}} \quad (6)$$

The UWB factor above includes two residuals: a ranging measurement residual and a virtual relative transformation measurement residual. γ_r and γ_s are the weights for these two residuals. \hat{d}_t are the UWB ranging measurements, which are compared with the predicted ranging from the robot body frame t to the anchor \mathbf{P}_A^{iw} in the world frame. To avoid the ranging factor excessively affecting the optimizer and breaking the relative pose property from visual-inertial estimation, we introduce a virtual relative transformation measurement $\hat{\mathbf{z}}_{b_t}^{b_j}$ between frames t and $j \in (t, t+s]$, which is extracted from the short sliding window estimation result. We select $s = 3$ in our experiments to create two consecutive

links between poses, as can be seen for example with x_{p7} in Fig. 3. $\rho(\cdot)$ is a pseudo-Huber loss function defined as $\rho(q) = \delta^2(\sqrt{1 + (q/\delta)^2} - 1)$.

2) *IMU factor*: IMU measurements are critical for monocular visual odometry. As the frequency of the IMU is usually higher than the camera image frame rate, the IMU measurements are preintegrated between two consecutive image frames [2]. By referring to the last body frame motion, this technique avoid repeated IMU reintegration and reduces computation during optimization. Forster et al. [3] extend this approach to manifold structures of the rotation group $\mathcal{S}o3$ for higher accuracy and robustness. We follow the same method as [4] in quaternion form.

The preintegrated IMU measurements between two consecutive frames of b_k and b_{k+1} referring to frame b_k can be expressed as:

$$\alpha_{b_{k+1}}^{b_k} = \iint_{t \in [t_k, t_{k+1}]} \mathbf{R}_t^{b_k} (\hat{\mathbf{a}}_t - \mathbf{b}_{a_t}) dt^2 \\ \beta_{b_{k+1}}^{b_k} = \int_{t \in [t_k, t_{k+1}]} \mathbf{R}_t^{b_k} (\hat{\mathbf{a}}_t - \mathbf{b}_{a_t}) dt \\ \gamma_{b_{k+1}}^{b_k} = \int_{t \in [t_k, t_{k+1}]} \frac{1}{2} \Omega(\hat{\boldsymbol{\omega}}_t - \mathbf{b}_{w_t}) \gamma_t^{b_k} dt$$

where

$$\Omega(\boldsymbol{\omega}) = \begin{bmatrix} -[\boldsymbol{\omega}]_{\times} & \boldsymbol{\omega} \\ \boldsymbol{\omega}^T & 0 \end{bmatrix}, [\boldsymbol{\omega}]_{\times} = \begin{bmatrix} 0 & -\omega_z & \omega_y \\ \omega_z & 0 & -\omega_x \\ -\omega_y & \omega_x & 0 \end{bmatrix}$$

$\hat{\mathbf{a}}_t$ and $\hat{\boldsymbol{\omega}}_t$ are the accelerometer and gyroscope measurement vectors, respectively. These three formulas correspond to relative the motion changes of position, velocity and orientation to the local body frame of b_k .

The IMU factor is the residual between predicted motion and the preintegrated results referring to the body frame:

$$\mathbf{r}_B(\hat{\mathbf{z}}_{b_{k+1}}^{b_k}, \mathcal{X}) = \begin{bmatrix} \delta \alpha_{b_{k+1}}^{b_k} \\ \delta \beta_{b_{k+1}}^{b_k} \\ \delta \boldsymbol{\theta}_{b_{k+1}}^{b_k} \\ \delta \mathbf{b}_a \\ \delta \mathbf{b}_g \end{bmatrix} = \begin{bmatrix} \mathbf{p}^{b_k} - \hat{\boldsymbol{\alpha}}_{b_{k+1}}^{b_k} \\ \mathbf{v}^{b_k} - \hat{\boldsymbol{\beta}}_{b_{k+1}}^{b_k} \\ \boldsymbol{\theta}^{b_k} \otimes (\hat{\boldsymbol{\gamma}}_{b_{k+1}}^{b_k})^{-1} \\ \mathbf{b}_w^{b_{k+1}} - \mathbf{b}_w^{b_k} \\ \mathbf{b}_a^{b_{k+1}} - \mathbf{b}_a^{b_k} \end{bmatrix}$$

