# Sea You Later: Metadata-Guided Long-Term Re-Identification for UAV-Based Multi-Object Tracking

Cheng-Yen Yang[1*]    Hsiang-Wei Huang[1]    Zhongyu Jiang[1]    Heng-Cheng Kuo[2]
Jie Mei[1]    Chung-I Huang[3]    Jenq-Neng Hwang[1]

[1]Information Processing Lab, University of Washington, USA    [2]National Taiwan University, Taiwan
[3]National Center for High-performance Computing, Taiwan

## Abstract

*Re-identification (ReID) in multi-object tracking (MOT) for UAVs in maritime computer vision has been challenging for several reasons. More specifically, short-term re-identification (ReID) is difficult due to the nature of the characteristics of small targets and the sudden movement of the drone's gimbal. Long-term ReID suffers from the lack of useful appearance diversity. In response to these challenges, we present an adaptable motion-based MOT algorithm, called Metadata Guided MOT (MG-MOT). This algorithm effectively merges short-term tracking data into coherent long-term tracks, harnessing crucial metadata from UAVs, including GPS position, drone altitude, and camera orientations. Extensive experiments are conducted to validate the efficacy of our MOT algorithm. Utilizing the challenging SeaDroneSee tracking dataset, which encompasses the aforementioned scenarios, we achieve a much-improved performance in the latest edition of the UAV-based Maritime Object Tracking Challenge with a state-of-the-art HOTA of 69.5% and an IDF1 of 85.9% on the testing split.*

## 1. Introduction

In recent years, the field of computer vision has witnessed remarkable advancements in multi-object tracking (MOT) techniques. These advancements have enabled the development of systems capable of detecting and following objects in various scenarios. However, one challenging and pressing domain where MOT capabilities are sought is maritime computer vision, including a wide range of applications [1, 17, 18]. The unique characteristics of this environment, including the presence of small-sized objects, challenging visibility conditions due to waves and sun reflections, and the dynamic nature of objects caused by gimbal movements and altitude changes, pose formidable chal-

lenges for conventional MOT algorithms. To exacerbate these difficulties, partial occlusions frequently occur in maritime scenes. Addressing these challenges requires a holistic approach that not only detects and tracks objects but also ensures the long-term re-identification of targets that temporarily vanish and reappear.

The SeaDronesSee-MOT benchmark [13, 21] is specifically designed to assess the capabilities of computer vision algorithms in the maritime domain, emphasizing the detection and tracking of humans, boats, and other objects in open water. While several MOT benchmarks exist, SeaDronesSee-MOT introduces the novel aspect of long-term tracking, a challenge that requires the re-identification (ReID) of objects that temporarily disappear from the scene and subsequently reappear within the same video clip. This challenge is particularly demanding for objects such as boats and swimmers, which may share similar visual characteristics. To address this complex task, we exploit the wealth of drone metadata accompanying each frame, including altitude, viewing angles, and gimbal information, among others. These metadata serve as a valuable resource for accurately associating objects over time, offering a promising avenue for enhancing the robustness and effectiveness of maritime MOT systems. In this paper, we detail our innovative approach, highlighting how the integration of drone metadata enables us to not only excel in multi-object tracking but also excel in long-term ReID, marking a significant step forward in the realm of maritime search and rescue missions.

There are two main challenges for the UAV-based tracking for the maritime:

1) First, the tracking performance is highly dependent on detection qualities, therefore, the nature of the objects in maritime tracking serves as one of the challenges. Similar to the scenario in UAV detection, the scale of the objects we try to detect is highly variant due to the height of the UAVs. Therefore, it is important to reconsider whether using one uniform detector is enough or not.
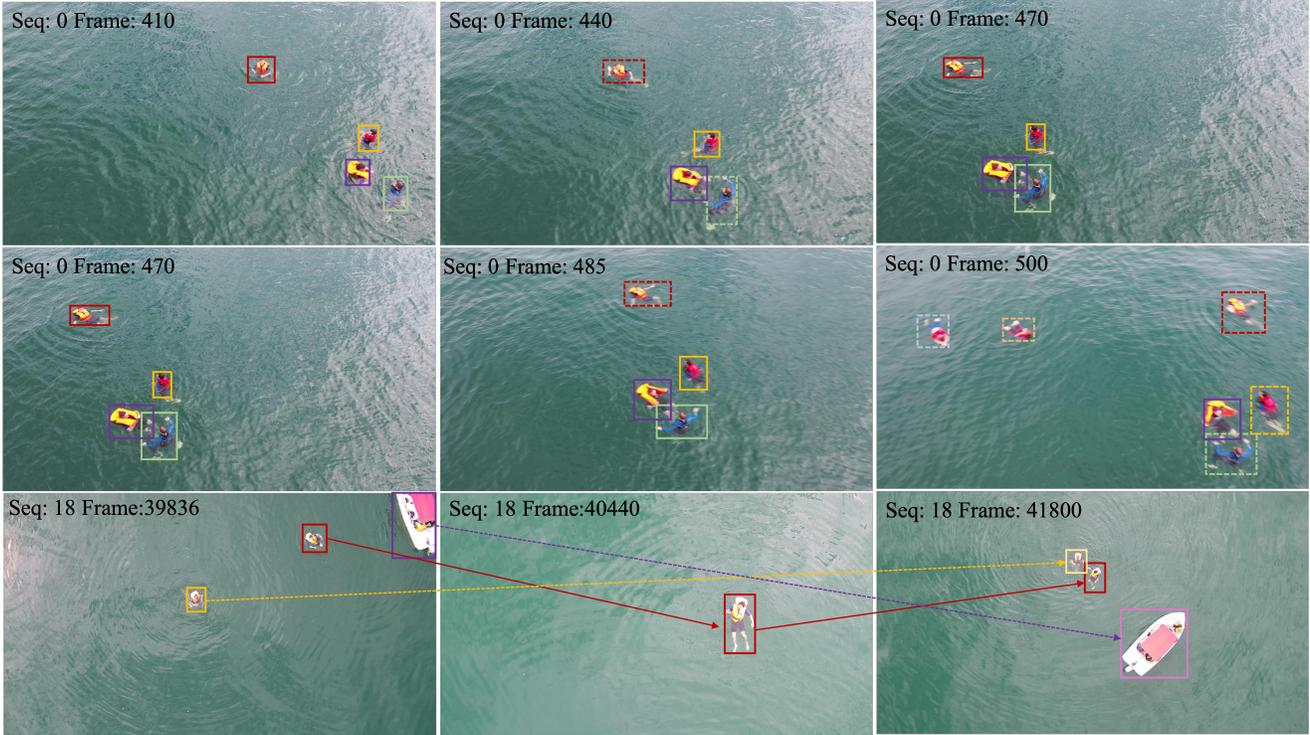
---

