**ChatGPT**

# Temporal Modeling for Reliable Drone Video Object Detection and Tracking

Modern drone applications increasingly leverage video-based object detection and tracking. Relying on single-frame detectors alone often leads to inconsistent or spurious predictions due to motion blur, occlusions, and changing viewpoints. Recent research has therefore explored how **temporal information** can be explicitly used to improve robustness and trustworthiness in drone vision. Short-term cues (frame-to-frame consistency) help smooth detections, long-term patterns (object trajectories, motion context) enable recovery of missed targets, and evolving scene context aids filtering of false alarms. Combining these temporal signals can greatly reduce false positives/negatives and produce more stable outputs.

## Short-Term Temporal Cues

A common strategy is to **aggregate information across adjacent frames**. For example, Xiao *et al.* (2023) introduced a *Temporal Feature Aggregation Module* (TFAM) that warps features from previous frames into the current frame via learned offsets (using deformable convolutions) and fuses them [1] [2]. TFAM effectively "remembers" objects that were briefly occluded or blurred, rescuing them in the current detection. In their UAV multi-object tracker, TFAM integration led to a ~2–3% absolute gain in MOTA on benchmarks (VisDrone2019, UAVDT) by recovering missed detections [3] [2]. Similarly, attention-augmented recurrent units can accumulate frame signals: Zhou *et al.* (2023) designed a **Temporal Attention GRU** (TA-GRU) add-on to YOLOv7 (see below). TA-GRU uses gated RNN layers with spatial-temporal attention to aggregate neighboring-frame features, which boosted VisDrone-VID mAP by about 5.9% [4] (with negligible speed cost).

*Figure: Temporal Attention GRU (TA-GRU) module from Zhou* et al. *(2023) [41]. The RNN-based TA-GRU combines features from adjacent frames via attention and gating, enhancing YOLO detection in video. (Image omitted credit.)*

Optical flow and motion-based linking also enforce short-term consistency. Many trackers perform *tracking-by-detection*, associating objects across frames to enforce temporal coherence. Guanxiong *et al.* (ECCV 2022) exploited "gradual change" priors: a **Location Prior Network (LPN)** restricts search to regions near previous detections, and a **Size Prior Network (SPN)** skips pyramid levels for objects expected to stay a similar size [5] [6]. By assuming objects move slowly frame-to-frame, LPN/SPN cut unnecessary computation and helped one-stage detectors focus on true object regions. In practice, such temporal priors sharpen detection by suppressing improbable new detections in the same area. Heuristic filters also use frame-to-frame context: Pi *et al.* (ICCVW 2019) analyzed false detections by consistency. They showed that *"objects in a video should be strongly correlated"* – for each detected box, they count how many similar detections appear in nearby frames [7]. If an object has too few temporal neighbors (e.g. under a confidence threshold or 3 neighbors in adjacent frames), it is likely a false alarm and can be dropped [8]. This kind of "temporal voting" removes sporadic mis-detections that do not persist over time.

Overall, **short-term temporal smoothing and linking** can greatly stabilize predictions. By warping and fusing features (via optical flow or learned offsets), gating recurrent units, or enforcing neighbor-count consistency, modern methods notably reduce false negatives and isolate false positives. For example, applying temporal consistency filters improved detector precision and recall without costly new annotations [1] [9].

## Long-Term Motion Patterns

Beyond adjacent frames, long-term trajectory context further improves robustness. Cores *et al.* (2021) proposed a two-stage *spatio-temporal detector* tailored to UAV videos [10]. In their system, object proposals are *linked across frames* (forming short tubelets) and then an attention module reasons over **distant-frame proposals**. Specifically, proposals sharing the same anchor in consecutive frames are connected, and an attention network compares proposal features across a long temporal window using updated trajectory positions [11]. This lets the detector enforce consistency over several seconds of motion, improving recall of objects that briefly disappear. The authors report that this combination of short-term linking and long-term attention achieved the best accuracy on benchmark UAV video datasets.

