

# Omni-swarm: A Decentralized Omnidirectional Visual-Inertial-UWB State Estimation System for Aerial Swarms

Hao Xu, Yichen Zhang, Boyu Zhou, Luqi Wang, Xinjie Yao, Guotao Meng, Shaojie Shen

**Abstract**—Decentralized state estimation is one of the most fundamental components of autonomous aerial swarm systems in GPS-denied areas yet it still remains a highly challenging research topic. Omni-swarm, a decentralized omnidirectional visual-inertial-UWB state estimation system for aerial swarms, is proposed in this paper to address this research niche. To solve the issues of observability, complicated initialization, insufficient accuracy, and lack of global consistency, we introduce an omnidirectional perception front-end in Omni-swarm. It consists of stereo wide-FoV cameras and ultra-wideband sensors, visual-inertial odometry, multi-drone map-based localization, and visual drone tracking algorithms. The measurements from the front-end are fused with graph-based optimization in the back-end. The proposed method achieves centimeter-level relative state estimation accuracy while guaranteeing global consistency in the aerial swarm, as evidenced by the experimental results. Moreover, supported by Omni-swarm, inter-drone collision avoidance can be accomplished without any external devices, demonstrating the potential of Omni-swarm as the foundation of autonomous aerial swarms.

**Index Terms**—Swarms, aerial systems: perception and autonomy, multi-robot systems, sensor fusion

## I. INTRODUCTION

FOR any aerial robotics system, the estimation of states, including positions and attitudes, is crucial. The estimation system lays a solid foundation for higher-level functions, such as path planning [1] and mapping [2]. The state estimation problems for single-drone systems are currently well-addressed through approaches such as visual-inertial odometry (VIO) [3]–[5] and LiDAR odometry [6], [7]. However, when we look beyond a single-drone to multiple drones working as an aerial swarm, the problem becomes much more complicated. In a swarm, each drone needs to estimate its ego state and also obtain the relative poses of other drones.

To date, the vast majority of aerial swarm researchers have adopted external devices, such as motion capture systems [8], ultra-wideband (UWB) systems with anchors [9] and GPS [10], and RTK-GPS [11] systems, to provide state estimations,

Manuscript received March 31, 2021; revised December 12, 2021 and March 25, 2022; accepted May 9, 2022. This paper was recommended for publication by Editor-in-Chief L. Kevin and Editor R. Paolo upon evaluation of the reviewers' comments. This work was supported by the HKUST Postgraduate Studentship, UGC RMGS20EG20 and HDJI lab. (Corresponding author: Hao Xu.)

All authors are with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong, China. {hxubc, yzhangec, bzhouai, lxwang, gmeng, xyaoab}@connect.ust.hk, eeshaojie@ust.hk

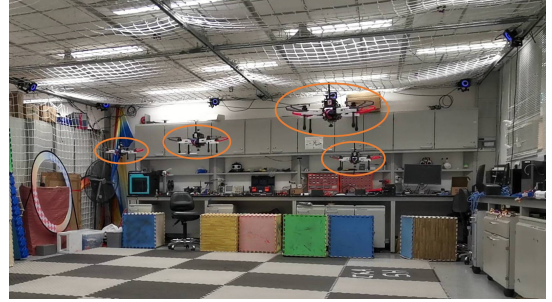


Fig. 1: Indoor aerial swarm formation flight with four drones. The customized drone platforms are circled.

which significantly limits the application of aerial swarms in the real world. Although motion capture systems and UWB modules with anchors can work in indoor environments with decent accuracy, they are centralized systems and require bulky external devices, meaning that they are susceptible to losing the central devices and are challenging to deploy. A requirement of a practical swarm application is simple deployment, which such approaches typically fail to meet.

The decentralized scheme of swarm robots is becoming popular in swarm robotics research [12], [12]–[17] because of its significant advantages. A robot swarm with this scheme does not require all robots to have stable communication with a central computer, which makes it more flexible in real-world environments where communication is limited. Additionally, each robot can act largely independently from the remainder of the team, making the whole system more fault-tolerant to single-point failure. To build a fully autonomous decentralized multi-robot system (multi-aerial-robot swarm, also known as an aerial swarm), a key problem is how to achieve relative state estimation in a decentralized fashion. The primary motivation of this paper is to solve this fundamental problem with additional global consistency of the estimated states, laying a solid foundation for a decentralized aerial swarm.

Recently, researchers have started to develop approaches to perform decentralized relative state estimation on aerial swarms. One of the most straightforward ideas is to utilize visual object detection to detect the drones in the aerial swarm to estimate the relative state [18]–[21]. Fusing distance measurements from the UWB and odometry (usually VIO) [18]–[23] is another viable approach, along with estimating relative states from common environment features captured by the drones [24], [25] or extracting relative states from sparse maps built by the aerial swarms [16], [17], [26], [27].

However, the practicability of these methods is limited by some serious issues, as follows,

- 1) Observability issue caused by a restricted field of view (FoV). For visual-drone-detection-based methods [18], [22], [23], the relative states are observable only when others drones are in the drone’s FoV.
- 2) Complicated initialization. UWB-odometry fusion methods [18]–[21] require large motions to initialize the system. The complex initialization procedure may cause severe safety issues and even crashes.
- 3) Insufficient accuracy. The estimated position errors in previous works [19]–[21] are generally around 20 cm–50 cm, meaning that these swarm systems can scarcely be adopted for use in confined indoor spaces or close formation scenarios.
- 4) Lack of global consistency. For all the current relative state estimation methods, the estimated poses drift and the ego state estimation cause consistency issues. Global consistency becomes especially important when we expect to build global maps based on state estimation.

To address the challenges, in this paper, we extend our previous method [18] and propose **Omni-swarm: a decentralized omnidirectional visual-inertial-UWB state estimation system for aerial swarms**, which combines the advantages of the UWB-odometry fusion method [18]–[23], visual-object-detection-based methods [18], [22], [23], and map-based methods [16], [17], [24], [25], [28], [29].

The most important contribution to address the aforementioned issues is the introduction of an omnidirectional perception front-end. The front-end includes the hardware capturing omnidirectional visual information and the algorithms processing this information. Specifically, we use two fisheye cameras with a wide-FoV (up to  $235^\circ$ ) to achieve omnidirectional observation of the surrounding area. We develop VINS-Fisheye as the ego-motion estimator of the front-end, which is a VIO state estimator using the measurements from stereo wide-FoV fisheye cameras for ego-motion estimation. Multi-drone map-based localization (MDML) based on real-time-generated sparse maps is further introduced to ensure the global consistency of state estimation and to achieve fast initialization. Finally, a visual drone tracking module detects and tracks the other drones to provide accurate relative pose estimation of the tracked targets. In addition to omnidirectional visual information, we also use UWB sensors to measure the relative distance between drones. This measurement can also be considered as omnidirectional.

As the back-end of the state estimation, we adopt the graph-based optimization method, which fuses the measurements from the front-end in real-time to estimate the states of the swarm with high accuracy. In the back-end, we adopt the state-of-the-art outlier rejection method to reduce the errors from the front-end to achieve high-accuracy and robustness.

Omni-swarm is designed to be decentralized, it runs on each drone’s onboard computer individually instead of using a central server. Differing from previous work on relative-state estimation, which only work in a line-of-sight situation [18]–[22], Omni-swarm can estimate the state of the swarm when another drone is in non-line-of-sight if the same place

has been visited. This capability is given by MDML. Another advantage over previous works [18]–[22] that MDML brings the global consistency; the long-term drifting of the ego state is eliminated by the map-based localization and thus guarantee global consistency of the estimation results. Omni-swarm also inherits the high accuracy of our previous work [18], which has a high relative localization accuracy compared to other related works [19]–[22].

In addition, the complex initialization problem is solved in Omni-swarm by introducing various system initialization methods, including map-based initialization and visual drone detection tracking initialization. Another Omni-swarm improvement that helps with practical applications is its plug-and-play feature, based on our newly introduced initialization methods, Omni-swarm allows the temporary joining or exiting of drones. Finally, Omni-swarm has redundancy in the information shared among the swarm and computations performed by the drones, which brings robustness to possible temporary signal loss and partial sensor failures.

To verify the above features of Omni-swarm, we perform comprehensive experiments in simulation and real-world experiments. Moreover, we design an inter-drone collision avoidance experiment to verify Omni-swarm under realistic conditions. The main contributions of this paper are:

- 1) Omni-swarm, a decentralized omnidirectional visual-inertial-UWB swarm state estimation system. Extensive experiments are conducted to validate Omni-swarm.
- 2) Open-source releases of the software and the custom datasets have been made public<sup>1</sup>.

In our previous work [18], we proposed a two-stage visual-inertial-UWB fusion method for relative state estimation. The method has the advantages of both visual-object-detection-based relative state estimation [22], [23] and UWB-odometry fusion relative state estimation [19], [20]. However, global consistency is absent. Also, the method still suffers from the same complicated initialization issue and the observability issue caused by the restricted FoV, as with other related methods. Although the previous method works in a non-line-of-sight case, the relative estimation accuracy will deteriorate. In extreme situations, the relative state can become unobservable. For example, drones fly parallelly side-by-side. In this paper, the initialization and global consistency issues are addressed by introducing multi-drone map-based localization, and the observability is fixed by using omnidirectional cameras.

In this paper, related works are discussed in Sect. II. Omni-swarm is briefly introduced in Sect. III, and a clear definition of the state estimation problem is defined in Sect. III-B. The front-end of Omni-swarm is presented in Sect. IV, while the back-end is introduced in Sect. V. Experimental results are discussed in Sect. VII. Finally, we conclude the paper in Sect. VIII and introduce potential future works.

## II. RELATED WORKS

### A. Visual inertial odometry on Aerial Swarm

In order to overcome the state estimation issues in multiple environments, including GPS-denied areas, visual-inertial si-

<sup>1</sup><https://github.com/HKUST-Aerial-Robotics/Omni-swarm>

multaneous localization and mapping (visual-inertial SLAM) [30]–[32] and visual-inertial odometry (VIO) [3], [5], [33], [34] are widely adopted on single drone systems. These visual inertial systems fuse the visual images and the inertial measurement unit (IMU) data together to estimate ego-motion without requiring external devices. Although VIO systems can be effortlessly adopted in various environments without cumbersome external devices and they are directly utilized by researchers on swarms for formation flights [35], [36], the state estimations suffer from severe drift, which can be detrimental to multi-drone systems. Specifically, the different drifts of individual drones in an aerial swarm induce distinct position estimations at the same location, possibly causing fatal crashes in collaborative missions in the absence of other correction methods.

### B. Relative Swarm State Estimation

One of the most effective methods to compensate for the drift of VIO in aerial swarms is to incorporate the relative state estimation between the drones. The existing relative state estimation methods can generally be divided into the following three categories: 1) UWB-odometry fusion relative state estimation methods [18]–[21], [37], 2) visual-object-detection-based relative state estimation methods [18], [22], [23], and 3) environment-feature-based relative state estimation methods [24], [25].

However, all of these methods have drawbacks in real application scenarios. The previous UWB-odometry based methods [19]–[21] fuse the ego-motion estimated by VIO and UWB distance measurements to achieve relative state estimation. These methods can only estimate relative localization in line-of-sight situations and provide merely a meter to decimeter level of accuracy. Even in line-of-sight situations, the relative localization can become unobservable in some cases, e.g., parallel flight tasks [38]. Furthermore, the initialization requires a certain amount of motion, usually several meters. Both of these drawbacks substantially limit the practicability of the UWB-visual fusion methods in feature-rich narrow environments. A viable idea to address observability issue was proposed by Nguyen et al. [38]. Briefly, a drone in the swarm keeps an irregular motion to guarantee observability of the relative state estimation, which may present safety issues and limit the aerial swarm’s cooperation. In Omniswarm, the observability issue can be solved by visual drone tracking and map-based localization while the formation is flying. In addition, by introducing map-based localization, we can achieve fast initialization without motion.

Visual-object-detection-based methods [18], [22], [23] are capable of delivering centimeter-level accuracy. However, the accuracy of these methods in a swarm of drones is highly dependent on the distance the drones are from each other. Data association from the detection results to the corresponding drone ID is also an issue when all the drones appear identical. A coupled-probabilistic-data-association-filter (CPDAF)-based approach is proposed in [23] to associate the detection result and estimate the relative state. Nevertheless, the visibility requirement between the drones and the limited FoV of their

cameras restrict the swarm formation and the distance between drones for stable state estimation. In this paper, the FoV issue is addressed by introducing an omnidirectional front-end. Due to the UWB measurements, the proposed method can still estimate the relative state when the distances between the drones are too far for them to detect each other.

Finally, all the aforementioned methods are limited to relative localization, and global consistency of the estimated states cannot be guaranteed. In this paper, global consistency is guaranteed by introducing map-based localization. This will also help us with global mapping in our future work.

### C. Environmental-Features-based Method and Collaborative Simultaneous Localization and Mapping

Environment-feature-based methods [24], [25] rely on multiple drones’ common environmental texture features to estimate the relative poses. The methods require sufficient overlapped features between the view of the drones, limiting the swarm formation and heading, while only working in feature-rich environments.

Collaborative simultaneous localization and mapping (CSLAM) methods [14]–[17], [26]–[29] focus on sparse and dense mapping utilizing the sensors on the multi-robot system, and can estimate the states of aerial swarms. DDF-SAM [26], [27] presents a landmark-based back-end implementation without providing the front-end and data association between features among the swarm. Choudhary et al. [16] utilized objects as landmarks, requiring known objects and their 3D models in the scene, which limits the method’s real-world practicability. DOOR-SLAM [17], [28], [29] achieves localization of robots using VIO incorporating loop closure detection (also known as map-based localization in this paper). [16], [17], [28], [29] estimate the state with the pose graph optimization (PGO). However, the accuracy of relative localization is crucial for cooperation inside the aerial swarm, and this has not been featured in these papers. Generally speaking, the PGO methods are not accurate because they do not directly use the features for relative localization. Zhu et al. [15] propose a distributed visual-inertial fusion for cooperative localization and reach centimeter-level accuracy. In [15] the environmental features are tightly coupled in the cooperative localization.

Due to the nature of map-based localization, the above CSLAM methods are limited to feature-rich environments and are also limited by the camera FoV. In this paper, with an omnidirectional front-end and UWB measurements, Omniswarm will not be limited by the FoV and can still estimate relative state without rich common environmental features.

### D. Visual-inertial-UWB Fusion with Global Consistency

Some of the works mentioned as being UWB-odometry fusion [18], [20], [21] can be considered to be visual-inertial-UWB fusion (also known as visual-inertial-range fusion, or VIR fusion). However, global consistency is absent in these works. Unlike these relative localization methods, a class of works [18], [20], [21] uses VIR fusion to achieve global localization by placing a fixed UWB anchor in the environment. With the help of this anchor, these methods can

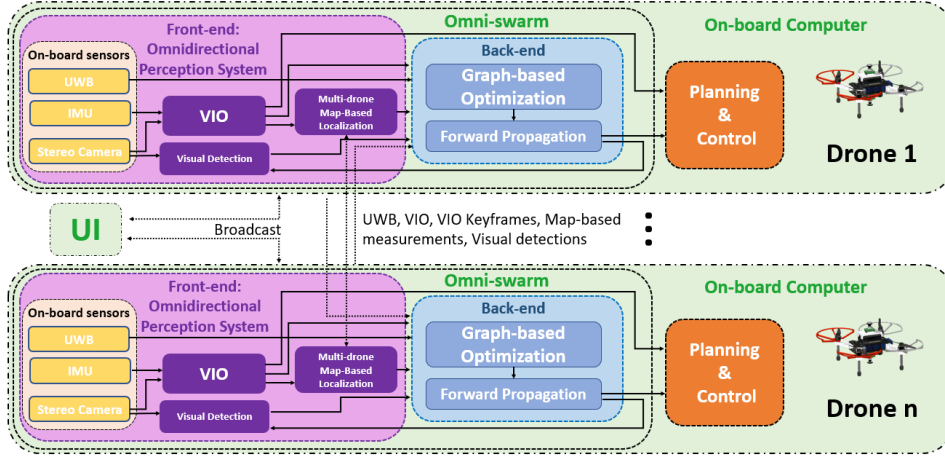


Fig. 2: The system architecture of Omni-swarm. The data from the onboard sensors are processed, and then broadcast to all the other drones. The swarm state estimation on each of the onboard computers collects both onboard and broadcast information, including the relative distance from the UWB modules, the VIO, map-based measurements, and the detection results, and performs optimization prediction to obtain real-time relative state estimations. The estimation results are then sent back to facilitate the matching procedure of detection and tracking and meanwhile serve the planning and control. The ground station obtains the drones' real-time information to monitor the flight status and concurrently sends the commands to the drones. All the communications between the devices are through UWB broadcast.

effectively eliminate the drift of odometry. The drawback of these methods is that they require additional infrastructure to be placed on the ground and they maintain global consistency only when drones are in the line of the sight of the anchors, limiting their real-world practical value. In this paper, the proposed Omni-swarm guarantees global consistency without any external infrastructure, which is more flexible in real-world applications.

