

A framework for detecting and tracking elephants in drone videos

Chaim Chai Elchik ^a, Serge Wich^b, and André Burger ^c

^aLancaster Environment Centre, Lancaster University, Lancaster, Lancashire, United Kingdom; ^bSchool of Biological and Environmental Sciences, Liverpool John Moores University, Liverpool, Merseyside, United Kingdom; ^cWelgevonden Game Reserve, Vaalwater, Limpopo Province, South Africa

Corresponding author: Chaim Chai Elchik (email: c.elchik123@gmail.com)

Abstract

The escalating global biodiversity crisis requires innovative and scalable solutions to monitor wildlife populations. Recent developments in remote sensing and deep learning offer promising avenues for improving the conservation of large mammals, including African elephants. This paper introduces a framework that utilizes drone video streams and integrates state-of-the-art object detection (YOLOv11) and tracking (BoT-SORT) methods, which are significantly enhanced by a custom post-track re-identification algorithm, to capture temporal dynamics and track individual elephants over time. The framework facilitates automated video analysis and elephant counting, generating key metrics such as individual elephant movement speed, group movement patterns, and Elephant Cluster Statistics. By automating aspects of data processing and analyses, this approach provides valuable insights that contribute to more efficient and data-driven decision-making in wildlife research.

Key words: object detection, object tracking, re-identification, drone videos, wildlife conservation, YOLO

1. Introduction

The escalating global biodiversity crisis requires innovative and scalable solutions to monitor wildlife populations to support conservation management (Kissling et al. 2024). The scale of the crisis is illustrated by the 2024 Living Planet Index report that showed a 73% average decline in wildlife population size for the set of species and populations they measured from 1970 to 2020 (WWF 2024). The IPBES (2019) further noted that major land-based habitats have declined by at least 20%, with over 40% of amphibian species, nearly 33% of reef-forming corals, and more than a third of all marine mammals now threatened. Additionally, a 2023 study (Finn et al. 2023), which analyzed population trend data for over 71 000 animal species across all 5 vertebrate groups, revealed a widespread global erosion of biodiversity, with 48% of species experiencing population declines.

The conservation of large mammals, such as African elephants (*Loxodonta africana*), can be significantly advanced by recent developments in conservation technology such as remote sensing and deep learning as seen in works by Wich and Piel (2021), Lamba et al. (2019) and Berger-Tal and Lahoz-Monfort (2018). Traditional wildlife monitoring is often costly, labor-intensive, and risky for researchers, particularly when studying elusive or dangerous species in remote areas highlighted in earlier research by Hodgson et al. (2016), McEvoy et al. (2016), Vermeulen et al. (2013), and Pedrazzi et al. (2025). One of these conservation technologies,

offers a potentially cost-effective, and less intrusive alternative for data acquisition than ground-based surveys, particularly if integrated with deep learning to (semi) automate analyses (López and Mulero-Pázmány 2019; Hamilton et al. 2020; Wich and Piel 2021). Early studies demonstrated the potential of drones for wildlife surveys. For example, Vermeulen et al. (2013) explored the use of drones to survey large mammals in Burkina Faso. They found that elephants were easily visible in drone images, with no observed reaction from the animals when the drone flew at 100 m. However, smaller mammals were harder to detect. The study concluded that drones could be a valuable tool for elephant enumeration, though the limited flight duration of the drones was a constraint.

Hodgson et al. (2016) further demonstrated the precision of UAVs for wildlife monitoring in various environments, showing that UAV-derived counts of nesting birds were more precise than traditional ground counts. This highlights the potential of UAVs to improve the accuracy and efficiency of wildlife monitoring.

Research has also addressed the potential impact of drones on wildlife. McEvoy et al. (2016) assessed the disturbance effects of UAVs on waterfowl, finding little to no disturbance when drones were flown at sufficient altitudes (60 m for fixed-wing, 40 m for multirotor). Further research by Mulero-Pázmány et al. (2017) and Afridi et al. (2025) also demonstrate how drones can disturb wildlife and what steps must be taken

to prevent this from happening as these findings are crucial for developing responsible drone-based monitoring practices.

The application of deep learning to drone imagery has been a key area of development. [Kellenberger et al. \(2018\)](#) tackled the challenges of mammal detection in drone images with imbalanced datasets, providing recommendations for scaling convolutional neural networks (CNNs) to large-scale wildlife census tasks. [Barbedo et al. \(2019\)](#) focused on cattle detection in drone images using deep learning, evaluating CNN architectures and image resolution. [Guirado et al. \(2019\)](#) developed a CNN-based system for automated whale detection and counting in satellite and aerial images, showcasing the potential of deep learning for marine mammal monitoring.

Previous research by [Delplanque et al. \(2021\)](#) first harnessed ultra-high-resolution (38–50 cm) panchromatic and true-color satellite imagery, pairing a U-Net segmentation network with K-means clustering to automatically localize and count sprawling mammal herds. Building on this, [Delplanque et al. \(2023a\)](#) swapped in oblique aerial RGB photographs—acquired via fixed-wing aircraft—and introduced HerdNet, a point-based CNN that outputs density maps to tally camels, donkeys, sheep, and goats more accurately than manual counts, though it omits both satellite data and elephant surveys. Subsequent research by [Delplanque et al. \(2023b\)](#) presented a semi-automated deep-learning pipeline that embedded the pretrained HerdNet model to slash human verification time by over 70% (and up to 98% in some surveys), while still mandating human quality checks to navigate shadows, occlusions, and species overlap. Concurrent field trials by [Delplanque et al. \(2024\)](#) further revealed that variable lighting, terrain heterogeneity, and mixed-species groupings can still hamper count precision, underscoring the imperative for richer, site-specific annotations rather than off-the-shelf detectors like Faster R-CNN or RetinaNet.

Other recent studies have further expanded the application of drones and deep learning in wildlife monitoring. [Rančić et al. \(2023\)](#) explored CNNs for animal detection and counting from drone images, [Koger et al. \(2023\)](#) presented a system for quantifying animal movement, behavior, and environmental context using drones and computer vision, and [Brickson et al. \(2023\)](#) reviewed the role of AI in elephant monitoring. Datasets specifically designed for wildlife detection in drone imagery, such as WAID ([Mou et al. 2023](#)), are also contributing to the advancement of the field.

Furthermore, [Mpouziotas et al. \(2024\)](#) presented methods for tracking wild birds from drone footage, [Alsaïdi et al. \(2024\)](#) detailed deep learning for tracking beluga whales in aerial video, and [Shukla et al. \(2024\)](#) explored estimating 3D poses and shapes of animals from drone imagery. Collectively, these studies illustrate a clear progression from labor-intensive manual approaches to advanced, automated monitoring systems based on high-resolution imaging and deep learning. Distinct from these static image-based approaches, [Pedrazzi et al. \(2025\)](#) provide a comprehensive review highlighting the transformative impact of drone technology on animal behaviour research, with a particular emphasis on the

role of automated data analysis. Their work underscores how rapid advancements in image-tracking technologies and AI, including deep-learning algorithms like CNNs, are enabling automated processes for species identification, counting, tracking, and behaviour recognition from drone-acquired data. While they acknowledge the use of these techniques for tracking and quantifying interactions to create activity budgets and association patterns, the broader literature, as implied by their review, has seen a stronger emphasis on the automation of animal detection rather than the fine-grained automation of dynamic behavioural analysis, such as movement speed within groups.

The recent evolution of object detection technology—exemplified by single-stage detectors such as You Only Look Once (YOLO)—has significantly pushed the boundaries of both detection accuracy and real-time performance ([Wang and Liao 2024](#)). While early versions of YOLO demonstrated powerful detection capabilities, subsequent refinements culminating in YOLOv11 ([Khanam and Hussain 2024](#)) have markedly improved small object detection, robustness under challenging environmental conditions, and frame-rate processing speeds. Such enhancements are critical for dynamic, real-time scenarios, particularly when processing high-resolution drone video feeds that directly influence effective conservation efforts.

Complementing these detection advances, breakthroughs in multiobject tracking have transformed real-time monitoring capabilities. The BoT-SORT framework ([Aharon et al. 2022](#)) exemplifies this progress by overcoming challenges related to rapidly moving objects and occlusions. Leveraging robust appearance-based re-identification along with refined motion association techniques, BoT-SORT integrates predictive filtering with dynamic feature matching to maintain consistent tracking even amidst erratic movements or partial obstructions. This level of robustness is vital in conservation applications, ensuring that individual elephants can be continuously tracked through complex and ever-changing scenes.

Building upon this foundation and motivated by the recent advancements in detection and tracking, our work specifically addresses the need for more automated approaches to analyze complex group behaviors, focusing on movement patterns, speeds, and group formations that are recently being studied with drones instead of from the ground ([Dai et al. 2007](#)) and facilitate our understanding of animal movement behaviour as well as the impact of the drone on movement itself ([Inoue et al. 2019](#); [Koger et al. 2023](#); [Schad and Fischer 2023](#)). Our methodology leverages drone video streams, integrating state-of-the-art detection (YOLOv11) and tracking (BoT-SORT), which are significantly enhanced by a custom post-track re-identification algorithm. This novel step, which is a core contribution of this work, is specifically designed to mitigate identity switching in complex drone video scenarios. This enables the derivation of movement dynamics and group patterns in an automated manner. This integrated approach mitigates challenges in real-time monitoring and behavioral analysis, providing finer temporal resolution and more robust conservation insights.