where $\mathbf{p}^{b_k} = \mathbf{R}_w^{b_k} (\mathbf{p}_{b_{k+1}}^w - \mathbf{p}_{b_k}^w + \frac{1}{2} \mathbf{g}^w \Delta t_k^2 - \mathbf{v}_{b_k}^w \Delta t_k)$, $\mathbf{v}^{b_k} = \mathbf{R}_w^{b_k} (\mathbf{v}_{b_{k+1}}^w + \mathbf{g}^w \Delta t_k - \mathbf{v}_{b_k}^w)$ and $\boldsymbol{\theta}^{b_k} = \mathbf{q}_w^{b_k} \otimes \mathbf{q}_{b_{k+1}}^w$ are the estimated position, velocity and orientation referred to the local body frame of b_k (see [4] for details).

3) *Vision Factor*: The vision factor consists of the re-projection error for the tracked features. We use a strategy similar to [4], comparing the reprojection of all features in the current frame with their first observations. We define the visual residual as

$$\mathbf{r}_C(\hat{\mathbf{z}}_l^{c_j}, \mathcal{X}) = \|\mathbf{u}_{l_k}^{c_j} - \pi(\mathbf{R}_{c_i}^{c_j}, \mathbf{T}_{c_i}^{c_j}, \mathbf{P}_{l_k}^{c_i})\|$$

where $\mathbf{u}_{l_k}^{c_j}$ is the coordinate of feature l_k in the image of the camera frame j , $\pi(\cdot)$ is the projection that converts homogeneous coordinates into image coordinates, $\mathbf{R}_{c_i}^{c_j}, \mathbf{T}_{c_i}^{c_j}$ represent the frame transformations (rotation and translation)

from the camera frame i to j , which are inferred from state poses, and $\mathbf{P}_k^{c_i}$ is the 3D position of the k th feature in the first observation frame i . The vision factor iterates through all the frames and all the tracked features in the estimated state.

4) *Anchor position estimation*: In the above discussion, we assume the anchor coordinates are available for the optimizer. As the anchor position is initially unknown, it needs to be estimated. Therefore, the state vector to estimate at the start becomes $\mathcal{X} = [\mathbf{P}_A^{iw}, \mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_s, \mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n, l_0, l_1, \dots, l_m]$. The cost function and factors are kept unchanged. The optimization result of \mathbf{P}_A^{iw} after the initialization stage is saved as a fixed value. We fixed the anchor position with two considerations: a) The anchor is static in practice, and b) usually, the initialization phase can be controlled and the robot moves in proximity of the anchor. Although the distance measurement error is not correlated with the distance value, the estimation of the anchor position depends on the distance. In other words, the ratio of trajectory length in the initialization with the distance affects the estimation. So we treat the initial estimate as fixed value.

C. Distributed Collaborative SLAM

Another significant benefit of our ranging-aided system is that, with a common anchor, robots can directly estimate inter-robot transformations when they rendezvous. Robots simply need to send their current position and anchor position while ranging their neighbors. After receiving this information twice, a robot can calculate the transformation matrix between itself and the sender. Once the transformation matrix is correctly estimated, all information received from neighbors can be correctly placed in the robot's frame.

Estimating the transformation between robots is a critical requirement for multi-robot systems. We mark the transformation of coordinate systems from robot i to j as $\mathbf{T}_i^j = [\mathbf{R}_i^j, \mathbf{t}_i^j]$, where \mathbf{R}_i^j is 3×3 matrix representing rotation and \mathbf{t}_i^j is 3×1 vector of translation. With the help of accelerometers and gyroscopes, we can establish the direction of gravity and define the same z axis for all robots. VINS-Mono [4] uses the same strategy: the z axis is aligned with the opposite of gravity when creating the coordinate system. Therefore, only the yaw angle θ and 3D offset \mathbf{t}_i^j between two coordinate systems need to be estimated. The transformation \mathbf{T}_i^j consists of:

$$\mathbf{R}_i^j = \begin{bmatrix} \cos(\theta) & -\sin(\theta) & 0 \\ \sin(\theta) & \cos(\theta) & 0 \\ 0 & 0 & 1 \end{bmatrix}, \mathbf{t}_i^j = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$$

Using the estimation of the anchor position for two robots, we can easily get $t_z = (\mathbf{P}_A^{jw} - \mathbf{P}_A^{iw})_z$, which is the projection of the difference vector on z axis. This leaves three parameters (θ, t_x, t_z) to be estimated.

