

[['Distractor suppression Siamese network with task-aware attention for visual tracking',
<https://www.sciencedirect.com/science/article/pii/S1077314225003303>,
'',
'Zhigang Liu, Fuyuan Xing, Hao Huang, Kexin Wang, Yuxuan Shao',
'1',
'2026',
'1'],
['Transformer tracking with high-low frequency attention',
<https://www.sciencedirect.com/science/article/pii/S1077314225002863>,
'',
'Zhi Chen, Zhen Yu',
'1',
'2026',
'1'],
['PlanarTrack: A high-quality and challenging benchmark for large-scale planar object tracking',
<https://www.sciencedirect.com/science/article/pii/S1077314225002760>,
'',
'Yifan Jiao, Xinran Liu, Xiaoqiong Liu, Xiaohui Yuan, Heng Fan, Libo Zhang',
'1',
'2025',
'12'],
['Curvi-Tracker: Curvilinear structure segmentation refinement by iterative tracking',
<https://www.sciencedirect.com/science/article/pii/S0031320325014608>,
'',
'Zhan Heng, Maurice Pagnucco, Erik Meijering, Yang Song',
'2',
'2026',
'5'],
['Hybrid-stage association with dynamicity adaptation and enhanced cues for multi-object tracking and segmentation',
<https://www.sciencedirect.com/science/article/pii/S0031320325014669>,
'',
'Longtao Chen, Guoxing Liao, Yifan Shi, Jing Lou, Fenglei Xu, Huanqiang Zeng',
'2',
'2026',
'5'],
['Spherical Vision Transformers for Audio-Visual Saliency Prediction in 360 $^{\circ}$ Videos',
<http://ieeexplore.ieee.org/document/11144923>,
'Omnidirectional videos (ODVs) are redefining viewer experiences in virtual reality (VR) by offering an unprecedented full field-of-view (FOV). This study extends the domain of saliency prediction to 360 $^{\circ}$ environments, addressing the complexities of spherical distortion and the integration of spatial audio. Contextually, ODVs have transformed user experience by adding a spatial audio dimension that aligns sound direction with the viewer's perspective in spherical scenes. Motivated by the lack of comprehensive datasets for 360 $^{\circ}$ audio-visual saliency prediction, our study curates YT360-EyeTracking, a new dataset of 81 ODVs, each observed under varying audio-visual conditions. Our goal is to explore how to utilize audio-visual cues to effectively predict visual saliency in 360 $^{\circ}$ videos. Towards this aim, we propose two novel saliency prediction models: SalViT360, a vision-transformer-based framework for ODVs equipped with spherical geometry-aware spatio-temporal attention layers, and SalViT360-AV, which further incorporates transformer adapters conditioned on audio input. Our results on a number of benchmark datasets, including our YT360-EyeTracking, demonstrate that SalViT360 and SalViT360-AV significantly outperform existing methods in predicting viewer attention in 360 $^{\circ}$ scenes. Interpreting these results, we suggest that integrating spatial audio cues in the model architecture is crucial for accurate saliency prediction in omnidirectional videos.',
'Mert Cokelek, Halit Ozsoy, Nevrez Imamoglu, Cagri Ozcinar, Inci Ayhan, Erkut Erdem, Aykut Erdem',

'3',
'2026',
'12'],

['PMGT-VR: A Decentralized Proximal-Gradient Algorithmic Framework With Variance Reduction',

'<http://ieeexplore.ieee.org/document/11152599>',

'This article considers the decentralized composite optimization problem. We propose a novel decentralized variance-reduction proximal-gradient algorithmic framework, called PMGT-VR, which combines several techniques, including multi-consensus, gradient tracking, and variance reduction. The proposed framework imitates centralized algorithms and algorithms under this framework achieve convergence rates similar to that of their centralized counterparts. We also describe and analyze two representative algorithms, PMGT-SAGA and PMGT-LSVRG, and compare them to existing state-of-the-art proximal algorithms. To the best of our knowledge, PMGT-VR is the first linearly convergent decentralized stochastic algorithm that can solve decentralized composite optimization problems. Numerical experiments are provided to demonstrate the effectiveness of the proposed algorithms.'