2. Methodology

2.1. Experimental setup

2.1.1. Dataset

The original dataset consisted of eight MP4 video files captured using a DJI *Mavic 3 Pro—Hasselblad camera—Drone* (see Appendix Table A1 for full technical details) flying over the Welgevonden Game Reserve in South Africa. All videos were recorded on the same day and in the same general area within the reserve under consistent atmospheric conditions. The elephants were located with the help of wildlife guides and trackers using cars or buggies. All videos were recorded in 4K resolution (3840×2160) at 30 frames per second (fps) and at varying altitudes and distances from the elephant subjects. The same herd of elephants was tracked and filmed in all eight videos. The total duration of the videos is 24 min and 1 s, with an average duration of 3 min.

To generate a robust dataset that can be used to train an object detection model, the video sequences were decomposed into individual frames. A sampling strategy of one fps was implemented to prevent overfitting and reduce annotation time. Splitting the video at its original rate of 30 fps would create many nearly identical frames. This could bias the model toward redundant features and increase the annotation workload. Therefore, we adopted a subsampling approach, selecting 1 frame every 30 frames.

Frame extraction was automated using a *Python* script. For each video, frames were extracted and saved if $f \equiv 0 \pmod{30}$, where f is the frame number. This process resulted in a dataset comprising 1441 representative frames.

2.1.2. Dataset annotation

To efficiently annotate the dataset with bounding boxes, we employed semi-automated techniques using Roboflow, a platform that provides a graphical interface to simplify manual data annotation. Roboflow leverages pre-trained object detection models, to generate initial bounding box annotations. These automatically suggested boxes can be accepted, rejected, or adjusted by the user, streamlining the annotation process.

The implemented workflow consisted of several steps. The first step was to use the pre-trained models to generate initial bounding box predictions, this provided a starting point for annotation. To enhance efficiency, Roboflow's box prompting feature then suggested bounding boxes based on user-provided annotations over time, enabling quick and accurate modifications through an easy-to-use interface. This was followed by manual reviewing of the proposed annotations and manually adjusting them as needed before adding the labels. Finally, the annotated dataset was used to retrain the detection model in a feedback loop, progressively improving its accuracy as it learned from newly labeled data (Roboflow 2025).

This process significantly accelerated the speed at which annotation could be made but the annotations were not flawless. The generated bounding boxes were often too large, too small, or entirely false positives. In some cases, elephant sub-

jects received multiple bounding box suggestions, splitting them up into several detections. The interface allowed for easy manual correction. Additionally, some frames were entirely rejected due to issues such as excessive camera motion, absence of elephants, or extreme zoom-ins/outs. After these adjustments, a total of 1337 frames were successfully annotated.

The dataset was then partitioned into training, validation, and testing sets comprising of 70%, 20%, and 10% of the frames. The training frames were then augmented by creating versions of them that were randomly rotated between -15° and $+15^\circ$, increasing the total number of training frames to 2367. This then increased the total amount of frames to 2705. The new ratios between training, validation, and testing sets therefore changed to 87.5% (2367 frames), 8.4% (225 frames), and 4.1% (113 frames), respectively. The dataset was then exported from Roboflow and included separate folders for each subset, along with corresponding annotation files in the required format for object detection model training. Each annotation file contained the object label (in this dataset, 0) and the x_{\min} , x_{\max} , y_{\min} , and y_{\max} coordinates.

2.1.3. Model selection

Traditionally, two-stage object detection models, such as Faster R-CNN (Ren et al. 2016), have demonstrated superior accuracy when compared to single-stage detection models. However, this has changed with the emergence of single-stage detection models such as the YOLO model series (Redmon et al. 2016; Wang and Liao 2024) and the Single Shot MultiBox Detector (Liu et al. 2016), which have closed the performance gap. This has led to computationally efficient single-stage models being able to be deployed where two-stage models were traditionally required. The YOLO series currently leads in both performance and inference speed, with YOLOv11 representing the latest advancement at the time of writing (Khanam and Hussain 2024; Wang and Liao 2024).

YOLOv11 builds upon the previous iterations of the YOLO series. Most notably, it integrates an optimized backbone network and improved anchor box strategies, which enhance object localization capabilities. This is an essential feature for detecting elephant subjects at varying distances and under diverse lighting conditions. Additionally, YOLOv11 leverages advanced transfer learning techniques, enabling efficient adaptation of pre-trained models to domain-specific datasets with limited or highly variable training samples. This ensures both rapid convergence and high detection accuracy (Khanam and Hussain 2024).

Three important design improvements contribute to YOLOv11's enhanced performance. The *C3K2 Block* utilizes smaller kernel sizes to optimize feature extraction, improving computational efficiency without compromising accuracy. Building on this, the Spatial Pyramid Pooling (SPP) Fusion Module, an evolution of the traditional SPP module, captures multiscale features, enhancing the model's ability to detect objects of varying sizes—an essential capability for processing aerial imagery. Additionally, the Cross Stage Partial

with Spatial Attention Block incorporates spatial attention mechanisms, allowing the model to focus on critical regions within an image, which is particularly beneficial for detecting partially occluded or overlapping objects (Khanam and Hussain 2024).

These innovative changes allow YOLOv11 to maintain real-time inference speeds while achieving higher mean average precision (mAP) than previous versions. Furthermore, its more streamlined processing pipeline minimizes latency. The enhanced nonmaximum suppression techniques also further refine object detection by reducing redundant bounding boxes and improving localization precision. Due to these improvements YOLOv11 is able to perform state of the art scalability and generalization which makes it a well-suited model for detecting elephants in drone images (Khanam and Hussain 2024).

The demonstrated success of YOLOv8 in challenging detection scenarios, particularly those involving complex motion and low-contrast subjects (Dave et al. 2023; Fang et al. 2024; Varghese and Sambath 2024; Yaseen 2024) highlights the ongoing evolution of the YOLO models. YOLOv11 builds upon the strengths of YOLOv8, which allows for real-time detection capabilities but with higher accuracy and robustness. These qualities are critical for the proposed framework, where timely and precise object detection serves as the foundation for effective post-track re-identification.

2.1.4. Tracker selection

Selecting a tracking algorithm that performs well with the complexity drone videos present is essential for ensuring that the proposed framework is robust and reliable. Although various tracking methodologies such as ByteTrack, DeepSORT, and BoT-SORT have been presented in recent literature, the BoT-SORT algorithm distinctly emerges as the most suited for drone-captured imagery. BoT-SORT capitalizes on robust association strategies that adeptly mitigate challenges inherent to aerial monitoring, including rapid target motion, pronounced scale variations, and frequent occlusions (Aharon et al. 2022).

BoT-SORT's architecture introduces several key improvements over traditional tracking methods. *Robust detection association* sets it apart from DeepSORT and its derivatives, which primarily rely on rudimentary motion models. Instead, BoT-SORT incorporates a sophisticated association mechanism that merges detection confidence with motion prediction, ensuring sustained object tracks even in cases of partial occlusion or abrupt motion changes (Wojke et al. 2017; Aharon et al. 2022; Zhao et al. 2024). Expanding on this, its *enhanced appearance modeling* refines re-identification processes by embedding improved appearance features, a crucial enhancement for distinguishing animals in drone videos, especially when dealing with overlapping trajectories and varying illumination conditions (Wojke et al. 2017; Zhao et al. 2024). Furthermore, the *adaptability to complex backgrounds* allows BoT-SORT to handle heterogeneous, cluttered drone imagery while mitigating false associations and ensuring precise object localization, outperforming alternative meth-

Table 1. Final YOLOv11x training hyperparameters.

Hyperparameter	Value
Model	YOLOv11x
Epochs	150
Imgsz	640 × 640
Lr0	0.01
Lrf	0.1
Batch	8
Weight_decay	0.0005
Save_period	10 epochs

ods like ByteTrack in maintaining detection accuracy with consistent tracking (Aharon et al. 2022; Zhang et al. 2022; Zhao et al. 2024). Finally, despite its intricate association strategy, its real-time performance is preserved while maintaining computational efficiency critical for real-time applications. This balance between precision and processing speed makes BoT-SORT well suited for animal tracking scenarios (Aharon et al. 2022; Zhao et al. 2024).

The combination of these architectural and algorithmic features makes it clear that BoT-SORT is the best choice for tracking animals in drone videos. Its proficiency in persistently associating detections across successive frames ensures that transient occlusions and rapid target movements do not result in track fragmentation. Furthermore, the algorithm's integrated utilization of both appearance-based and motion-based cues offers a comprehensive and adaptable solution tailored to the multifaceted nature of aerial surveillance imagery (Zhao et al. 2024).

2.1.5. Detection and tracking model training and fine tuning

The YOLOv11x detection model was trained iteratively with varying hyperparameters, leading to several configurations that were evaluated to determine the optimal settings. Table 1 outlines the final selected hyperparameters. The largest variant, YOLOv11x, was chosen to maximize performance, as smaller models like YOLOv11n yielded lower detection scores. Training was conducted for 150 epochs to ensure robust generalization across varying perspectives, lighting conditions, and object scales in drone imagery. The input image size (*imgsz*) was set to 640 × 640 pixels, balancing detail preservation with computational efficiency. A high initial learning rate (*lr0*) of 0.01 facilitated rapid convergence, while a final learning rate (*lrf*) of 0.1 ensured refined weight adjustments in later epochs. The batch size of eight was selected based on GPU memory constraints, optimizing computational feasibility and gradient updates. Weight decay was set to 0.0005 to prevent overfitting, and model checkpoints were saved every 10 epochs to allow for rollback in case of instability.