* Corresponding author: *cycyang@uw.edu*.

Figure 1. Three types of detection and tracking errors caused by fast drone movement (Top), rapid camera gimbal changes (Mid), and target re-entries (Bottom). Dashed bounding boxes represent the missing detections in our baseline model, while the dashed lines across image frames represent the IDs due to not being able to associate correctly.

2) Secondly, long-term and short-term ReIDs are also challenging. For short-term ReID, the difficulties come from the quick movement of either the drone or the camera, such as quick rotating of pitch or yaw can result in unsatisfying bounding box tracking results. In the case of long-term ReID, which is a significant focus of the SeaDronesSee-MOT benchmark, traditional appearance features, often effective in pedestrian tracking scenarios, may not perform well for object tracking in maritime environments. The challenge arises from the characteristics of maritime objects, including boats, which may share similar visual appearances.

Therefore, the paper argues that leveraging drone metadata, such as altitude, viewing angles, and gimbal information, can help overcome these challenges and significantly enhance the robustness and effectiveness of maritime MOT systems. The integration of the metadata marks a significant step forward in the field of maritime search and rescue missions, enabling more accurate and reliable tracking and ReID of objects over extended time frames. Our proposed Metadata-Guided MOT (MG-MOT) effectively merges short-term tracking data into coherent long-term tracks, harnessing crucial metadata from UAVs, including GPS position, drone altitude, and camera orientations, which achieves first place in the UAV-based Multi-Object

Tracking with Re-Identification Track in the latest edition of Workshop on Maritime Computer Vision (MaCVi).

The paper is organized as follows: we will first introduce some related prior works towards MOT in Sec 2. Then our main proposed method used in the challenge will be described in Sec 3. Sec 4 and 5 will cover the implementation details and the experiments. Finally, we will have the conclusion in Sec 6.

## 2. Related Work

**Multi-Object Tracking.** Multi-object tracking algorithms have been significantly improved by the advancement in deep learning-based detectors. The recent tracking algorithms usually follow the tracking by detection paradigm and utilize an object detector and an association algorithm to conduct tracking. Several popular tracking algorithms include DeepSORT [22], ByteTrack [25], and BoTSORT [2]. These methods usually focus on the tracking of pedestrians and vehicles, in which the target objects and cameras usually demonstrate simple movement. Several popular existing MOT datasets include MOT [19] and BDD [24]. However, with the recent increased popularity of maritime computer vision applications, more and more research starts to focus on the MOT task in maritime environments [13, 23].

**Tracking with Moving Cameras.** Strong camera movement can cause failure in object motion modeling and object detection during multi-object tracking tasks. Camera movements exist in multiple benchmarks [4, 13, 19]. To increase the robustness of tracking and reduce the negative effect of camera motion, BoTSORT [2] uses the global motion compensation (GMC) technique, allowing the tracker to estimate the background motion and thus produce a more accurate object motion prediction for the association progress. StrongSORT [5] incorporates an enhanced correlation coefficient maximization (ECC) model [6] for camera motion compensation and helps the tracker's estimatation of the global rotation as well as the translation of adjacent frames. While some other algorithms leverage different association methods to conduct tracking under large camera motion, e.g., bounding box distance [11], different IoU association method [9]. These methods aim to generate higher spatial similarity of the same object in different time stamps even if they share no IoU in adjacent frames, thus providing more robustness during the tracking process.

**Tracking with Metadata.** Several previous works leverage the information from metadata to conduct multi-object tracking. The metadata might include useful information like the target object's information, camera-related information (e.g., intrinsic and extrinsic parameters), etc. [7, 8] utilizes vehicle metadata and vehicle travel distance to increase the ReID accuracy and multi-camera vehicle tracking performance in urban surveillance scenarios. Huang et al. [10] estimates the camera calibration by selecting corresponding points in the bird-eye-view map and camera frames using the PnP method in [20] to further improves the association between camera-views in an indoor scene. Kiefer et al. [14] show that metadata from UAVs can create a memory map of object locations in actual world coordinates, improving the representation of object locations, which has proven to enhance the downstream tasks, e.g., object detection, multi-object tracking, and video anomaly detection.