### III. SYSTEM OVERVIEW

#### A. Notation

To aid understanding of the proposed system, the notations are defined below.

#### NOMENCLATURE

$\hat{(\cdot)}$	The estimated state.
$d_{i,j}^t$	Distance between drone $i$ and drone $j$ at time $t$ .
$\mathcal{F}_k^t$	The keyframe of the drone $k$ at time $t$ , which contains the 4D pose ${}^{v_k}\hat{\mathbf{P}}_i^t$ to be estimated, a few virtual camera keyframes ${}^c\mathcal{K}_k^t$ and other essential information of the drone.
${}^c\mathcal{K}_k^t$	The keyframe of drone $k$ 's virtual camera $c$ at time $t$ , which contains the global descriptor, local features, virtual camera's extrinsic and other essential information of the drone. The virtual camera $c$ is cropped from the raw fisheye camera.
$\mathcal{SF}_k^t$	Equal to $[\mathcal{F}_1^t \mathcal{F}_2^t \dots \mathcal{F}_n^t]$ . The swarm keyframe of the drone $k$ at time $t$ , which contains $n$ keyframes.
$\mathcal{G}_k$	The graph built on drone $k$ for state estimation.
$(\cdot)_R$	The rotation part of the transformation matrix.
$(\cdot)_P$	The corresponding 4-DoF transformation matrix.
$(\cdot)_T$	The corresponding 6-DoF transformation matrix.
$(\cdot)_X$	The translation part of the transformation matrix.
$\mathbf{z}(\cdot)^t$	The measurement data at time $t$ .
$(\cdot)_\psi$	The yaw angle of the rotation matrix.
$(\mathcal{F})_f$	The global descriptor of the keyframe $\mathcal{F}$ .
$(\mathcal{K})_{\mathcal{F}}$	The corresponding keyframe $\mathcal{F}_k^t$ of virtual camera keyframe $\mathcal{K}$ .

$(\mathcal{F})_{lf}$	The local descriptors of the features of the keyframe $\mathcal{F}$ .
$\ (\cdot)\ $	Euclidean norm of $(\cdot)$ if $(\cdot)$ is a vector or matrix; otherwise if $\cdot$ is a set, $\ (\cdot)\ $ is its size.
$\ (\cdot)\ _{\Sigma}$	Mahalanobis norm of $\cdot$ .
$\mathcal{D}$	The set of all existing drones, including the currently unavailable drones due to loss of communication, user poweroff and accident.
$\mathcal{D}_a^k$	The set of all available drones for drone $k$ .
$\mathcal{D}_e^k$	The set of all estimated drones of drone $k$ 's state estimation.
$\mathcal{D}_u^k$	The set of all uninitialized drones of drone $k$ 's state estimation, where $\mathcal{D}_u^k = \mathcal{D}_a^k - \mathcal{D}_e^k$ .
${}^{b_k}(\cdot)_i^t$	State of drone $i$ in drone $k$ 's body frame. For simplicity, the pose in the body frame is defined as a 4-DoF pose, i.e., ${}^{b_k}(\cdot)_i^t = ({}^{v_k}\mathbf{P}_k^t)^{-1} {}^{v_k}(\cdot)_i^t$ .
$D_i$	The $i$ th drone.
${}^{v_k}(\cdot)_i^t$	State of drone $i$ in drone $i$ 's local frame.
${}^{v_k}\mathbf{P}_i^t$	Equal to $\begin{bmatrix} \mathbf{R}_z({}^{v_k}\psi_i^t) & {}^{v_k}\mathbf{X}_i^t \\ 0 & 1 \end{bmatrix}$ . The pose of drone $i$ in drone $k$ 's local frame at time $t$ . For simplicity, the notation of $\mathbf{P}_k^t$ represents ${}^{v_k}\mathbf{P}_k^t$ . ${}^{v_k}\mathbf{R}_z({}^{v_k}\psi_i^t)$ represents the rotation matrix rotated over the z axis with angle ${}^{v_k}\psi_i^t = ({}^{v_k}\mathbf{R}_i^t)_\psi$ .
${}^{v_k}\mathbf{T}_i^t$	Equal to $\begin{bmatrix} {}^{v_k}\mathbf{R}_i^t & {}^{v_k}\mathbf{X}_i^t \\ 0 & 1 \end{bmatrix}$ . The 6-DoF pose of drone $i$ in drone $k$ 's local frame at time $t$ . ${}^{v_k}\mathbf{R}_i^t$ represents the rotation matrix
$\tilde{\mathbf{P}}_k^t, \tilde{\mathbf{T}}_k^t$	The 4-DoF and 6-DoF pose, respectively, of drone $k$ in its local frame, as estimated by VIO, which drifts through time. This pose is initialized to an identity matrix after the drone start-up, which also gives the origin and the axis of the local frame.
${}^{v_k}\mathbf{X}_i^t$	Equal to $[{}^{v_k}x_i^t, {}^{v_k}y_i^t, {}^{v_k}z_i^t]^T$ . The translation part of ${}^{v_k}\mathbf{P}_i^t$ .
$\delta\mathbf{P}_i^t$	The transformation matrix from time $t-1$ to $t$ of drone $i$ from the VIO result, i.e., $\mathbf{P}_i^t = \mathbf{P}_i^{t-1} \delta\mathbf{P}_i^t$ .

Suppose our aerial swarm contains up to  $n$  drones. A

drone with ID  $i$ , will be denoted as drone  $i$  or  $D_i$ , where  $i \in \mathcal{D}$ ,  $\mathcal{D} = \{1, 2, 3, \dots, n\}$ . Although Omni-swarm runs independently on each drone, to simplify the discussion, we discuss drone  $k$ 's Omni-swarm by default in the following unless otherwise stated. We say a drone  $i$  is available to drone  $k$  when it satisfies following conditions: 1) drone  $i$ 's Omni-swarm and ego-motion state estimation (VIO) work properly; 2) drone  $i$  and drone  $k$  have a stable network connection. The set of all available drones for drone  $k$  is denoted as  $\mathcal{D}_a^k$ , which also contains the drone itself.

### B. State Estimation Problem of Aerial Swarm

For an aerial swarm that contains a maximum of  $n$  homogeneous drones, the state estimation problem can be represented as follows. For every drone  $k \in \mathcal{D}$ , estimate the 6-DOF pose  ${}^{v_k}\mathbf{T}_i^t$  for every drone  $i \in \mathcal{D}_a^k$  at time  $t$  in drone  $k$ 's local frame. The state estimation problem for drone  $k$  can be split into two parts:

- 1) Estimating the ego-motion state of drone  $k$  in a local frame, i.e.,  $\hat{\mathbf{T}}_k^t$ .
- 2) Estimating the state of any other arbitrary drone  $i$ , i.e.,  ${}^{v_k}\hat{\mathbf{T}}_i^t$ , and 4-DoF relative state  ${}^{b_k}\hat{\mathbf{P}}_i^t$ .

The VIO utilizes the gravitational acceleration measured by the IMU to help extract the roll and pitch angles in the attitude. Because the gravity acceleration is consistent among drones, with the estimation of the 4-DoF pose  ${}^{v_k}\hat{\mathbf{P}}_i^t$  and relative pose  ${}^{b_k}\hat{\mathbf{P}}_i^t$ , we are able to combine the drone's own VIO  $\tilde{\mathbf{T}}_i^t$  to obtain the 6-DoF pose:

$${}^{v_k}\hat{\mathbf{T}}_i^t = \begin{bmatrix} \mathbf{R}_z \left( ({}^{v_k}\hat{\mathbf{P}}_i^t)_\psi - (\tilde{\mathbf{T}}_i^t)_\psi \right) \tilde{\mathbf{R}}_i^t & ({}^{v_k}\hat{\mathbf{P}}_i^t)_X \\ 0 & 1 \end{bmatrix}, \quad (1)$$

where  $\mathbf{R}_z \left( ({}^{v_k}\hat{\mathbf{P}}_i^t)_\psi - (\tilde{\mathbf{T}}_i^t)_\psi \right)$  eliminates the yaw drift of the rotation  $\tilde{\mathbf{R}}_i^t$  estimated by the VIO, where  $\mathbf{R}_z(\psi)$  is the rotation matrix rotated over the z axis with angle  $\psi$ .

### C. Global Consistency of the State Estimation

In SLAM research [39], [40], global consistency of state estimation represents that the estimate results are drift-free, e.g., the  ${}^{b_k}\hat{\mathbf{P}}_i^t$  do not drift along with the robot move.

In our previous work [18], we focused on estimating  ${}^{b_k}\hat{\mathbf{P}}_i^t$ , which is the relative state estimation for the aerial swarm. The ego-motion was directly estimated by the VIO; i.e.,  $\hat{\mathbf{T}}_k^t = \tilde{\mathbf{T}}_k^t$  was assumed. However, the VIO is concerned with local accuracy and suffers from long-term drifting. In this paper, by adopting the map-based localization method, we can estimate the state of the aerial swarm with guaranteed global consistency, which means both  ${}^{v_k}\hat{\mathbf{T}}_k^t$  and  ${}^{b_k}\hat{\mathbf{P}}_i^t$  are estimated by the proposed state estimator for the aerial swarm.

### D. Observability and Initialization

For all state estimation systems, initialization is critical, especially for Omni-swarm. The necessary condition for initializing the state estimation is that observability must be satisfied. Initialization with improper states and insufficient

observability are likely to direct Omni-swarm to incorrect state estimation.

Compared to previous works, a major highlight of the system proposed in this paper is the extended observability, which is important for its practical application. First, ego-motion estimation, namely the VIO in our system, is always assumed to be available for drones in an aerial swarm since it is the fundamental module for stable flight. Therefore, the Omni-swarm observability problem is focused on the relative pose. The observability of drone  $i$ 's states estimated by drone  $k$  can be classified into two levels: 1) 3-DoF observable: The position  ${}^{v_k}\mathbf{X}_i^t$  is observable, and 2) 6-DoF observable: The 6-DoF pose  ${}^{v_k}\mathbf{T}_i^t$ . When the 4-DoF pose  ${}^{v_k}\mathbf{P}_i^t$  and VIO  $\tilde{\mathbf{T}}_i^t$  are observable,  ${}^{v_k}\mathbf{T}_i^t$  will also be observable, following Eq. (1).

However, the observability of each drone in a swarm running Omni-swarm may be different. Thus, to provide the ability of plug-and-play, we track the observability of each available drone in the swarm and only initialize and estimate the state of observable drones. We use set  $\mathcal{D}_u^k$  to describe all the uninitialized drones and the set  $\mathcal{D}_e^k$  to describe all the initialized (estimated) drones. The details of the observability and initialization of Omni-swarm will be stated in Sect. V-D and Sect. V-E.

### E. System Architecture

As shown in Fig. 2, our proposed method is divided into a front-end and a back-end, and Omni-swarm independently runs on each drone in the swarm. In the front-end, the raw measurements are pre-processed by an ego-motion estimator (VIO), visual drone tracking module (VDT), and multi-drone map-based localization module (MDML). In the back-end, graph-based optimization is utilized for state estimation. After this, we propagate the real-time state of the aerial swarm with the latest VIO and the estimated states.

### F. Graph-based Optimization

We adopt a 4-DoF graph-based optimization for the swarm state estimation in the back-end of Omni-swarm. Suppose  $\mathcal{SF}^t$  denotes the swarm keyframe, which contains the keyframes  $\{\mathcal{F}_1^t \dots \mathcal{F}_n^t\}$  of  $n$  drones at time  $t$ .

The graph  $\mathcal{G}$  contains  $m$  swarm keyframes, the keyframes  $\{\mathcal{F}_1^t \dots \mathcal{F}_m^t\}$  of these swarm keyframes serve as the vertices of  $\mathcal{G}$  and the measurements serve as edges to connect these keyframes:

$$\mathcal{G} = \{\mathcal{SF}^1, \mathcal{SF}^2, \mathcal{SF}^3 \dots \mathcal{SF}^m, \\ \dots \mathbf{z}_{\mathcal{L}_{i \rightarrow j}}^{t_0 \rightarrow t_1} \dots \mathbf{z}_{\mathcal{D}_{k \rightarrow l}}^{t_2} \dots \mathbf{z}_{\mathcal{d}_{p,q}}^{t_3}\}.$$

An illustration of  $\mathcal{G}$  can be found in Fig. 3a.

In our graph-based optimization, we use maximum a posteriori (MAP) inference for the factor graph based on non-linear least-squares optimization [41] for swarm state estimation. A factor graph  $\mathcal{G}_f$  is constructed from graph  $\mathcal{G}$  in our back-end after the outlier rejection module filters out the outlier measurements in the graph  $\mathcal{G}$ . As shown in Fig. 3b, the poses of the vertices in  $\mathcal{G}$  serve as the variables in  $\mathcal{G}_f$ , and the

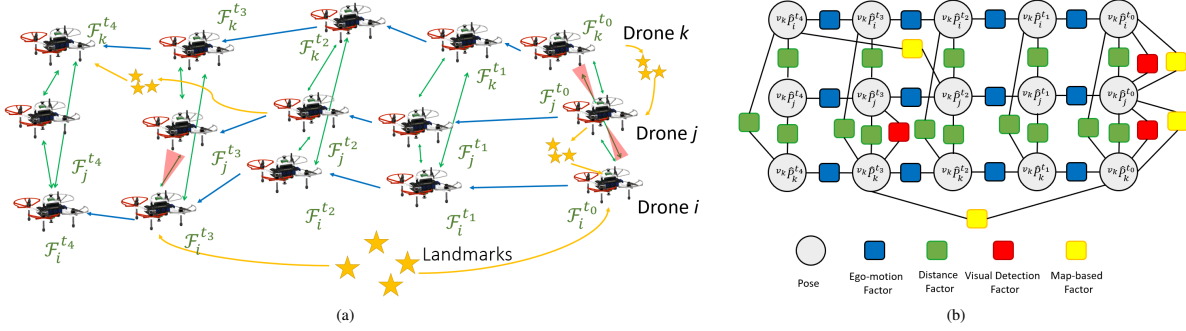


Fig. 3: A demonstration of measurements and corresponding factor graph in the graph-based optimization. a) Measurements involved in the graph-based optimization, including the UWB measurements (green), the VIO measurements (blue arrow) and the map-based factor (yellow). The yellow stars represent the landmarks of the sparse map, and the map-based factors are detected from these landmarks. b) Factor graph for swarm state estimation. The poses are connected by ego-motion factors, distance factors, visual detection factors and map-based factors.

variables are connected by four types of factors, which are built from the measurements of  $\mathcal{G}$ :

- Ego-motion factor: Each pose is connected to the previous pose of the same drone by an ego-motion factor, which represents the 4-DoF relative pose  $\mathbf{z}_{\delta P_j}^t$  from the previous keyframe. This factor smoothes the state estimation and provides local accuracy.
- Map-based factor: The poses are connected by the corresponding map-based factors, which represent the 4-DoF relative poses  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}}^t$  from the keyframe  $\mathcal{F}_i^{t_0}$  to  $\mathcal{F}_j^{t_1}$  to ensure global consistency and provide relative state estimation.
- Distance factor: The poses at one timestamp are connected to each other by distance factors, which are the distance measurement  $\mathbf{z}_{d_{i,j}}^t$  measured by UWB for relative state estimation.
- Visual detection factor: The poses are connected by a visual detection factor  $\mathbf{z}_{D_{k \rightarrow j}^{t_0}}^t$  if one drone successfully detects another drone in the visual drone tracking module. This factor stands for accurate relative state estimation.

In the following discussions, we will first introduce these measurements and factors and then give the optimization problem for state estimation from the factor graph  $\mathcal{G}_f$ .