The final trained model demonstrated strong performance across multiple evaluation metrics. The preprocessing time was 0.3 ms, inference time was 13.9 ms, and postprocessing time was 4.2 ms, ensuring real-time detection capabilities. In

Table 2. Final BoT-SORT hyperparameters.

Hyperparameter	Value
Track_high_thresh	0.20
Track_low_thresh	0.05
New_track_thresh	0.75
Track_buffer	90
Match_thresh	0.85
Fuse_score	True
Gmc_method	sparseOptFlow
Proximity_thresh	0.5
Appearance_thresh (with re-id)	0.25
Size_ratio_thresh	0.8
Iou_thresh	0.5

terms of complexity, the model contained 464 layers, 56.8 million parameters, and had a computational cost of 194.4 GFLOPS. These results indicate that YOLOv11x achieves high accuracy while maintaining efficiency suitable for real-time applications.

The BoT-SORT tracking model was fine-tuned for tracking elephants in drone video footage by iteratively adjusting its hyperparameters via the YAML configuration file. Table 2 summarizes the final selected hyperparameters. Given the challenges posed by aerial views, a lower *track_high_thresh* of 0.20 was chosen to allow associations even when detection confidence was reduced due to partial occlusions. Additionally, a *track_low_thresh* of 0.05 enabled a secondary matching stage for borderline detections. To minimize false tracks, *new_track_thresh* was set to 0.75, ensuring that only highly confident detections initiated new tracks. A *track_buffer* of 90 frames allowed tracks to persist through temporary detection lapses, which are common in drone footage due to motion blur or occlusions. A high *match_thresh* of 0.85 was used to enforce strict spatial and appearance-based correspondence between detections and tracks, reducing false associations. The *fuse_score* parameter was enabled to integrate raw detection confidence into the matching process, enhancing robustness. Given the significant camera motion in drone footage, *gmc_method* was set to *sparseOptFlow* for efficient global motion compensation. To ensure spatial consistency, *proximity_thresh* was set to 0.5, allowing detections to be associated only if they were sufficiently close. With re-identification enabled, *appearance_thresh* was set to 0.25, enforcing strict similarity requirements to accurately track visually similar elephants even after occlusions. To prevent erroneous associations caused by scale variations, *size_ratio_thresh* was set to 0.8. Finally, an *iou_thresh* of 0.5 was maintained to balance strictness and leniency in spatial alignment between detections and existing tracks.

2.1.6. Post-track re-identification algorithm

Preliminary model outputs revealed a significant ID switching issue: elephant objects that temporarily “disappeared” due to occlusions—whether by moving behind other elephants, exiting the frame, or becoming obstructed by struc-

tural elements—were later “reappearing” with new IDs. This problem stemmed from the BoT-SORT tracking model’s inability to match objects when the disappearance persisted for an extended period or when the reappearing elephant’s orientation had substantially changed (e.g., shifting from upward to downward or from leftward to rightward). To address this limitation, a post-track re-identification algorithm was implemented. This algorithm detects instances of ID switching and reassigns the original IDs, thereby ensuring a more accurate count of unique elephant objects and enhancing overall tracking performance.

The algorithm begins by identifying potential disappearances by scanning each elephant object ID across all video frames. If an elephant object ID is absent from one or more frames, it is flagged as a potential disappearance and recorded for further analysis. Once disappearances are identified, the algorithm searches for potential reappearances, examining frames following the last recorded occurrence of each disappeared ID. Any new elephant object ID appearing in these frames is considered a candidate for reassignment, forming a list of potential reappearances.

Next, the algorithm constructs candidate matches by associating each disappeared elephant object ID with one or more potential reappearing IDs. These candidate pairs undergo evaluation based on three key conditions: edge, distance, and similarity. The edge condition ensures that if an elephant disappears near the frame’s edge, its reappearance must also occur near the same edge, within a defined Euclidean distance relative to its bounding box size (see Fig. 1). If the last known frame of the elephant is not near an edge, this condition is disregarded. The distance condition estimates the maximum travel distance of the disappeared elephant based on its observed speed and compares it to the normalized Euclidean distance between the last known position and the first detected position of the candidate reappearance (see Fig. 2). If the estimated travel range does not align with the actual observed movement, the match is rejected. The similarity condition further refines the matching process by analyzing the visual similarity between the last recorded frame of the disappeared elephant and the first frame of the candidate reappearance, assigning a similarity score accordingly.

Following the evaluation, the algorithm selects the best match for each disappeared ID by identifying the candidate with the highest combined distance and similarity scores. Not all disappearances yield valid matches, meaning that the final list of confirmed re-identifications may be shorter than the initial set of candidate pairs. Finally, the identified elephant object IDs are updated, replacing the disappeared ID with the matched reappearing ID. This ID correction process supports chain reactions; for example, if ID 3 is matched with ID 4, and ID 4 is later matched with ID 5, the correction propagates through the entire sequence to maintain consistency.

2.1.7. Data analysis

The CSV output generated by the post-track re-identification algorithm serves as the foundation for a comprehensive data analysis, enabling the creation of rel-

Fig. 1. Edge condition for ID matching video 0395. The left image shows the frame before the camera pans to the right, while the right image shows the frame after the camera pans back to the left. The elephant with ID 3 in the left image is reassigned a new ID, 8, after the panning motion. These two IDs correspond to the same elephant.



Fig. 2. Distance condition for ID matching video 0406. The left image shows the last frame where ID 4 is visible, and the right image shows the first frame where ID 7 appears. Both IDs belong to the same elephant. The black dotted circle indicates the maximum range the elephant could have traveled. The red dot marks ID 4's last location, the purple dot marks ID 7's first location, and the black line represents the normalized Euclidean distance between the two.



evant statistical summaries and visualizations for further ecological research. This analysis aims to highlight key segments of the videos that may warrant manual review, facilitating the identification of significant behavioral patterns and ecological events. By automating aspects of data processing and visualization, this analysis reduces the workload of ecologists while providing valuable insights that contribute to more efficient and data-driven decision-making in wildlife research.

It is important to note that all spatial metrics described in this section—including movement speed, trajectories, and travel distance—are calculated in pixel units. This was a deliberate methodological choice. The framework is designed for broad accessibility and ease of use, allowing researchers to apply it to any standard drone video without requiring complex camera calibration or the integration of drone telemetry data. This approach ensures the tool remains practical for field conditions where such setups are often infeasible, as will be expanded upon in the **Discussion** section.

2.1.7.1. Individual Elephant Movement Speed Plot

The *Individual Elephant Movement Speed Plot* analysis visualizes the movement speed of an elephant by analyzing changes in its central position over time using the Euclidean distance between consecutive center points recorded every 30 frames. This approach provides a measure of the elephant's speed fluctuations per second. Specifically, the func-

tion extracts the coordinates of the elephant's center at 30-frame intervals and computes the distance between these points. A greater distance corresponds to a higher movement speed within that time frame.

However, due to the movement of the drone capturing the footage, abrupt changes in speed may occasionally occur as a result of sudden shifts in the drone's position rather than the elephant's movement. The primary objective of this analysis is to offer insights into individual elephants' movement speed patterns by assessing their speed variations over time. Notably, sharp spikes in velocity may indicate significant moments in the footage, potentially highlighting behaviors or external influences that require further investigation.

2.1.7.2. Average Elephant Movement Speed Plot

The *Average Elephant Movement Speed Plot* analysis does the same as the individual elephant distance analysis but averages the changes in central positions of the elephants to create a plot that visualizes the average movement speed of the entire group of elephants. This makes it easier to highlight moments that trigger a shift in movement speed across the entire group of detected elephants.

2.1.7.3. Elephant Movement Trajectories Plot

The *Elephant Movement Trajectories Plot* analysis visualizes the movement trajectories of the elephants by plotting the sequence of their center (x, y) coordinates over time. This

is done by constructing a trajectory for each unique elephant by connecting the center points from the bounding boxes for each frame. These trajectories provide a spatial representation of how each elephant moves through the duration of the video. As these trajectories are rendered in pixel coordinates, they reflect apparent movement within the frame and are not compensated for the drone's own motion. This analysis visualizes the spatial distribution of elephants by generating a Kernel Density Estimate (KDE) heatmap of their detected positions. This provides insights into the areas where elephants are most frequently observed throughout the video. The KDE is computed based on the (x, y) coordinates of the elephants, using Seaborn's *kdeplot* to estimate the density of their locations. The heatmaps can be affected by the motion of the drone however in videos that contain significant drone motion.

2.1.7.4. Visual Appearance Statistics

The Visual Appearance Statistics analysis calculates the statistics regarding how long each elephant is detected in the video and the total average among all elephants. This is done in exact frames and seconds, by summing up the amount of frames that each elephant is detected in for the frame count and dividing this by 30 to calculate the corresponding amount of seconds. These statistics offer insights into the persistence and visibility of each elephant within the video, potentially highlighting which elephants may be more interesting for further evaluation based on their visual presence in the video.