# 3. Proposed MG-MOT Method

Due to the nature of the characteristics of small and similar targets and the sudden movement of the drone's gimbal, appearance-based and 2D motion-based ReID or tracking methods can easily fail in this maritime MOT task. Taking advantage of rich metadata from the drone, including latitude, longitude, altitude, pitch and yaw angles, we are able to construct the camera model and build our long-term ReID method based on 3D geometry.

## 3.1. Estimated Camera Model

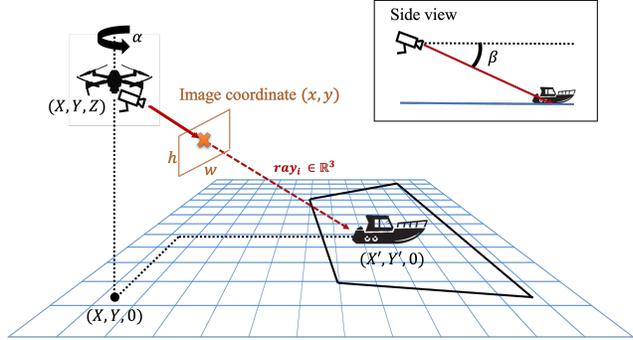More specifically, the camera's intrinsic parameters can be obtained by field-of-view as



Figure 2. Illustration of our estimated camera model used to project the location of the object from image coordinates $(x, y)$ to world coordinates $(X', Y', 0)$.

$$
K = \begin{bmatrix} w/2 \tan(fov/2) & 0 & w/2 \\ 0 & w/2 \tan(fov/2) & h/2 \\ 0 & 0 & 0 \end{bmatrix},
$$
(1)

where $w$ and $h$ are the width and height of images, and $fov$ is the field-of-view of the camera.

As shown in Fig 2, for each frame, with the metadata from the drone, the rotation matrix between the camera coordinate and the world coordinate can be established by:

$$
R(\alpha, \beta, 0) = \begin{bmatrix} \cos\alpha\cos\beta & -\sin\alpha & \cos\alpha\sin\beta \\ \sin\alpha\cos\beta & \cos\alpha & \sin\alpha\sin\beta \\ -\sin\beta & 0 & \cos\beta \end{bmatrix}, \quad (2)
$$

where $\alpha$ and $\beta$ are the yaw and pitch angles of the drone gimbal. Specifically, $\alpha = 90° - gimbal\_heading$ and $\beta = gimbal\_pitch$, where gimbal_heading and gimbal_pitch are given by metadata. Combining $K$ and $R$, any detected object $X_i^{2D}$ and its corresponding image coordinates $x_i$ can be projected to the world coordinates using the 3D directional vector:

$$
ray_i = R \cdot K^{-1} \cdot \frac{x_i}{||x_i||} \in \mathbb{R}^3. \quad (3)
$$

Finally, we use the altitude information $Z_{drone}$ to compute the intersection of our $ray_i$ with the sea surface, which we assume is a plane of $z = 0$:

$$
X_i^{3D} = Loc_{object} = Loc_{drone} + \frac{Z_{drone}}{Z_{ray_i}} \cdot ray_i. \quad (4)
$$

## 3.2. Metadata-Guided Re-Identification

**Short-Term ReID.** For short-term ReID, we focus on efficiently associating the broken tracklets which do not exit the camera view, a scenario that typically occurs within a

shorter time window. This often results from sudden drone movements or rapid gimbal adjustments. Given the limited time frame in such situations, we adopt a direct approach by computing the world coordinates directly. This method allows us to track more effectively and enhance our ability to maintain continuous object tracking even during these dynamic and challenging conditions.

For each detection at frame $t$, the world coordinate $X_t^{3d}$ is computed using the projection $H(\cdot)$ derived from the metadata $m_t$ as mentioned in Sec 3.1:

$$X_t^{3D} = H_{m_t}\big(\hat{X}_t^{2D}\big), \qquad (5)$$

where $\hat{X}^{2D}$ represents the center point of the bounding box $X^{2D}$. For short-term ReID, we directly use Hungarian assignment to match the tracks using the world coordinate distance:

$$\arg\min C_{\text{dis}}(T_i, T_j) = \left\| X_{t_i^{\text{exit}}}^{3D} - X_{t_j^{\text{entry}}}^{3D} \right\|_2 \qquad (6)$$

where $T_i$ represents the tracks left in the memory bank waiting to be associated, while $T_j$ represents the entering or new tracks waiting to be matched. An additional $\tau_{match}$ is used to constraint the matching; if the cost is greater than such threshold, we will treat the entering or new tracks as a new one.