Similarly, Telegrap and Kyrkou (2024) showed that a YOLO detector "enhanced with temporal dynamics" can dramatically outperform single-frame models for UAV traffic monitoring [12]. They created a new **Spatio-Temporal Vehicle Detection (STVD)** dataset of consecutive drone frames, and modified YOLO to process frame pairs or streams. Their best spatiotemporal model *outperformed* the static baseline by **16.2% mean AP** [12]. Notably, they also found that adding *attention mechanisms* in the spatiotemporal model gave further gains, indicating that long-range dependencies (which objects follow which trajectory) are important. [12]

Recurrent and memory-based networks likewise capture long-term trends. Methods using LSTMs, GRUs or transformers can propagate information over many frames. For example, Jiang *et al.* (Drones 2023) applied a convolutional GRU on top of object features to smooth detection scores and box coordinates over time, making outputs more stable. (They reported improved mAP on VisDrone video.) While we focus on broad approaches, this line of work converges on the idea that learning an explicit temporal model (RNNs or attention) yields more consistent tracking across long sequences, as evidenced by large mAP or MOTA gains [12] [11].

## Evolving Scene Context

Drones often survey large, changing scenes (e.g. shifting background in a moving aerial view), so modeling *context* can boost trust. Some approaches integrate spatial context: for instance, Xiao *et al.*'s tracker includes a **Topology-Integrated Embedding Module (TIEM)** that embeds objects along with global scene layout [1]. By learning relationships between objects and background (e.g. cars typically drive on roads), the model improves re-identification consistency and suppresses implausible associations. In practice this means a car is less likely to be suddenly re-detected in the sky or a person on the road if it violates the learned topology. These context features, though not purely temporal, evolve over time and help maintain consistent identities.

Moreover, unsupervised motion segmentation exploits scene context: Fan *et al.* (2025) devised an algorithm that enforces **foreground sparsity and spatial–temporal consistency** as prior knowledge in an optical-flow based detector [13]. Their loss penalizes detections inconsistent with global scene motion, so spurious

moving-background clutters (like tree shadows or parallax) are ignored. They show that adding these UAV-specific priors into a 3D U-Net detector *significantly reduces false positives* (e.g. a 28% IoU gain by sparsity constraint alone) [14] . In short, the system treats a true moving object as a cluster of coherent motion over space and time, thereby filtering out random noise. The result is near-continuous, stable detection of moving targets, even under occlusion [15] .

Finally, techniques like *sliding window context* use the fact that some object categories should not appear randomly. Pi *et al.* (2019) exploited drone camera motion patterns: by inferring the flight direction from object size changes across frames, they could eliminate detections of implausible size/distribution (e.g. a tiny truck on a pedestrian sidewalk) [16] . They also noted that if a purported object rarely appears (say in <10% of frames), it is likely a false alarm [16] . Such temporal-statistical filters use evolving context to adapt the detector's output, improving reliability.

## Enhancing Robustness and Trust

Across these techniques, the goal is to **improve robustness and reduce false alarms** in real-world drone streams. Metrics of technical trustworthiness – consistency, stability, low false positive rate – are explicitly optimized by temporal methods. For example, Tung *et al.* (2022) argued that detector accuracy alone is insufficient: detectors must be *consistent* on similar consecutive frames [17] . They measured consistency of modern detectors at ~83–97% across video and showed that simple pre-processing (de-noising, unsharp masking) could raise consistency by a few percent without accuracy loss [17] . This underlines that temporal smoothness is an orthogonal but critical reliability criterion.

Quantitatively, nearly every spatiotemporal approach reports gains in standard metrics: TA-GRU gave +5.9% mAP [18] , TFAM+TIEM gave +2–3% MOTA [3] , the STVD YOLO work +16% mAP [12] , and the FairMOT-based tracker of Lin *et al.* (2022) achieved +4.9% MOTA with fewer false/missed detections [19] . In all cases, temporal modeling not only boosted accuracy but also *tightly reduced false and missed detections* by leveraging continuity. For instance, in the unsupervised motion detector [50], foreground-sparsity loss alone cut false positives dramatically (28% Jaccard gain) [14] . In short, by grounding predictions in temporal evidence, these methods make drones' vision systems more reliable and trustworthy.

## Representative Models and Techniques

Key modern methods illustrate these principles:

- **Feature Warping/Aggregation:** TFAM [1] and similar approaches (e.g. FGFA, DFF in the literature) warp features from neighbor frames to augment the current frame, boosting recall of small/blurred objects.
- **Attention/Transformer-based:** Cores *et al.* (2021) and Telegraph *et al.* (2024) use attention modules across time to relate distant-frame proposals [10] [12] . TA-GRU [4] combines RNNs with attention to focus on relevant temporal cues, essentially an efficient spatiotemporal transformer.
- **Motion Prior Networks:** Guanxiong *et al.* (2022) introduced LPN/SPN to encode the idea that objects move and change size slowly [5] [6] . This prunes irrelevant search space, indirectly reducing false positives by ignoring unlikely detections.

