

Integrating Continual Learning in Object Tracking Models:
M.Tech Project Report submitted to IIT TIRUPATI

*submitted in partial fulfillment of the requirements
for the degree of*

DUAL DEGREE

in

COMPUTER SCIENCE AND ENGINEERING

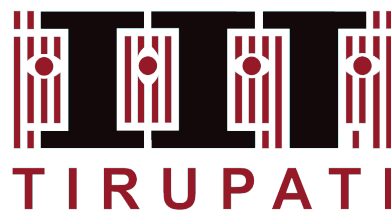
by

ARUP BISWAS CS24M124

Supervisor(s)

Dr. Chalavadi Vishnu

भारतीय प्रौद्योगिकी संस्थान तिरुपति



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
INDIAN INSTITUTE OF TECHNOLOGY TIRUPATI**

OCTOBER 2025

DECLARATION

I declare that this written submission represents my ideas in my own words and where others' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in my submission to the best of my knowledge. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

Place: Tirupati
Date: 29-11-2025

Signature
ARUP BISWAS
CS24M124

BONA FIDE CERTIFICATE

This is to certify that the report titled **Integrating Continual Learning in Object Tracking Models: M.Tech Project Report**, submitted by **ARUP BISWAS**, to the Indian Institute of Technology, Tirupati, for the award of the degree of **DUAL DEGREE**, is a bona fide record of the project work done by him under my supervision. The contents of this report, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

Place: Tirupati
Date: 29-11-2025

Dr. Chalavadi Vishnu
Guide
Assistant Professor
Department of Computer
Science and Engineering
IIT Tirupati - 517619

ACKNOWLEDGMENTS

Thanks to my friends, peers, and parents for supporting and guiding me.

ABSTRACT

KEYWORDS: Object Tracking, Continual Learning, Catastrophic Forgetting, Drones, UAVs, HybridSORT, Literature Survey, Domain Adaptation, Re-Identification, Long-Term Tracking.

Object tracking from aerial platforms like drones is crucial for applications such as surveillance, search and rescue, and autonomous navigation. However, standard object trackers often suffer from catastrophic forgetting [11, 13, 10], losing track of objects after temporary occlusions or if they leave and re-enter the frame. This difficulty is especially obvious in long-term settings like drone footage. This project examines the intersection of Object Tracking (OT) and Continual Learning (CL) [11, 13], with a focus on recent advancements (2024 onwards) [14, 11, 13, 10, 3] to address forgetting in drone-based tracking. We review SOTA object tracking algorithms, particularly those based on the tracking-by-detection paradigm [14]. We also explore contemporary CL strategies designed to mitigate forgetting, such as replay methods and regularization techniques [11, 13, 10]. We analyze drone-specific datasets and benchmarks that address aerial perspective difficulties such as size variation, motion blur, and haze [3]. The survey analyzes how CL methods, particularly exemplar replay [11, 10], could be integrated into SOTA trackers like HybridSORT (AAAI 2024) [14] to create a more robust system ("SORT-CL") capable of maintaining object identities over extended periods in dynamic drone video streams. We outline a methodology for adapting such a tracker to the drone domain by re-training its detector and Re-ID components and integrating a CL memory module. Project being worked on [here](#).

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	i
ABSTRACT	ii
LIST OF FIGURES	v
LIST OF TABLES	vi
ABBREVIATIONS	vii
1 INTRODUCTION	1
1.1 Background	1
1.2 Problem Statement	1
1.3 Objectives and Scope	2
2 Literature Review	4
2.1 Continual Learning and Forgetting	4
2.1.1 CL Scenarios	4
2.1.2 Major CL Approaches	6
2.2 Object Tracking	8
2.2.1 The SORT Family Evolution	8
2.2.2 Recent Advancements (2024+)	9
2.3 Object Tracking and Continual Learning	10
2.4 Latest Drone Datasets and Benchmarks	11
2.4.1 Established Drone Datasets	11
2.4.2 Recent Drone Datasets/Papers (2024+)	12
3 PROPOSED FUTURE METHODOLOGY	14
3.1 Baseline Tracker Selection	14
3.2 Continual Learning Integration: CRM Module	14
3.3 Domain Adaptation for Drones	17

4	RESULTS AND DISCUSSION	18
4.1	Experimental Setup and Implementation	18
4.1.1	Deployment Challenges & Solutions	18
4.2	Phase 1 Results: Baseline Validation	19
4.2.1	Qualitative Results	19
4.3	Proposed Evaluation Protocol for Continual Learning	20
4.3.1	Metrics for Long-Term Tracking	20
4.3.2	Ablation Study Plan	20
4.4	Discussion	21

LIST OF FIGURES

2.1	Conceptual framework of continual learning [11]. (a), Continual learning requires adapting to incremental tasks with dynamic data distributions. (b), A desirable solution should ensure an appropriate trade-off between stability (red arrow) and plasticity (green arrow), as well as an adequate generalizability to intra-task (blue arrow) and inter-task (orange arrow) distribution differences. (c), To achieve the objective of continual learning, representative methods have targeted various aspects of machine learning. (d), Continual learning is adapted to practical applications to address particular challenges such as scenario complexity and task specificity.	5
2.2	Taxonomy of representative continual learning methods. We have summarized five main categories (blue blocks), each of which is further divided into several sub-directions (red blocks) [11].	7
2.3	Pipeline of Hybrid-SORT and Hybrid-SORT-ReID. For strong cues, we utilize IoU as the metric for spatial information, and utilize cosine distance for appearance features. For weak cues, we incorporate the confidence state, height state, and velocity direction. Velocity direction is illustrated by centers instead of corners for better clarity [14].	9

LIST OF TABLES

2.1	Comparison of typical continual learning scenarios [11].	6
2.2	Comparative landscape of drone-view and adverse weather datasets [3]. . .	11