**Long-term Re-ID.** Different from short-term ReID, long-term ReID is somewhat more challenging because we need to closely monitor the potential movement of tracked objects which are out of the image view. Therefore, on top of the standard world coordinate distance matching, we add in two crucial components, i.e. **Bi-directional Movement Extrapolation (BiME)** and **Matching Threshold Expansion (MTE)**. Bi-Directional Movement Extrapolation is a naive method that extrapolates the world coordinates of the tracks that exit or enter the image. Given some tracks $T_i = \{X_{t_i^{enter}}, \dots, X_{t_i^{exit}}\}$ that enter and exit the image during frame $t^{enter}$ and $t^{exit}$. The world coordinate of $T_i$ at $t'$ will be forward-extrapolating as:

$$T_i^{3D} = X_{t_i^{exit}}^{3D} + \Delta t \cdot V_{T_i}^{exit}, \quad \text{if } \Delta t < \tau_{memory}, \qquad (7)$$

where $\Delta t = t' - t^{exit}$ and the track last-seen exit velocity can be estimated as $V_{T_i}$:

$$V_{T_i}^{exit} = \left( \frac{X_{t_i^{exit}}^{3D} - X_{t_i^{exit}-w}^{3D}}{w} \right). \qquad (8)$$

using the window size $w$ as a constant. Note that we will not extrapolate further beyond a certain duration $\tau_{memory}$ since the world coordinates after such duration are often unreliable and may be confused with some new tracks. We also backward-extrapolating the tracks using similar logistics in Eq. 7 and 8 by substituting the exit terms with the entry

| Dataset | # of Seq. | # of Frame | # of Bbox | Longest Seq. |
|---|---|---|---|---|
| Training | 20 | 27,259 | 160,470 | 6,296 |
| Validation | 17 | 8,584 | 47,678 | 2,069 |
| Testing | 19 | 18,253 | - | 4,138 |
| Total | 21 | 54,096 | 207,938 | - |

Table 1. The overall statistics of the SeaDroneSee-MOT dataset.

| Metadata | | Data Split (Mean/Min/Max) | |
|---|---|---|---|
| | | Train+Val | Test |
| **P** | gps_latitude | 0.072/0.070/0.074 | 0.072/0.071/0.074 |
| | gps_longitude | 0.069/0.067/0.071 | 0.069/0.067/0.071 |
| | altitude | 34.29/4.90/140.39 | 42.55/4.60/149.49 |
| **D** | gimbal_pitch | 42.5/-2.5/90.0 | 57.333/6.8/90.0 |
| | gimbal_heading | 201.6/0.1/359.8 | 184.437/0/359.3 |
| **S** | x_speed | 1.24/0/13.799 | 1.331/0/10.799 |
| | y_speed | 0.953/0/12.799 | 1.046/0/9.399 |
| | z_speed | 0.129/0/3.999 | 0.251/0/5.199 |

Table 2. The frame-level drone metadata provided in the SeaDroneSee-MOT dataset and its corresponding statistics according to the train/val/test splits. **P** stands for positions (latitude and longitude are reported with a relative center at (N47.6° E9.2°)), **D** stands for directions, and **S** stands for absolute speed.

terms. In contrast to short-term re-identification, Matching Threshold Expansion is a strategy to expand the matching space as the tracks disappear or reappear from the image:

$$\tau'_{match} = \lambda \cdot (\Delta t_{exit} + \Delta t_{entry}) \cdot \tau_{match}. \qquad (9)$$

**Class-wise Re-ID.** Recognizing the distinct characteristics of boat and swimmer movement patterns and appearances, we conduct short-term and long-term Re-identification (Re-ID) in a class-specific manner. We tailor the thresholds mentioned earlier for each class, such as setting a larger $\tau_{match}$ for boats due to their potential for higher relative velocity. The association steps are performed individually, enabling us to significantly lower the chances of associating tracks from different classes.

## 4. Dataset: SeaDroneSee-MOT

The SeaDronesSee-MOT dataset consists of 21 clips in the training set, 17 clips in the validation set, and 19 clips in the testing set with a total of 54,105 frames and 403,192 annotated instances as in Table 1. Metadata for the drone are also provided for all training, validation, and testing split as we summarized it in Table 2. Detail studies and analyses on the characteristics of each sequence (e.g., challenging scenario, per sequence HOTA performance, distribution of metadata) are being carefully investigated and reported in [13]. Note that the training and validation sets do not contain long-term tracking labels, i.e., objects that have gone missing are assigned new IDs when they reappear.

| Team | HOTA↑ | MOTA↑ | IDF1↑ | MOTP↓ | MT↑ | ML↓ | FP↓ | FN↓ | Rec↑ | Pre↑ | IDs↓ | Frag↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Baseline (MaCVi) | 37.5 | 43.8 | 38.8 | 22.9 | 69 | 75 | 13566 | 38858 | 59.4 | 80.7 | 1340 | 2556 |
| Team 400 | 43.4 | 53.2 | 46.9 | 22.4 | 86 | 64 | 10587 | 33075 | 65.4 | 85.5 | 1110 | 2233 |
| NCKU ACVLab (Team 403) | 49.9 | 32.0 | 57.9 | -1.000 | 86 | 97 | 23601 | 41253 | 56.9 | 69.7 | 161 | 646 |
| MI-SIT (Team 399) | 53.1 | 62.5 | 58.8 | **19.8** | 116 | 54 | **9068** | 26650 | 72.1 | 88.4 | 123 | _896_ |
| Team 412 | 55.4 | 65.3 | 61.5 | 20.9 | 134 | 31 | 13603 | 19426 | 79.7 | 84.9 | 137 | 1162 |
| Lenovo (Team 395) | 61.5 | 76.8 | 70.4 | 20.6 | 145 | 31 | 10155 | 11982 | 87.5 | _89.2_ | 100 | 1350 |
| Franunhofer IOSB (Team 220) | _69.3_ | **78.0** | _84.4_ | _20.5_ | **165** | **20** | 10643 | **10391** | 89.1 | 88.9 | **16** | 984 |
| Ours (Team 198) | **69.5** | _77.9_ | **85.9** | 20.7 | _158_ | _22_ | 9700 | _11425_ | 88.1 | **89.7** | _18_ | _784_ |