Fig.4 shows how the remaining parameters are estimated. Let us assume two robots i and j moving independently in 2D for simplicity (without loss of generality, as we already

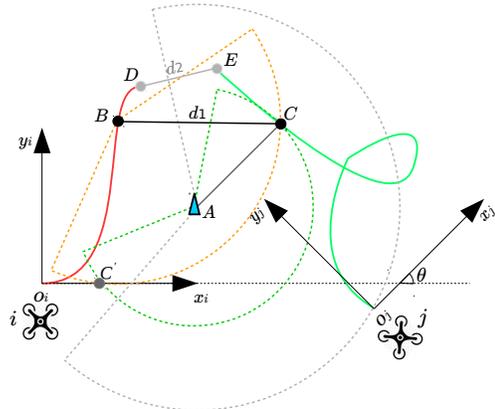


Fig. 4. Transformation matrix estimation for multiple robots scenario.

know the z axis offset t_z). Both robots have their own coordinate systems, $o_i x_i y_i$ and $o_j x_j y_j$. They also have the anchor A position after initialization, and have their own trajectory tracked in their respective coordinate systems. When the robot i passes point B and robot j passes point C , they enter in communication range and they can both obtain a reciprocal distance measurement d_1 . Simultaneously, they transmit their current position and A 's position (which are in their own coordinate system) to each other.

We take the view of robot i to illustrate our solution. When robot i receives the position of the anchor and the current position of robot j (point C) expressed in j 's frame, $\mathbf{P}_{b_A}^{jw}$ and $\mathbf{P}_{b_C}^{jw}$, respectively. Robot i then knows the origin of $o_j x_j y_j$ must be on the gray circle around the anchor with radius $|A o_j|$ (known from \mathbf{P}_A^{jw}). In addition, robot i calculates $|AC|$ from \mathbf{P}_A^{jw} and $\mathbf{P}_{b_C}^{jw}$ to find robot j 's position C must lie on a green circle around A with radius $|AC|$. As we know, the distance between point B and C is the range measurement d_1 , and point C also must lie on the orange circle around its current position B with radius d_1 .

Therefore, the current position lies at one of the intersections between the green and orange circles, C and C' . To find which intersection is the correct one, we take an additional measurement d_2 between two following points D and E . Following the same procedure, we can find the candidates of E . For clarity, we did not draw these circles for D and E . By comparing the relation between E, E', C and C' with the true motion from C to E , we can find C and determine the transformation \mathbf{T}_i^j . This is the same as solving:

$$\begin{cases} \mathbf{P}_{b_A}^{jw} = \mathbf{T}_i^j \cdot \mathbf{P}_{b_A}^{iw} \\ d_1 = \|\mathbf{P}_{b_B}^{jw} - \mathbf{T}_i^j \cdot \mathbf{P}_{b_C}^{iw}\| \\ d_2 = \|\mathbf{P}_{b_D}^{jw} - \mathbf{T}_i^j \cdot \mathbf{P}_{b_E}^{iw}\| \end{cases} \quad (7)$$

where $\mathbf{P}_{b_Q}^{rw}$ represent the 3D position vector of points Q in world coordinate frame of robot r .