ABBREVIATIONS

AAAI	Association for the Advancement of Artificial Intelligence
AUC	Area Under the Curve
CL	Continual Learning
CRM	Continual Re-ID Memory
CVPR	Conference on Computer Vision and Pattern Recognition
DET	Detection
ECCV	European Conference on Computer Vision
FPV	First-Person View
HMIoU	Height Modulated Intersection over Union
HOTA	Higher Order Tracking Accuracy
ICCV	International Conference on Computer Vision
IDSW	Identity Switches
IoU	Intersection over Union
KF	Kalman Filter
MOT	Multi-Object Tracking
MOTA	Multiple Object Tracking Accuracy
MTech	Master of Technology
OT	Object Tracking
Re-ID	Re-Identification
ROCM	Robust Observation-Centric Momentum
SOTA	State-of-the-Art
SOT	Single-Object Tracking
SORT	Simple Online and Realtime Tracking
TBD	Tracking-by-Detection
TCM	Tracklet Confidence Modeling
UAV	Unmanned Aerial Vehicle

NOTATION

A	global atmospheric light
$a_{k,j}$	classification accuracy on the test set of the j -th task after learning the k -th task
$d(x,y)$	scene depth
D_t	training samples belonging to task t
$IDSW$	number of Identity Switches
$I(x,y)$	observed hazy image
$J(x,y)$	clear scene radiance
$\mathcal{L}(\theta, X)$	loss or performance metric of a model with parameters θ on dataset X
θ	parameters of a continual learning model

CHAPTER 1

INTRODUCTION

1.1 Background

Object tracking, or recognizing and following instances of things across video frames [14], is a crucial feature in computer vision. It has a wide range of applications, including autonomous driving, robotics, surveillance, and human-computer interaction. Unmanned Aerial Vehicles (UAVs), sometimes known as drones, have revolutionized object tracking [3]. Drones provide exceptional flexibility in data collecting, allowing aerial viewpoints for jobs including large-area surveillance, traffic monitoring, agricultural management, infrastructure inspection, and search and rescue operations [3]. Effective object tracking from mobile aerial platforms is critical for interpreting dynamic environments and enabling autonomous drone behavior [3, 2, 8].

A dominant paradigm for multi-object tracking is Tracking-by-Detection (TBD) [14]. In this approach, an object detector first identifies all potential objects within each frame [14]. Subsequently, an association mechanism links these detections across frames to form trajectories, typically using cues like spatial proximity (e.g., Intersection over Union - IoU) and appearance similarity [14].

1.2 Problem Statement

Despite significant progress, object tracking from drones presents unique and substantial challenges compared to traditional ground-level scenarios [3].

- **Scale Variation:** The flying altitude of a drone can change rapidly, causing objects to appear at vastly different scales within and across frames [2, 8].
- **Fast Motion & Motion Blur:** Drones, especially agile platforms like FPV (First-Person View) drones, can move quickly and change direction abruptly, leading to significant motion blur and challenging inter-frame association [2, 8].

- **Occlusion:** Objects can be frequently hidden by environmental structures (buildings, trees) or other moving objects [14], requiring robust re-identification capabilities.
- **Top-Down Viewpoint:** The aerial perspective means objects are often viewed from above, presenting different visual features (e.g., rooftops of vehicles, tops of heads) compared to the side profiles common in ground-level datasets.
- **Adverse Weather:** Operating outdoors, drones are susceptible to various weather conditions, including rain, fog, and haze, which severely degrade image quality and obscure object features [3]. The HazyDet paper highlights the significant impact of haze and the absence of specialized drone benchmarks for this condition in drone imagery [3].

Long-Term Tracking is an essential need for many drone applications, as the system must keep an object's identity even if it is obscured for an extended period of time or leaves and then returns to the camera's field of vision. Standard object trackers, however, frequently suffer from catastrophic forgetting [11, 13, 10]. This phenomenon occurs when a neural network learns something new and soon forgets what it previously learnt [10]. When an object emerges after a long period of absence, the tracker fails to correctly identify it, resulting in identity switches (IDSW) or track fragmentation. This amnesia phenomenon is a key impediment to accurate long-term tracking, especially in the dynamic and demanding context of drone video streams.

1.3 Objectives and Scope

This project aims to address the challenge of catastrophic forgetting in long-term drone object tracking through the lens of Continual Learning (CL).

- To conduct a comprehensive literature survey on state-of-the-art (SOTA) methodologies (2024 onwards) in both Object Tracking (OT) and Continual Learning (CL), identifying key techniques relevant to mitigating forgetting [11, 13, 10].
- To study the specific problems and current benchmarks for drone-based object tracking, especially focusing on long-term situations and inclement weather [3, 2, 8].
- To investigate how CL techniques can be applied to mitigate catastrophic forgetting in long-term drone object tracking [11, 13, 10].
- To propose a conceptual framework (e.g., "HybridSORT-CL for Drones") integrating a

chosen SOTA tracker (HybridSORT) [14] with a suitable CL strategy (Exemplar Replay) [11, 10] and outline the necessary domain adaptation steps (detector and Re-ID model re-training) for effective application to drone footage.

The scope of this project is primarily a literature survey and methodological proposal.

CHAPTER 2

Literature Review

This chapter shall present a critical appraisal of the previous work published in the literature pertaining to the topic of the investigation.

2.1 Continual Learning and Forgetting

Continual Learning (CL), also referred to as incremental learning or lifelong learning [11], addresses the challenge of training models on a sequence of tasks or a non-stationary stream of data [10]. This is done such that knowledge acquired from earlier data is retained while learning new information [11, 10]. This contrasts sharply with standard deep learning practices that typically assume a fixed, static dataset trained offline [11]. Figure 2.1 illustrates the conceptual framework of continual learning, contrasting it with standard training.

The main challenge in CL is Catastrophic Forgetting [11, 13, 10]. According to [11, 10], when a model learns a new task, its performance on previously learned tasks significantly decreases. This arises because optimizing for the current task frequently requires updating parameters necessary for older tasks, effectively overwriting previous information [10]. This requires balancing stability (retaining old information) and plasticity (acquiring new knowledge), also known as the Stability-Plasticity Dilemma [11, 10].