#### IV. FRONT-END: OMNIDIRECTIONAL PERCEPTION SYSTEM AND MEASUREMENT MODELING

##### A. Omnidirectional Visual Inertial Odometry

To achieve omnidirectional visual perception, Omni-swarm adopts two 235-degree FoV fisheye cameras to cover all the surrounding directions, as shown in Fig. 4c. Based on previous works on omnidirectional VIO [3], [42], [43], we develop VINS-Fisheye<sup>2</sup>, which is an omnidirectional visual-inertial navigation system derived from VINS-Fusion<sup>3</sup>. VINS-Fisheye uses an IMU and previously mentioned stereo fisheye cameras to estimate ego-motion.

Because of the massive distortion of the fisheye cameras, it is hard to directly apply the existing visual algorithms to the raw image data produced by them. As an alternative, we reproject the raw image captured by a fisheye camera into five distortion-free images for later algorithms, which follows the

pipeline proposed in [42] and [43]. An example of a raw image and processed distortion-free images is shown in Fig. 4d. After reprojecting the raw fisheye images, VIO will be extracted based on these distortion-free images to provide real-time local pose and velocity estimation. Due to the long-term drifting of the VIO, instead of directly fusing the original odometry, we fuse the 4-DoF relative pose extracted from the VIO in the back-end, which can be modeled as

$$\mathbf{z}_{\delta P_i}^t = (\tilde{\mathbf{P}}_i^{t-1})^{-1}(\tilde{\mathbf{P}}_i^t) = (\mathbf{P}_i^{t-1})^{-1}(\mathbf{P}_i^t) + \mathbf{n}_{vio}, \quad (2)$$

where the noise of the relative pose is assumed as Gaussian,  $\mathbf{n}_{vio} \sim \mathcal{N}(0, \sigma_{vio}^2)$ .

The VIO keyframes contain the distortion-free images generated from the fisheye stereo cameras and the external parameters of the camera, and the real-time pose estimation is utilized for further processing to avoid redundant computation.

##### B. Visual Drone Tracking Module

In Omni-swarm, we introduce a visual drone tracking (VDT) module for visual tracking and relative pose estimation of other drones in the swarm. VDT uses visual object tracking techniques to keep tracking all the drone targets it detects. A drone target tracked by VDT is recorded as a visual target file  $V_{TF_{i^*}}$  in the module, where  $i^* \in \{1^*, 2^*, 3^* \dots\}$  is its ID inside VDT, and  $V_{TF_{i^*}}.d$  is the drone that tracks this visual target file. A visual target file  $V_{TF_{i^*}}$  can be anonymous, i.e., its drone ID is unknown, or identified, i.e., its drone ID has been successfully associated.

Alg. 1 demonstrates the algorithm of VDT, which will be described in detail in this section. VDT starts by detecting the drones in the image using a visual object detector. A drone detected by this detector creates a visual target candidate  $V_{TC_{i^*}}$ , which may create a visual target file or update a visual target file in subsequent algorithms. VDT will attempt to associate the visual target candidates with the drones observed by the state estimation and visual target files to identify its drone ID and VDT ID. However, this process may not be successful, either because the system is not properly initialized or because the  $V_{TC_{i^*}}$  is a false target. Visual target candidates that are not correctly matched to a drone will create anonymous visual target files in the subsequent algorithms. Finally, we use a 6-DoF pose estimator for relative pose estimation of these visual target files.

<sup>2</sup><https://github.com/HKUST-Aerial-Robotics/VINS-Fisheye>

<sup>3</sup><https://github.com/HKUST-Aerial-Robotics/VINS-Fusion>

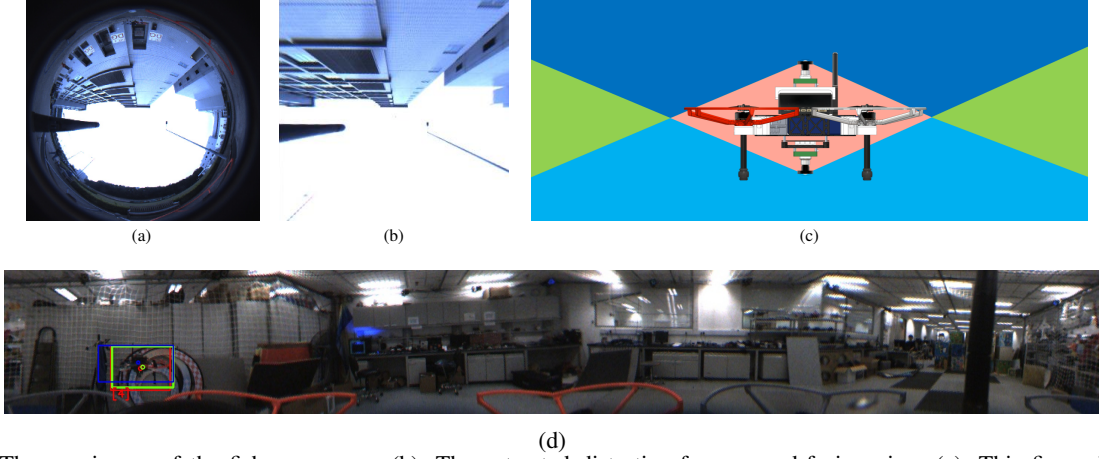


Fig. 4: (a): The raw image of the fisheye camera. (b): The extracted distortion-free upward-facing view. (c): This figure illuminates the available FoV for Omni-swarm. The blue and green areas denote the fisheye cameras' coverage, which is omnidirectional. The pink area denotes the dead zone. The green area denotes the stereo coverage, which is utilized for map-based localization. (d): The concatenation of distortion-free images. The red bounding box and circle are extracted by the visual object detection. The green bounding box and circle are tracked by the visual object tracking. The blue bounding box and circle are reprojected estimated states of the drone. The circles are the center of the bounding box.

---

#### Algorithm 1: Visual Target Tracking Module

---

**Data:** Set of estimated drones  $\mathcal{D}_e$ , set of visual target files  $\mathcal{V}_{TF}$

**Input:**  $\mathcal{I}$ , the set of distortion-free images.

**Output:**  $\mathcal{Z}_D$ , the set of visual target measurements.

```

1 Function VisualTargetTracking ( $\mathcal{I}$ )
2    $\mathcal{V}_{TC} \leftarrow \text{YOLO}(\mathcal{I})$ 
3    $\mathcal{V}_{TF} \leftarrow \text{MOSSETrackers}(\mathcal{I}, \mathcal{V}_{TF})$ 
4   for  $V_{TCi} \in \mathcal{V}_{TC}$  do
5     for  $V_{TFj} \in \mathcal{V}_{TF}$  do
6        $C_{i,j} \leftarrow c_{i,j}(V_{TCi}, V_{TFj})$ 
7       for  $j \in \mathcal{D}_e$  do
8          $C_{i,j+\|\mathcal{V}_{TF}\|} \leftarrow c_{i,j}(V_{TCi}, D_i)$ 
9    $\mathcal{M} \leftarrow \text{Hungarian}(C)$ 
10   $\mathcal{V}_{TF}^u \leftarrow \emptyset$ 
11  for  $V_{TCi} \in \mathcal{V}_{TC}$  do
12    if  $V_{TCi} \in \mathcal{M}$  then
13       $j \leftarrow \mathcal{M}(i)$ 
14      if  $j < \|\mathcal{V}_{TF}\|$  then
15         $V_{TFj}.T \leftarrow \text{NewMOSSE}(\mathcal{I}, V_{TCi}.B)$ 
16      else
17         $\mathcal{V}_{TF}^u \leftarrow \mathcal{V}_{TF}^u \cup \text{NewVTF}(V_{TCi}, D_i)$ 
18    else
19       $\mathcal{V}_{TF}^u \leftarrow \mathcal{V}_{TF}^u \cup \text{NewAnonVTF}(V_{TCi})$ 
20   $\mathcal{Z}_D \leftarrow \emptyset$ 
21  for  $V_{TFi} \in \mathcal{V}_{TF}^u$  do
22     $\mathcal{Z}_D \leftarrow \mathcal{Z}_D \cup \text{PoseEstimation}(V_{TFi})$ 
23   $\mathcal{V}_{TF} \leftarrow \mathcal{V}_{TF} \cup \mathcal{V}_{TF}^u$ 
24  return  $\mathcal{Z}_D$ 

```

---

1) *Visual Drone Detection:* We adopt YOLOv4-tiny [44], [45], one of the state-of-the-art visual object detection approaches based on a convolutional neural network (CNN), for detecting the 2D bounding boxes of the drones on the distortion-free images extracted from raw fisheye images. In practice, limited to the computational resources, only the upward-facing camera is used. The network is trained with

our custom data to efficiently detect our custom drones. A demonstration of the detected bounding boxes is shown in Fig. 4d.

As the YOLO function in Alg. 1 shows, the result of the visual object tracking is a set of target candidates  $\mathcal{V}_{TC} = \{V_{TC1^*}, V_{TC2^*} \dots\}$ . A target candidate  $V_{TCi^*}$  detected by the visual object detector can be represented by a bounding box  $\mathbf{B}_{i^*} = [\mathbf{B}_{i^*}^L, \mathbf{B}_{i^*}^R]$ , where  $\mathbf{B}_{i^*}^L$  is the coordinates of the left-upper corner of bounding box and  $\mathbf{B}_{i^*}^R$  is the coordinates of the right-bottom corner.

2) *Visual Object Tracking of Drones:* In VDT, we use MOSSE [46] as the object tracking technique to keep track of visual target files. MOSSE is a visual tracking approach that combines robustness and efficiency. The images  $\mathcal{I}$  extracted from the original image by VINS-Fisheye are also utilized for visual object tracking, as the function **MOSSETrackers** in Alg. 1.

If a new visual target file  $V_{TCi^*}$  is successfully associated with an existing visual target file  $V_{TFj^*}$ , then the visual tracker of  $V_{TFj^*}$  will be replaced by a new visual object tracker initialized by  $\mathbf{B}_{i^*}$  (line 15 of Alg. 1). Otherwise, we use  $\mathbf{B}_{i^*}$  to create a new visual target file in the function **NewVTF** and function **NewAnonVTF**.

In practice, since visual object detection consumes far more computational resources than visual object tracking, we run the former at a lower frequency (1 Hz in practice) than the latter (10 Hz in practice).

3) *Data Association:* For a homogeneous swarm, one issue is the association between the detected targets and the drones observed by state estimation, but this is essential for subsequent state estimation. The problem of data association for visual drone detection in Omni-swarm is divided into two sub-problems. The first is to initialize the state estimation, and the second is to perform global-nearest-neighbor (GNN) matching on the target based on the existing estimation. The state estimation of Omni-swarm can be initialized by a variety of methods, which will be described in detail later in Sect. V-E, including an approach using only anonymous visual target

files. On the other hand, GNN has been widely adopted in multi-target tracking algorithms [47]–[49] for data association.

Here, we propose a unified framework to associate the visual target candidates to visual target files and to the estimated drones. The latter two are collectively referred to as objects. This method does not require the system to be fully initialized, considering that it also has the responsibility of tracking of anonymous targets. We project all the drones estimated by Omni-swarm onto the image plane, as shown by the blue bounding box in Fig. 4d. These projected bounding boxes are treated the same as the bounding box of the visual target files (the green bounding box in Fig. 4d), as the objects in GNN. The corresponding cost between candidates (the red bounding box in Fig. 4d) and objects is defined by

$$c_{i^*,j} = \begin{cases} 1 - o_{i^*,j}, & o_{i^*,j} > 0 \\ +\infty, & o_{i^*,j} = 0, \end{cases} \quad (3)$$

where  $o_{i^*,j}$  is defined as the overlap of the two bounding boxes, which is defined as

$$o_{i^*,j} = A_{i^*,j}^i / \max(A_j, A_{i^*}),$$

where  $A_{i^*,j}^i$  is the intersection area of bounding box  $\mathbf{B}_{i^*}$  and  $\mathbf{B}_j$ , and  $A_j$  and  $A_{i^*}$  are the areas of two the bounding boxes.

With the corresponding cost defined here, we create a corresponding cost matrix  $\mathbf{C}$  of the candidates and objects. In addition, we add dummy objects and dummy candidates with a  $+\infty$  distance to all normal candidates and objects to represent missing detections and false alarms. Finally, Omni-swarm uses the Hungarian algorithm [50], [51] to solve the GNN problem defined by Betke et al. [47] with the corresponding cost matrix  $\mathbf{C}$ . Alg. 1, line 9, shows this procedure. The output of the Hungarian algorithm is a corresponding dictionary in which the candidates assigned to dummy objects are not present.

After the Hungarian algorithm is performed, for each candidate  $V_{TFi^*}$ , if  $V_{TFi^*}$  is assigned to a visual target file, the old visual target files' visual object tracker will be updated in Alg. 1, line 15. If a  $V_{TFi^*}$  has been assigned to a drone observed by state estimation, it will create a named visual target; otherwise, it creates a new anonymous visual target file. This is shown in Alg. 1, line 17 and line 19.

4) *6-DoF Pose Estimator*: The VDT adopts a CNN-based 6-DoF pose estimator [52], [53] to extract the accurate 6-DoF relative pose estimation of the visual target files as the function **PoseEstimation**. Specifically, this approach performs relative pose estimation by using CNNs to extract semantic feature points of objects in images and build optimization problems. The network is trained with the distortion-free images collected from the onboard fisheye cameras for efficiently detecting our custom drones. Before solving the full-perspective pose estimation problem, an additional outlier rejection on features using EPnP with RANSAC [54] is performed.

The 6-DoF relative poses extracted by VDT from drone  $k$  to drone  $j$  will be adopted as visual detection measurements, which can be modeled as

$$\mathbf{z}_{Dk \rightarrow j}^{t_0} = ({}^{v_i} \mathbf{P}_k^{t_0})^{-1} ({}^{v_i} \mathbf{P}_j^{t_0}) + \mathbf{n}_D, \quad (4)$$

where this relationship stands with any reference frame  $v_i$  and  $\mathbf{n}_D$  is Gaussian noise. We call this result the visual detection

---

### Algorithm 2: Multi-drone Map-based Localization Algorithm for Drone $k$

---

**Data:** Local Visual Database  $\mathcal{D}_l$ , Remote Visual Database  $\mathcal{D}_r$ , Local Drone  $k$

**Input:**  $\mathcal{F}_j^{t_1}$

```

1 Function KF_QUERY ( $\mathcal{F}_j^{t_1}, \mathcal{D}$ )
2    $\mathcal{C} \leftarrow \emptyset$ 
3   for  ${}^c \mathcal{K}_j^{t_1} \in \mathcal{F}_j^{t_1}$  do
4     if  $j = k$  then
5        $\mathcal{C} \leftarrow \mathcal{C} \cup \text{KNN\_SEARCH}({}^c \mathcal{K}_j^{t_1}, \mathcal{D}_r, \tau_{fl})$ 
6        $\mathcal{C} \leftarrow \mathcal{C} \cup \text{KNN\_SEARCH}({}^c \mathcal{K}_j^{t_1}, \mathcal{D}_l, \tau_{fl})$ 
7     ADD ( $\mathcal{K}_j^{t_1}, j = k, \mathcal{D}_l, \mathcal{D}_r$ )
8      $\mathcal{K} \leftarrow \text{argmin}_{({}^c \mathcal{K})_f} \|({}^c \mathcal{K})_f - ({}^c \mathcal{K}_j^{t_1})_f\|$ 
9     return  $\mathcal{K}.\mathcal{F}$ 
10 Function G_CHECK ( $V_j \hat{\mathbf{R}}_i^{t_0}, \tilde{\mathbf{R}}_i^{t_0}, \tilde{\mathbf{R}}_j^{t_1}$ )
11    $\delta \hat{\mathbf{R}}_{i \rightarrow j}^{t_0 \rightarrow t_1} \leftarrow (V_j \hat{\mathbf{R}}_i^{t_0})^{-1} \tilde{\mathbf{R}}_i^{t_0}$ 
12    $V_i \hat{\mathbf{R}}_j^{t_1} \leftarrow \tilde{\mathbf{R}}_i^{t_0} \delta \hat{\mathbf{R}}_{i \rightarrow j}^{t_0 \rightarrow t_1}$ 
13    $\delta \psi \leftarrow (\tilde{\mathbf{R}}_j^{t_1})_\psi - (V_i \hat{\mathbf{R}}_j^{t_1})_\psi$ 
14    $\Delta \mathbf{R} \leftarrow (\mathbf{R}_z(\delta \psi) V_i \hat{\mathbf{R}}_j^{t_1})^T \tilde{\mathbf{R}}_j^{t_1}$ 
15   return  $\|\Delta \mathbf{R}\| > \tau_{rot}$ 
16 Function LOOP_DETECTION ( $\mathcal{F}_j^{t_1}$ )
17    $\mathcal{F}_i^{t_0} \leftarrow \text{KF\_QUERY}(\mathcal{F}_j^{t_1}, \mathcal{D})$ 
18   if  $\mathcal{F}_i^{t_0} \neq \emptyset$  then
19      $\mathcal{P}_{2d_i}^{t_0}, \mathcal{P}_{3d_j}^{t_1} \leftarrow \text{BF\_MATCHER}(\mathcal{F}_i^{t_0}, \mathcal{F}_j^{t_1})$ 
20      $\text{inliers}, V_j \hat{\mathbf{T}}_i^{t_0} \leftarrow \text{PNP\_RANSAC}(\mathcal{P}_{2d_i}^{t_0}, \mathcal{P}_{3d_j}^{t_1})$ 
21     if  $\text{inliers} \geq \tau_{in}$  then
22       if G_CHECK ( $(V_j \hat{\mathbf{T}}_i^{t_0})_R, (\mathcal{F}_i^{t_0})_R, (\mathcal{F}_j^{t_1})_R$ )
23         then
24            $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}} \leftarrow ((V_j \hat{\mathbf{T}}_i^{t_0})^{-1} \tilde{\mathbf{T}}_j^{t_1})_P$ 
25           return  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}}$ 
25   return  $\emptyset$ 

```

---

measurement. We only perform pose estimation on newly created or updated VTFs to save computational resources.