Table 3. Leaderboard Results of the 2024 MacVi SeaDronesSee Multi-Object Tracking with Re-Identification. The **bold** numbers indicate the best while the underlined numbers indicate the second to the best among all participants. Our proposed method achieved the SOTA in terms of HOTA and IDF1 on the testing split with a noticeable amount of FP, ID switches, and track fragments being reduced.

| Method | Ensemble | Re-ID | HOTA↑ | MOTA↑ | IDF1↑ | MOTP↓ | FP↓ | FN↓ | Rec↑ | Pre↑ | IDs↓ | Frag↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Baseline | | | 60.3 | 77.6 | 68.2 | 20.9 | 9932 | 11381 | 88.1 | 89.5 | 84 | 753 |
| (BoT-SORT) | ✓ | | 61.2 | 77.7 | 69.9 | 20.7 | 9956 | 11338 | 88.1 | 89.4 | 78 | 769 |
| | | Short | 64.0 | 76.8 | 76.6 | 20.8 | 10543 | 11621 | 87.9 | 88.9 | 84 | 758 |
| Ours | | Long | 69.0 | 77.2 | 85.4 | 20.7 | 10082 | 11674 | 87.8 | 89.3 | 26 | 809 |
| (MG-MOT) | ✓ | Long | 69.4 | **77.9** | **85.9** | 20.7 | 9626 | 11476 | 88.0 | 89.7 | 20 | 787 |
| | ✓ | Short+Long | **69.5** | **77.9** | 85.0 | 20.7 | 9700 | 11425 | 88.1 | 89.7 | 18 | 784 |

Table 4. Our implemented method using different types of detector backbones and different Re-ID method combinations.

## 5. Experiment Results

### 5.1. Implementation Details

**Detector.** We use the YOLOv8-x [12], pretrained weights on COCO [15] and with an additional p2 head, to predict tiny object. We also find the need to train a multi-class detector instead of a single-class one since the characteristics of the boat and swimmer classes are totally different in terms of movement, appearance, and others. We experiment with several settings of input image size along with the input data source (e.g., jpg or png) and record the detection results. All detectors are trained using similar hyperparameters for 100 epochs with an initial learning rate of 0.01 and a decay of 0.05. Eventually, we settle for the image size of 1280. Our detector performs poorly when detecting swimmers while the drone flew low to the sea surface with a considerable heading angle. This is due to a lack of similar data in the training dataset. Therefore, we adopt a pre-trained YOLOv8-x detector to perform an ensemble of detections for the sake in the challenge. The yolov8-p2 detectors are trained and inferenced using Nvidia Tesla V100 GPUs.

**Tracking.** Before applying the proposed meta-guided ReID method, we obtained the initial tracking results using BoT-SORT [2] with sparse optical flow as the Generalized Motion Compensation (GMC) method. We use the same tracking hyper-parameters throughout all 18 testing sequences for a generic tracking method. The thresholding for the high confidence detection is set to 0.5, while the low confidence detection is set to 0.1. We initialize a new track with a confidence higher than 0.2 if the detection does not match any existing track and use a buffer of 100 frames to remove those unmatched tracks.

**Camera Parameters Calibration.** We manually pick several segments of sequences from training and validation, with non-moving targeted objects and a continuously moving or turning camera. Then we try to minimize the standard deviation of the locations of the same target with different field-of-view angles. Finally, we select a field-of-view of $70 \deg$. The intrinsic matrix is obtained using such field-of-view angle and image size being $3840 \times 2160$ to obtain the world coordinates of each object from image coordinates.

**Post-Processing.** We employ a simple linear interpolation strategy to recover missing detections, whether caused by poor lighting conditions or sudden camera movements. We have two iterations of interpolations, one before and one after the ReID step. Additionally, in some sequences with lower altitudes, we may encounter overlapping detections, leading to a degradation in detection performance. To address this issue, we apply non-maximum suppression to filter out overlapping detections. It is important to note that, unlike traditional non-maximum-suppression, we prioritize detections with larger bounding box areas while filtering out smaller boxes, as in most situations, the latter are false positive detections.

### 5.2. Results on SeaDroneSee-MOT

**Evaluation Metrics.** The results of long-term multi-object tracking are evaluated using HOTA [16] and CLEAR MOT metrics [3] including MOTA, IDF1, MOTP, MT, ML, FP, FN, Recall, Precision, ID Switches, Fragments. Note that despite the multi-classes labels of the tracks being available in the training and validation data, the testing only considers all tracks as the same classes during evaluation.
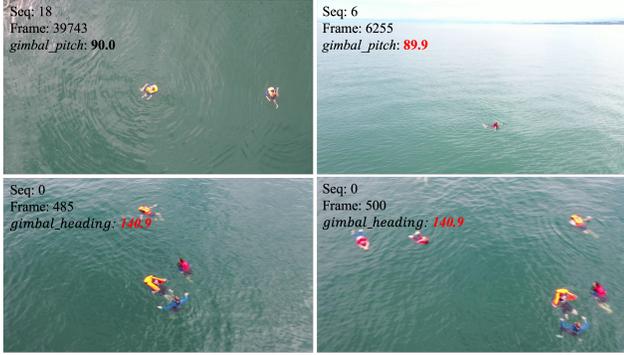
Figure 3. Two types of error in our final results. Top row: Example of gimbal pitch error in the metadata found in the testing split. The two frames have nearly identical pitch, but the drone's directions are clearly different. Bottom row: Example of the metadata is not synchronized with the actual drone gimbal.