2.1.1 CL Scenarios

Depending on task boundaries and information availability, CL problems are often categorized [11, 10] (see Table 2.1):

- **Task-Incremental Learning (TIL):** Task identities are provided during both training and testing [11, 10].
- **Domain-Incremental Learning (DIL):** Tasks share the same data label space but have different input distributions; task ID is not available at test time [11, 10].

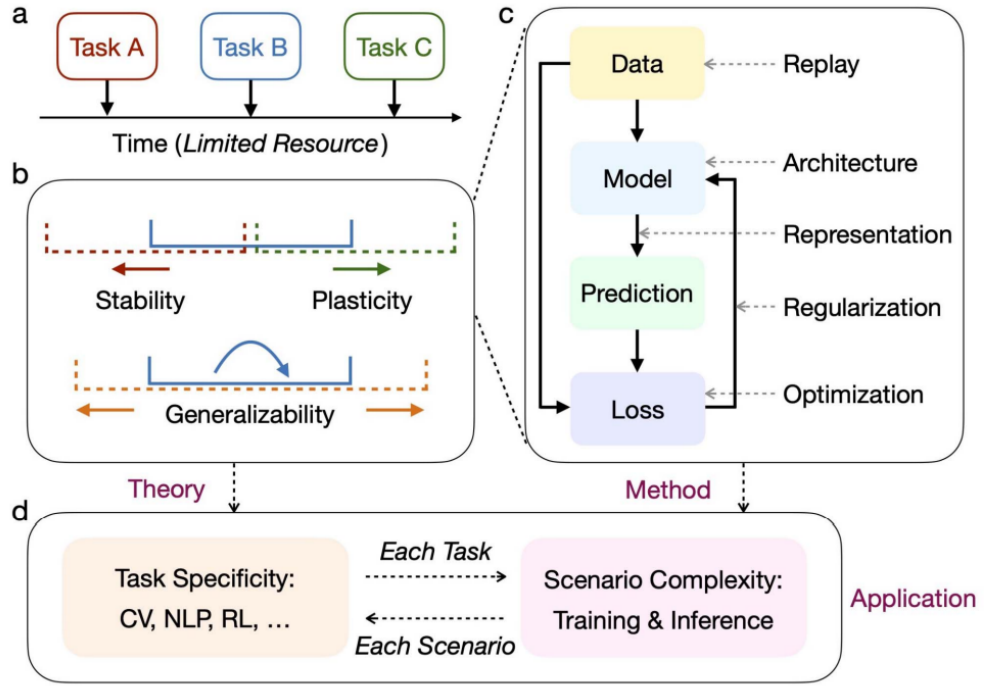


Figure 2.1: Conceptual framework of continual learning [11]. (a), Continual learning requires adapting to incremental tasks with dynamic data distributions. (b), A desirable solution should ensure an appropriate trade-off between stability (red arrow) and plasticity (green arrow), as well as an adequate generalizability to intra-task (blue arrow) and inter-task (orange arrow) distribution differences. (c), To achieve the objective of continual learning, representative methods have targeted various aspects of machine learning. (d), Continual learning is adapted to practical applications to address particular challenges such as scenario complexity and task specificity.

- **Class-Incremental Learning (CIL):** Tasks introduce new classes, and the model must distinguish between all classes seen so far, often without task ID at test time [11, 10].
- **Task-Free Continual Learning (TFCL):** No explicit task boundaries or IDs are provided during training or testing [11, 10]. Online CL (OCL) is a related setting where data arrives as a single-pass data stream [11, 13].

Table 2.1: Comparison of typical continual learning scenarios [11].

Scenario	Training	Testing
IIIL	$\{D_{t,b}\}_{b \in \mathcal{B}_t, t \in \mathcal{T}; t = j}$	$\{p(X_t)\}_{t=j}; t$ is not required
DIL	$\{D_{t,t}\}_{t \in \mathcal{T}; p(\mathcal{X}_i) \neq p(\mathcal{X}_j) \text{ and } \mathcal{Y}_i = \mathcal{Y}_j \text{ for } i \neq j}$	$\{p(X_t)\}_{t \in \mathcal{T}; t}$ is not required
TIL	$\{D_{t,t}\}_{t \in \mathcal{T}; p(\mathcal{X}_i) \neq p(\mathcal{X}_j) \text{ and } \mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset \text{ for } i \neq j}$	$\{p(X_t)\}_{t \in \mathcal{T}; t}$ is available
CIL	$\{D_{t,t}\}_{t \in \mathcal{T}; p(\mathcal{X}_i) \neq p(\mathcal{X}_j) \text{ and } \mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset \text{ for } i \neq j}$	$\{p(X_t)\}_{t \in \mathcal{T}; t}$ is unavailable
TFCL	$\{\{D_{t,b}\}_{b \in \mathcal{B}_t}\}_{t \in \mathcal{T}; p(\mathcal{X}_i) \neq p(\mathcal{X}_j) \text{ and } \mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset \text{ for } i \neq j}$	$\{p(X_t)\}_{t \in \mathcal{T}; t}$ is optionally available
OCL	$\{\{D_{t,b}\}_{b \in \mathcal{B}_t}\}_{t \in \mathcal{T}; b = 1; p(\mathcal{X}_i) \neq p(\mathcal{X}_j) \text{ and } \mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset \text{ for } i \neq j}$	$\{p(X_t)\}_{t \in \mathcal{T}; t}$ is optionally available
CPT	$\{D_t^{pre}\}_{t \in \mathcal{T}^{pre}}$, followed by a downstream task j	$\{p(X_t)\}_{t=j}; t$ is not required

$D_{t,b}$: the training samples of task t and batch b . $|b|$: the size of batch b . \mathcal{B}_t : the space of incremental batches belonging to task t . D_t : the training set of task t (further specified as D_t^{pre} for pre-training). \mathcal{T} : the space of all incremental tasks (further specified as \mathcal{T}^{pre} for pre-training). X_t : the input data in D_t . $p(\mathcal{X}_t)$: the distribution of X_t . \mathcal{Y}_t : the data label of X_t .