### C. Multi-drone Map-based Localization Module

Similar to the approach used in VINS-Fisheye, we also perform visual object detection on distortion-free images. The multi-drone map-based localization (MDML) module performs relative localization and eliminates the drifting of the VIO by identifying the locations visited by all the drones of the aerial swarm. Sparse maps, which contain landmarks and keyframes, are simultaneously generated by local and remote measurements on every drone. Beyond optimizing the sparse maps, we utilize them to extract relative poses among drones in the swarm by a loop closure detection procedure. The MDML module is decentralized, running on each drone separately.

1) *Multi-drone map-based Localization Procedure*: When the MDML module receives the VIO keyframes, it uses MobileNetVLAD [55], [56] to extract the global features and uses SuperPoint [57] to extract the landmarks and the corresponding descriptor. Correspondences between landmarks from the upward-facing and downward-facing cameras are established by performing feature matching, and the matched landmarks are triangulated for estimating their 3D positions in



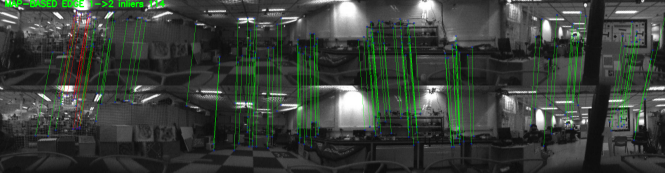


Fig. 5: A demonstration of detected map-based measurements between different drones. Green lines denote the inlier landmarks correspondences, and red lines denote the outliers.

the local frame. The global descriptor and landmarks, together with the odometry and extrinsic, are packed into keyframe  $\mathcal{F}_i^t$ , which will be broadcast to the entire swarm later. It is worth noting that we broadcast only the feature points with the 3D position successfully estimated (some of the features will fail to establish correspondence) to reduce the amount of communication.

To store and retrieve these keyframes, we build databases based on Faiss [58], which is a vector similarity retrieval database. The keyframes  $\mathcal{F}_j^{t_1}$ , which are indexed by the corresponding global descriptors, are saved in the databases as maps. After the module receives a keyframe remotely or locally, loop closure detection is adopted to extract the relative pose. There are two separate visual databases on each drone; a remote database  $\mathcal{D}_r$  that stores keyframes from remote drones, and a local database  $\mathcal{D}_l$  that stores keyframes from the local drone. Extracting the map-based measurement for a pair of keyframes from  $\mathcal{D}_r$  is avoided to save computational resources since all the extracted map-based measurements are broadcast to the whole swarm.

The procedure of loop closure detection is shown in Alg. 2. Suppose a drone  $k$  receives a keyframe  $\mathcal{F}_j^{t_1}$  from drone  $j$ , where if  $j = k$ , the keyframe is generated by the drone itself. The most similar keyframe  $\mathcal{F}_i^{t_0}$  to  $\mathcal{F}_j^{t_1}$  in the databases is retrieved by function **KF\_QUERY**. We use function **KNN\_SEARCH** to retrieve the  $K$  nearest-neighbor of the descriptors from the Faiss database, where  $K$  is set to 5 in practice. When the search is completed, the new keyframes are added to the local or remote database by function **ADD**. Finally, function **KF\_QUERY** returns the corresponding keyframe  $\mathcal{K}.\mathcal{F}$  of the nearest-neighbor distortion-frame image keyframe  $\mathcal{K}$ .

2) *Relative pose extraction*: Once the keyframe  $\mathcal{F}_i^{t_0}$  is returned from the visual database, we establish the 2D-3D matches from the landmarks of  $\mathcal{F}_i^{t_0}$  to landmarks of the  $\mathcal{F}_j^{t_1}$  by using the landmark descriptors with a brute-force matcher **BF\_MATCHER**. The brute-force matcher finds the correspondences of features with a minimum L2 distance, and uses a cross check to reduce outlier correspondence<sup>4</sup>. Either keyframe retrieved from the database or brute-force matching may bring abnormal results. To remove the outliers, the map-based measurement is verified with two methods in the map-based localization module.

- Homography test [59] and Perspective-n-point (PnP) test with RANSAC [54] between the 2D features of the incoming keyframe and the 3D positions of the queried

keyframe. These tests are performed in the function **PNP\_RANSAC**.

- Geometric test performed by **G\_CHECK**. In this test, the consistency of gravity is checked with the 6-DoF pose extracted by PnP RANSAC.

If enough inliers are found in **PNP\_RANSAC** and the geometric test is passed, the map-based measurement  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}}^{t_0 \rightarrow t_1}$  is considered valid. In addition,  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}}^{t_0 \rightarrow t_1}$  can be modeled as:

$$\mathbf{z}_{\mathcal{L}_{i \rightarrow j}}^{t_0 \rightarrow t_1} = ({}^{v_k} \mathbf{P}_i^{t_0})^{-1} ({}^{v_k} \mathbf{P}_j^{t_1}) + \mathbf{n}_{\mathcal{L}}, \quad (5)$$

where this relationship stands with any reference frame  $v_k$  and  $\mathbf{n}_{\mathcal{L}}$  is Gaussian noise.

When  $i \neq j$ , the map-based measurement  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}}^{t_0 \rightarrow t_1}$  provides sufficient observability of the relative pose of drone  $i$  and drone  $j$ . This makes map-based measurements essential in the observability verification and initialization of state estimation. If  $t_0 \neq t_1$ , the map-based measurement represents the relative pose of the drones that visit the same place at different times, which eliminates the accumulated drifting error of the VIO.

#### D. UWB Measurement

The distance measurements from the UWB module can be modeled as

$$\mathbf{z}_{d_{i,j}}^t = \| {}^{v_k} \mathbf{X}_i^t - {}^{v_k} \mathbf{X}_j^t \|_2 + \mathbf{n}_d, \quad (6)$$

where  $\mathbf{n}_d \sim \mathcal{N}(0, \sigma_d^2)$  is the Gaussian noise of the distance measurement. The installation length of the antenna relative to the IMU is ignored in our model.

## V. BACK-END: GRAPH-BASED OPTIMIZATION FOR STATE ESTIMATION

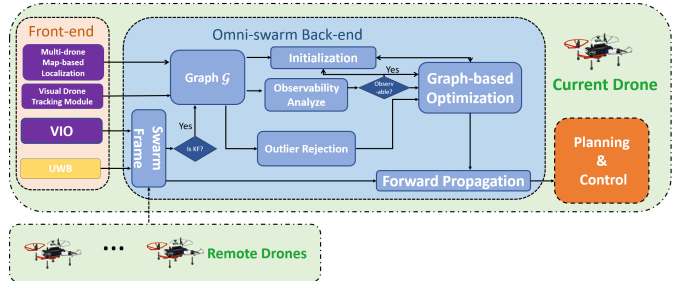


Fig. 6: The structure of the back-end of Omni-swarm. The swarm frames  $\mathcal{S}\mathcal{F}^t$  composed of VIO and UWB measurements are first judged to be swarm keyframes, and the non-keyframes are only utilized for forward propagation, while the swarm frames that qualify as swarm keyframes are added to the graph. The system will evaluate whether the measurements in the graph meet the observability requirements. When the requirements are met, the system solves the graph-based optimization for swarm state estimation. Finally, we use the optimization results to perform forward propagation to obtain the real-time state of the swarm.

#### A. Swarm Keyframe Limitation

Similar to our previous work [18], a swarm frame  $\mathcal{S}\mathcal{F}^t$  is judged to be a swarm keyframe is by the motion of the drones in this frame relative to the last swarm keyframe,

<sup>4</sup>A more detailed introduction to the brute-force matcher can be found at [https://docs.opencv.org/4.x/dc/dc3/tutorial\\_py\\_matcher.html](https://docs.opencv.org/4.x/dc/dc3/tutorial_py_matcher.html)

or a new drone is discovered in  $\mathcal{SF}^t$ . To limit the computational resources utilized by the back-end, we keep the numbers of swarm keyframes in  $\mathcal{G}$  within a preset maximum number  $m_{max}$ . In contrast to the sliding window of the swarm keyframes employed in our previous work [18], here we use a random deletion mechanism to achieve better global consistency; once the graph's swarm keyframe in the graph grows to more than  $m_{max}$ , we randomly delete one swarm keyframe. The intuition of this approach is that the earlier swarm keyframes in  $\mathcal{G}$  have less impact on the relative estimation. With this approach, the earlier swarm keyframes in  $\mathcal{G}$  are not immediately discarded but become more sparse, which allows us to use early swarm keyframes for optimization without affecting the computational speed due to too many keyframes. This leads to less drifting and better global consistency. In practice,  $m_{max}$  is set to 100 as a tradeoff between the performance and accuracy. Taken together, compared to the previous sliding window method, this method does not affect the accuracy of relative estimation, but obtains better global consistency, which is verified in Sect. VII-C4.

### B. Outlier Rejection

One challenge for achieving robust swarm state estimation is the variety of factors that can cause measurement outliers. Although we perform several outlier rejections in the front-end, some outliers will always be transferred to the back-end. Here we introduce the outlier rejection techniques adopted by the back-end of Omni-swarm.

1) *Unified outlier rejection for visual target measurements and map-based measurements*: Pairwise-consistency measurement set maximization (PCM), proposed by Mangelson et al. [60], is the state-of-the-art outlier rejection technique for loop closure (also called map-based measurements in the paper) and has been verified on real-world SLAM systems [17] [61]. In Omni-swarm, we introduce a unified outlier rejection module for the visual target measurements and the map-based measurements (collectively referred to as relative pose measurements) based on PCM. This method is applicable to visual detection measurement because the visual target measurements are also 6-DoF relative pose measurements. Our PCM module combines the inter-drone and intra-drone PCM outlier rejection methods [17], [60], [61].

The core of the PCM outlier rejection is consistency graphs [60]. A consistency graph is a graph  $\mathbf{G}_{i,j}^{PCM} = \{V_{i,j}^{PCM}, \mathcal{E}_{i,j}^{PCM}\}$ , where vertex  $v \in V_{i,j}^{PCM}$  represents a relative pose measurement between drone  $i$  and drone  $j$ , and each edge  $e \in \mathcal{E}_{i,j}^{PCM}$  represents the consistency of the measurements. Only the consistent pairs-of-vertices are connected with an edge. We maintain a set of consistency graphs  $\{\mathbf{G}_{i,j}^{PCM} | i \in \mathcal{D}, j \in \mathcal{D}\}$  of each pair of drones. When  $i \neq j$ ,  $\mathbf{G}_{i,j}^{PCM}$  denotes the inter-drone measurements, and vice versa for the intra-drone measurements.

Prior to graph-based optimization, we will first update the consistency graphs  $\{\mathbf{G}_{i,j}^{PCM} | i \in \mathcal{D}, j \in \mathcal{D}\}$  with the newly received relative pose measurements in  $\mathcal{G}$ . Similar to [61], the consistency graph is updated incrementally to save the computational power. For every updated  $\mathbf{G}_{i,j}^{PCM}$ , we adopt

the fast maximum clique method proposed by Pattabiraman et al. [62] to solve the PCM problem defined in [60] and thus determine the inner measurements. Each drone only solves the PCM problem with itself involved to save computational power and shares the inlier measurements with the other drones in the aerial swarm.

2) *Outlier Rejection for Distance Measurements*: In our experiments, we find that the UWB generates significant outliers when two drones have a large relative elevation angle. This is because of the occlusion of the drones' airframe, and we decide to flag out the measurement  $z_{d_{i,j}}^t$ , which satisfies

$$\left| \arcsin \left( z_{d_{i,j}}^t / (v_k z_i^t - v_k z_j^t) \right) \right| > \tau_{ele}, \quad (7)$$

where  $v_k z_i^t$  and  $v_k z_j^t$  is the estimated height of drone  $i$  and  $j$ , respectively, and  $\tau_{ele}$  is the elevation threshold, which is set to  $37^\circ$  in practice.  $\tau_{ele}$  is chosen regarding the angle at which the drone fuselage obscures the UWB. Finally, we utilize the residual defined in Eq. (11) without the Huber loss to test if the distance measurement is consistent with the current estimation, and  $r_d(\mathbf{z}_{d_{i,j}}^t, \mathcal{X}_k) > \tau_d$  will be flagged as an outlier, where  $\tau_d$  is the distance outlier threshold, which is set to 0.3 m in practice. The expected error of UWB measurements is no more than 10 cm, and we choose it three times as  $\tau_d$  to ensure that the correct values are not overly filtered out. When  $\tau_d$  is too large, the outlier will interfere with the estimate estimation, and when  $\tau_d$  is too small, the correct measurement will be rejected due to inaccurate initialization.

### C. Optimization Problem

After the outlier rejection, the factor graph  $\mathcal{G}_f$  will be constructed from the graph  $\mathcal{G}$  with poses as its variables and inlier measurements as the factors. Once we set up the factor graph, a non-linear least-squares optimization problem is built to solve the maximum a posteriori (MAP) inference of the factor graph [41] for swarm state estimation. For a drone  $k$ , the full state vector  $\mathcal{X}_k$  of the swarm state estimation problem is defined as:

$$\begin{aligned} & [v_k \hat{\mathbf{X}}_0^{t_0 T}, v_k \hat{\psi}_0^{t_0} \dots v_k \hat{\mathbf{X}}_0^{t_{m-1} T}, v_k \hat{\psi}_0^{t_{m-1}}, \\ & \dots v_k \hat{\mathbf{X}}_1^{t_0 T}, v_k \hat{\psi}_1^{t_0} \dots v_k \hat{\mathbf{X}}_{n-1}^{t_{m-1} T}, v_k \hat{\psi}_{n-1}^{t_{m-1}}]^T, \end{aligned} \quad (8)$$

where  $[v_k \hat{\mathbf{X}}_i^t, v_k \hat{\psi}_i^t]^T$  is the state vector of 4-DoF pose  $v_k \hat{\mathbf{P}}_i^t$ ,  $n$  is the number of drones in the swarm system, and  $m$  is the number of keyframes in the graph. The non-linear least-squares optimization problem for MAP inference is expressed as the following formulation:

$$\begin{aligned} \min_{\mathcal{X}_k} & \left\{ \sum_{(i,t) \in \mathcal{S}} \|\mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{\delta \mathbf{P}_i}^t, \mathcal{X}_k)\|_{\Sigma}^2 + \right. \\ & \sum_{(i,j,t) \in \mathcal{U}} \rho \left( \left\| \mathbf{r}_d(\mathbf{z}_{d_{i,j}}^t, \mathcal{X}_k) \right\|_{\Sigma}^2 \right) \\ & + \sum_{(i,j,t) \in \mathcal{VD}} \rho \left( \left\| \mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{D_{i \rightarrow j}}^t, \mathcal{X}_k) \right\|_{\Sigma}^2 \right) \\ & \left. + \sum_{\mathcal{L}_{k \rightarrow j}^{t_0 \rightarrow t_1} \in \mathcal{L}} \rho \left( \left\| \mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}}, \mathcal{X}_k) \right\|_{\Sigma}^2 \right) \right\}, \end{aligned} \quad (9)$$

where  $\mathcal{S}$  is the set of all odometry factors;  $\mathcal{VD}$  is the set of all visual detection factors,  $\mathcal{U}$  is the set of all distance factors,  $\mathcal{L}$  is the set of map-based factors,  $\mathbf{r}_d(\mathbf{z}_{d_{ij}}^t, \mathcal{X}_k)$  is the residual of the distance factor, where  $\mathbf{r}_{\mathcal{RP}}(\cdot, \mathcal{X}_k)$  represents the residual of the relative pose, which is applicable to the ego-motion factors, map-based factors and visual detection factors.  $\mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{\delta \mathbf{P}_i}^t, \mathcal{X}_k)$  represents the residual of the ego-motion factor, ensuring the local consistency of drone  $i$ 's state,  $\mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{D_{i \rightarrow j}}^t, \mathcal{X}_k)$  is the residual of the visual detection factor, where  $(i, j)$  presents the drone  $j$  detected by drone  $i$ ,  $\mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}}, \mathcal{X}_k)$  represents the residual of the map-based factor, which ensures the global consistency and observability of the relative state. Since some outlier measurements may be generated in distance measurements, map-based factor measurements, and visual detection measurements, we adopt the Huber norm  $\rho(s)$  [63] to reduce the effect of possible outlier factors.