2.1.7.5. Elephant Overlap Statistics

The Elephant Overlap Statistics analysis identifies instances where the bounding boxes of different elephants overlap within the same frame, potentially indicating social interactions or close proximity. For each frame, all detected elephants are compared to determine if their bounding boxes overlap. An overlap is identified when the bounding boxes intersect along both the x and y axes using the formula below.

$$(1) \quad x_{\max}^{(i)} > x_{\min}^{(j)} \quad \text{and} \quad x_{\min}^{(i)} < x_{\max}^{(j)}$$

The overlapping pairs, along with their corresponding frame numbers, are recorded. The percentage of frames in which each elephant is involved in an overlap is then computed, and the results are saved as a CSV file for further analysis.

2.1.7.6. Elephant Cluster Statistics

The Elephant Cluster Statistics analysis identifies clusters of elephants that are spatially close to one another and tracks how these clusters persist over time. For each frame, the diagonal length of each elephant's bounding box is computed as a reference for spatial proximity. An average diagonal length per frame is calculated, and a threshold is set at 1.5 times this average. Elephants whose Euclidean distance falls below this threshold are grouped into clusters using a depth-first search algorithm. The continuity of these clusters is then tracked across consecutive frames to determine the time periods during which specific clustering patterns persist. The results, including the frame ranges of detected clusters, are saved as

a CSV file. This allows for automated detection of potential herds, sub herds which with further investigation can be used to find mother calf pairs or other insightful herds and dynamics.

2.1.7.7. Elephant Travel Distance Statistics

The Elephant Travel Distance Statistics analysis calculates the total Euclidean distance traveled by each elephant over the duration of the video. For each detected elephant, the analysis aggregates the total traveled distance by grouping the data by ID and summing the calculated Euclidean distance values per frame. The final summary, listing the total movement for each elephant in pixels, is then saved as a CSV file. The total travel distance serves as an important metric for assessing elephant movement. However, in videos with significant drone motion, the computed distances in pixels may reflect both the elephants' movement and the movement of the camera. As mentioned, this is a trade-off to ensure the framework's accessibility, and it should be considered when interpreting the results.

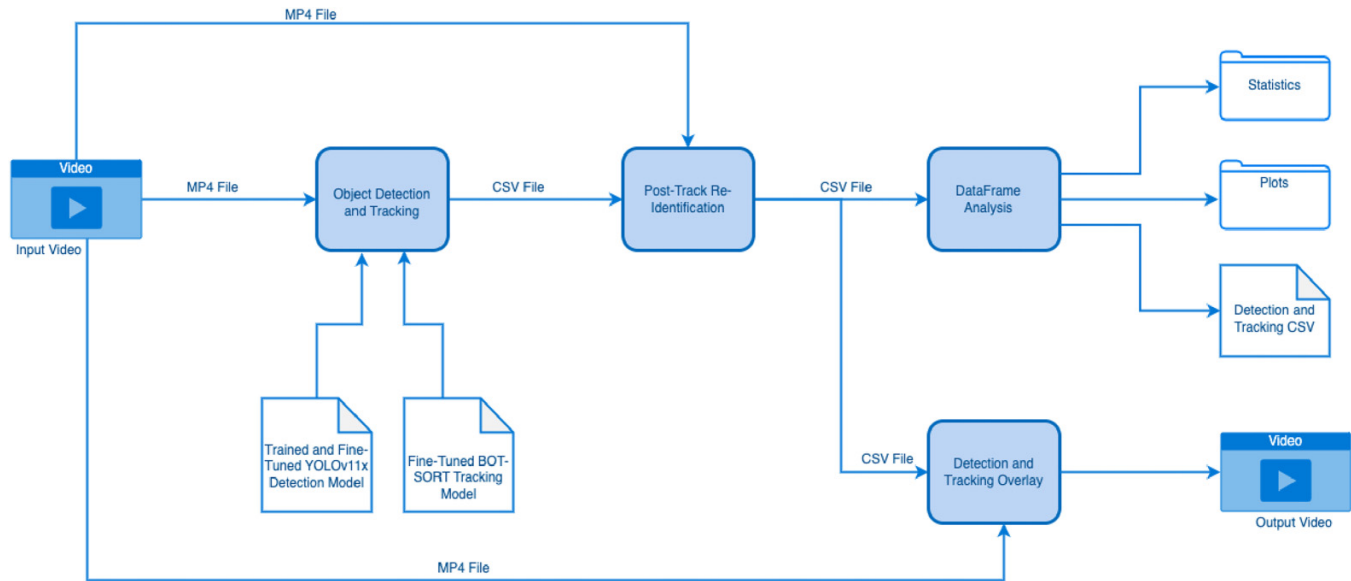
2.1.8. Framework overview

The framework consists of a pipeline that processes a single video input through multiple steps to generate a video output, featuring bounding boxes with unique IDs overlaid on the original video, along with detection and tracking data. This output includes various visualizations, such as speed analysis plots for each elephant, an aggregated plot showing the average speed across all elephants, a trajectory plot illustrating the movement paths of the elephants in one graph, and a density plot highlighting the locations where elephants spend the most time. Additionally, the framework provides statistical data, including the number of frames and seconds each elephant is visible, the percentage of frames in which an elephant overlaps with others, and the total distance traveled by each elephant (in pixels). The output also includes an analysis of group/herd dynamics, identifying which elephants remain together in groups or herds, and the frames during which this occurs. Finally, a CSV file is generated, containing tracking data for each frame.

This framework facilitates automated video analysis and elephant counting through its diverse outputs, significantly accelerating the work of ecologists and enabling the extraction of new insights from drone footage. An overview of the steps that make up the pipeline is provided below, with additional details illustrated in Fig. 3.

The first step in the pipeline involves object detection and tracking, where the trained and fine-tuned YOLOv11x model, in combination with the BoT-SORT tracker, processes the video input using *Python*. This step generates two outputs: a copy of the original video with detection and tracking results overlaid and a CSV file containing detailed detection data for each frame. The CSV file includes a row for each detected elephant, capturing the frame number, elephant ID, bounding box coordinates $(x_{\min}, x_{\max}, y_{\min}, y_{\max})$, and confidence score. While the video file is saved for visualization purposes, only the CSV file is used in subsequent steps.

Fig. 3. Framework pipeline schematic visualization.



Next, the CSV file is processed by the post-track re-identification algorithm, which updates the tracking information to refine the association of elephants across frames. The revised CSV file produced in this step is then passed to the next stage of the pipeline. Following this, the data undergoes detailed analysis using *Python*, generating meaningful plots and statistics related to elephant movement and behavior. These outputs, include individual elephant movements speed plots, an Average Elephant Movement Speed Plot, a combined Elephant Movement Trajectories Plot, Visual Appearance Statistics, overlap statistics, individual and average travel distance statistics, and finally cluster/herd statistics are saved in corresponding folders alongside the processed CSV file.

Finally, the updated CSV file is used to create a new visual overlay on the original video, integrating bounding boxes with corresponding IDs and detection confidence scores. This visualization is generated using a *Python* script that reconstructs the detection and tracking data, ensuring a comprehensive representation of elephant movements in the footage. The resulting video output, along with the various analytical outputs, provides a robust tool for understanding elephant behavior and movement patterns in drone footage.

2.2. Evaluation

We evaluated the elephant object detection performance using precision, recall, mAP50, mAP50-95, and F1 Score, on the train, test, and validation set frames (Powers 2020). The elephant tracking was evaluated using the association accuracy (AssA) metric, a standard measure in multiobject detection and tracking tasks that quantifies association consistency and, consequently, the effectiveness of the tracking component (Bernardin and Stiefelhagen 2008; Ristani et al. 2016; Luiten et al. 2020; Yu et al. 2023; Gao and Wang 2024). The post-track re-identification algorithm was also evaluated by comparing AssA and amount of unique elephant IDs per

video with the results before and after the implementation of the post-track re-identification algorithm. To do this for each video file a ground truth tracking file was created by hand by using the annotated bounding boxes data and adding unique IDs to each unique elephant.

mAP50:

$$(2) \quad mAP_{50} = \frac{1}{C} \sum_{c=1}^C AP_{50}^{(c)}$$

where

- C is the total number of object classes.
- $AP_{50}^{(c)}$ (average precision for class c at an Intersection over Union (IoU) threshold of 50%) is defined as

$$(3) \quad AP_{50}^{(c)} = \int_0^1 p_{50}^{(c)}(r) dr$$

- $p_{50}^{(c)}(r)$ denotes the precision as a function of recall r for class c when using an IoU threshold of 50% (Khanam and Hussain 2024; Ultralytics 2025).

mAP50-95:

$$(4) \quad mAP_{50-95} = \frac{1}{10} \sum_{k=1}^{10} AP_{t_k}$$

where

- $t_k = 0.5 + 0.05 \times (k - 1)$ for $k = 1, 2, \dots, 10$ represents the set of IoU thresholds from 50% to 95%.
- AP_{t_k} (average precision at IoU threshold t_k) is defined as

$$(5) \quad AP_{t_k} = \int_0^1 p_{t_k}(r) dr$$

Table 3. Detection evaluation metrics.

Set	Precision	Recall	mAP50	mAP50-95	F1
Validation	0.967	0.965	0.982	0.827	0.966
Training	0.973	0.987	0.989	0.882	0.980
Test	0.960	0.971	0.988	0.839	0.966
Average	0.967	0.974	0.986	0.849	0.971

Note: mAP, mean average precision.

Table 4. Tracking results without post-track re-identification algorithm.