IV. EXPERIMENTS

We validate our algorithms using public datasets and real hardware experiments. We choose VINS-Mono¹ and

¹<https://github.com/HKUST-Aerial-Robotics/VINS-Mono>

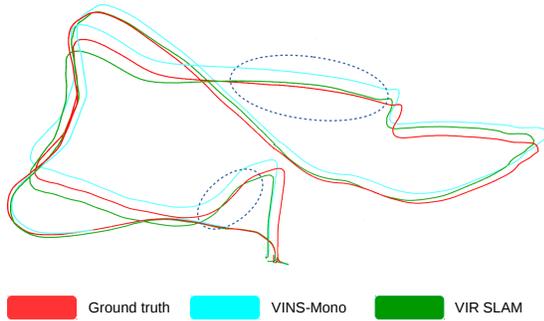


Fig. 5. VIR results in EuRoC MH05 Dataset. ATE error of our method is 0.291m, compared to VINS-Mono MONO (0.388m). We set our method is closer to the ground truth in most sections, which testify the capability to correct the drift error.

TABLE I
ERROR COMPARISON BETWEEN VINS-MONO AND OUR VIR SLAM.

ATE Error (m)	VINS-Mono	VIR SLAM	Improvement
EuRoC MH_01	0.186	0.178	4.33 %
EuRoC MH_02	0.240	0.188	21.76 %
EuRoC MH_03	0.271	0.260	4.04 %
EuRoC MH_04	0.402	0.366	9.0 %
EuRoC MH_05	0.388	0.291	24.84 %
Lab Seq1	0.266	0.193	27.81 %
Lab Seq2	0.195	0.169	15.08 %

compare with it using the TUM evaluation tool². We compute the absolute trajectory error (ATE) when ground truth is available. For our large area experiments we do not have ground truth, and we start and end the experiment in the same location and then calculate the start-to-end error.

A. EuRoC Dataset Experiments

We test our single robot tracking system on the EuRoC [15] dataset. We simulate the UWB ranging measurements from ground truth data. The static anchor is assumed in the origin of the frame created during robot initialization. We add Gaussian white noise $\mathcal{N}(0, 0.05)$ to model the error of our UWB sensor. The sliding window size of VINS-Mono is set to 10. For our system, we tested with a short window size of 10 and a long window size of 100.

Table I shows our VIR-SLAM outperform VINS-Mono in terms of ATE. As an example, we show the trajectory of the comparison for the EuRoC MH-05 sequence in Fig.5. The green trajectory of our method is closer to the red ground truth. From the ellipse indicator, we can see our estimation can correct the error even when drifting.

B. Single Robot experiments

We also test our system with a Spiri robot with real UWB sensors, shown in Fig. 6. The system has a D435 Realsense camera but we use it as a monocular camera. The IMU measurements come from the Pixracer flight controller, and the robot carries an Nvidia TX1 as an on-board computer. We

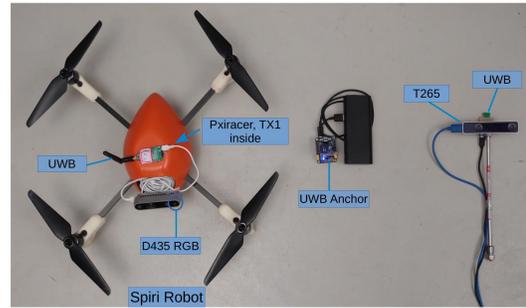


Fig. 6. For single robot experiments, we use the Spiri robot (w/ Pixracer and NVIDIA TX1). We add a Realsense D435 camera (although we only use RGB) and UWB sensor module. The IMU measurements are obtained from the Pixracer. The static is in the middle, while the Realsense T265 on the right simulates a second robot in our multi-robot experiment.

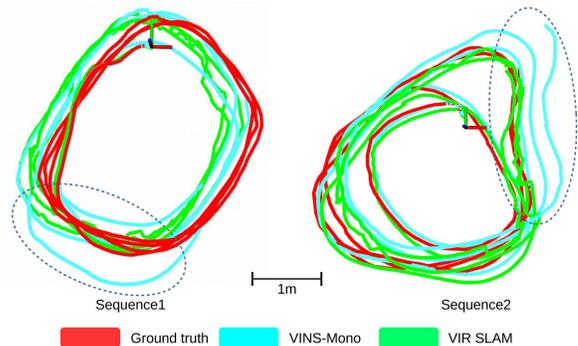


Fig. 7. VIR indoor experiments with OptiTrack ground truth. The left and right figures are two sequences. All trajectories start at the origin and are not aligned. Our estimator does not present excessive drift.

test the system in our lab with an OptiTrack motion capture system as ground truth.