According to the measurement models stated in (2),(4) and (5), the residual of the relative pose is defined as

$$\mathbf{r}_{\mathcal{RP}}(\mathbf{z}_{\mathcal{RP}_{i \rightarrow j}^{t_0 \rightarrow t_1}}, \mathcal{X}_k) = (\mathbf{z}_{\mathcal{RP}_{i \rightarrow j}^{t_0 \rightarrow t_1}})^{-1} \left( ({}^{v_k} \hat{\mathbf{P}}_i^{t_0})^{-1} {}^{v_k} \hat{\mathbf{P}}_j^{t_1} \right), \quad (10)$$

where  $\mathbf{z}_{\mathcal{RP}}$  represents the relative pose measurements including the ego-motion measurement, map-based measurement, and visual target measurement. Referring to the measurement model stated in Eq. (6), the residual of the distance is

$$\mathbf{r}_d(\mathbf{z}_{d_{ij}}^t, \mathcal{X}_k) = \mathbf{z}_{d_{ij}}^t - \| {}^{v_k} \hat{\mathbf{x}}_i^t - {}^{v_k} \hat{\mathbf{x}}_j^t \|. \quad (11)$$

Finally, the Ceres Solver, an open-source C++ library developed by Google for modeling and solving non-linear least-squares optimization problems [64], is adopted to solve the optimization problem with a trust region method using a sparse normal Cholesky decomposition is chosen as the optimization algorithm. Meanwhile, since the planner or controller may require high-frequency real-time poses, the poses  ${}^{v_k} \hat{\mathbf{P}}_i^t$  and  ${}^{v_k} \hat{\mathbf{T}}_i^t$  are forward propagated based on the latest IMU propagated VIO  $\tilde{\mathbf{P}}_i^t$  and  $\tilde{\mathbf{T}}_{i1}^t$  to predict the swarm state at a higher rate of 100 Hz following (1) and

$${}^{v_k} \hat{\mathbf{P}}_i^{t_1} = {}^{v_k} \hat{\mathbf{P}}_i^t (\tilde{\mathbf{P}}_i^t)^{-1} \tilde{\mathbf{P}}_i^{t_1}. \quad (12)$$

#### D. Observability Analysis

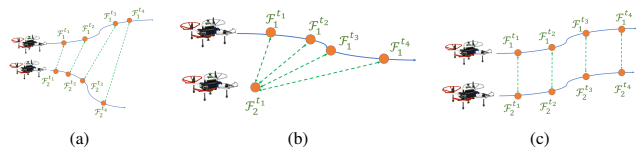


Fig. 7: Observability of UWB-odometry fusion: a) Two drones fly separately and the flight direction is not parallel. At this time, the two drones are observable to each other due to the presence of motion. b) Drone 1 is flying and drone 2 is in place. At this point, drone 1 can observe the position of drone 2, but the yaw of 2 is not observable. Drone 1 is not observable from drone 2. c) The two drones are flying parallel to each other and at this point, the two drones are unobservable to each other.

One of the key features of this system is the fusion of several sensors with completely different characteristics, and this multi-sensor fusion can effectively improve the observability

TABLE I: Observability with typical measurements combinations in different scenarios. T in the table means the measurement exists. F in the table means the measurement is absent. T/F in the table indicates that the existence of the measurement does not change the observability.

Typical Scenarios	Motion $k$	Motion $i$	UWB	Detection $k \leftrightarrow i$	Map-based $k \leftrightarrow i$	Observability $k \leftrightarrow i$
Feature-rich	T/F	T/F	T/F	T/F	T	6-DoF
	F	F	T/F	F	F	No
Feature-poor	T/F	T/F	T/F	T	T/F	6-DoF
	T	F	T	F	F	3-DoF
	T	T	T	T/F	T/F	6-DoF

in various environments. In feature-rich scenarios (usually narrow indoor environment), the rich features allow us to take full advantage of the map-based localization for state estimation. The inter-drone map-based measurements and visual detection measurements provide direct 6-DoF observability between drones, as the first and third rows in Table I show.

However, map-based localization is challenging to perform in feature-poor scenarios, e.g., open environments such as stadiums, lawns, and farmland, where there are few features, and visual detection measurements are limited by the visual detection range. In such scenarios, UWB-odometry fusion can be utilized for state estimation instead. The necessary condition for the observability of the relative state estimated by UWB-odometry fusion is the presence of relative motion, i.e., an aerial swarm stationary on the ground or flying in parallel does not satisfy this condition. As shown in Fig. 7, when drone 1 in the factor graph has sufficient motion and drone 2 is fixed, drone 1 can estimate the position of drone 2. However, the yaw angle of drone 2 cannot be estimated. Only when both drone 1 and 2 in the graph have motion and there is relative motion can drone 1 and 2 estimate each other's 6-DoF pose with UWB-odometry fusion. In practice, this brings a considerable level of annoyance. For example, aerial swarms in which case relative motion does not exist. In Omni-swarm, this issue is solved by fusing the relative pose given by visual drone tracking.

In summary, Table I shows the observability of Omni-swarm in various scenarios. We find that, by incorporating multiple sensors and measurements, Omni-swarm is capable of maintaining observability in various complex environments.

#### E. Initialization

Omni-swarm independently tracks the observability of each drone in the  $\mathcal{D}_u^k$  for initialization. Depending on the observability, as given in Table I, Omni-swarm has various ways of performing initialization: 1) UWB-odometry initialization with sufficient motion; 2) map-based measurements for initialization; and 3) anonymous visual detection measurements for initialization. A special case is drone  $k$  itself, whose state is initialized by its VIO.

1) *Map-based measurements for initialization*: One of the biggest advantages of map-based localization is that only one measurement is needed to provide sufficient observability, as shown in the first row of Table I. In this case, when a map-based measurement  $\mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}}$ , where  $i \in \mathcal{D}_e^k$  is an estimated drone and  $j \in \mathcal{D}_u^k$  is an uninitialized drone, is detected by

the map-based localization module, the system immediately initializes the drone, and we use this map-based measurement and the ego-motion of drone  $i$  to initialize its poses in  $\mathcal{G}$ :

$${}^{v_k}\hat{\mathbf{P}}_j^t \leftarrow {}^{v_k}\hat{\mathbf{P}}_i^{t_0} \mathbf{z}_{\mathcal{L}_{i \rightarrow j}^{t_0 \rightarrow t_1}} \left( (\tilde{\mathbf{P}}_j^{t_1})^{-1} \tilde{\mathbf{P}}_j^t \right),$$

where  ${}^{v_k}\hat{\mathbf{P}}_j^t$  is the state to be initialized in  $\mathcal{G}$ .

2) *Anonymous Visual Detection Measurements for Initialization*: Similar to map-based measurements, the visual detection measurements also contain 6-DoF information, which provides sufficient observability to initialize the state estimation, as shown in the third row of Table I. The difference is that the target ID of these visual detection measurements is anonymous when the system is not fully initialized. Therefore, we need to associate the anonymous IDs with the IDs of the uninitialized drones  $\mathcal{D}_u^k$  first, and then initialize the state estimates of the matched drone using the same method as in Sect. V-E1.

The data association method with anonymous visual detection measurements in Omni-swarm is a well-pruned depth-first search (DFS) [23], [65]. This method was first proposed by Zhou et al. [65] for multi-target tracking and was extended by Nguyen et al. [23] to the case where multiple homogeneous drones in an aerial swarm detect each other with anonymous relative measurements. We do not expand the details of this DFS algorithm here due to the space limitation.

3) *UWB-odometry for Initialization*: Unlike the previous two cases, initialization with the UWB-odometry requires a combination of measurements in multiple frames to provide sufficient observability, as shown in the fourth and fifth rows of Table I. The initialization condition for drone  $i$  is that there exists enough motion of drone  $i$  in the swarm keyframes where both ego-motion and UWB measurements of drone  $i$  and drone  $k$  exist, i.e., drone  $i$  has at least 3-DoF observability by drone  $k$ . In practice, we require the motion of drone  $i$  to be larger than a bounding box of [1.5, 1.5, 0.8] m. This bounding box is chosen so that the system has a sufficient baseline length.

In contrast to the previous cases, we cannot obtain an initial estimate of the states directly from the UWB-odometry measurement. The initialization with UWB-odometry fusion is accomplished by solving Eq. (9). To avoid the problem of the local optimal in this process, we will set three different random values for initializing the state of drone  $i$  and choose the one with the lowest cost as result of the optimization. Specifically, the poses of the other drones are initialized with the values of  $x$ ,  $y$  positions randomly distributed in  $[-5, 5]$  m and  $z$  position randomly distributed in  $[-1, 1]$  m. The initialized value of yaw is equal to the yaw estimated by VIO. This scale corresponds to the possible distribution of an aerial swarm of 10-20 drones in the takeoff state. This procedure can be accomplished within within 15 ms on a PC and 100 ms on the onboard computer.

A point worth noting is that this UWB-odometry fusion cannot avoid the unobservability in the case in Fig. 7c; i.e., if the flight paths of the two drones are parallel during the UWB-odometry initialization, the method may lead to incorrect results. Nevertheless, with other measurements in the Omni-swarm, this problem is very unlikely to occur. This is one advantage of Omni-swarm over other UWB-odometry fusion methods.

4) *Initialization Procedure*: When a new drone starts up and starts sending its measurements, Omni-swarm adds it to the set of available drones  $\mathcal{D}_a^k$  and set of uninitialized drones  $\mathcal{D}_u^k$ . The observability condition shown in Table I of each drone in  $\mathcal{D}_u^k$  is checked before each optimization. Once a drone  $i$  in  $\mathcal{D}_u^k$  is found to satisfy at least 3-DoF observability, the initialization approaches corresponding to the available measurements will be adopted to initialize its states.

## VI. SYSTEM IMPLEMENTATION

To fully demonstrate the Omni-swarm method proposed in this paper in real-world experiments, we present an aerial swarm system containing drones, algorithms, software management, a communication network, and 3D user interface. We also develop an aerial swarm trajectory planning method for verifying the practical value of Omni-swarm in inter-drone collision avoidance experiments.

### A. Aerial Platform

Our aerial swarm contains a few homogeneous custom drones, as shown in Fig. 8. The drone is equipped with a DJI N3 flight controller, two Pointgrey cameras with fisheye lenses, a NoopLoop UWB module, a DJI Manifold 2-G onboard computer with an Nvidia TX2 Module, and a WiFi module for communication. The front-end of VINS-Fisheye is accelerated with CUDA, and the various CNNs used for state estimation are accelerated with TensorRT. On the TX2, the TensorRT accelerated CNNs are two times faster than using TensorFlow and PyTorch and it takes up less memory. The computational performance will be detailed in Sect. VII-F

To facilitate the distribution of all the algorithms, including state estimation and planning, we divide the software running on the onboard computer into two layers: the boot layer running on the operating system and the algorithm layer running in the docker container. Omni-swarm together with the trajectory planning method, is deployed by a docker image and individually runs on each drone. The docker images are distributed through our wireless ad hoc network and each drone automatically pulls updates to improve the development efficiency. An additional computer serves as the maintenance server of the aerial swarm to facilitate the distribution of the docker image, which is not essential for flight. After we update the code, we may quickly get every drone's docker image up to date by simply booting up the drone within the operating range of the maintenance server's wireless ad hoc network.

### B. Frequency and Synchronization

Different measurements have different frequencies, which generally depend on the nature of the sensor and its computing power. In practice, we set the acquisition frequency of the camera to 20 Hz (for collecting datasets) or 16 Hz (for inter-drone collision avoidance experiments), IMU to 400 Hz and UWB to 100 Hz for distance measurements. VIO uses the IMU to predict the current ego-motion, so its output is also 400 Hz.

Timestamp synchronization is another vexing problem in robot swarms. The UWB module completes mutual timestamp



Fig. 8: One of the aerial platforms in the swarm system, which is equipped with stereo fisheye cameras, a DJI N3 flight controller, a Nooploop UWB module, and a DJI Manifold2-G onboard computer with an Nvidia TX-2 chip.

synchronization during ranging and communication, and sends the synchronized timestamps to the onboard computer via the serial port. Therefore, we choose the UWB timestamp as the time reference to obtain a swarm-wide synchronized timestamp.

To cope with the different frequencies and the time differences between different measurements, we take the timestamps of UWB measurements and convert swarm keyframes and all the other measurements to these timestamps. For VIO, this conversion process is achieved by finding the nearest UWB timestamp. Since VIO has 400 Hz output, the error caused by the difference between the VIO timestamp and the nearest UWB timestamp is small. For visual target measurement and map-based measurement, we use ego-motion to convert their relative pose by

$$\mathbf{z}_{\mathcal{R}\mathcal{P}_{i \rightarrow j}}^{t_a \rightarrow t_b} \leftarrow \left( (\tilde{\mathbf{P}}_i^{t_a})^{-1} \tilde{\mathbf{P}}_i^{t_0} \right) \mathbf{z}_{\mathcal{R}\mathcal{P}_{i \rightarrow j}}^{t_0 \rightarrow t_1} \left( (\tilde{\mathbf{P}}_j^{t_1})^{-1} \tilde{\mathbf{P}}_j^{t_b} \right), \quad (13)$$

where  $\mathbf{z}_{\mathcal{R}\mathcal{P}_{i \rightarrow j}}^{t_0 \rightarrow t_1}$  is the relative pose measurement,  $t_a$  and  $t_b$  are the timestamps of swarm keyframes, which are nearest to  $t_0$  and  $t_1$ . Eq. (13) converts the  $\mathbf{z}_{\mathcal{R}\mathcal{P}_{i \rightarrow j}}^{t_0 \rightarrow t_1}$  to a relative pose measurement between timestamp  $t_a$  and  $t_b$ . The accuracy loss caused by this conversion process is also negligible due to the good local accuracy of VIO.

Finally, Omni-swarm's graph-based optimization is solved at a speed of 1 Hz and is forward propagated through VIO at 100 Hz. Table II shows the computational frequency of the main components in Omni-swarm. Compared to on a PC, we reduce the computational frequency of the Omni-swarm on the onboard computer so that Omni-swarm has sufficient real-time performance in real-world applications. This is shown in the bottom row of Table II. In addition, on the onboard computer, the position control frequency is 50 Hz.

### C. Redundant Computations

The computations involved in the Omni-swarm can be classified into two types according to whether they are computed redundantly or not: 1) Distributed computations: VIO, visual target tracking, and map-based localization, which are not duplicated computed on different drones. 2) Redundant computations: Graph-based optimization of (Eq. (9)) is redundant when computed on different drones. We choose to

TABLE II: The computational frequency of each component on the PC and onboard computer. In the table, VINSF and VINSB is the frontend and back-end of VINS-Fisheye, respectively. Det. is the visual detection (YOLO) in VDT, and Trk. is the visual tracking (MOSSE). Desc. is the descriptors extraction (NetVLAD and SuperPoint) in MDML. LoopDet. is the loop closure detection in MDML, including database retrieval and brute-force matching. Opti. is the graph-based optimization.  $n$  in the table is the drone number of the swarm. The units in this table are Hz.