Video name	Amount of frames	Amount of IDs	Ground truth (GT) amount of IDs	AssA
0391	6868	12	6	0.819
0392	3967	11	5	0.761
0393	6864	24	11	0.633
0394	6872	47	23	0.683
0395	1853	9	7	0.909
0404	6825	37	15	0.745
0405	6870	20	9	0.898
0406	2849	6	5	0.999
Average	5371	20.700	10.125	0.806

Note: AssA, association accuracy.

- $p_{t_k}(r)$ denotes the precision as a function of recall r for a given IoU threshold t_k (Khanam and Hussain 2024; Ultralytics 2025).

AssA:

$$(6) \quad \text{AssA} = \frac{\text{Correctly Associated Pairs}}{\text{Total Number of Associations}}$$

where

- Correctly associated pairs are pairs of detections that are correctly identified as the same object across consecutive frames.
- Total number of associations (TNAs) is the TNAs that the tracking algorithm makes, including both correct and incorrect associations (Luiten et al. 2020).

F1:

$$(7) \quad \text{F1} = \frac{\text{TP}}{\text{TP} + \frac{1}{2}(\text{FP} + \text{FN})}$$

where

- TP is correctly assigned objects.
- FP is incorrectly assigned objects.
- FN is missed objects.

(Ristani et al. 2016)

2.3. Ethics statement

Drone flights and data collection were conducted with approval from the Welgevonden Game Reserve management. All procedures fell under the general permission to fly drones

for animal observation at Liverpool John Moores University and were performed in accordance with its institutional animal care and ethics policies.

3. Results

3.1. Detection results

The fine-tuned YOLOv11x detection model demonstrated strong performance in detecting elephants from drone imagery. As shown in Table 3, on the validation set the model achieved a precision of 0.967, indicating high accuracy in identifying elephants with minimal false positives. The recall of 0.965 suggests that the model successfully detects nearly all actual instances, ensuring reliable detection. This is further supported by the F1 score of 0.966, reflecting a balanced trade-off between precision and recall. Additionally, the model attained a mAP50 of 0.982, confirming its ability to localize elephants effectively under moderate overlap conditions. The mAP50-95 score of 0.827 further demonstrates the model's robustness under stricter localization criteria.

3.2. Tracking results

Table 4 shows the tracking results for our fine-tuned BoT-SORT tracker, the tracker achieved an average AssA score of 0.806 across all evaluated video sequences. This score reflects a strong capacity for identity preservation, suggesting that BoT-SORT is well-suited for maintaining coherent object trajectories under relatively stable visual conditions.

Despite this promising overall performance, notable variation in tracking quality is observed between different videos. The lowest AssA scores are found in sequences containing frequent and abrupt scene transitions—particularly zoom-ins and zoom-outs—which significantly alter both the spatial and

Table 5. Tracking results with post-track re-identification algorithm.

Video name	Amount of frames	Amount of IDs	GT amount of IDs	AssA
0391	6868	10	6	0.875
0392	3967	9	5	0.762
0393	6864	13	11	0.958
0394	6872	31	23	0.772
0395	1853	7	7	1.000
0404	6825	21	15	0.971
0405	6870	11	9	0.959
0406	2849	6	5	0.999
Average	5371	13.500	10.125	0.912

Note: AssA, association accuracy.

visual characteristics of the scene. This can be seen in videos 0393, 0394, 0404, and 0392. These disruptions hinder the tracker's ability to maintain consistent object associations, a known limitation of conventional tracking models that lack mechanisms for robust adaptation to rapid changes in perspective or scale. Additionally, videos in which elephant subjects temporarily disappear—due to occlusion by vegetation or moving outside the frame—and reappear after extended gaps also exhibit decreased performance. In such cases, the tracker often fails to reassociate the reappearing elephant with its original ID, instead assigning a new ID and thereby inflating the apparent number of individuals. This leads to an average difference between the ground truth amount of elephants detected and model output of 10.575.

These results highlight three key findings. First, BoT-SORT demonstrates a strong baseline capability for tracking elephants in aerial drone footage, provided that the video remains relatively continuous and free from abrupt scene changes. Second, the tracking performance is highly sensitive to sudden camera movements, particularly zoom operations, which should be minimized in future data collection efforts to preserve tracking integrity. Third, extended occlusions—such as those caused by dense foliage or long absences from the frame—pose a significant challenge to identity continuity, underscoring the need for additional post-processing steps, such as re-identification algorithms, to recover lost associations and improve the overall reliability of elephant counting in ecological monitoring applications.

The integration of the custom post-track re-identification algorithm substantially enhances the performance of the BoT-SORT tracker, addressing several of its key limitations in standalone operation. As presented in Table 5, the average AssA score increases to 0.912 following the application of the re-identification step—a 10.6% improvement compared to the pre-processing results. This increase in AssA reflects a more consistent preservation of object identities across frames, reinforcing the algorithm's value in correcting erroneous ID switches. Notably, videos 0393 and 0404 exhibit significant improvements in tracking accuracy, with large reductions in the number of unique IDs detected. These values now more closely align with the ground truth, indicating a reduced incidence of ID fragmentation and a corresponding increase in tracking reliability.

Despite this overall improvement, the limitations imposed by abrupt scene transitions—particularly zoom-ins and zoom-outs—remain evident. Such transitions drastically alter the spatial and visual features leveraged by the tracker, introducing inconsistencies that even the re-identification algorithm struggles to resolve. Nevertheless, in videos that do not suffer from such disturbances, the benefits of the re-identification algorithm are striking. For instance, videos 0395 and 0406 achieve near-perfect or perfect tracking, with AssA scores of 1.000 and 0.999, respectively. These sequences feature smooth camera motion and limited occlusion, demonstrating that the algorithm performs exceptionally well under favorable recording conditions, even when elephants temporarily disappear behind foliage or move briefly out of frame due to gradual panning.

On average, the difference between the number of detected unique elephant IDs and the ground truth decreases to 3.3375 following the implementation of the re-identification step, an improvement compared to the pre-processing difference of 10.757. These results prove the impact of the post-track re-identification algorithm on tracking consistency and accuracy in the videos.

3.3. Analysis results

All analyses in this section were produced automatically by our framework's analysis module applied to Video 0395, a continuous 61.8 s (1854-frame) aerial recording of seven individually identified elephants (IDs 1, 2, 4, 5, 6, 8, 10). We report metrics on appearance duration, spatial overlap, cumulative travel distance, temporal clustering of group composition, instantaneous speed profiles, and spatial trajectories, along with summary statistics and parameter details to ensure full reproducibility.

In Table 6, we report each elephant's visibility expressed both in absolute frame count and in seconds. Elephants 1, 2, 5, and 6 are detected in every frame (1854 frames; 100%; 61.8 ± 0.0 s), demonstrating uninterrupted coverage. Elephant 4 exhibits only minor drop-outs, appearing in 1850 frames (99.8%; 61.7 ± 0.1 s). By contrast, elephant 8 is visible for 1779 frames (96.0%; 59.3 ± 1.5 s) and elephant 10 for 1618 frames (87.3%; 53.9 ± 3.0 s). Across all individuals, the mean visibility is 1809 frames (97.6%; 60.3 ± 2.1 s), with a standard deviation

Table 6. Visual Appearance Statistics video 0395.

Elephant ID	Frame count	Seconds
1.0	1854	61.80
2.0	1854	61.80
4.0	1850	61.67
5.0	1854	61.80
6.0	1854	61.80
8.0	1779	59.30
10.0	1618	53.93
Average	1809.00	60.30

Table 7. Elephant Overlap Statistics.

ID	Overlap percentage
10	84.574%
4	65.912%
6	52.643%
5	52.211%
2	35.922%
1	3.937%
8	0.000%
Average	42.186%

Table 8. Elephant Travel Distance Statistics.

ID	Distance
1	7710.265 px
2	7983.044 px
4	6798.674 px
5	7475.164 px
6	7628.081px
8	6260.099 px
10	8002.789 px
Average	7408.302 px

of 93 frames (5.0 s), indicating consistently high track retention throughout the recording.

Table 7 quantifies spatial overlap by calculating the proportion of each elephant's visible frames in which its bounding box intersects that of at least one other herd member. Elephant 10 displays the highest overlap rate at 84.6%, followed by elephants 4 and 6 at 65.9% and 52.6%, respectively. Elephant 8 registers no overlap (0%), confirming its peripheral positioning. The group mean overlap rate is 42.2% with a standard deviation of 28.5%, reflecting heterogeneous inter-individual spacing patterns.

In **Table 8**, cumulative travel distances are computed by summing the Euclidean displacement between successive frames for each individual, reported in pixel units. Elephant 10 traverses the greatest path length of 8002.8 px, while elephant 8 covers the shortest distance of 6260.1 px. The mean travel distance across all elephants is 7408.3 px (SD = 527.3 px), suggesting modest variability in movement magnitude that may derive from both behavioral differences and camera parallax.

Table 9 details the results of a frame-wise clustering analysis performed to detect stable herd compositions. Six elephants (excluding ID 8) form a core cluster during most intervals: frames 0–375, 797–1190, 1191–1496, 1497–1618, and 1619–1853. A transient reconfiguration occurs in frames 376–796, during which elephant 1 briefly joins elephant 8 in a secondary grouping. Interval durations vary between 122 and 421 frames, illustrating both prolonged cohesion and short-term fission events.