'Haishan Ye, Wei Xiong, Tong Zhang',

'3',
'2026',
'12'],

['ADA-Track++: End-to-End Multi-Camera 3D Multi-Object Tracking With Alternating Detection and Association',

'<http://ieeexplore.ieee.org/document/11175504>',

'Many query-based approaches for 3D Multi-Object Tracking (MOT) adopt the tracking-by-attention paradigm, utilizing track queries for identity-consistent detection and object queries for identity-agnostic track spawning. Tracking-by-attention, however, entangles detection and tracking queries in one embedding for both the detection and tracking task, which is sub-optimal. Other approaches resemble the tracking-by-detection paradigm and detect objects using decoupled track and detection queries followed by a subsequent association. These methods, however, do not leverage synergies between the detection and association task. Combining the strengths of both paradigms, we introduce ADA-Track++, a novel end-to-end framework for 3D MOT from multi-view cameras. We introduce a learnable data association module based on edge-augmented cross-attention, leveraging appearance and geometric features. We also propose an auxiliary token in this attention-based association module, which helps mitigate disproportionately high attention to incorrect association targets caused by attention normalization. Furthermore, we integrate this association module into the decoder layer of a DETR-based 3D detector, enabling simultaneous DETR-like query-to-image cross-attention for detection and query-to-query cross-attention for data association. By stacking these decoder layers, queries are refined for the detection and association task alternately, effectively harnessing the task dependencies. We evaluate our method on the nuScenes dataset and demonstrate the advantage of our approach compared to the two previous paradigms.'

'Shuxiao Ding, Lukas Schneider, Marius Cordts, Juergen Gall',

'3',
'2026',
'12'],

['SNNTracker: Online High-Speed Multi-Object Tracking With Spike Camera',

'<http://ieeexplore.ieee.org/document/11165142>',

'Multi-object tracking (MOT) is crucial for applications such as autonomous driving and robotics, yet traditional image-based methods struggle in high-speed scenarios due to motion blur and temporal gaps caused by low frame rates. Spike cameras, with their ability to continuously record spatiotemporal signals, overcome these limitations. However, existing spike-based methods often rely on intermediate image reconstruction or discrete clustering, limiting real-time performance and temporal continuity. To address this, we propose SNNTracker, the first fully spiking neural network (SNN)-based MOT algorithm tailored for spike cameras. SNNTracker integrates a dynamic neural field (DNF)-based attention mechanism for target detection and a winner-take-all (WTA)-based tracking module with online spike-timing-dependent plasticity (STDP) for adaptive learning of object trajectories. By directly processing spike streams without

reconstruction, SNNTracker reduces latency, computational overhead, and dependency on image quality, making it ideal for ultra-high-speed environments. It maintains robust, continuous tracking even under occlusions, severe lighting variations, or temporary object disappearance, by leveraging SNN-estimated motion predictions and long-term online clustering. We construct three types of spike-camera MOT datasets covering dense and sparse annotations across diverse real-world scenarios, including camera ego-motion, deformable and ultra-fast motion (up to 2600 RPM), occlusion, indoor/outdoor lighting changes, and low-visibility tracking. Extensive experiments demonstrate that SNNTracker consistently outperforms state-of-the-art MOT methods—both ANN- and SNN-based—achieving MOTA scores above 96% and up to 100% in many sequences. Our results highlight the advantages of spike-driven SNNs for low-latency, high-speed, and label-free multi-object tracking, advancing neuromorphic vision for real-time perception.'