Plat	VINSF	VINSB	Det.	Trk.	Desc.	Loop. Det.	Opti.
PC	20	10	2	20	0.3	$0.3 * n$	10
Onboard	16	8	1	10	0.3	$0.3 * n$	1

TABLE III: The broadcasting bandwidth requirements of Omni-swarm for each drone. Odom & UWB is the odometry and UWB distances measurements. Keyframe represents the keyframe information for map-based localization. Inliers are the broadcast PCM result for each other drone in the swarm.

Item	Freq	Size	Bandwidth
Odom & Dis	100 Hz	49 Bytes	4.90
Keyframe	0.3 Hz	180 kB	54.0 kB/s
Inliers	1 Hz	552 Bytes	0.47 kB/s
Sum		59.4 kB/s	

solve this optimization on each drone individually to allow the state estimation to be robust to temporary network failure or single-point failure. It also avoids the communication overhead induced by distributed optimization methods [16], [21].

### D. Network Setup

The latency, bandwidth, interference immunity and working range of the communication directly determine if the proposed method will work well in practice. In the real-world experiments, the drones need to exchange image descriptors and flight paths, requiring a much larger communication bandwidth. In this case, we use a wireless ad hoc network based on the independent basic service set (IBSS) configuration of IEEE 802.11 standard [66] to meet the communication requirements, which cover the task execution system and maintenance requirements. One advantage of the wireless ad hoc network is its decentralized nature. No central server or router is required for deploying the wireless ad hoc network. With a 5.2 Ghz wireless ad hoc network, our communication network can cover a sufficient area in our experiments and provide high bandwidth and low latency communication services. In our laboratory environment, this network reached a 4.9 MB/s bandwidth at a distance of 22.4 m.

We use lightweight communications and marshalling (LCM) [67] to achieve efficient data encoding and communication broadcast in the ad hoc wireless network. The communication broadcast is limited to one-hop neighbors, covering all the swarm members in our current experiment. Some of the transmissions with lower bandwidth requirements are simultaneously backed up by UWB. While this network architecture is not part of the proposed state estimation, it accelerates our development and debugging process.

### E. Communication Bandwidth Requirements

Table III summarizes the major bandwidth of data broadcasting to the swarm by Omni-swarm for each drone, where

TABLE IV: Packet loss rate when there are different numbers of drones in the field.

Number	2	3	4
Avg. Loss.	3.5%	4.7%	3.9%

keyframe takes the biggest proportion. The major bandwidth is consumed on broadcasting landmark descriptors generated by the SuperPoint. In Omni-swarm, we compress the landmark descriptors (averaging 696 per keyframe) from 256 dimensions to 64 dimensions using PCA dimensionality reduction, which has no significant impact on the feature matching but reduces the communication amount by 534.5 kB (three times to current) per keyframe. We send each landmark descriptor separately in broadcasting the keyframe for map-based localization. Losing some of landmark descriptors does not affect the loop closure detection.

Table IV shows the packet loss rate with different numbers of drones in the field in the actual test. This result indicates that the packet loss rate does not change significantly with the current number of drones, which shows that Omni-swarm is scalable in the network. According to the available bandwidth in the lab environment (4.9 MB/s), the ideal limitation of the size of the aerial swarm given by the ad hoc network is 83.

#### F. Trajectory Planning for Swarm

To demonstrate the capability of our state estimation method, we conduct fully autonomous formation flights and inter-drone collision avoidance tests. In these flights, trajectories of the drones are generated by the method extended from Fast-Planner [1], a real-time trajectory planner based on kinodynamic search and gradient-based optimization.

The original Fast-Planner [1] assumes a static environment and does not consider the interaction with other agents. Therefore, we adapt it to achieve decentralized collision avoidance. Among all the drones, trajectories are shared through the wireless ad hoc network. Each drone checks whether any conflict exists in its own trajectory and that of others. Whenever a conflict is found, it replans a new trajectory immediately. The replanning starts by searching a safe and dynamically feasible initial path, in which motion primitives conflicting with other drones' trajectories are pruned to avoid an inter-drone collision. Collision avoidance is achieved by properly designing the collision penalty function in the trajectory optimization, following recent literature [68], [69]. If a collision happens, the penalty term grows rapidly, which quickly dominates the overall optimization objective function. The process is repeated for each drone independently until it reaches the designated goal. In this process, all computations, including trajectory planning, state estimation, and flight control, are done on the onboard computer.

## VII. EXPERIMENTS

To validate Omni-swarm's feasibility, accuracy, and practical value, we run a series of experiments: 1) timing and scalability of Omni-swarm on a PC and TX2 onboard computer; 2) experiments to verify the features of Omni-swarm; 3) accuracy comparison with the ground truth on a custom dataset with various other approaches; 4) verification of the global

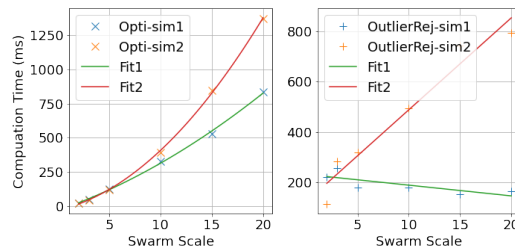


Fig. 9: A demonstration of the growth of computational time with different swarm scales in the simulation of two scenarios. The left figure shows the computation time of graph-based optimization (symbol 'x') and fitted quadratic splines. The right figure shows the computation time of PCM outlier rejection (symbol '+') and fitted lines. Suffix 1 and 2 in the legends correspond to Scenario 1 and Scenario 2.

consistency on an outdoor dataset; 5) real-time performance verification in an inter-drone collision avoidance experiment.

#### A. Scalability

It should be clear that the computation effort scale of the Omni-swarm's front-end is not significant with the number of drones. This is because most components in the front-end are running at a fixed frequency, which is not relevant to the scale of the drone. There are two special cases: 1) VDT's pose estimation is computed separately for each detected drone. It is related to the number of drones in the FoV. In practice, it is unlikely that there will be too many drones in the FoV, so the scalability issue here is not severe. 2) Loop closure detection requires the Faiss database to retrieve the keyframe. Faiss has the complexity of  $O(k \log(l/w))$  for KNN searching [70], where  $k = 5$  is the selection number each time,  $w = 32$  is a parameter of Faiss,  $l$  is the size of the keyframes in the database, which can be considered as linear to the swarm scale. In practice, the computation time that grows logarithmically with the scale of the swarm does not introduce scalability difficulties.

Next, we discuss the scalability of the Omni-swarm's back-end. Here we assume that measurements are consistent between any two drones. Two main computational components are involved in the back-end: 1) PCM outlier rejection: The drone only needs to compute the PCM problem with another drone using relative pose measurements making the complexity of this algorithm  $O(n)$ , where  $n$  is the swarm scale. 2) Graph-based optimization: The scale of Problem (9) grows quadratically with the swarm scale, so the computational effort of graph-based optimization is  $O(n^2)$ .

To verify the conclusions here, we use simulations to simulate the complexity of the back-end computation. Aerial swarms of scales from 2 to 20 are set up in the simulations and are tested under different scenarios: 1) The drones in the swarm are lined up, each drone can only generate relative pose measurements with several neighboring drones, but all drones have distance measurements between them. 2) The drones in the swarm are closely grouped in a grid formation, and all drones can generate relative pose measurements with each other.

The results of the simulation are shown in Fig. 9. It can be seen that the quadratic curves perfectly fit the computation

for graph-based optimization in both scenarios. While in Scenario1, the average computation time of PCM outlier rejection is not significantly related to the swarm scale because the consistency graph that needs to be computed for each drone in Scenario 1 only depends on the number of neighbors (2-4). In Scenario2, PCM outlier rejection takes linearly increasing time because all drones have relative pose measurements with each other.

## B. Features Verification

1) *Fast initialization and plug-and-play*: In this experiment, drone 1 and 2 are firstly booted up in the field, with drone 1 carrying a battery in poor condition to simulate a real emergency. Omni-swarm is properly initialized after the program is loaded and the two drones take off. Drone 3 is then placed in the field and powered on, at which point drone 1 is automatically landed by the onboard computer's control system due to its lack of power. This action does not affect Omni-swarm's estimation, and the program on drone 3 is then successfully loaded and correctly estimated by drone 1. Drone 3 establishes a correct estimation of 1 and 2 at the same time. Next, we launch drone 3 and then remove the battery of drone 1 and remove it from the field. We then add drone 4 in the field, and after it is estimated by drone 3, we launch drone 4 and let drone 2 go around drones 3 and 4 to the back of the field with the inter-drone collision avoidance algorithm. Finally, we land drone 2 and shut it down, and the state estimation of the remaining drones 3 and 4 works normally.

2) *Non-line-of-sight condition flight experiment*: In this experiment, the two drones are separated by a short wall at the beginning. Then, the two drones take-off and fly around the wall anti-clockwise. During the flight, the drones are not in each others' line-of-sight. We disable the UWB measurement in Omni-swarm since the UWB modules can still measure distances crossing the foam obstacles between the drones. Once a drone visits the other drone's starting place, Omni-swarm immediately initializes, and the state of the two drones is estimated correctly. This experiment verifies that the proposed method works even in totally non-line-of-sight cases if the same location has been revisited.

Recording of these two experiments can be found in our video<sup>5</sup>.

## C. Evaluation on Datasets

We collect indoor and outdoor datasets for accuracy comparison. The datasets include the raw sensor data and ground truth given by a motion capture system in the indoor laboratory and outdoor environment without ground truth. The indoor datasets include: 1) two double-drone-fixed-heading formation flight datasets; 2) a double-drone formation flight dataset with the drones heading towards the velocity direction; 3) a double-drone-random-target flight dataset. The details of the datasets are shown in Table V.

TABLE V: The datasets we use for accuracy comparison. Traj. Len. is the trajectory length of different drones in the aerial swarm. Takeoff time is the average time for aircraft takeoff in different datasets.

Dataset	Scale	Avg. Traj. Len. (m)	Takeoff Time (s)	Type	Envir.
<b>Parallel1</b>	2	93.0	21.5	Real-world	Indoor
<b>Parallel2</b>	2	106.4	41.3	Real-world	Indoor
<b>Parallel3</b>	2	130.7	39.1	Real-world	Indoor
<b>RandFlight</b>	2	45.5	57.8	Real-world	Indoor
<b>Outdoor1</b>	2	237.2	31	Real-world	Outdoor

1) *Metric Definitions*: In this paper, two metrics are involved, absolute trajectory error (ATE) and relative error (RE), which are used to measure the global consistency and relative estimation accuracy, respectively, of Omni-swarm. The definition of ATE used in this paper is the same as that in [71]. A point worth noting is that the drift of each trajectory and the relative state estimation error between trajectories estimated by the state estimation are included in the ATE, i.e., the ATE measures the global consistency and the relative localization accuracy of the swarm.

The RE is used to characterize the relative state estimation accuracy, which is defined as

$$\begin{aligned}
 RE_{rot_k}^i &= \left( \frac{1}{N} \sum_{t=0}^{N-1} \|\angle \left( v_k \hat{\mathbf{R}}_k^t T v_k \hat{\mathbf{R}}_i^t \right) \|^2 \right)^{\frac{1}{2}} \\
 RE_{pos_k}^i &= \left( \frac{1}{N} \sum_{t=0}^{N-1} \left\| \left( v_k \hat{\mathbf{R}}_k^t T (v_k \hat{\mathbf{T}}_i^t - v_k \hat{\mathbf{T}}_k^t) \right) \right\|^2 \right)^{\frac{1}{2}}, \\
 RE_{pos_{k_{axis}}}^i &= \left( \frac{1}{N} \sum_{t=0}^{N-1} \left| \left( v_k \hat{\mathbf{R}}_k^t T (v_k \hat{\mathbf{T}}_i^t - v_k \hat{\mathbf{T}}_k^t) \right)_{axis} \right|^2 \right)^{\frac{1}{2}}, \tag{14}
 \end{aligned}$$

where  $RE_{rot_k}^i$  and  $RE_{pos_k}^i$  are the rotation and translation relative error, respectively, of drone  $i$  estimated by drone  $k$ ,  $N$  is the total length of the trajectory, and  $\angle(\cdot)$  is the angle of the rotation error, which is defined in [71].  $axis \in \{x, y, z\}$ ,  $RE_{pos_{k_x}}^i, RE_{pos_{k_y}}^i, RE_{pos_{k_z}}^i$  are the translation relative errors on the  $x, y$  and  $z$  axis, which helps the reader to establish a more precise perception of the error. Due to the space limitation, we use  $RE_{pos_k}^i$  in the dataset comparison and only use per-axis errors for inter-drone collision avoidance in real-world environments. The ATE and RE in this paper are the averaged results over the dataset and the units are in meters and degrees, respectively.

2) *Accuracy Comparison with Ground Truth*: Table VI shows the accuracy comparison on recorded datasets (**Parallel1**, **Parallel2**, **RandFlight**). **Parallel1** and **Parallel2** are the parallel flight tasks and **RandFlight** is the double-drone-random-target flight datasets task. The table contains a comparison of four methods: 1) our previous method in [18] (named *Xu2020*). 2) An approach with only map-based measurements with VIO (named *PGO*), namely the pose graph optimization method for CSLAM [16], [17], [26], [27]. 3) A *VIO-only* method [36]. VIO does not have the ability to perform relative positioning, so we use the ground truth to align the coordinate system of different drones. 4) Our proposed method. In addition, the table also consists of four ablation studies: 1) remove UWB (named *Without UWB*), 2)

<sup>5</sup><https://www.youtube.com/watch?v=SMtJUKoza4>

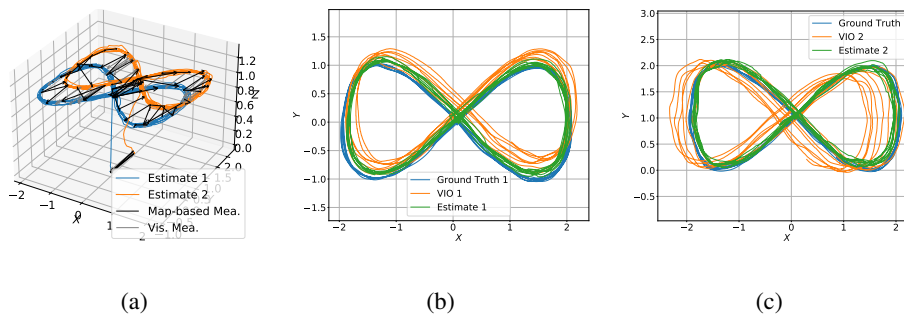


Fig. 10: The estimated trajectories of the two drones in an indoor dataset. (a): The black arrows are the map-based measurements. The gray arrows are the visual detection measurements. Only part of the measurements are shown in this figure for better visualization. (b)-(c): The comparison of the ground truth trajectories (blue), the VIO trajectories (orange), and the estimated trajectories (green). The VIO trajectories drift away from the ground truth trajectories, while the estimated trajectories closely follow the ground truth.

TABLE VI: Comparison of swarm state estimation methods on the indoor datasets. Init. time. is the average time for successful initialization of state estimation. The overall best results in ATE and RE are in bold.