Instantaneous speed for each elephant is calculated by dividing frame-to-frame displacement by the inter-frame interval (0.033 s). As shown in the top of **Fig. 4**, the herd's mean speed trace fluctuates around a baseline of 20 px/s, with two pronounced peaks reaching approximately 45 px/s at 15 and 45 s. The bottom of **Fig. 4** presents individual speed trajectories, which exhibit high temporal correlation with the group mean (mean cross-correlation $r = 0.92$), and indicate that elephants 2 and 10 lead acceleration events by 0.2–0.4 s.

Figure 5 overlays the two-dimensional spatial trajectories of all elephants in image coordinates. The predominant path follows a linear corridor from the lower-left to the upper-right portion of the frame, with lateral dispersion of ± 150 px around the central axis. Elephant 8's trajectory deviates by more than 200 px laterally, corroborating its peripheral role as evidenced by the overlap and distance metrics.

4. Discussion

Our end-to-end pipeline for aerial elephant monitoring integrates three key components—YOLOv11x for detection, BoT-SORT for tracking, and a bespoke post-track re-identification module—to deliver both high accuracy and robust identity continuity. In the detection stage, YOLOv11x attains precision ≥ 0.96 , recall ≥ 0.965 , and mAP50 ≥ 0.982 across all splits, corroborating recent advances in single-stage detectors for wildlife monitoring (Khanam and Husain 2024; Wang and Liao 2024). Compared to earlier findings that single-stage models can struggle with small or occluded targets (Kellenberger et al. 2018), our results suggest that YOLOv11x's enhanced attention mechanisms and neck design substantially mitigate these shortcomings.

The BoT-SORT tracker alone yields an AssA of 0.806, consistent with its performance in pedestrian domains (Ristani et al. 2016; Aharon et al. 2022). However, abrupt drone maneuvers and prolonged occlusions still induce fragmentation and identity switches, mirroring challenges reported in aerial bird tracking (Mpouziotas et al. 2024). Our post-track re-identification algorithm, which reunites fragmented tracks via spatial continuity and appearance similarity heuristics, raises mean AssA to 0.912 and cuts ID-count errors by 68%. This lightweight approach parallels deep-metric methods (Wojke et al. 2017) but avoids the heavy data and compute demands of end-to-end embedding training.

Beyond technical metrics, the framework has operational benefits for ecologists. Traditional manual annotation of herd videos is time-consuming and error-prone, often requiring frame-by-frame labelling (Delplanque et al. 2023a). By automating detection, tracking, and re-identification, our system dramatically reduces human labour, enabling broader

Table 9. Elephant Cluster Statistics.

Clusters	Frame ranges
[1,2,4,5,6,10] [8]	[0–375] [1662–1663] [1676–1697][1702–1712] [1716–1718]
[1] [2,4,5,6,10] [8]	[376–796] [818–1260] [1262–1264] [1270–1273] [1619–1666]
[1] [2,4,5,6,10]	[797–817]
[1] [2] [4,5,6] [8]	[1261–1261] [1265–1269] [1274–1440]
[1] [2] [5] [4,6] [8]	[1441–1490]
[1] [2] [10] [5] [4,6] [8]	[1491–1496]
[1] [2,4,6,10] [5] [8]	[1497–1618]
[1,2,4,5,6] [8]	[1667–1675] [1698–1701]
[1,2,5,6,10] [4] [8]	[1713–1715] [1719–1799]
[1,2,5,6,10] [4]	[1800–1853]

Fig. 4. Average Elephant Movement Speed Plot and Individual Movement Speed Plot.

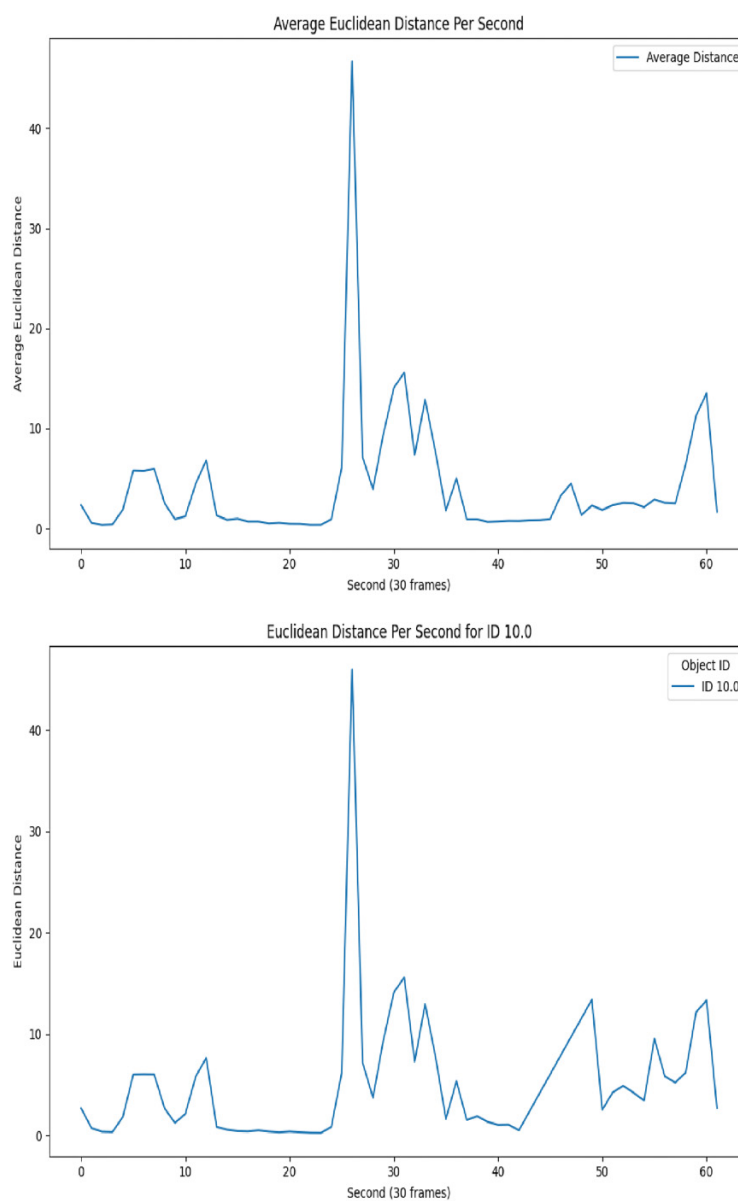
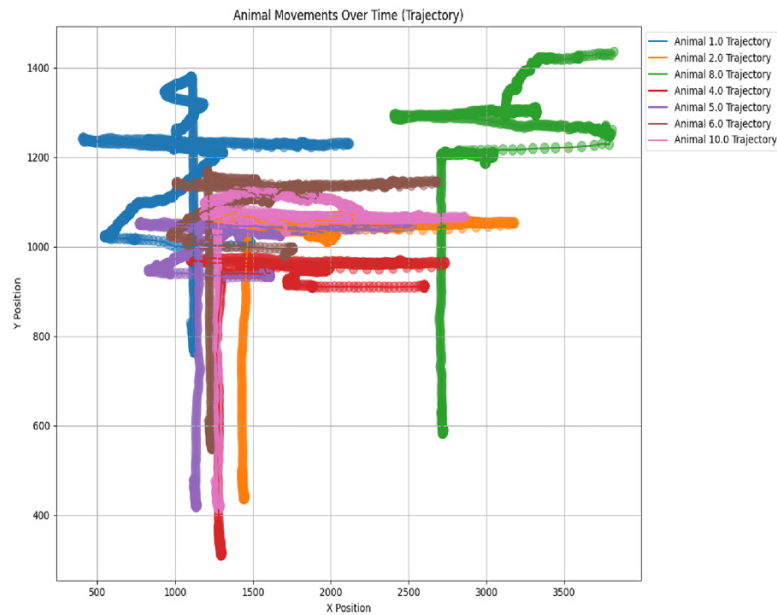


Fig. 5. Elephant Movement Trajectories Plot.



surveys, and more frequent sampling of elephant populations. For example, time savings on a single 1 h flight can translate into multiple additional flights per field season, allowing researchers to detect emergent behaviors such as sudden range shifts or drought-induced dispersals with minimal delay.

Moreover, standardized deployment of our pipeline across different reserves can facilitate multisite meta-analyses. As noted by Kellenberger et al. (2018), variability in model performance on imbalanced datasets hinders comparisons; our demonstration of YOLOv11x's robustness suggests that a unified detection-tracking framework could serve as a common baseline for inter-regional studies of movement ecology and social structure.

Our automated extraction of behavioral metrics opens new avenues in social and spatial ecology. Pairwise overlap and clustering analyses reveal fission–fusion dynamics and subgroup formation that might elude manual observation, while trajectory heatmaps identify preferred travel corridors akin to the habitat-use insights obtained from avian studies (Mpouziotas et al. 2024). Metrics such as distances between individuals, individual tracks, and travel speed of individuals and the herd are all useful to understand animal movement behaviour which is an important field of study (Boinski and Garber 2000) and for which ground observations have been used (Dai et al. 2007) in addition to VHF or satellite tracking for elephants (Tchamba et al. 1995). Recently, drones have started to be used to derive such metrics either manual or by using automated analyses (Inoue et al. 2019; Koger et al. 2023; Schad and Fischer 2023). Measuring animals' speed can also be used to determine the influence a drone might have on animals as it gets closer, and thus it would be useful as a way to measure animal disturbance by the drone through the images the drone itself obtains. Integrating heatmaps with habitat features—such as water sources or vegetation indices—

could further elucidate resource-driven movement patterns, informing targeted conservation interventions.