'Yajing Zheng, Chengen Li, Jiyuan Zhang, Zhao Fei Yu, Tiejun Huang',

'3',

'2026',

'12'],

['LVOS: A Benchmark for Large-Scale Long-Term Video Object Segmentation',

'<http://ieeexplore.ieee.org/document/11168273>',

'Video object segmentation (VOS) aims to distinguish and track target objects in a video. Despite the excellent performance achieved by off-the-shelf VOS models, part of the existing VOS benchmarks mainly focuses on short-term videos, where objects remain visible most of the time. However, these benchmarks may not fully capture challenges encountered in practical applications, and the absence of long-term datasets restricts further investigation of VOS in realistic scenarios. Thus, we propose a novel benchmark named LVOS, comprising 720 videos with 296,401 frames and 407,945 high-quality annotations. Videos in LVOS last 1.14 minutes on average. Each video includes various attributes, especially challenges encountered in the wild, such as long-term reappearing and cross-temporal similar objects. Compared to previous benchmarks, our LVOS better reflects VOS models' performance in real scenarios. Based on LVOS, we evaluate 15 existing VOS models under 3 different settings and conduct a comprehensive analysis. On LVOS, these models suffer a large performance drop, highlighting the challenge of achieving precise tracking and segmentation in real-world scenarios. Attribute-based analysis indicates that one of the significant factors contributing to accuracy decline is the increased video length, interacting with complex challenges such as long-term reappearance, cross-temporal confusion, and occlusion, which emphasize LVOS's crucial role. We hope our LVOS can advance development of VOS in real scenes.'

'Lingyi Hong, Zhongying Liu, Wenchao Chen, Chenzhi Tan, Yuang Feng, Xinyu Zhou, Pinxue Guo, Jinglun Li, Zhaoyu Chen, Shuyong Gao, Wei Zhang, Wenqiang Zhang',

'3',

'2026',

'12'],

['SUIT: Spatial-Spectral Union-Intersection Interaction Network for Hyperspectral Object Tracking',

'<http://ieeexplore.ieee.org/document/11267013>',

'Hyperspectral videos (HSVs), with their inherent spatial-spectral-temporal structure, offer distinct advantages in challenging tracking scenarios such as cluttered backgrounds and small objects. However, existing methods primarily focus on spatial interactions between the template and search regions, often overlooking spectral interactions, leading to suboptimal performance. To address this issue, this paper investigates spectral interactions from both the architectural and training perspectives. At the architectural level, we first establish band-wise long-range spatial relationships between the template and search regions using Transformers. We then model spectral interactions using the inclusion-exclusion principle from set theory, treating them as the union of spatial interactions across all bands. This enables the effective integration of both shared and band-specific spatial cues. At the training level, we introduce a spectral loss to enforce material distribution alignment between the template and predicted regions, enhancing robustness to shape deformation and appearance variations. Extensive experiments demonstrate that our tracker achieves state-

of-the-art tracking performance. The source code, trained models and results will be publicly available via <https://github.com/bearshng/suit> to support reproducibility',

'Fengchao Xiong, Zhenxing Wu, Jun Zhou, Sen Jia, Yuntao Qian',

'4',

'2025',

'12'],

['Quality-Aware Spatio-Temporal Transformer Network for RGBT Tracking',

'<http://ieeexplore.ieee.org/document/11270003>',

'Transformer-based RGBT tracking has attracted much attention due to the strong modeling capacity of self attention and cross attention mechanisms. These attention mechanisms utilize the correlations among tokens to construct powerful feature representations, but are easily affected by low-quality tokens. To address this issue, we propose a novel Quality-aware Spatio-temporal Transformer Network (QSTNet), which calculates the quality weights of tokens in search regions based on the correlation with multimodal template tokens to suppress the negative effects of low-quality tokens in spatio-temporal feature representations, for robust RGBT tracking. In particular, we argue that the correlation between search tokens of one modality and multimodal template tokens could reflect the quality of these search tokens, and thus design the Quality-aware Token Weighting Module (QTWM) based on the correlation matrix of search and template tokens to suppress the negative effects of low-quality tokens. Specifically, we calculate the difference matrix derived from the attention matrices of the search tokens from both modalities and the multimodal template tokens, and then assign the quality weight for each search token based on the difference matrix, which reflects the relative correlation of search tokens from different modalities to multimodal template tokens. In addition, we propose the Prompt-based Spatio-temporal Encoder Module (PSEM) to utilize spatio-temporal multimodal information while alleviating the impact of low-quality spatio-temporal features. Extensive experiments on four RGBT benchmark datasets demonstrate that the proposed QSTNet exhibits superior performance compared to other state-of-the-art tracking methods. Our code and supplementary video are now available: <https://zhaodongah.github.io/QSTNet>',