Dataset	Metrics	<i>Xu2020</i>		<i>PGO</i>		<i>VIO Only</i>		<b>Proposed</b>		<i>Without UWB</i>		<i>Without Tracking</i>		<i>Without Map-based</i>		<i>Without Outlier Rej.</i>	
		Pos	Rot	Pos	Rot	Pos	Rot	Pos	Rot	Pos	Rot	Pos	Rot	Pos	Rot	Pos	Rot
<b>Parallel1</b>	ATE	0.600	8.6°	0.149	2.5°	0.228	2.5°	0.127	2.2°	<b>0.126</b>	<b>2.2°</b>	0.149	2.3°	0.160	2.4°	0.360	2.6°
	RE	0.959	17.7°	0.114	3.3°	0.261	2.3°	<b>0.062</b>	<b>2.3°</b>	0.063	2.7°	0.100	3.7°	0.062	2.6°	0.064	3.2°
<b>Parallel2</b>	ATE	0.351	3.2°	0.113	3.2°	0.225	3.0°	<b>0.086</b>	<b>3.0°</b>	0.094	3.1°	0.103	3.0°	0.225	3.1°	0.119	3.2°
	RE	0.241	4.7°	0.124	4.3°	0.230	4.3°	<b>0.063</b>	<b>4.3°</b>	0.083	4.0°	0.096	4.3°	0.069	4.2°	0.069	4.3°
<b>Parallel3</b>	ATE	0.470	2.6°	0.153	1.8°	0.281	3.1°	<b>0.119</b>	<b>1.8°</b>	0.128	1.8°	0.134	1.7°	0.242	3.0°	0.535	5.1°
	RE	0.566	1.8°	0.130	1.7°	0.225	4.7°	<b>0.072</b>	<b>1.8°</b>	0.082	1.8°	0.081	1.8°	0.075	1.9°	0.143	4.0°
<b>RandFlight</b>	ATE	0.197	2.7°	0.153	1.8°	0.125	2.6°	0.088	2.7°	0.092	2.7°	0.101	2.7°	<b>0.079</b>	<b>2.7°</b>	0.165	5.9°
	RE	0.191	2.1°	0.130	1.7°	0.160	1.1°	<b>0.069</b>	<b>1.3°</b>	0.071	1.2°	0.078	1.4°	0.074	1.4°	0.147	9.8°

TABLE VII: Comparison of initialization and observability of different methods. In the table, Init. Time is the time when the system is initialized after boot-up, Init. Trials is the number of times the system attempted the initialization optimization, Succ. Rate is the success rate of the initialization, and RE is the relative measurement error. '-' in the table means the method is not applicable to the dataset.

Dataset	<i>Xu2020</i>					<i>PGO</i>					<b>Proposed</b>				
	Init. Time (s)	Init. Trials	Succ. Rate	$RE_{pos}$	$RE_{rot}$	Init. Time (s)	Init. Trials	Succ. Rate	$RE_{pos}$	$RE_{rot}$	Init. Time (s)	Init. Trials	Succ. Rate	$RE_{pos}$	$RE_{rot}$
<b>Parallel1</b>	52.1	3	0%	0.959	17.7°	1.36	1	100%	0.114	3.3°	1.29	1	100%	0.062	2.3°
<b>RandFlight1</b>	84.4	3	100%	0.191	2.1°	1.36	1	100%	0.130	1.7°	1.43	1	100%	0.069	1.3°
<b>Parallel1 Outdoor</b>	52.3	3	0%	1.047	3.2°	-	-	-	-	-	1.31	1	100%	0.057	1.1°
<b>RandFlight1 Outdoor</b>	83.1	3	100%	0.260	4.5°	-	-	-	-	-	1.38	1	100%	0.111	3.0°

remove visual drone tracking (named *Without Tracking*), 3) remove map-based localization (named *Without Map-based*), and 4) remove the outlier rejection and Huber norm (named *Without Outlier Rej.*).

The proposed method has the best ATE and RE overall, where the estimated trajectories of the proposed method on **Parallel1** are shown in Fig. 10. Our method also shows the best performance on the simulated outdoor dataset in Table VII (the details will be discussed in Sect. VII-C3), which proves that Omni-swarm has the best adaptation to the environment. *Xu2020* nearly fails on the **Parallel1** and **Parallel3** datasets because it lacks observability in the parallel flight task. The reason that relative state estimation of *Xu2020* does not fail on **Parallel2** is that the drone position control errors in **Parallel2** cause 15 cm of relative motion between the drones, which causes the *Xu2020* approach to obtain an inaccurate estimate with relative motion. However, the visual detection of *Xu2020* does not function on this dataset due to the FoV limitation. This makes the final relative localization accuracy only reach 0.241 m. The same situation applies to the **RandFlight** dataset, where UWB-odometry in *Xu2020* establishes a rough relative estimate, but the method does not yield accurate relative

measurements due to the FoV limitation. As a comparison, *Without Map-based* is similar to *Xu2020* but has the new omnidirectional VDT module and initialization, and it has much better relative localization accuracy than *Xu2020*. This illustrates the significance of the omnidirectional front-end and initialization in this paper.

On all datasets, **PGO** obtains accuracy in the order of decimeters, but we will show the drawbacks of **PGO** in outdoor environments later. The *VIO-only* method is not as accurate as **PGO** and the proposed methods, and for this method we need to manually align the coordinate system of different drones at the very beginning.

From the ablation studies *Without UWB* and *Without Tracking*, we can see the contribution of these two types of factors (distance and visual detection factors) to system accuracy. In addition, good relative localization accuracy can be obtained without using the map-based localization (*Without Map-based*), which is especially important for feature-poor outdoor environments. However, the global consistency (ATE) deteriorates to *VIO-Only* in this case. This comparison demonstrates the importance of map-based localization for global consistency. Finally, *Without Outlier Rej.* is two to four



times less accurate than the proposed method in some datasets, demonstrating the benefits of introducing the outlier rejection module and Huber norm.

3) *Observability analysis and initialization*: We conduct comparative experiments to verify the improvement in the observability of Omni-swarm relative to previous works. The experiments give a comparison of Omni-swarm with the *PGO* and *Xu2020* methods in terms of initialization performance under different environments. The comparison shows the initialization of the indoor environment and the outdoor environment under parallel and random flight. The comparison results are shown in Table VII.

We argue that open outdoor scenes do not have sufficient feature points for map-based localization. While VIO still works in this case, it is not as accurate as it is in indoor environments. On this basis, we use the **Parallel1** and **RandFlight** datasets to simulate outdoor datasets with ground truth, named **Parallel1Out** and **RandFlightOut** in Table VII. We restrict VINS-Fisheye to use ground features and limit surrounding features (no more than 20), and turn off the map-based localization to simulate feature-poor outdoor environments.

As shown in Table VII, our proposed method has accurate relative localization on all indoor and outdoor datasets, which means that it is observable in all cases. In indoor environments, map-based localization can be used for initialization (Sect. V-E1), while visual detection initialization (Sect. V-E2) is used for indoor and outdoor environments. The *PGO* method also has fair relative localization accuracy indoors but is not applicable to open outdoor environments. The *Xu2020* method has a fair relative localization accuracy on the dataset with relative motion (**RandFlight**) in indoor and outdoor cases due to its use of UWB-odometry fusion, which is observable in this case. However, on the parallel flight dataset (neither indoor or outdoor), the relative localization of this method is considered as having failed with a huge relative localization error.

Table VII also shows the initialization quantitatively. The proposed method can be initialized quickly before the aerial swarm takes off in all cases and requires only one optimization. The *PGO* method can also be quickly initialized in indoor environments with the help of surrounding feature points, but it does not work in outdoor environments. The *Xu2020* method, on the other hand, does not succeed in initializing the correct state estimate on the parallel flight dataset. On the **RandFlight** dataset, *Xu2020* is able to initialize successfully, but it requires a period after the aerial swarm has flown to complete the process. This period brings the issue of flight safety. As a reference, the take-off time of these datasets can be found in Table V.

4) *Impact of the random deletion mechanism*: Table VIII shows the comparison between the random deletion proposed in Sect. V-A (**RandDel** in the table) and the other two deletion mechanisms, an approach without any deletion mechanism (*NoDel* in the table), and the sliding window mechanism used in [18] (*SldWin* in the table), which deletes the first frame when swarm keyframes in  $\mathcal{G}$  are more than  $m_{max}$ . In the comparison,  $m_{max}$  is set to 100. The **Parallel1** and **Parallel2** dataset generates 255 and 321 swarm keyframes, respectively.

In the table, all three approaches have similar relative

TABLE VIII: Accuracy comparison of different deletion mechanisms. Avg. Time is the averaged solving time of the graph-based optimization.

Dataset	Metrics	NoDel		SldWin		RandDel	
		Pos	Rot	Pos	Rot	Pos	Rot
Parallel1	ATE	0.091	3.2°	0.153	3.3°	0.086	3.0°
	RE	0.066	4.4°	0.067	4.4°	0.063	4.3°
	Avg Time.	21.00ms		12.11ms		14.36ms	
Parallel2	ATE	0.115	1.8°	0.175	2.3°	0.119	1.8°
	RE	0.072	1.8°	0.073	1.8°	0.072	1.8°
	Avg Time.	21.51ms		11.62ms		13.86ms	

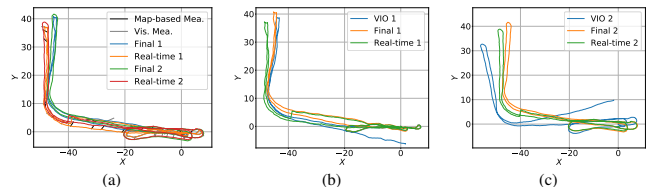


Fig. 11: The estimated trajectories and the VIO trajectories of the two drones in an outdoor dataset. For better visualization of the global consistency, the real-time estimated trajectories and final estimated trajectories are shown in the figures. (a) The estimated trajectories of two drones with the map-based measurements and visual detection measurements. The black lines are the map-based measurements. (b)-(c) The detailed comparison of the estimated trajectories (blue) and the VIO trajectories (orange).

measurement accuracy (measured using RE), showing the deletion mechanism has no effect on the relative state estimates. Nevertheless, the *NoDel* approach has the best global consistency (measured using ATE), **RandDel** is slightly worse than *NoDel*, and the worst approach, *SldWin*, is close to the level of VIO-only in Table VI. Concluding, **RandDel** has the same optimization speed compared to *SldWin* but is 50% faster than *NoDel*.

#### D. Evaluation on Outdoor Dataset

We also evaluate Omni-swarm on an outdoor dataset to verify that Omni-swarm is flexible in many scenarios. Fig. 11 shows the estimation trajectories on a double-drone outdoor dataset. The real-time trajectory in the figure is the output of state estimation in real-time with forward propagation (Eq. 12). The Final trajectory is the final estimate poses of keyframes in the  $\mathcal{G}$  at the end of the task.

As a result of a 235-meter-run in the outdoor environment, the estimated states drift only 1.9 m from the start point, which is 0.8% of the trajectory length, showing that the global consistency is ensured. As a comparison, the original VIO trajectories' drifts averaged 6.35 m from the start point, which is 3.3 times bigger compared to our proposed method. The real-time trajectories in the figure have a noticeable jump, while the final trajectories do not. This is because Omni-swarm eliminates the drifting error of VIO at this instant after detecting a map-based measurement, so the real-time estimation jumps to a more global consistency estimate. On the other hand, the final trajectory is smoother without this jump because it is a whole trajectory estimated result.

#### E. Inter-drone collision avoidance experiments

In the inter-drone collision avoidance experiments, each drone performs independent swarm state estimation in real-

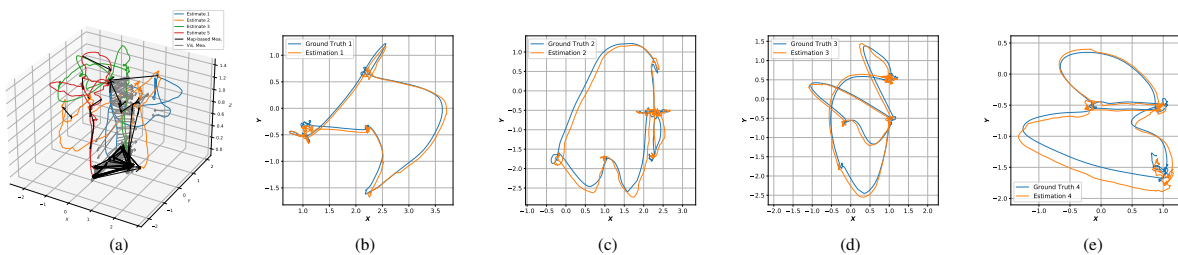


Fig. 12: These figures show the estimated trajectories, and ground truth trajectories a four-drone inter-drone collision avoidance experiment. (a): The estimated trajectories of the drones. The black lines are the map-based measurements. The gray lines are the detection measurements. Only part of the measurements is shown in this figure for better visualization. (b)-(e): The detailed comparison of the estimated trajectories (blue) and the ground truth trajectories (orange).

TABLE IX: Accuracy comparison on three-drone and four-drone inter-drone collision avoidance experiments. Traj. Len. is the averaged trajectory length of the experiment.

Task		Proposed Method			VIO-Only		
		x	y	z	x	y	z
Three Drone	Traj. Len.	27.8					
	$RE_{x,y,z}$	0.076	0.074	0.083	0.064	0.061	0.030
	$RE_{rot}$	1.20°			0.8°		
	$ATE_{pos}$	0.103			0.078		
	$ATE_{rot}$	0.427°			0.096°		
Four Drone	Traj. Len.	19.9					
	$RE_{x,y,z}$	0.065	0.056	0.088	0.070	0.131	0.064
	$RE_{rot}$	1.46°			1.2°		
	$ATE_{pos}$	0.103			0.091		
	$ATE_{rot}$	0.399°			0.160°		

time and uses the estimated results for inter-drone collision avoidance and planning, which follows Sect. VI-F. Table IX shows the accuracy of real-time estimated states in the inter-drone collision avoidance experiments, and Fig. 12 shows the estimated trajectories of a four-drone inter-drone collision avoidance experiment. This experiment successfully validates the practical value of Omni-swarm in a real-world environment. One may notice that VIO has better performance on some data. This is because in short flight tasks (19-30 m as listed in the table), VIO drifts less and therefore has higher accuracy. However, our proposed method outperforms VIO on all datasets in the previous comparison, where drones have longer trajectories. Additionally, VIO-only requires the aerial swarm to be initialized using ground truth or in using predefined starting points, and its error increases with flight range, limiting the practical value of this approach.

#### F. Computation Time

Table X shows the computation time of the main components of Omni-swarm in experiments. The first P&P scenario is the four drone plug-and-play experiment in Sect. VII-B1. The second C&3 scenario is the three drone inter-drone collision avoidance experiment in Sect. VII-E. The C&3 scenario's back-end computation time is much longer than that of P&P because it is a more complex task. The table shows that Omni-swarm is able to achieve real-time performance with the computational frequency shown in Table II.

### VIII. CONCLUSION AND FUTURE WORK

This paper introduced Omni-swarm, a decentralized omnidirectional visual-inertial-UWB state estimation system for aerial swarms. Compared to previous works, the proposed

TABLE X: The average computation time on each component in real-world experiments on the onboard computer. The units in this table are milliseconds. The abbreviations VINSF, VINSB, Det, Trk, LoopDet., and Opti. are the same as in Table II. PCM is the PCM outlier rejection.

Scenario	VINSF	VINSB	Det & Trk	PoseEsti.	Desc.	Loop Det.	Init.	PCM	Opti.
P&P	12.1	43.8	32.8	24.1	273.7	48.9	8.6	0.8	6.2
C&3	19.5	82.1	41.4	16.6	262.6	28.1	23.7	110.5	155.5

approach addresses the complicated initialization issue, observability issue caused by a restricted FoV, and global consistency issue. To demonstrate the feasibility and effectiveness of Omni-swarm, we tested it by extensive aerial swarm flight experiments. Compared with ground truth data from the motion capture system, the state estimation results reach centimeter-level accuracy on relative estimation while ensuring global consistency. With Omni-swarm, formation flights in various complex environments are no longer impossible. Inter-drone collision avoidance is successfully demonstrated based on Omni-swarm. We believe that Omni-swarm can be widely adopted in a variety of scenarios and on multiple scales.

However, our system is also inevitably characterized by a number of shortcomings that need to be improved in future work, including: 1) The dependence on camera intrinsic and extrinsic calibration, which is a common problem for all visual SLAM systems. In the future, we will try to develop online fault detection and calibration to further improve the robustness. 2) Scalability. Our current back-end algorithm is  $O(n^2)$  for the growth of the scale of the swarm, which makes it difficult to apply Omni-swarm to larger-scale swarms (more than 100 drones). In the future, we will try to extend Omni-swarm to large aerial swarms. 3) Communication range. Current work is focused on state estimation, and the network setup currently limits our aerial swarm to work within a short distance (22.4 m) from each other. In the future, we will extend the communication distance by deploying routing algorithms, e.g., Batman-adv [72], and upgrading communication devices.

In the future, we will deploy Omni-swarm to real-world applications and also use it to build the dense global map to exploit the advantages of Omni-swarm fully.

### REFERENCES

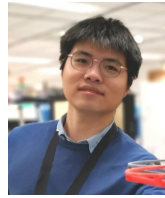
- [1] B. Zhou, F. Gao, L. Wang, C. Liu, and S. Shen, "Robust and efficient quadrotor trajectory generation for fast autonomous flight," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3529–3536, 2019.
- [2] K. Wang, F. Gao, and S. Shen, "Real-time scalable dense surfel mapping," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*, 2019.