We manually control the robot and collect two sequences: the results are listed in Table I. To clearly see the difference, we show the trajectories from same origin. This is different from aligning the two trajectory to compute the ATE. As Fig. 7 shows, our estimator has less accumulated error than VINS-Mono.

We also test our system performance in large atrium. The environment in the atrium is quite challenging: featureless walls, low light conditions, glass walls with reflections, etc. We move the robot around the atrium and make sure the robot ends at the same point as it starts. Then we compare the start-to-end error. As Fig. 8 and Table II shows, the VINS-Mono estimation has significant drift, more than 5.4m. Although the loop closure version can correct the drift as the loop closure is detected, there is still a clear error in the trajectory. However, our VIR SLAM works correctly without any loop closure. The start and end point have a small translation error (0.148m) in 2D, but our system introduces a bigger accumulated error on the z axis. We believe this is due to the anchor being close to the horizontal plane of the robot, and the small difference in measurement does not help correct the z coordinate. We manually checked the UWB ranging information along the top and bottom edges and we can confirm that our estimation is at the right position on the

²<https://vision.in.tum.de/data/datasets/rgbd-dataset/tools>

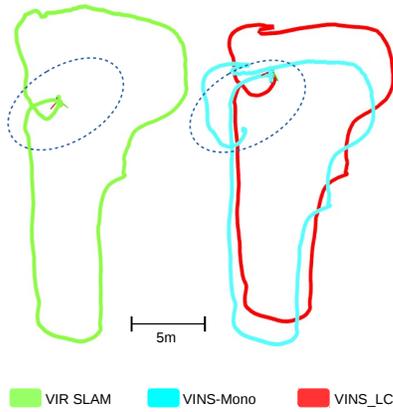


Fig. 8. VIR results in the atrium, which is very visually challenging. The robot starts and ends at the same point and we compare the start-to-end error. As the figure shows, the start and end point of VIR-SLAM are overlapped. Although the VINS-Mono with loop closure (VINS.LC) can close the loop, it does not eliminate the accumulated error from the VIO.

boundaries.

TABLE II
START-TO-END ERROR COMPARISON

	End point	2D Error(m)	3D Error(m)
VINS-Mono	(4.29, 3.35, 0.52)	5.439	5.465
VIR-SLAM	(0.18, 0.43, -0.38)	0.464	0.602
VINS-LC	(-0.04, 0.14, 0.84)	0.148	0.853

C. Multi-Robot Experiments

We also test our multi-robot technique, as shown in Fig. 1. We manually and independently control two robots to make a trajectory similar to the "MIST", the name of our lab. We place one static UWB anchor in the environment. The first robot is imply a Realsense T265 and an UWB module, show in Fig. 6 (right). We move it to form an "M" shape. The other robot, Spiri, moves along a "IST" and it controlled independently. Their trajectories in each robot's frame are shown at the bottom of Fig. 1. With simply two inter-robot measurements (with data exchange), robot 1 can estimate the transformation of robot 2 and map robot 2's trajectory in its own frame as shown at the top of Fig. 1.

V. CONCLUSIONS AND DISCUSSIONS

In this paper, we propose VIR-SLAM, a novel SLAM paradigm combining vision, inertial, and UWB sensors. By arbitrarily setting up a static UWB anchor in the environment, robots can have drift-free state estimation and collaboration on SLAM. Our solution combines the accurate relative pose estimation from VIO and enhances it using ranging to correct the accumulated error. As our experiments show, the introduced static anchor can help correct the drift effectively to improve localization accuracy. We also show an example with two independent robots mapping each other's trajectory to their own frame after obtaining two range measurements. This technique allows the robots to find the inter-robot transformation, which is extremely useful for multi-robot SLAM.