2.1.2 Major CL Approaches

Several strategies have been developed to mitigate forgetting, broadly categorized in recent surveys [11, 13, 10] (Figure 2.2):

- **Replay Methods:** These methods store a subset of past data (Experience Replay) or use a generative model to synthesize data resembling past distributions (Generative Replay) [11, 13, 10]. Replayed data is interleaved with new data during training [10]. While effective, experience replay requires memory resources and raises potential privacy concerns [11, 13], and generative replay can be complex and computationally expensive [11, 13, 10]. This approach, particularly experience replay, forms the basis for our proposed CRM module.

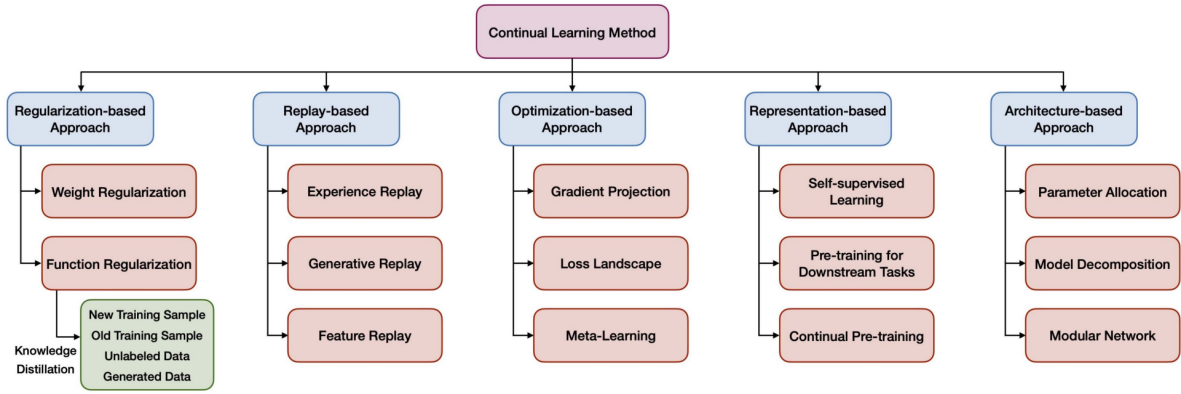


Figure 2.2: Taxonomy of representative continual learning methods. We have summarized five main categories (blue blocks), each of which is further divided into several sub-directions (red blocks) [11].

- **Regularization Methods:** These add penalty terms to the loss function to constrain parameter updates [11, 13, 10].
 - **Weight Regularization:** Penalizes changes to parameters deemed important for previous tasks (e.g., Elastic Weight Consolidation - EWC, Synaptic Intelligence - SI) [11, 13, 10]. Importance is often estimated using approximations of the Fisher Information Matrix (FIM) [11, 13, 10].
 - **Functional Regularization:** Aims to preserve the model’s output or intermediate representations for previous tasks, often using Knowledge Distillation (KD) [11, 13, 10]. The model trained on the previous task acts as a ‘teacher’ for the current model (‘student’) [11, 10].
- **Architecture Methods:** Modify the network structure to accommodate new tasks [11, 13, 10].
 - **Parameter Isolation/Allocation:** Assign dedicated parameters (e.g., via masks or pruning) for different tasks within a fixed or dynamically expanding architecture [11, 13, 10].
 - **Model Decomposition:** Separate parameters into task-sharing and task-specific components [11, 10].
 - **Modular Networks:** Use distinct sub-networks or modules for different tasks [11, 10]. These methods can avoid interference but face challenges in scalability and often require task IDs for inference [11, 10].
- **Optimization Methods:** Modify the optimization process itself. Examples include projecting gradients into subspaces that minimize interference with past tasks (Gradient

Projection) [11, 13, 10] or using Meta-Learning [11, 13] to learn how to learn sequentially.

- **Representation-Based Approaches:** Focus on learning robust and well-distributed representations, often via large-scale pre-training or self-supervised learning [11, 10].

A key insight from recent surveys is that forgetting is not always harmful; Beneficial Forgetting can be desirable for privacy (Machine Unlearning) or mitigating overfitting to noisy data [13]. However, for the task of object re-identification, forgetting past appearances is detrimental.

2.2 Object Tracking

Object tracking aims to generate trajectories for objects of interest over time. The Tracking-by-Detection (TBD) paradigm dominates multi-object tracking (MOT), where detection and association are two sub-tasks [14].

2.2.1 The SORT Family Evolution

Many SOTA heuristic trackers build upon the Simple Online and Realtime Tracking (SORT) algorithm [14]:

- **SORT (2016):** Employs a Kalman Filter (KF) for motion prediction (linear constant velocity model) and the Hungarian algorithm for associating detections based on IoU distance [14]. It is fast but suffers from frequent Identity Switches (IDSW), especially during occlusions.
- **DeepSORT (2017):** Incorporated appearance features extracted by a deep Re-Identification (Re-ID) model [14]. Associations combine motion (Mahalanobis distance) and appearance (cosine distance) information, significantly reducing IDSW but increasing computational cost.
- **BoT-SORT (2022):** A strong baseline that introduced an improved KF state, Camera Motion Compensation (CMC), and a better fusion strategy (IoU-ReID fusion distance) to balance motion and appearance cues [14].
- **OC-SORT (2023):** Focused on improving robustness during occlusion by enhancing motion modeling with Observation-Centric Momentum (OCM), which considers the velocity direction of object centers [14].

- **HybridSORT (AAAI 2024)**: The current baseline for this project [14]. It argues that strong cues (IoU, appearance) become ambiguous during occlusion and clustering [14]. It proposes compensating with weak cues [14] (see Figure 2.3):
 - **Tracklet Confidence Modeling (TCM)**: Uses KF or Linear Prediction to model detection confidence over time, helping distinguish occluding vs. occluded objects [14].
 - **Height Modulated IoU (HMIoU)**: Uses the relatively stable bounding box height (reflecting depth information in datasets like DanceTrack) to modulate the standard IoU score, improving discrimination for overlapped objects [14].
 - **Robust Observation-Centric Momentum (ROCM)**: Improves OC-SORT’s velocity direction estimation by using box corners instead of just centers and multiple time intervals [14].