Despite these strengths, several limitations remain. First, without drone pose or GPS/IMU data, our movement estimates are in pixel units and can overestimate true displacement when the camera itself moves (Zhao et al. 2024). While established techniques for motion compensation exist, they were deliberately excluded to maintain the framework's accessibility and ease of use. Typically, this is achieved through visual-based methods, like optical flow, which track how static background elements move between frames to model the camera's motion, or through sensor-based methods that use the drone's own telemetry (GPS and IMU data) for a direct measurement of its movement.

However, integrating these techniques would introduce the significant technical barriers we sought to avoid. Mandating camera calibration for visual methods or the integration and validation of telemetry data would limit the framework's versatility, as it requires complex setups like Ground Control Points and detailed terrain maps. Such requirements make the process less practical for researchers in the field and would prevent the framework from being a generalizable, “plug-and-play” tool. Furthermore, the telemetry from many consumer-grade drones lacks the accuracy needed for reliable real-world speed calculations, and relying on it could create a false sense of accuracy. Our current approach is therefore a deliberate trade-off, ensuring the framework remains a practical tool for a broader user base.

Second, extreme viewpoint shifts or extended occlusions can still fragment tracks; future incorporation of transformer-based memory modules may enhance long-term appearance retention (Gao and Wang 2024). Third, our current focus on localization and tracking leaves fine-grained behavior recognition—such as foraging, social interactions, or stress indicators—as a topic for further study, poten-

tially leveraging 3D pose estimation from oblique drone imagery (Shukla et al. 2024).

Looking forward, integrating nonconsumer-grade drone-mounted inertial/GPS sensors will yield georeferenced tracks for absolute movement metrics and home-range estimation (Zhao et al. 2024). To further enhance identity continuity, future iterations of our re-identification module could draw on adaptive appearance-model management strategies such as those proposed by Cho and Kim (2023). By dynamically updating per-target appearance galleries and incorporating confidence-weighted template selection, such an approach would better handle gradual appearance changes and mitigate drift during long occlusions. Embedding these concepts into our lightweight post-track re-identification stage could reduce residual ID fragmentation without imposing significant computational overhead. In the longer term, extending the framework to fine-grained behavior recognition and 3D pose estimation from oblique imagery will enable automated classification of foraging, social interactions, and stress behaviors (Shukla et al. 2024).

By harnessing advances in detection, tracking, and lightweight re-identification, this pipeline turns drone footage into actionable intelligence—empowering wildlife stewards to count, monitor, and protect elephant populations at a reduced manual labor cost.

Acknowledgements

The authors thank Carmen Warmenhove and Jonathan Swart for their invaluable support during the drone flights at the Welgevonden Game Reserve. The authors also acknowledge the use of Gemini 2.5 Pro for assistance in checking for grammar and spelling mistakes and for reformulation of the text in the preparation of this manuscript.

Article information

History dates

Received: 6 June 2025

Accepted: 30 September 2025

Accepted manuscript online: 7 October 2025

Version of record online: 25 November 2025

Copyright

© 2025 The Authors. This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/) (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

Data availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

Author information

Author ORCIDs

Chaim Chai Elchik <https://orcid.org/0009-0005-0254-5244>

André Burger <https://orcid.org/0000-0003-3954-5174>

Author notes

Serge Wich served as an Associate Editor at the time of manuscript review and acceptance; peer review and editorial decisions regarding this manuscript were handled by another Editorial Board Member.

Author contributions

Conceptualization: CCE

Data curation: CCE

Formal analysis: CCE

Funding acquisition: SW

Investigation: CCE

Methodology: CCE

Project administration: SW

Resources: AB, SW

Software: CCE

Supervision: SW

Validation: CCE

Visualization: CCE

Writing – original draft: CCE

Writing – review & editing: CCE, AB, SW

Competing interests

The authors declare that there are no competing interests.

Funding information

This research was supported by Liverpool John Moores University (LJMU).

References

- Afridi, S., Laporte-Devyllder, L., Maalouf, G., Kline, J.M., Penny, S.G., Hlebowicz, K., et al. 2025. Impact of drone disturbances on wildlife: a review. *Drones*, **9**(4): 311. doi:10.3390/drones9040311.
- Aharon, N., Orfaig, R., and Bobrovsky, B.Z. 2022. BoT-SORT: robust associations multi-pedestrian tracking. arXiv:2206.14651 [cs.CV]. Available from <https://arxiv.org/abs/2206.14651> [accessed 15 April 2025].
- Alsaidi, M., Al-Jassani, M.G., Bang, C., O'Corry-Crowe, G.M., Watt, C., Ghazal, M., and Zhuang, H. 2024. Localization and tracking of beluga whales in aerial video using deep learning. *Front. Mar. Sci.* **11**: 1445698. doi:10.3389/fmars.2024.1445698.
- Barbedo, J.G.A., Koenigkan, L.V., Santos, T.T., and Santos, P.M. 2019. A study on the detection of cattle in UAV images using deep learning. *Sensors*, **19**(24): 5436. doi:10.3390/s19245436.
- Berger-Tal, O., and Lahoz-Monfort, J.J. 2018. Conservation technology: the next generation. *Conserv. Lett.* **11**(6): e12458. doi:10.1111/conl.12458.
- Bernardin, K., and Stiefelhagen, R. 2008. Evaluating multiple object tracking performance: the clear MOT metrics. *EURASIP J. Image Video Proc.* **2008**: 1–10.
- Boinski, S., and Garber, P.A. 2000. *On the move: how and why animals travel in groups*. University of Chicago Press.
- Brickson, L., Zhang, L., Vollrath, F., Douglas-Hamilton, I., and Titus, A.J. 2023. Elephants and algorithms: a review of the current and future role of AI in elephant monitoring. *J. R. Soc. Interface*, **20**(197): 20230367. doi:10.1098/rsif.2023.0367.
- Cho, Y.J., and Kim, D. 2023. Rethinking multi-object tracking based on re-identification and appearance model management. *IEEE Access*, **11**: 54337–54351. doi:10.1109/ACCESS.2023.3274662.
- Dai, X., Shannon, G., Slotow, R., Page, B., and Duffy, K.J. 2007. Short-duration daytime movements of a cow herd of African elephants. *J. Mammal.* **88**(1): 151–157. doi:10.1644/06-MAMM-A-035R1.1.
- Dave, B., Mori, M., Bathani, A., and Goel, P. 2023. Wild animal detection using YOLOv8. *Procedia Comput. Sci.* **230**: 100–111. doi:10.1016/j.procs.2023.12.065.