REFERENCES

- [1] J. Delmerico and D. Scaramuzza, "A benchmark comparison of monocular visual-inertial odometry algorithms for flying robots," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, pp. 2502–2509.
- [2] T. Lupton and S. Sukkarieh, "Preintegration-visual-Inertial-Aided Navigation for High-Dynamic Motion in Built Environments Without Initial Conditions," *IEEE Transactions on Robotics*, vol. 28, no. 1, pp. 61–76, Feb. 2012.
- [3] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-Manifold Preintegration for Real-Time Visual-Inertial Odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, Feb. 2017.
- [4] T. Qin, P. Li, and S. Shen, "VINS-Mono: A Robust and Versatile Monocular Visual-Inertial State Estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, Aug. 2018.
- [5] "Loco Positioning system | Bitcraze." [Online]. Available: <https://www.bitcraze.io/loco-pos-system/>
- [6] S. Saeedi, M. Trentini, M. Seto, and H. Li, "multi robot slam review Multiple-Robot Simultaneous Localization and Mapping: A Review," *Journal of Field Robotics*, vol. 33, no. 1, pp. 3–46, 2016.
- [7] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [8] J. Engel, V. Koltun, and D. Cremers, "Direct Sparse Odometry," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 3, pp. 611–625, Mar. 2018.
- [9] A. Prorok and A. Martinoli, "Accurate indoor localization with ultra-wideband using spatial models and collaboration," *The International Journal of Robotics Research*, vol. 33, no. 4, pp. 547–568, 2014, 05.
- [10] M. W. Mueller, M. Hamer, and R. D'Andrea, "Fusing ultra-wideband range measurements with accelerometers and rate gyroscopes for quadcopter state estimation," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. Seattle, WA, USA: IEEE, May 2015, pp. 1730–1736.
- [11] X. Fang, C. Wang, T.-M. Nguyen, and L. Xie, "Graph Optimization Approach to Localization with Range Measurements," *arXiv:1802.10276 [cs]*, Feb. 2018, arXiv: 1802.10276.
- [12] J. Tiemann, A. Ramsey, and C. Wietfeld, "Enhanced UAV Indoor Navigation through SLAM-Augmented UWB Localization," in *2018 IEEE International Conference on Communications Workshops (ICC Workshops)*. Kansas City, MO: IEEE, May 2018, pp. 1–6.
- [13] C. Wang, H. Zhang, T.-M. Nguyen, and L. Xie, "Ultra-wideband aided fast localization and mapping system," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1602–1609.
- [14] Q. Shi, X. Cui, W. Li, Y. Xia, and M. Lu, "Visual-UWB Navigation System for Unknown Environments," p. 11.
- [15] 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, Sep. 2016.
- [16] S. Saeedi, M. Trentini, M. Seto, and H. Li, "Multiple-Robot Simultaneous Localization and Mapping: A Review," *Journal of Field Robotics*, vol. 33, no. 1, pp. 3–46, 2016.
- [17] P. Schmuck and M. Chli, "CCM-SLAM: Robust and efficient centralized collaborative monocular simultaneous localization and mapping for robotic teams," *Journal of Field Robotics*, vol. 36, no. 4, pp. 763–781, 2019.
- [18] S. Choudhary, L. Carlone, C. Nieto, J. Rogers, H. I. Christensen, and F. Dellaert, "Distributed Mapping with Privacy and Communication Constraints: Lightweight Algorithms and Object-based Models," *arXiv:1702.03435 [cs]*, Feb. 2017, arXiv: 1702.03435.
- [19] J. G. Mangelson, D. Dominic, R. M. Eustice, and R. Vasudevan, "Pairwise consistent measurement set maximization for robust multi-robot map merging," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 2916–2923.
- [20] N. Trawny, X. S. Zhou, and S. I. Roumeliotis, "3d relative pose estimation from six distances," 2009.
- [21] F. M. Martel, J. Sidorenko, C. Bodensteiner, M. Arens, and U. Hugenobler, "Unique 4-DOF Relative Pose Estimation with Six Distances for UWB/V-SLAM-Based Devices," *Sensors*, vol. 19, no. 20, p. 4366, Oct. 2019.