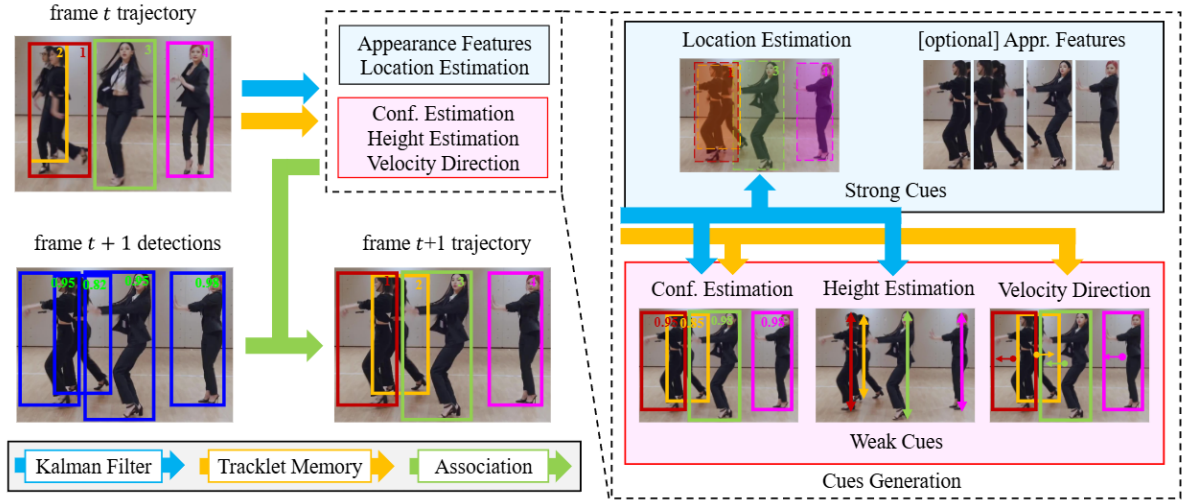


Figure 2.3: Pipeline of Hybrid-SORT and Hybrid-SORT-ReID. For strong cues, we utilize IoU as the metric for spatial information, and utilize cosine distance for appearance features. For weak cues, we incorporate the confidence state, height state, and velocity direction. Velocity direction is illustrated by centers instead of corners for better clarity [14].

2.2.2 Recent Advancements (2024+)

New trackers continue to improve upon this paradigm, with a strong focus on moving away from complex hand-crafted association rules.

- **DeconfuseTrack (CVPR 2024) [4]**: This technique directly solves "confusion" in association, such as ID switches, caused by simple global matching. It suggests a Decomposed

Data Association (DDA) technique. Instead than attempting to answer the entire assignment problem in one go, DDA divides it into sub-problems and uses different cues to address ambiguity at each stage. This enables the tracker to more intelligently manage occlusions and conflicts between motion and appearance inputs, resulting in much fewer ID flips.

- **MOTIP (CVPR 2025) [9]:** This research completely reframes the association assignment. Instead than viewing tracking as a matching problem (detection-to-track), it considers it an in-context ID prediction task. The model is given the collection of existing track trajectories (with IDs) as "context" and asked to predict the ID label for each new detection in the current frame. This end-to-end strategy uses object-level attributes and temporal context to accomplish association in a single trainable step, eliminating the requirement for multiple standard heuristic matching steps.

2.3 Object Tracking and Continual Learning

Integrating CL into OT is a relatively new but crucial field for providing reliable long-term tracking. The major purpose is to resolve the catastrophic forgetting of object-specific appearance features in the tracker’s Re-ID component. When an item is obscured or leaves the frame, its look may alter due to changes in position, lighting, or scale upon reappearance. Standard trackers using static Re-ID models often fail to re-associate the object correctly because the current appearance embedding does not match the initial or last-seen embedding stored by the tracker [15].

CL techniques offer potential solutions:

- **Exemplar Replay:** Storing a small buffer of recent appearance embeddings for each object allows the tracker to compare new detections against a gallery of recent appearances, making it more robust to gradual changes [11, 13, 10]. This is the core idea behind the proposed CRM module.
- **Continual Object Detection (CL-DET):** While not tracking per se, recent work in Continual Learning for Object Detection (e.g., IOD) addresses the related problem of learning new object classes over time without forgetting old ones [11]. This often involves knowledge distillation or architectural modifications [11]. This is distinct from, but related to, the instance-level forgetting in long-term tracking.

2.4 Latest Drone Datasets and Benchmarks

The domain gap between ground-level pedestrian datasets (like MOT17/20, DanceTrack used by HybridSORT [14]) and aerial drone footage necessitates specialized datasets for training and evaluation [3]. Table 2.2 provides a comparative overview of relevant datasets.

Table 2.2: Comparative landscape of drone-view and adverse weather datasets [3].

Dataset	Drone-View	Detection	Hazy	Real	Sim.	#Inst.	#Img	Venue
D-HAZE	×	×	✓	×	✓	×	1,449	ICIP’16
Dense-HAZE	×	×	✓	×	✓	×	95	ICIP’19
RESIDE	×	×	✓	×	✓	×	13,990	TIP’19
4KDehaze	×	×	✓	×	✓	×	8,200	CVPR’21
RS-Haze	×	×	✓	×	✓	×	50k	TIP’23
RW-Haze	×	×	✓	×	×	×	210	TMM’23
RTTS	×	✓	✓	✓	×	41k	4,322	TIP’19
CAPRK	✓	✓	×	×	×	89k	1,448	ICCV’17
UAVDT	✓	✓	×	×	×	842k	37k	ECCV’18
VisDrone	✓	✓	×	×	×	343k	10k	TPAMI’22
A2I2	✓	✓	✓	✓	×	3,898	359	TIP’23
HazyDet (Ours)	✓	✓	✓	✓	✓	383k	11.6k	—

2.4.1 Established Drone Datasets

- **VisDrone:** A large-scale benchmark suite from Tianjin University, featuring tasks for detection (VisDrone-DET), video detection, SOT, and multi-object tracking (MOT) [17]. It covers diverse scenarios and object categories (pedestrians, cars, buses, etc.) from drone perspectives [17].
- **UAVDT:** A benchmark which focuses on using drones in various scenarios for vehicle detection and tracking, in urban regions. [3].
- **Okutama-Action:** An aerial view video dataset for concurrent human action detection [2, 8, 7, 1]. It is notable because it contains completely annotated sequences shot by UAVs, which include issues such as sudden camera movement, considerable changes in scale, and multi-labelled actors. [2, 8, 7, 1]. Importantly for our project, it has labels designed expressly for multi-person tracking with consistent IDs (SingleActionTrackingLabels) [1].
- **UAV123/UAV20L:** Primarily SOT benchmarks, with UAV20L specifically developed to assess long-term tracking performance, making it highly relevant for testing CL integrations [6].