Figure 4. An example of splitted tracks, in the earlier frame, the swimmers are wearing lifejackets, however, in the later frame, after the lifejackets are taken off, new track ids emerges.

**Final Results.** In the results section of the competition, as illustrated in Table 3, the performance of various methods and teams in the 2024 MacVi SeaDronesSee Multi-Object Tracking with Reidentification Challenge is rigorously evaluated across several critical metrics (see https://macvi.org/leaderboard/airborne/seadronessee/multi-object-tracking-reid). Notably, our proposed MG-MOT (Team 198) demonstrates exceptional tracking accuracy, achieving the highest HOTA score of 69.5%, while also securing the highest IDF1 score, an impressive 85.9%. This method's performance showcases the ability to track objects in challenging scenarios accurately. Furthermore, Franunfofer IOSB (Team 220) emerges as a noteworthy contender, securing the highest MOTA at 78.0%, highlighting their prowess in the detection results as they also have the fewest FN detection. These results underline the significant advancements in maritime MOT capabilities and affirm the strong potential of our proposed MG-MOT method for real-world applications in object tracking.

## 5.3. Ablation Studies

**Re-ID Modules.** The effectiveness of our MG-MOT components, short- and long-term re-identification are being analyzes in Table 4. Our MG-MOT is compared to the base-

| Backbone | Source | Input Size | MOTA↑ | FP↓ | FN↓ |
|---|---|---|---|---|---|
| yolov8x | jpg | 640 | 64.3 | 8998 | 25063 |
| yolov8x-p2 | jpg | 640 | 68.0 | **8010** | 22448 |
| yolov8x-p2 | jpg | 960 | 70.4 | 8971 | 19249 |
| yolov8x-p2 | png | 960 | 76.0 | 10052 | 12817 |
| yolov8x-p2 | png | 1280 | **77.3** | 9701 | **11930** |

Table 5. Ablation studies on the detector performance.

line using the BoT-SORT, which also serves as our initial tracking results (in Sec 5.1). We found that both the short- and long-term re-identification strategies can improve the IDF1 of the tracking therefore boosting the HOTA performance based on the same detections.

**Detection.** We report the detection result using different YOLOv8 backbones, input sources, and input sizes in Table 5 with the same tracking hyper-parameters. The performance is obtained by submitting the results to the evaluation server. Due to a significant number of tiny detections in the dataset, the $p2$ head combined with a greater input size of the backbone model resulted in better MOTA, FP, and FN.

**Tracking Parameters.** We conduct ablation studies on various tracking parameters' impact in the inference stage. To assess performance on annotated data, a new model is trained based on the same architecture as the previously reported best one, using the training set and evaluating on the validation set. Notably, interpolation and ReID are not adopted. The effect of the IoU threshold on tracking performance, tuning it from 0.1 to 0.9. Results in Table 6 show that setting the IoU threshold to 0.5 achieves a more balanced performance. Detecting swimmers prove to be more challenging than detecting boats due to small bounding boxes and interactions, resulting in box overlaps. The impact of the matching threshold by varying it from 0.1 to 0.9. Results illustrate that a higher matching threshold improves performance. In videos with multiple swimmers, they tend to gather as a group and follow similar trajectories, making a higher matching threshold valuable in preventing incorrect matches.