'Zhaodong Ding, Chenglong Li, Tao Wang, Futian Wang',

'4',

'2025',

'12'],

['Self-Adaptive Vision-Language Tracking With Context Prompting',

'<http://ieeexplore.ieee.org/document/11284903>',

'Due to the substantial gap between vision and language modalities, along with the mismatch problem between fixed language descriptions and dynamic visual information, existing vision-language tracking methods exhibit performance on par with or slightly worse than vision-only tracking. Effectively exploiting the rich semantics of language to enhance tracking robustness remains an open challenge. To address these issues, we propose a self-adaptive vision-language tracking framework that leverages the pre-trained multi-modal CLIP model to obtain well-aligned visual-language representations. A novel context-aware prompting mechanism is introduced to dynamically adapt linguistic cues based on the evolving visual context during tracking. Specifically, our context prompter extracts dynamic visual features from the current search image and integrates them into the text encoding process, enabling self-updating language embeddings. Furthermore, our framework employs a unified one-stream Transformer architecture, supporting joint training for both vision-only and vision-language tracking scenarios. Our method not only bridges the modality gap but also enhances robustness by allowing language features to evolve with visual context. Extensive experiments on four vision-language tracking benchmarks demonstrate that our method effectively leverages the advantages of language to enhance visual tracking. Our large model can obtain 55.0% AUC on $\text{\texttt{\{LaSOT\}_\text{\texttt{\{EXT\}}}}$ and 69.0% AUC on TNL2K. Additionally, our language-only tracking model achieves performance comparable to that of state-of-the-art vision-only tracking methods on TNL2K. Code is available at <https://github.com/zj5559/SAVLT>',

'Jie Zhao, Xin Chen, Shengming Li, Chunjuan Bo, Dong Wang, Huchuan Lu',

'4',

framework designed for ship features.

ABSTRACT

Multiobject tracking of ships is crucial for various applications, such as maritime security and the development of ship autopilot systems. However, existing ship visual datasets primarily focus on ship detection tasks, lacking a fully open-source dataset for multiobject tracking research. Furthermore, current methods often struggle with extracting appearance features under complex sea conditions, varying scales and different ship types, affecting tracking precision. To address these issues, we propose ShipsMOT, a new benchmark dataset containing 121 video sequences with an average of 15.45s per sequence, covering 15 distinct ship types and a total of 237,999 annotated bounding boxes. Additionally, we propose JDR-CSTrack, a ship multiobject tracking framework that improves feature extraction at different scales by optimising a joint detection and Re-ID network. JDR-CSTrack utilises the fusion of appearance and motion features for multilevel data association, thereby minimising track loss and ID switches. Experimental results confirm that ShipsMOT can serve as a benchmark for future research in ship multiobject tracking and validate the superiority of the proposed JDR-CSTrack framework. The dataset and code can be found on <https://github.com/jpj0916/ShipsMOT>.

Fang Luo, Pengju Jiang, George To Sum Ho, Wenjing Zeng,

'6',

'2025',

'10'],

['A Review of Multi-Object Tracking in Recent Times',

<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cvi2.70010?af=R>,

This paper discusses many recent deep-learning MOT methods. Moreover, to highlight their contributions, these methods are categorised into four main groups: detection-based, SOT-based, and segmentation-based methods according to the integrated core technologies.

ABSTRACT

Multi-object tracking (MOT) is a fundamental problem in computer vision that involves tracing the trajectories of foreground targets throughout a video sequence while establishing correspondences for identical objects across frames. With the advancement of deep learning techniques, methods based on deep learning have significantly improved accuracy and efficiency in MOT. This paper reviews several recent deep learning-based MOT methods and categorises them into three main groups: detection-based, single-object tracking (SOT)-based, and segmentation-based methods, according to their core technologies. Additionally, this paper discusses the metrics and datasets used for evaluating MOT performance, the challenges faced in the field, and future directions for research.