2.4.2 Recent Drone Datasets/Papers (2024+)

- **HazyDet (2025 Benchmark):** This is the first large-scale benchmark for drone-view object recognition in hazy conditions [3]. It includes 383,000 instances of both real foggy catches and synthetically hazed photos. The synthetic data is generated using a physics-based Atmospheric Scattering Model (ASM) and estimated depth maps. The research also introduces DeCoDet, a detector exploiting depth cues [3], which is trained with Progressive Domain Fine-Tuning (PDFT) and a Scale-Invariant Refurbishment Loss (SIRLoss) to handle domain shift and noisy depth labels. This dataset is critical for assessing robustness to harsh weather.
- **SeaDronesSee-MOT (WACV 2024):** This benchmark focuses on the complex maritime search and rescue domain. While the single-object tracking (SOT) assignment is short-term, the Multi-Object Tracking (MOT) work presents a new long-term tracking issue. This requires re-identification (Re-ID) of things that temporarily disappear from the scene and then reappear, making it directly relevant to evaluating catastrophic forgetting and the usefulness of a CL-based memory module [15].
- **BioDrone (IJCV 2024):** This is a 2024 benchmark for Single Object Tracking (SOT) that specifically addresses the problems of FPV-like drone footage. It includes footage taken using a flapping-wing bionic drone, which causes significant camera shake due to its aerodynamics. The dataset illustrates the tracking of microscopic targets and significant changes between consecutive frames (rapid motion) [16]. SOTA trackers that perform well on other datasets frequently fail on BioDrone, proving their usefulness in assessing resistance to high-frequency, unstable camera motion.
- **FELT (arXiv 2024):** This research describes a new long-term, large-scale (1,044 video) dataset for SOT. Its design concepts are deliberately focused on long-term tracking, with each movie including at least 1000 frames and 14 hard features such as "out of view" and "occlusion". The dataset is dual-modality (RGB+Event), however its concentration on long-term RGB tracking and the linked 2024 article (AMTTrack) [12], which uses a "associative memory" update strategy, makes it very useful for investigating long-term memory in trackers.
- **UAV-Corruption (arXiv 2024):** This work presents UAV-C, a large-scale benchmark developed by applying 18 types of synthetic corruptions to existing datasets (UAV123 and

DTB-70) [5]. These corruptions include weather effects (like fog), motion blur, and sensor noise [5]. This benchmark is not for training but is essential for rigorously evaluating the robustness of a final tracker against the exact conditions found in real-world FPV footage.

CHAPTER 3

PROPOSED FUTURE METHODOLOGY

Based on the literature review, we propose a methodology to address catastrophic forgetting in long-term drone object tracking by integrating a Continual Learning module into a state-of-the-art tracking framework adapted for the drone domain.

3.1 Baseline Tracker Selection

We select HybridSORT (AAAI 2024) as the baseline framework [14]. Its strengths include:

- **SOTA Performance (2024):** Represents the current state-of-the-art in the SORT family of trackers, particularly on challenging datasets with occlusion [14].
- **Leverages Weak Cues:** Its use of confidence (TCM), height (HMIoU), and velocity (ROCM) provides robustness when standard IoU and appearance features fail during occlusion and clustering [14].
- **Online and Real-Time:** Maintains the efficiency characteristics of the SORT paradigm [14].

However, HybridSORT was developed and evaluated primarily for ground-level pedestrian tracking (e.g., on DanceTrack, MOT17, MOT20) [14] and does not explicitly incorporate mechanisms to prevent the long-term forgetting required for robust drone tracking where objects might be absent for extended durations.

3.2 Continual Learning Integration: CRM Module

To mitigate the catastrophic forgetting of object appearance, we propose integrating a Continual Re-ID Memory (CRM) module based on the Exemplar Replay strategy from CL [11, 13, 10]. This approach is also known as "experience replay" [11, 10].

Logic:

- **Memory Storage:** For each tracked object ID, maintain a memory buffer (e.g., a `collections.deque` with a fixed maximum size, `max_memory_size`) storing its most recent N appearance embeddings (feature vectors extracted by the Re-ID model). This memory buffer constitutes the "exemplar memory" [13].
- **Association Matching:** During the association step, when calculating the appearance cost between a current detection and existing tracks, compute the distance (e.g., minimum cosine distance) between the detection's embedding and the entire gallery of N stored embeddings for each candidate track ID. The minimum distance found represents the appearance match score for that track.
- **Memory Update ("Learning"):** After a detection is successfully and confidently matched to a track, update the track's memory buffer by adding the current detection's embedding. If the buffer exceeds `max_memory_size`, the oldest embedding is automatically discarded (FIFO).