## 5.4. Limitations

We observe certain limitations in our proposed method and the standard algorithm for UAV-based MOT. One notable challenge is ensuring the synchronization and accuracy of metadata with the captured frames. However, this issue could potentially be resolved by cross-referencing with the presence of the sea surface in the images. Moreover, when objects temporarily vanish from the camera frame due to drone movement or camera repositioning, there is a risk of long-term ReID failure, as these objects may reappear after an unusually prolonged interval. For this, we can employ an expansion of the search and rematch zone or simply readopt the appearance-based ReID for further filtering of possible matches.
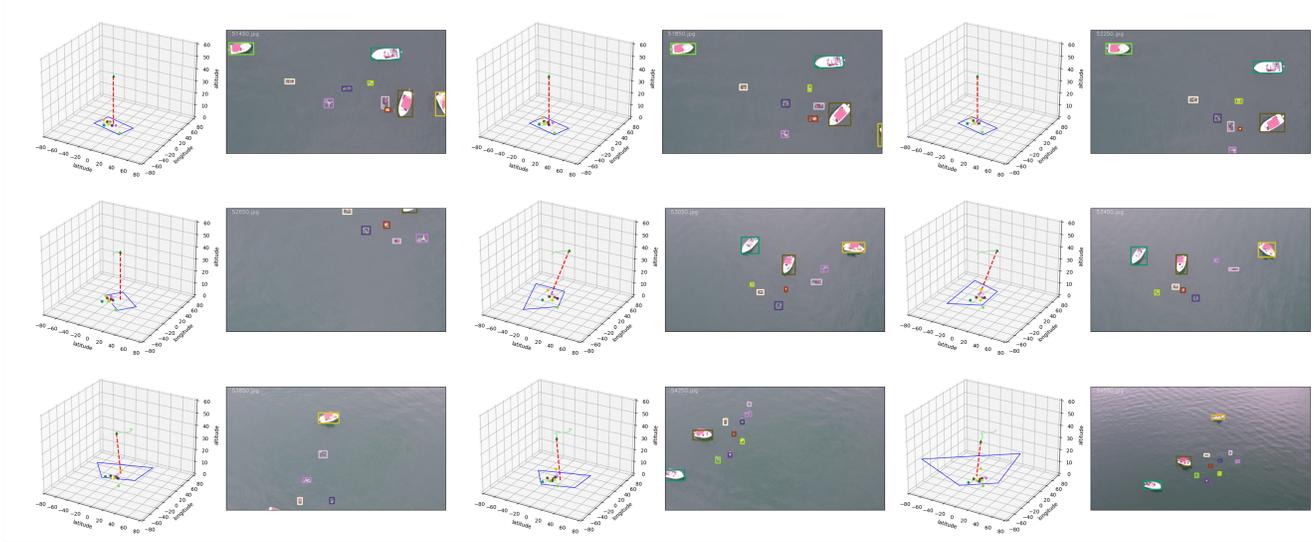
Figure 5. Visualization of the result of MG-MOT on the testing sequence 21 with world coordinates (best view in color).

| Tracking Parameter | Metric | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 |
|---|---|---|---|---|---|---|---|---|---|---|
| IoU Threshold | HOTA↑ | 69.4 | 69.4 | 69.6 | 69.7 | 69.9 | 69.6 | 69.6 | 69.1 | 67.3 |
| | MOTA↑ | 80.3 | 80.6 | 81.4 | 81.3 | 81.3 | 81.4 | 81.4 | 81.1 | 80.5 |
| | IDF1↑ | 80.8 | 80.7 | 80.6 | 80.5 | 81.0 | 80.6 | 80.6 | 79.2 | 75.7 |
| | FP↓ | 7078 | 7262 | 7732 | 7575 | 7340 | 7732 | 7732 | 8225 | 9998 |
| | FN↓ | 10573 | 10526 | 10282 | 10467 | 10225 | 10282 | 10282 | 11005 | 12529 |
| Matching Threshold | HOTA↑ | 33.6 | 47.5 | 69.6 | 62.4 | 66.1 | 69.6 | 69.6 | 69.6 | 69.9 |
| | MOTA↑ | 21.4 | 55.0 | 81.4 | 77.9 | 80.3 | 81.4 | 81.4 | 81.4 | 81.4 |
| | IDF1↑ | 27.6 | 41.2 | 80.6 | 66.9 | 72.2 | 80.6 | 80.6 | 80.6 | 81.4 |
| | FP↓ | 7917 | 7917 | 7732 | 7917 | 7917 | 7732 | 7732 | 6769 | 4953 |
| | FN↓ | 38711 | 31167 | 10282 | 17455 | 14575 | 10282 | 10282 | 10282 | 9894 |

Table 6. Ablation studies on IoU threshold and matching threshold.

**Error in Metadata.** We find the metadata of the drone is not that reliable. As shown in Fig 3, although the gimbal pitch angles are almost the same, the image from Seq 6 shows an entirely different view, which is more similar to the pitch angle around $45°$.

**Error in Splitted Tracks.** In the testing set, there are some cases when a swimmer takes off his/her life jacket during the tracking process, which causes a "split tracks" situation to happen. Since our tracking algorithm does not consider object class during the association stage, this might result in ID switch during the tracking process and lead to a degradation in tracking performance.

## 6. Conclusion

In summary, our motion-based multi-object tracking algorithm, MG-MOT, enhanced by the inclusion of UAV metadata, represents a significant performance stride in the domain of maritime computer vision. We've made substan-

tial progress by addressing the challenges of short-term and long-term ReID with the help of UAV metadata. The effectiveness of our algorithm, as demonstrated through the SeaDroneSee dataset and the UAV-based Maritime Multi-Object Tracking Challenge, is a testament to its practicality. We achieve a much-improved performance in the latest edition of the UAV-based Maritime Object Tracking Challenge with a state-of-the-art HOTA of 69.5% and an IDF1 of 85.9% on the testing split.