- [3] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Trans. Robot. (TRO)*, vol. 34, no. 4, pp. 1004–1020, 2018.
- [4] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras," *IEEE Trans. Robot. (TRO)*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [5] P. Geneva, K. Eickenhoff, W. Lee, Y. Yang, and G. Huang, "Openvins: A research platform for visual-inertial estimation," in *Proc. of the IEEE International Conference on Robotics and Automation*, Paris, France, 2020.
- [6] J. Zhang and S. Singh, "LOAM: Lidar odometry and mapping in real-time," in *Robotics: Science and Systems*, vol. 2, no. 9. Berkeley, CA, 2014, pp. 1–9.
- [7] J. Lin and F. Zhang, "Loam livox: A fast, robust, high-precision lidar odometry and mapping package for lidars of small fov," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2020, pp. 3126–3131.
- [8] J. A. Preiss, W. Honig, G. S. Sukhatme, and N. Ayanian, "Crazyswarm: A large nano-quadcopter swarm," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2017, pp. 3299–3304.
- [9] A. Ledergerber, M. Hamer, and R. D'Andrea, "A robot self-localization system using one-way ultra-wideband communication," in *Proc. of the IEEE/RSJ Intl. Conf. on Intell. Robots and Syst.(IROS)*. IEEE, 2015, pp. 3131–3137.
- [10] A. Jaimes, S. Kota, and J. Gomez, "An approach to surveillance an area using swarm of fixed wing and quad-rotor unmanned aerial vehicles uav (s)," in *2008 IEEE International Conference on System of Systems Engineering*. IEEE, 2008, pp. 1–6.
- [11] S. Moon, Y. Choi, D. Kim, M. Seung, and H. Gong, "Outdoor swarm flight system based on RTK-GPS," *Journal of KIISE*, vol. 43, no. 12, pp. 1315–1324, 2016.
- [12] P. Petráček, V. Krátký, M. Petrлік, T. Báča, R. Kratochvíl, and M. Saska, "Large-scale exploration of cave environments by unmanned aerial vehicles," vol. 6, no. 4, pp. 7596–7603, 2021.
- [13] A. J. Smith and G. A. Hollinger, "Distributed inference-based multi-robot exploration," *Auton. Robots*, vol. 42, no. 8, pp. 1651–1668, 2018.
- [14] P. Zhu, Y. Yang, W. Ren, and G. Huang, "Cooperative visual-inertial odometry," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*, 2021.
- [15] P. Zhu, P. Geneva, W. Ren, and G. Huang, "Distributed visual-inertial cooperative localization," *Proc. of the IEEE/RSJ Intl. Conf. on Intell. Robots and Syst.(IROS)*, 2021.
- [16] 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," *Intl. J. Robot. Research (IJRR)*, vol. 36, no. 12, pp. 1286–1311, 2017.
- [17] P.-Y. Lajoie, B. Ramtoula, Y. Chang, L. Carlone, and G. Beltrame, "DOOR-SLAM: Distributed, online, and outlier resilient slam for robotic teams," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 2, pp. 1656–1663, 2020.
- [18] H. Xu, L. Wang, Y. Zhang, K. Qiu, and S. Shen, "Decentralized visual-inertial-UWB fusion for relative state estimation of aerial swarm," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2020, pp. 8776–8782.
- [19] K. Guo, Z. Qiu, W. Meng, L. Xie, and R. Teo, "Ultra-wideband based cooperative relative localization algorithm and experiments for multiple unmanned aerial vehicles in gps denied environments," *International Journal of Micro Air Vehicles*, vol. 9, no. 3, pp. 169–186, 2017.
- [20] K. Guo, X. Li, and L. Xie, "Ultra-wideband and odometry-based cooperative relative localization with application to multi-uav formation control," *IEEE Transactions on Cybernetics*, 2019.
- [21] T. Ziegler, M. Karrer, P. Schmuck, and M. Chli, "Distributed Formation Estimation via Pairwise Distance Measurements," *IEEE Robotics and Automation Letters*, 2021.
- [22] V. Walter, N. Staub, A. Franchi, and M. Saska, "Uvdar system for visual relative localization with application to leader-follower formations of multirotor uavs," *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 2637–2644, 2019.
- [23] T. Nguyen, K. Mohta, C. J. Taylor, and V. Kumar, "Vision-based Multi-MAV Localization with Anonymous Relative Measurements Using Coupled Probabilistic Data Association Filter," 2020.
- [24] N. Piasco, J. Marzat, and M. Sanfourche, "Collaborative localization and formation flying using distributed stereo-vision," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2016, pp. 1202–1207.
- [25] M. W. Achtelik, S. Weiss, M. Chli, F. Dellaert, and R. Siegwart, "Collaborative stereo," in *Proc. of the IEEE/RSJ Intl. Conf. on Intell. Robots and Syst.(IROS)*. IEEE, 2011, pp. 2242–2248.
- [26] A. Cunningham, M. Paluri, and F. Dellaert, "DDF-SAM: Fully distributed SLAM using constrained factor graphs," in *Proc. of the IEEE/RSJ Intl. Conf. on Intell. Robots and Syst.(IROS)*. IEEE, 2010, pp. 3025–3030.
- [27] A. Cunningham, V. Indelman, and F. Dellaert, "DDF-SAM 2.0: Consistent distributed smoothing and mapping," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2013, pp. 5220–5227.
- [28] L. Wang, D. Cheng, F. Gao, F. Cai, J. Guo, M. Lin, and S. Shen, "A Collaborative Aerial-Ground Robotic System for Fast Exploration," in *Proc. of the Intl. Sym. on Exp. Robot. (ISER)*, 2018, pp. 59–71.
- [29] N. Michael, S. Shen, K. Mohta, V. Kumar, K. Nagatani, Y. Okada, S. Kiribayashi, K. Otake, K. Yoshida, K. Ohno *et al.*, "Collaborative mapping of an earthquake damaged building via ground and aerial robots," in *Field and Service Robotics*. Springer, 2014, pp. 33–47.
- [30] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, "ORB-SLAM3: An accurate open-source library for visual, visual-inertial, and multimap slam," *IEEE Trans. Robot. (TRO)*, 2021.
- [31] R. Mur-Artal and J. D. Tardós, "Visual-inertial monocular slam with map reuse," *IEEE Robotics and Automation Letters (RA-L)*, vol. 2, no. 2, pp. 796–803, 2017.
- [32] V. Usenko, N. Demmel, D. Schubert, J. Stückler, and D. Cremers, "Visual-inertial mapping with non-linear factor recovery," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 422–429, 2019.
- [33] T. Qin, S. Cao, J. Pan, and S. Shen, "A general optimization-based framework for global pose estimation with multiple sensors," *arXiv preprint arXiv:1901.03642*, 2019.
- [34] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct ekf-based approach," in *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2015, pp. 298–304.
- [35] A. Weinstein, A. Cho, G. Loianno, and V. Kumar, "Visual inertial odometry swarm: An autonomous swarm of vision-based quadrotors," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1801–1807, 2018.
- [36] P. C. Lusk, X. Cai, S. Wadhwan, A. Paris, K. Fathian, and J. P. How, "A Distributed Pipeline for Scalable, Deconflicted Formation Flying," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5213–5220, 2020.
- [37] T.-M. Nguyen, Z. Qiu, T. H. Nguyen, M. Cao, and L. Xie, "Persistently excited adaptive relative localization and time-varying formation of robot swarms," *IEEE Trans. Robot. (TRO)*, vol. 36, no. 2, pp. 553–560, 2019.
- [38] T.-M. Nguyen, Z. Qiu, T. H. Nguyen, and *et al.*, "Distance-Based Cooperative Relative Localization for Leader-Following Control of MAVs," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3641–3648, 2019.
- [39] T. Qin, P. Li, and S. Shen, "Relocalization, global optimization and map merging for monocular visual-inertial slam," in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2018, pp. 1197–1204.
- [40] K. Qiu, T. Liu, and S. Shen, "Model-based global localization for aerial robots using edge alignment," *IEEE Robotics and Automation Letters*, vol. 2, no. 3, pp. 1256–1263, 2017.
- [41] F. Dellaert, M. Kaess *et al.*, "Factor Graphs for Robot Perception," *Foundations and Trends® in Robotics*, vol. 6, no. 1-2, pp. 1–139, 2017.
- [42] W. Gao and S. Shen, "Dual-fisheye omnidirectional stereo," in *Proc. of the IEEE/RSJ Intl. Conf. on Intell. Robots and Syst.(IROS)*, 2017, pp. 6715–6722.
- [43] W. Gao, K. Wang, W. Ding, F. Gao, T. Qin, and S. Shen, "Autonomous aerial robot using dual-fisheye cameras," *J. Field Robot. (JFR)*, vol. 37, no. 4, pp. 497–514, 2020.
- [44] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv*, 2018.
- [45] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [46] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in *Proc. of the IEEE Intl. Conf. on Pattern Recognition*. IEEE, 2010, pp. 2544–2550.
- [47] M. Betke and Z. Wu, "Data Association for Multi-Object Visual Tracking," *Synthesis Lectures on Computer Vision*, vol. 6, no. 2, pp. 9–11, 2016.
- [48] P. Konstantinova, A. Udvarev, and T. Semerdjiev, "A Study of a Target Tracking Algorithm Using Global Nearest Neighbor Approach," in *Proceedings of the International Conference on Computer Systems and Technologies (CompSysTech'03)*, 2003, pp. 290–295.
- [49] A. Sinha, Z. Ding, T. Kirubarajan, and M. Farooq, "Track Quality

Based Multitarget Tracking Approach for Global Nearest-Neighbor Association,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. 48, no. 2, pp. 1179–1191, 2012.

- [50] H. W. Kuhn, “The Hungarian method for the assignment problem,” *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [51] D. P. Bertsekas and D. A. Castanon, “A Forward/Reverse Auction Algorithm for Asymmetric Assignment Problems,” *Computational Optimization and Applications*, vol. 1, no. 3, pp. 277–297, 1992.
- [52] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, “6-DoF Object Pose from Semantic Keypoints,” in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2017, pp. 2011–2018.
- [53] M. Pavliv, F. Schiano, C. Reardon, D. Floreano, and G. Loianno, “Tracking and Relative Localization of Drone Swarms With a Vision-Based Headset,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1455–1462, 2021.
- [54] V. Lepetit, F. Moreno-Noguer, and P. Fua, “Eppn: An accurate o (n) solution to the pnp problem,” *International journal of computer vision*, vol. 81, no. 2, p. 155, 2009.
- [55] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “NetVLAD: CNN architecture for weakly supervised place recognition,” in *Proc. of the IEEE Intl. Conf. on Pattern Recognition*, 2016, pp. 5297–5307.
- [56] P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk, “From coarse to fine: Robust hierarchical localization at large scale,” in *Proc. of the IEEE Intl. Conf. on Pattern Recognition*, 2019, pp. 12716–12725.
- [57] D. DeTone, T. Malisiewicz, and A. Rabinovich, “Superpoint: Self-supervised interest point detection and description,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 224–236.
- [58] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [59] Z. Zhang and A. R. Hanson, “3d reconstruction based on homography mapping,” *Proc. ARPA96*, pp. 1007–1012, 1996.
- [60] J. G. Mangelson, D. Dominic, R. M. Eustice, and R. Vasudevan, “Pairwise consistent measurement set maximization for robust multi-robot map merging,” in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2018, pp. 2916–2923.
- [61] A. Rosinol, M. Abate, Y. Chang, and L. Carlone, “Kimera: an Open-Source Library for Real-Time Metric-Semantic Localization and Mapping,” in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2020, pp. 1689–1696.
- [62] B. Pattabiraman, M. M. A. Patwary, A. H. Gebremedhin, W.-k. Liao, and A. Choudhary, “Fast algorithms for the maximum clique problem on massive graphs with applications to overlapping community detection,” *Internet Mathematics*, vol. 11, no. 4-5, pp. 421–448, 2015.
- [63] P. J. Huber, “Robust estimation of a location parameter,” in *Breakthroughs in Statistics*. Springer, 1992, pp. 492–518.
- [64] S. Agarwal, K. Mierle, and T. C. S. Team, “Ceres Solver,” 3 2022. [Online]. Available: <https://github.com/ceres-solver/ceres-solver>
- [65] B. Zhou and N. Bose, “Multitarget tracking in clutter: Fast algorithms for data association,” *IEEE Transactions on aerospace and electronic systems*, vol. 29, no. 2, pp. 352–363, 1993.
- [66] B. P. Crow, I. Widjaja, J. G. Kim, and P. T. Sakai, “IEEE 802.11 Wireless Local Area Networks,” *IEEE Communications magazine*, vol. 35, no. 9, pp. 116–126, 1997.
- [67] A. S. Huang, E. Olson, and D. C. Moore, “LCM: Lightweight communications and marshalling,” in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2010, pp. 4057–4062.
- [68] X. Zhou, J. Zhu, H. Zhou, C. Xu, and F. Gao, “Ego-swarm: A fully autonomous and decentralized quadrotor swarm system in cluttered environments,” in *Proc. of the IEEE Intl. Conf. on Robot. and Autom. (ICRA)*. IEEE, 2021, pp. 4101–4107.
- [69] B. Zhou, J. Pan, F. Gao, and S. Shen, “Raptor: Robust and perception-aware trajectory replanning for quadrotor fast flight,” *IEEE Trans. Robot. (TRO)*, 2021.
- [70] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *arXiv preprint arXiv:1702.08734*, 2017.
- [71] Z. Zhang and D. Scaramuzza, “A Tutorial on Quantitative Trajectory Evaluation for Visual(-Inertial) Odometry,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 7244–7251.
- [72] D. Seither, A. König, and M. Hollick, “Routing performance of Wireless Mesh Networks: A practical evaluation of BATMAN advanced,” in *2011 IEEE 36th Conference on Local Computer Networks*. IEEE, 2011, pp.

897–904.



**Hao Xu** received the B.Sc. degree in Physics from the University of Science and Technology of China, Hefei, China, in 2016. He is currently working toward the Ph.D. degree with the Hong Kong University of Science and Technology, Hong Kong, under the supervision of Prof. Shaojie Shen. His research interests include unmanned aerial vehicles, aerial swarm, state estimation, sensor fusion, localization and mapping.



**Yichen Zhang** received the B.Eng. degree in computer engineering and B.B.A in general business management in 2020 from the Hong Kong University of Science and Technology, Hong Kong, where he is currently working toward the Ph.D. degree in electronic and computer engineering under the supervision of Prof. S. Shen. His research interests include motion planning, dense mapping and active SLAM for autonomous robots.



**Boyu Zhou** received the B.Eng. degree in mechanical engineering from Shanghai Jiao Tong University, Shanghai, China, in 2018. He is currently working toward the Ph.D. degree in electronic and computer engineering with the Hong Kong University of Science and Technology, Hong Kong. His research interests include aerial robots, autonomous navigation, motion planning, dense mapping, exploration and swarm.



**Luqi Wang** received the B.Eng. degree in computer engineering and aerospace engineering in 2018 from the Hong Kong University of Science and Technology, Hong Kong, where he is currently working toward the Ph.D. degree in electronic and computer engineering. His research interests include control, navigation, and path planning for autonomous robots.



**Xinjie Yao** received her B.Eng in Computer Engineering from the Hong Kong University of Science and Technology in 2020. She received her M.S. in Robotics in 2022 from Carnegie Mellon University. Her research interests are in robotics, focusing on navigation, autonomy, and human-robot interaction.



**Guotao Meng** received his B.Eng. degree in Automation from the School of Electronic and Information Engineering of Xi'an Jiaotong University in 2018. He is working towards his Ph.D. degree at the Visual Intelligence Laboratory, the Department of Electronic and Computer Engineering, the Hong Kong University of Science and Technology. His research interests include image processing, deep learning, computer vision and robotics.



**Shaojie Shen** received the B.Eng. degree in electronic engineering from the Hong Kong University of Science and Technology, Hong Kong, in 2009, and the M.S. degree in robotics and the Ph.D. degree in electrical and systems engineering from the University of Pennsylvania, Philadelphia, PA, USA, in 2011 and 2014, respectively.

He was with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology in September 2014 as an Assistant Professor, and was promoted to an Associate Professor in 2020. His research interests include the areas of robotics and unmanned aerial vehicles, with focus on state estimation, sensor fusion, computer vision, localization and mapping, and autonomous navigation in complex environments.