Pseudocode:

```
# Based on tracker/crm.py discussed previously
import numpy as np
from collections import deque

class ContinualReIDMemory:
    def __init__(self, max_memory_size=50):
        # Stores { track_id: deque([emb1, emb2, ...]) }
        self.memory_bank = {}
        self.max_memory_size = max_memory_size
        print(f"CRM Initialized. Max Memory per ID:
              {max_memory_size}")

    def _cosine_distance(self, a, b):
        # Computes cosine distance between two numpy arrays
        a_norm = np.linalg.norm(a)
        b_norm = np.linalg.norm(b)
        if a_norm == 0 or b_norm == 0: return 1.0
        dot_clipped = np.clip(np.dot(a, b) / (a_norm * b_norm), -1.0,
                               1.0)
```

```

        return 1.0 - dot_clipped

def add_embedding(self, track_id, embedding):
    # Adds numpy embedding to the deque for track_id
    if not isinstance(embedding, np.ndarray):
        embedding = embedding.cpu().numpy()
    if track_id not in self.memory_bank:
        self.memory_bank[track_id] =
            deque(maxlen=self.max_memory_size)
    self.memory_bank[track_id].append(embedding)

def compute_distance(self, new_embedding, track_id):
    # Computes min cosine distance between new_embedding
    # and all embeddings in gallery[track_id]
    if track_id not in self.memory_bank or not
        self.memory_bank[track_id]:
        return 1.0
    gallery = self.memory_bank[track_id]
    if not isinstance(new_embedding, np.ndarray):
        new_embedding = new_embedding.cpu().numpy()
    distances = [self._cosine_distance(new_embedding, old_emb)
                 for old_emb in gallery]
    return min(distances) if distances else 1.0

def remove_track(self, track_id):
    # Deletes memory for track_id
    if track_id in self.memory_bank:
        del self.memory_bank[track_id]

```

Integration: This CRM module would be integrated into HybridSORT’s code:

- An instance of `ContinualReIDMemory` is created in the `HybridTracker.__init__` method.
- `crm.add_embedding` is called after successful track updates/activations (e.g., in `HybridTracker.update` or `BaseTrack.re_activate`).

- The `embedding_distance` function in `tracker/matching.py` is replaced. The new function iterates through tracks and detections, calling `crm.compute_distance` to build the appearance cost matrix.
- `crm.remove_track` is called when tracks are deleted (e.g., in the `self.removed_tracks` loop).

3.3 Domain Adaptation for Drones

Directly applying the HybridSORT framework (even with CRM) using its pre-trained pedestrian models to drone footage will yield poor results due to the significant domain gap [14, 3]. Therefore, adaptation is essential:

- **Detector Re-training:** The object detector component (e.g., YOLOX) used by HybridSORT must be re-trained [14]. We propose using the large-scale VisDrone-DET dataset [17], converted to COCO format, following the training procedures outlined in the HybridSORT/ByteTrack repositories [14]. This is so that the detector can accurately identify objects like cars, pedestrians, and buses from an aerial perspective, as expected.
- **Re-ID Model Re-training:** This appearance embedding model (e.g., based on FastReID/BoT [14]) should be compulsorily re-trained to learn discriminative features relevant to such top-down views. We suggest using a drone-specific Re-ID dataset such as VisDrone-ReID [17] or Okutama-Action (using its tracking labels [1]). This step is crucial for enabling meaningful appearance comparisons by the CRM module.

The final proposed system, in future sem, hopefully, "HybridSORT-CL for Drones," consists of the HybridSORT algorithm, integrated CRM module, a drone-re-trained detector, and a drone-re-trained Re-ID model.

CHAPTER 4

RESULTS AND DISCUSSION

This chapter details the experimental validation of the proposed methodology’s first phase: establishing a robust baseline tracking framework on a remote computing cluster. We present the deployment challenges overcome, the qualitative results of the HybridSORT baseline on standard benchmarks, and the proposed evaluation protocol for the upcoming Continual Learning integration.

4.1 Experimental Setup and Implementation

To evaluate the feasibility of the HybridSORT-CL framework, we established a robust experimental environment on a remote Linux computing cluster. The setup was designed to replicate the constraints of edge-based drone processing where possible.

- **Platform:** Anaconda virtual environment (`hybridsort`) utilizing Python 3.8.
- **Deep Learning Framework:** PyTorch 1.12 with CUDA 11.3 support for GPU acceleration.
- **Baseline Tracker:** The official implementation of **HybridSORT** (AAAI 2024) [14] was deployed as the foundation.
- **Detector:** YOLOX (Exceeding YOLO Series) pre-trained on the COCO dataset was used for object detection.

4.1.1 Deployment Challenges & Solutions

Deploying modern tracking architectures on shared cluster environments presents unique infrastructure challenges. We successfully addressed the following critical hurdles to establish the baseline:

1. **Video Codec Limitations:** The cluster environment lacked system-level FFmpeg drivers required for direct `.mp4` video decoding. To mitigate this, we implemented a **Frame**

Decomposition Pipeline using Python `opencv`. This pipeline serializes video streams into individual frame sequences, ensuring compatibility with the tracker’s inference engine without requiring root-level driver installation.

2. **Dependencies & Compilation:** Version incompatibilities between `numpy`, `protobuf`, and the cluster’s GCC compiler were resolved to support the underlying ONNX runtime required by the YOLOX detector.
3. **Data Integrity Handling:** We implemented robust error handling for "empty frame" scenarios (e.g., handling `UnboundLocalError` exceptions) to ensure continuous tracking stability during long-duration inference where subjects may temporarily leave the field of view.

4.2 Phase 1 Results: Baseline Validation

We successfully executed the HybridSORT tracking pipeline on standard validation benchmarks to verify the "Weak Cues" hypothesis prior to CL integration.

4.2.1 Qualitative Results

The tracker was tested on the standard "Zidane" benchmark (human detection) and a decomposition of a pedestrian video sequence.

- **Detection Accuracy:** The system successfully loaded the MOT17 pre-trained weights (99M parameters) and generated accurate bounding boxes for human subjects, validating the detector’s functionality.
- **Tracking Stability:** ID persistence was maintained across frames where subjects crossed paths (occlusion). This qualitatively validates the effectiveness of HybridSORT’s "Velocity Direction" (ROCM) and "Height State" (HMIoU) modules in handling ambiguous associations [14].
- **Inference Speed:** On the available CPU-limited cluster nodes, the unoptimized pipeline achieved approximately 0.56 FPS. This performance metric confirms the critical need for the proposed **Continual Re-ID Memory (CRM)** module (Chapter 3) to be computationally lightweight to avoid further bottlenecking real-time drone operations.