## 7. Acknowledgment

## References

[1] Bharat Sharma Acharya, Mahendra Bhandari, Filippo Bandini, Alonso Pizarro, Matthew Perks, Deepak Raj Joshi,

Sheng Wang, Toby Dogwiler, Ram L Ray, Gehendra Kharel, et al. Unmanned aerial vehicles in hydrology and water management: Applications, challenges, and perspectives. *Water Resources Research*, 57(11):e2021WR029925, 2021. 1

[2] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Botsort: Robust associations multi-pedestrian tracking. *arXiv preprint arXiv:2206.14651*, 2022. 2, 3, 5

[3] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: The clear mot metrics. *J. Image Video Process.*, 2008, jan 2008. 5

[4] Yutao Cui, Chenkai Zeng, Xiaoyu Zhao, Yichun Yang, Gangshan Wu, and Limin Wang. Sportsmot: A large multi-object tracking dataset in multiple sports scenes. *arXiv preprint arXiv:2304.05170*, 2023. 3

[5] Yunhao Du, Zhicheng Zhao, Yang Song, Yanyun Zhao, Fei Su, Tao Gong, and Hongying Meng. Strongsort: Make deepsort great again. *IEEE Transactions on Multimedia*, 2023. 3

[6] Georgios D Evangelidis and Emmanouil Z Psarakis. Parametric image alignment using enhanced correlation coefficient maximization. *IEEE transactions on pattern analysis and machine intelligence*, 30(10):1858–1865, 2008. 3

[7] Hung-Min Hsu, Jiarui Cai, Yizhou Wang, Jenq-Neng Hwang, and Kwang-Ju Kim. Multi-target multi-camera tracking of vehicles using metadata-aided re-id and trajectory-based camera link model. *IEEE Transactions on Image Processing*, 30:5198–5210, 2021. 3

[8] Hsiang-Wei Huang and Jenq-Neng Hwang. Multi-target multi-camera vehicle tracking using transformer-based camera link model and spatial-temporal information. *arXiv preprint arXiv:2301.07805*, 2023. 3

[9] Hsiang-Wei Huang, Cheng-Yen Yang, Jenq-Neng Hwang, and Chung-I Huang. Iterative scale-up expansioniou and deep features association for multi-object tracking in sports. *arXiv preprint arXiv:2306.13074*, 2023. 3

[10] Hsiang-Wei Huang, Cheng-Yen Yang, Zhongyu Jiang, Pyong-Kun Kim, Kyoungoh Lee, Kwangju Kim, Samartha Ramkumar, Chaitanya Mullapudi, In-Su Jang, Chung-I Huang, and Jenq-Neng Hwang. Enhancing multi-camera people tracking with anchor-guided clustering and spatio-temporal consistency id re-assignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2023. 3

[11] Hsiang-Wei Huang, Cheng-Yen Yang, Samartha Ramkumar, Chung-I Huang, Jenq-Neng Hwang, Pyong-Kun Kim, Kyoungoh Lee, and Kwangju Kim. Observation centric and central distance recovery for athlete tracking. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 454–460, 2023. 3

[12] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics yolov8, 2023. 5

[13] Benjamin Kiefer, Matej Kristan, Janez Perš, Lojze Žust, Fabio Poiesi, Fabio Augusto de Alcantara Andrade, Alexandre Bernardino, Matthew Dawkins, Jenni Raitoharju, Yitong Quan, et al. 1st workshop on maritime computer vision (macvi) 2023: Challenge results. *arXiv preprint arXiv:2211.13508*, 2022. 1, 2, 3, 4

[14] Benjamin Kiefer, Yitong Quan, and Andreas Zell. Memory maps for video object detection and tracking on uavs. *arXiv preprint arXiv:2303.03508*, 2023. 3

[15] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. 5

[16] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. Hota: A higher order metric for evaluating multi-object tracking. *International journal of computer vision*, 129:548–578, 2021. 5

[17] Eleftherios Lygouras, Nicholas Santavas, Anastasios Taitzoglou, Konstantinos Tarchanidis, Athanasios Mitropoulos, and Antonios Gasteratos. Unsupervised human detection with an embedded vision system on a fully autonomous uav for search and rescue operations. *Sensors*, 19(16):3542, 2019. 1

[18] Mario Monteiro Marques, Pedro Dias, Nuno Pessanha Santos, Vitor Lobo, Ricardo Batista, D Salgueiro, A Aguiar, M Costa, J Estrela Da Silva, A Sérgio Ferreira, et al. Unmanned aircraft systems in maritime operations: Challenges addressed in the scope of the seagull project. In *OCEANS 2015-Genova*, pages 1–6. IEEE, 2015. 1

[19] A. Milan, L. Leal-Taixé, I. Reid, S. Roth, and K. Schindler. MOT16: A benchmark for multi-object tracking. *arXiv:1603.00831 [cs]*, Mar. 2016. arXiv: 1603.00831. 2, 3

[20] Zheng Tang, Gaoang Wang, Hao Xiao, Aotian Zheng, and Jenq-Neng Hwang. Single-camera and inter-camera vehicle tracking and 3D speed estimation based on fusion of visual and semantic features. In *Proc. CVPR Workshops*, pages 108–115, 2018. 3

[21] Leon Amadeus Varga, Benjamin Kiefer, Martin Messmer, and Andreas Zell. Seadronessee: A maritime benchmark for detecting humans in open water. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2260–2270, 2022. 1

[22] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *2017 IEEE international conference on image processing (ICIP)*, pages 3645–3649. IEEE, 2017. 2

[23] Cheng-Yen Yang, Alan Yu Shyang Tan, Melanie J Underwood, Charlotte Bodie, Zhongyu Jiang, Steve George, Karl Warr, Jenq-Neng Hwang, and Emma Jones. Multi-object tracking by iteratively associating detections with uniform appearance for trawl-based fishing bycatch monitoring. *arXiv preprint arXiv:2304.04816*, 2023. 2

[24] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020. 2

[25] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. Bytetrack: Multi-object tracking by associating every detection box. 2022. 2