- **Spatiotemporal Tracking:** FairMOT variants (e.g. Lin *et al.* 2022) separate detection and Re-ID while adding *temporal embedding branches* to reinforce track consistency [19]. ByteTrack and CenterTrack (not cited here) similarly rely on temporal continuity for ID consistency.
- **Contextual/Heuristic Filters:** Pi *et al.* (2019) used statistical context of drone motion to filter false alarms [9]. Unsupervised motion methods (Fan *et al.* 2025) integrate scene priors to eliminate spurious motion [14].

These methods often combine multiple cues. For example, Zhou *et al.*'s TA-GRU explicitly gates out background interference while fusing temporal features [4], and the unsupervised network [50] combines sparse-motion and spatial-consistency losses. Together, they exemplify how *leveraging every possible temporal clue* yields much more reliable output than per-frame detectors.

## Datasets and Benchmarks

Evaluation and benchmarking are crucial. Standard **UAV video datasets** where temporal methods are tested include VisDrone-VID (video object detection track) [20], UAVDT, UAV123/UAV20L, Stanford Drone, and specialized sets like UAV-Person-in-Crowd. VisDrone (2018/2019) provides annotated videos (96 clips, tens of thousands of frames) in diverse urban scenes [20]. UAV123 (2016) and its long-term subset UAV20L contain >100K frames of low-altitude aerial footage with various objects [21]. Recent datasets like the STVD (Telegrap 2024) and UAVMD (Fan 2025) explicitly offer sequential frames for spatiotemporal training [22] [13]. Benchmarks frequently report video AP and tracking MOTA. For instance, STVD allowed evaluating the YOLO-based spatiotemporal detector that outperformed static baselines by >16% mAP [12]. The VisDrone video DET&TRK challenge continues to be a standard benchmark for UAV detection/tracking accuracy and consistency.

## Conclusion

In summary, **temporal modeling** has emerged as a powerful tool for enhancing the reliability of AI-driven detection/tracking in drone video. By exploiting both short-term cues (frame-to-frame smoothing, motion links) and long-term patterns (trajectories, global attention), and by integrating evolving scene context (spatial layout, motion priors), recent research has achieved more robust, consistent predictions. These approaches demonstrably lower false alarm rates and maintain object continuity, thus increasing technical trustworthiness. Notable methods – from deformable feature aggregation [1] to transformer/RNN hybrids [11] [4], from motion-prior networks [5] to unsupervised motion segmentation [14] – highlight that temporal coherence is key. As drone vision systems become increasingly critical, leveraging temporal information will remain essential to ensure their outputs are both accurate and dependable.

[1] [2] [3] Enhancing Online UAV Multi-Object Tracking with Temporal Context and Spatial Topological Relationships
https://www.mdpi.com/2504-446X/7/6/389

[4] [18] Object Detection in Drone Video with Temporal Attention Gated Recurrent Unit Based on Transformer
https://www.mdpi.com/2504-446X/7/7/466

[5] [6] ecva.net
https://www.ecva.net/papers/eccv_2022/papers_ECCV/papers/136950001.pdf

[7] [8] [9] [16] [20] A Novel Spatial and Temporal Context-Aware Approach for Drone-Based Video Object Detection
https://openaccess.thecvf.com/content_ICCVW_2019/papers/VISDrone/Pi_A_Novel_Spatial_and_Temporal_Context-Aware_Approach_for_Drone-Based_Video_ICCVW_2019_paper.pdf

[10] [11] persoal.citius.usc.es
http://persoal.citius.usc.es/manuel.mucientes/pubs/Cores21_caip.pdf

[12] [22] Spatiotemporal Object Detection for Improved Aerial Vehicle Detection in Traffic Monitoring
https://arxiv.org/html/2410.13616v1

[13] [14] [15] An Unsupervised Moving Object Detection Network for UAV Videos
https://www.mdpi.com/2504-446X/9/2/150

[17] arxiv.org
https://arxiv.org/pdf/2207.13890

[19] Multiple Object Tracking of Drone Videos by a Temporal-Association Network with Separated-Tasks Structure
https://www.mdpi.com/2072-4292/14/16/3862

[21] A Survey of Object Detection for UAVs Based on Deep Learning
https://www.mdpi.com/2072-4292/16/1/149