Suya Li, Hengyi Ren, Xin Xie, Ying Cao,

'6',

'2025',

'3'],

['A New Large-Scale Dataset for Marine Vessel Re-Identification Based on Swin Transformer Network in Ocean Surveillance Scenario',

<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cvi2.70007?af=R>,

A new large-scale marine vessel dataset with well-annotated vessel orientation, vessel colour, and vessel type labels has been collected and created in a real marine environment for vessel Re-ID research. (2) A side information embedding module is introduced through the learnable embedding layer to encode more kinds of information, including marine vessel orientation, type and colour. (3) A deep neural network framework based on Swin Transformer for marine vessel Re-ID task is proposed to learn and extract discriminative features, and achieves SOTA performance on vessel, vehicle and person Re-ID benchmark datasets.

ABSTRACT

In recent years, there has been an upward trend that marine vessels, an important object category in marine monitoring, have gradually become a research focal point in the field of computer vision, such as detection, tracking, and classification. Among them, marine vessel re-identification (Re-ID) emerges as a significant frontier research topics, which not only faces the dual challenge of huge intra-class and small inter-class differences, but also has complex environmental interference in the port monitoring scenarios. To propel advancements in marine vessel Re-ID technology, SwinTransReID, a framework grounded in the Swin Transformer for marine vessel Re-ID, is introduced. Specifically, the project initially encodes

the triplet images separately as a sequence of blocks and construct a baseline model leveraging the Swin Transformer, achieving better performance on the Re-ID benchmark dataset in comparison to convolution neural network (CNN)-based approaches. And it introduces side information embedding (SIE) to further enhance the robust feature-learning capabilities of Swin Transformer, thus, integrating non-visual cues (orientation and type of vessel) and other auxiliary information (hull colour) through the insertion of learnable embedding modules. Additionally, the project presents VesselReID-1656, the first annotated large-scale benchmark dataset for vessel Re-ID in real-world ocean surveillance, comprising 135,866 images of 1656 vessels along with 5 orientations, 12 types, and 17 colours. The proposed method achieves 87.1% mAP and 96.1% Rank-1 accuracy on the newly-labelled challenging dataset, which surpasses the state-of-the-art (SOTA) method by 1.9% mAP regarding to performance. Moreover, extensive empirical results demonstrate the superiority of the proposed SwinTransReID on the person Market-1501 dataset, vehicle VeRi-776 dataset, and Boat Re-ID vessel dataset.',

'Zhi Lu, Liguu Sun, Pin Lv, Jiuwu Hao, Bo Tang, Xuanzhen Chen',

'6',

'2025',

'3'],

['Unlocking the power of multi-modal fusion in 3D object tracking',

'<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cvi2.12335?af=R>',

'3D Single Object Tracking plays a vital role in autonomous driving and robotics, yet traditional approaches have predominantly focused on using pure LiDAR-based point cloud data, often neglecting the benefits of integrating image modalities. To address this gap, we propose a novel Multi-modal Image-LiDAR Tracker (MILT) designed to overcome the limitations of single-modality methods by effectively combining RGB and point cloud data. Our key contribution is a dual-branch architecture that separately extracts geometric features from LiDAR and texture features from images. These features are then fused in a BEV perspective to achieve a comprehensive representation of the tracked object. A significant innovation in our approach is the Image-to-LiDAR Adapter module, which transfers the rich feature representation capabilities of the image modality to the 3D tracking task, and the BEV-Fusion module, which facilitates the interactive fusion of geometry and texture features. By validating MILT on public datasets, we demonstrate substantial performance improvements over traditional methods, effectively showcasing the advantages of our multi-modal fusion strategy. This work advances the state-of-the-art in SOT by integrating complementary information from RGB and LiDAR modalities, resulting in enhanced tracking accuracy and robustness.',

'Yue Hu',

'6',

'2025',

'2'],

['Representation alignment contrastive regularisation for multi-object tracking',