- Delplanque, A., Foucher, S., Lejeune, P., Linchant, J., and Théau, J. 2021. Multispecies detection and identification of African mammals in aerial imagery using convolutional neural networks. *Remote Sens. Ecol. Conserv.* **8**. doi:10.1002/rse2.234.
- Delplanque, A., Foucher, S., Théau, J., Bussière, E., Vermeulen, C., and Lejeune, P. 2023a. From crowd to herd counting: how to precisely detect and count African mammals using aerial imagery and deep learning? *ISPRS J. Photogramm. Remote Sens.* **197**, 167–180. doi:10.1016/j.isprsjprs.2023.01.025.
- Delplanque, A., Lamprey, R., Foucher, S., Théau, J., and Lejeune, P. 2023b. Surveying wildlife and livestock in Uganda with aerial cameras: deep learning reduces the workload of human interpretation by over 70%. *Front. Ecol. Evol.* **11**. doi:10.3389/fevo.2023.1270857.
- Delplanque, A., Linchant, J., Vincke, X., Lamprey, R., Théau, J., Vermeulen, C., et al. 2024. Will artificial intelligence revolutionize aerial surveys? A first large-scale semi-automated survey of African wildlife using oblique imagery and deep learning. *Ecol. Inform.* **82**, 102679. doi:10.1016/j.ecoinf.2024.102679.
- Fang, C., Li, C., Yang, P., Kong, S., Han, Y., Huang, X., and Niu, J. 2024. Enhancing livestock detection: an efficient model based on YOLOv8. *Appl. Sci.* **14**(11). doi:10.3390/app14114809.
- Finn, C., Grattarola, F., and Pincheira-Donoso, D. 2023. More losers than winners: investigating Anthropocene defaunation through the diversity of population trends. *Biol. Rev. Cambridge Philos. Soc.* **98**. doi:10.1111/brv.12974.
- Gao, R., and Wang, L. 2024. Memotr: long-term memory-augmented transformer for multi-object tracking. arXiv:2307.15700 [cs.CV]. Available from <https://arxiv.org/abs/2307.15700> [accessed 10 April 2025].
- Guirado, E., Tabik, S., Rivas, M.L., Alcaraz-Segura, D., and Herrera, F. 2019. Whale counting in satellite and aerial images with deep learning. *Sci. Rep.* **9**(1): 14271. doi:10.1038/s41598-019-50795-9.
- Hamilton, G., Corcoran, E., Denman, S., Hennekam, M.E., and Koh, L.P. 2020. When you can't see the koalas for the trees: using drones and machine learning in complex environments. *Biol. Conserv.* **247**: 108598. doi:10.1016/j.biocon.2020.108598.
- Hodgson, J.C., Baylis, S.M., Mott, R., Herrod, A., and Clarke, R.H. 2016. Precision wildlife monitoring using unmanned aerial vehicles. *Sci. Rep.* **6**: 22574. doi:10.1038/srep22574.
- Inoue, S., Yamamoto, S., Ringhofer, M., Mendonça, R.S., Pereira, C., and Hirata, S. 2019. Spatial positioning of individuals in a group of feral horses: a case study using drone technology. *Mammal Res.* **64**(2): 249–259.
- IPBES. 2019. Global assessment report on biodiversity and ecosystem services of the Intergovernmental Science-Policy Platform on Biodiversity and Ecosystem Services. Version 1. doi:10.5281/zenodo.6417333.
- Kellenberger, B., Marcos, D., and Tuia, D. 2018. Detecting mammals in UAV images: best practices to address a substantially imbalanced dataset with deep learning. **216**: 139–153. doi:10.1016/j.rse.2018.06.028. www.sciencedirect.com/science/article/pii/S0034425718303067.
- Khanam, R., and Hussain, M. 2024. YOLOv11: an overview of the key architectural enhancements. arXiv:2410.17725 [cs.CV]. Available from <https://arxiv.org/abs/2410.17725> [accessed 10 April 2025].
- Kissling, W.D., Evans, J.C., Zilber, R., Breeze, T.D., Shinneman, S., Schneider, L.C., et al. 2024. Development of a cost-efficient automated wildlife camera network in a European Natura 2000 site. *Basic Appl. Ecol.* **79**: 141–152. doi:10.1016/j.baec.2024.06.006.
- Koger, B., Deshpande, A., Kerby, J.T., Graving, J.M., Costelloe, B.R., and Couzin, I.D. 2023. Quantifying the movement, behaviour and environmental context of group-living animals using drones and computer vision. *J. Anim. Ecol.* **92**(7): 1357–1371. doi:10.1111/1365-2656.13904.
- Lamba, A., Cassey, P., Segaran, R.R., and Koh, L.P. 2019. Deep learning for environmental conservation. *Curr. Biol.* **29**(19): R977–R982. doi:10.1016/j.cub.2019.08.016.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., and Berg, A.C. 2016. SSD: Single Shot Multibox Detector. In *Computer Vision - ECCV 2016*. Springer International Publishing. pp. 21–37. doi:10.1007/978-3-319-46448-0_2.
- López, J.J., and Mulero-Pázmány, M. 2019. Drones for conservation in protected areas: present and future. *Drones*, **3**(1): 10. doi:10.3390/drones3010010.
- Luiten, J., Osep, A., Dendorfer, P., Torr, P., Geiger, A., Leal-Taixé, L., and Leibe, B. 2020. HOTA: a higher order metric for evaluating multi-object tracking. *Int. J. Comput. Vision*, **129**(2): 548–578. doi:10.1007/s11263-020-01375-2.
- McEvoy, J.F., Hall, G.P., and McDonald, P.G. 2016. Evaluation of unmanned aerial vehicle shape, flight path and camera type for waterfowl surveys: disturbance effects and species recognition. *PeerJ*, **4**: e1831. doi:10.7717/peerj.1831.
- Mou, C., Liu, T., Zhu, C., and Cui, X. 2023. WAID: a large-scale dataset for wildlife detection with drones. *Appl. Sci.* **13**(18): 10397. doi:10.3390/app131810397.
- Mpouziotas, D., Karvelis, P., and Stylios, C. 2024. Advanced computer vision methods for tracking wild birds from drone footage. *Drones*, **8**(6): 259. doi:10.3390/drones8060259.
- Mulero-Pázmány, M., Jenni-Eiermann, S., Strebel, N., Sattler, T., Negro, J.J., and Tablado, Z. 2017. Unmanned aircraft systems as a new source of disturbance for wildlife: a systematic review. *PLoS One*, **12**(6): e0178448. doi:10.1371/journal.pone.0178448.
- Pedrazzi, L., Naik, H., Sandbrook, C., Lurgi, M., Fürtbauer, I., and King, A.J. 2025. Advancing animal behaviour research using drone technology. *Anim. Behav.* **222**: 123147. doi:10.1016/j.anbehav.2025.123147.
- Powers, D.M. 2020. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. <https://arxiv.org/abs/2010.16061>.
- Rančić, K., Blagojević, B., Bezdan, A., Ivošević, B., Tubić, B., Vranešević, M., et al. 2023. Animal detection and counting from UAV images using convolutional neural networks. *Drones*, **7**(3): 179. doi:10.3390/drones7030179.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. 2016. You Only Look Once: unified, real-time object detection. arXiv:1506.02640 [cs.CV]. <https://arxiv.org/abs/1506.02640>.
- Ren, S., He, K., Girshick, R., and Sun, J. 2016. Faster R-CNN: towards real-time object detection with region proposal networks. arXiv:1506.01497 [cs.CV].
- Ristani, E., Solera, F., Zou, R.S., Cucchiara, R., and Tomasi, C. 2016. Performance measures and a data set for multi-target, multi-camera tracking. arXiv:1609.01775 [cs.CV].
- Roboflow. 2025. Documentation: build vision models with Roboflow. Available from <https://docs.roboflow.com/> [accessed 6 March 2025].
- Schad, L., and Fischer, J. 2023. Opportunities and risks in the use of drones for studying animal behaviour. *Methods Ecol. Evol.* **14**(8): 1864–1872.
- Shukla, V., Morelli, L., Remondino, F., Micheli, A., Tuia, D., and Risse, B. 2024. Towards estimation of 3D poses and shapes of animals from oblique drone imagery. In *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, Volume XLVIII-2-2024. pp. 379–386. doi:10.5194/isprs-archives-XLVIII-2-2024-379-2024.
- Tchamba, M., Bauer, H., and IONGH, H.D. 1995. Application of VHF-radio and satellite telemetry techniques on elephants in northern cameroon. *Afr. J. Ecol.* **33**(4): 335–346.
- Ultralytics. 2025. Ultralytics documentation: YOLO performance metrics. Available from <https://docs.ultralytics.com/guides/yolo-performance-metrics/> [accessed 6 March 2025].
- Varghese, R., and Sambath, M. 2024. YOLOv8: a novel object detection algorithm with enhanced performance and robustness. In *2024 International Conference on Advances in Data Engineering and Intelligent Computing Systems (ADICS)*. pp. 1–6. doi:10.1109/ADICS58448.2024.10533619.
- Vermeulen, C., Lejeune, P., Lisein, J., Sawadogo, P., and Bouché, P. 2013. Unmanned aerial survey of elephants. *PLoS One*, **8**(1): e54700. doi:10.1371/journal.pone.0054700.
- Wang, C.Y., and Liao, H.Y.M. 2024. YOLOv1 to YOLOv10: the fastest and most accurate real-time object detection systems. arXiv:2408.09332[cs.CV].
- Wich, S.A., and Piel, A.K. (Editors). 2021. *Conservation technology*. Oxford University Press, Oxford. Available from <https://global.oup.com/academic/product/conservation-technology-9780198850243> [accessed 11 April 2025].
- Wojke, N., Bewley, A., and Paulus, D. 2017. Simple online and realtime tracking with a deep association metric. arXiv:1703.07402 [cs.CV].
- WWF. 2024. *Living planet report 2024—a system in peril*. WWF, Gland, Switzerland.

Yaseen, M. 2024. What is YOLOv8: an in-depth exploration of the internal features of the next-generation object detector. arXiv:2408.15857 [cs.CV].

Yu, E., Wang, T., Li, Z., Zhang, Y., Zhang, X., Tao, W., et al. 2023. MOTRv3: release-fetch supervision for end-to-end multi-object tracking. arXiv preprint arXiv:2305.14298. arXiv:2305.14298 [cs.CV].

Zhang, Y., Sun, P., Jiang, Y., Yu, D., Weng, F., Yuan, Z., et al. 2022. ByteTrack: multi-object tracking by associating every detection box. arXiv:2110.06864 [cs.CV].

Zhao, X., Huang, X., Cheng, J., Xia, Z., and Tu, Z. 2024. A vision-based end-to-end reinforcement learning framework for drone target tracking. Drones, 8(11). doi:10.3390/drones8110628.

Appendix A

Table A1. Key technical specifications of the DJI Mavic 3 Pro.

Characteristic	Specification
Hasselblad camera	
Sensor	4/3 CMOS, 20 MP
Lens FOV	84°
Equivalent focal length	24 mm
Aperture	f/2.8 to f/11 (adjustable)
ISO range (video)	100–12 800
Shutter speed	8–1/8000 s
Video and imaging	
Max video resolution	5.1K: 5120 × 2700 @ 50fps DCI 4K: 4096 × 2160 @ 120fps 4K: 3840 × 2160 @ 120fps
Video formats	MP4/MOV (MPEG-4 AVC/H.264, HEVC/H.265) Apple ProRes 422 HQ, 422, 422 LT (Cine Model)
Color profiles	Normal, HLG, 10-bit D-Log M
Max video bitrate	H.264/H.265: 200 Mbps
Digital zoom	Hasselblad Camera: 1-3× Medium Tele Camera: 3-7× Tele Camera: 7-28×
Gimbal	
Stabilization	3-axis mechanical (tilt, roll, pan)
Mechanical range	Tilt: –135° to 100° Roll: –45° to 45° Pan: –27° to 27°
Controllable range	Tilt: –90° to 35°
Max control speed (tilt)	100°/s

Drone Syst. Appl. Downloaded from cdnsciencepub.com by LIVERPOOL JOHN MOORES UNIV on 06/10/26