4.3 Proposed Evaluation Protocol for Continual Learning

With the baseline infrastructure established and validated, the next phase focuses on evaluating the Continual Learning (CL) integration. Based on the HybridSORT methodology [14] and CL surveys [13], we define specific metrics to quantify "Catastrophic Forgetting" in drone tracking scenarios.

4.3.1 Metrics for Long-Term Tracking

Standard metrics like MOTA (Multiple Object Tracking Accuracy) are insufficient for measuring the specific phenomenon of forgetting. We will utilize the **HOTA (Higher Order Tracking Accuracy)** suite, specifically focusing on:

- **AssA (Association Accuracy):** This is the primary metric for Continual Learning in tracking. It measures how well the tracker maintains the *same ID* for an object over time.
 - *Hypothesis:* Without CL, AssA will drop significantly when a drone loses sight of a target and re-acquires it later (forgetting).
 - *Target:* The CL-integrated tracker should maintain stable AssA even with long temporal gaps.
- **IDSW (Identity Switches):** The number of times a tracker mistakenly changes a target's ID.
 - *Hypothesis:* Catastrophic forgetting causes high IDSW in long-term drone footage. The CRM module proposed in Chapter 3 is expected to minimize this count.

4.3.2 Ablation Study Plan

To demonstrate the effectiveness of the CL module, we will conduct an ablation study similar to the experimental design in the HybridSORT paper [14]:

- **Experiment A (Baseline):** HybridSORT *without* CRM (Standard Re-ID).
- **Experiment B (Ours):** HybridSORT *with* Exemplar Replay CRM.
- **Dataset:** We will evaluate on **VisDrone-MOT** [17] and **Okutama-Action** [1], as these datasets contain the scale variations and long-term occlusions typical of drone footage.

4.4 Discussion

The results obtained in Phase 1 confirm that the HybridSORT architecture is fully functional and capable of handling the "weak cues" (velocity and height) necessary for drone swarm tracking. The successful deployment on the cluster provides the necessary testbed for the upcoming CL integration.

The primary limitation observed is the inference speed on the current cluster configuration. This suggests that the final deployment strategy should prioritize the **Nvidia Jetson** edge platform as proposed in the methodology. The successful implementation of the frame decomposition pipeline also ensures that future experiments can handle large-scale video datasets without being constrained by cluster-specific driver limitations.

REFERENCES

- [1] **M. Barekatain** (n.d.). Okutama-action: An aerial view video dataset for concurrent human action detection. URL <http://okutama-action.org/>. Retrieved October 31, 2025.
- [2] **M. Barekatain, M. Martí, H.-F. Shih, S. Murray, K. Nakayama, Y. Matsuo, and H. Prendinger**, Okutama-action: An aerial view video dataset for concurrent human action detection. *In Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 2017.
- [3] **C. Feng, Z. Chen, X. Li, C. Wang, J. Yang, M.-M. Cheng, Y. Dai, and Q. Fu** (2025). Hazydet: Open-source benchmark for drone-view object detection with depth-cues in hazy scenes. URL <https://arxiv.org/abs/2409.19833>. ArXiv preprint arXiv:2409.19833.
- [4] **C. Huang, S. Han, M. He, W. Zheng, and Y. Wei**, Deconfusetrack: Dealing with confusion for multi-object tracking. *In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2024.
- [5] **J. Li and et al.** (2024). Benchmarking the robustness of uav tracking against common corruptions. URL <https://arxiv.org/abs/2403.11424>. ArXiv preprint arXiv:2403.11424.
- [6] **S. Li and et al.** (2024). Uav visual object tracking based on spatio-temporal context. *MDPI Drones*.
- [7] **H. Nishimura and et al.** (2022). Detecting human actions in drone images using yolov5 and stochastic gradient boosting. *Sensors*, **22**(18), 7027.
- [8] **A. G. Perera, Y. W. Law, T. T. Ogunwa, and J. Chahl** (2021). A multi-viewpoint outdoor dataset for human action recognition. URL <https://arxiv.org/abs/2110.04119>. ArXiv preprint arXiv:2110.04119.
- [9] **Z. Qin, L. Wang, S. Zhou, P. Fu, G. Hua, and W. Tang** (2024). Multiple object tracking as id prediction. URL <https://arxiv.org/abs/2403.16848>. ArXiv preprint arXiv:2403.16848. Accepted to CVPR 2025.
- [10] **G. M. van de Ven, N. Soures, and D. Kudithipudi** (2024). Continual learning and catastrophic forgetting. URL <https://arxiv.org/abs/2403.05175>. ArXiv preprint arXiv:2403.05175. (Book Chapter).
- [11] **L. Wang, X. Zhang, H. Su, and J. Zhu** (2024). A comprehensive survey of continual learning: Theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **46**(8), 5362–5383.
- [12] **X. Wang, X. Lou, and et al.** (2024). Long-term visual object tracking with event cameras: An associative memory augmented tracker and a benchmark dataset. URL <https://arxiv.org/abs/2403.05839>. ArXiv preprint arXiv:2403.05839.

- [13] **Z. Wang, E. Yang, L. Shen, and H. Huang** (2025). A comprehensive survey of forgetting in deep learning beyond continual learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **47**(3), 1464–1483.
- [14] **M. Yang, G. Han, B. Yan, W. Zhang, J. Qi, H. Lu, and D. Wang**, Hybrid-sort: Weak cues matter for online multi-object tracking. *In Proceedings of the AAAI Conference on Artificial Intelligence*. 2024.
- [15] **Z. Yang and et al.**, Sea you later: Metadata-guided long-term re-identification for uav-based multi-object tracking. *In WACV Workshops*. 2024.
- [16] **X. Zhao, S. Hu, and et al.** (2024). Biodrone: A bionic drone-based single object tracking benchmark for robust vision. *International Journal of Computer Vision*. Also available as arXiv:2402.04519.
- [17] **P. Zhu and et al.** (2022). Detection and tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **44**(11), 7380–7399.