'<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cvi2.12331?af=R>',

"Achieving high-performance in multi-object tracking algorithms heavily relies on modelling spatial-temporal relationships during the data association stage. Mainstream approaches encompass rule-based and deep learning-based methods for spatial-temporal relationship modelling. While the former relies on physical motion laws, offering wider applicability but yielding suboptimal results for complex object movements, the latter, though achieving high-performance, lacks interpretability and involves complex module designs. This work aims to simplify deep learning-based spatial-temporal relationship models and introduce interpretability into features for data association. Specifically, a lightweight single-layer transformer encoder is utilised to model spatial-temporal relationships. To make features more interpretative, two contrastive regularisation losses based on representation alignment are proposed, derived from spatial-temporal consistency rules. By applying weighted summation to affinity matrices, the aligned features can seamlessly integrate into the data association stage of the original tracking workflow. Experimental results showcase that our model enhances the majority of existing tracking networks' performance without excessive complexity, with minimal increase in training

coverhead and nearly negligible computational and storage costs.",
 'Shujie Chen, Zhonglin Liu, Jianfeng Dong, Xun Wang, Di Zhou',
 '6',
 '2025',
 '2'],

['Droplet Detection and Tracking in Complex Motions Based on YOLOv5s Network',
<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ipr2.70245?af=R>,
 'This study presents an enhanced model for droplet detection and tracking in solvent extraction operations, using a refined YOLOv5s framework integrated with the OC-SORT algorithm. The model accurately detects droplets and tracks their nonlinear motion trajectories, achieving a detection accuracy of 0.977 and a 40% increase in speed compared to the original model.\n\n\n\n\n\n\n\n\n\nABSTRACT\nIn the realm of chemical engineering, solvent extraction is pivotal for separating valuable components from mixtures. Precise understanding of droplet dynamics is crucial for enhancing extraction column efficiency. However, analysing dispersed-phase droplets during solvent extraction operations remains challenging due to their complex motions and interactions. This study proposes an enhanced model for droplet detection and tracking in complex motions. Building upon the YOLOv5s object detection framework, we refine the model to accurately detect droplets and recognise their size characteristics. Additionally, we integrate YOLOv5s with the OC-SORT algorithm to develop a robust droplet tracking model capable of tracing nonlinear droplet motion trajectories under both pulsed and non-pulsed conditions. Experimental results demonstrate that our approach achieves high accuracy and faster detection speeds, with a detection accuracy of 0.977 and a 40% increase in speed compared to the original model. This study not only validates the effectiveness of our model through manual labelling but also provides a solid foundation for further research into the complex fluid dynamics in extraction columns. The source code of the proposed method is available at <https://github.com/polangjushi/DDTA>.',
 'Yubo Zhang, Zhenning Su, Yong Wang, Boren Tan, Chen Zhao',
 '7',
 '2025',
 '11'],

['AI-Powered Human Activity Detection and Tracking in Dense Crowds Using YOLOv8-DeepSORT',
<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ipr2.70227?af=R>,
 'This study presents a real-time human activity detection and recognition system that combines YOLOv8l with deep simple online real-time tracking for enhanced performance in crowded and dynamic environments. Evaluated on over 10,000 frames from public and experimental datasets, the system achieves a mean average precision of 96.1%, with the Adam optimiser outperforming stochastic gradient descent. The model demonstrates robustness and scalability, offering a significant advancement in intelligent public surveillance.\n\n\n\n\n\n\n\n\n\nABSTRACT\nMany organisations face challenges in accurately and efficiently identifying human activities using manually operated CCTV cameras and small datasets, particularly in densely populated and dynamic environments. This paper presents a real-time human activity detection and recognition system that leverages You Only Look Once version 8 (YOLOv8) for object detection and integrates the deep simple online real-time tracking algorithm for associate frames, robust tracking and individual identity preservation. Evaluated on a combination of public and experimental datasets containing over 10,000 frames, the system achieves notable enhancements in efficiency along with accuracy. Through data preprocessing and algorithm optimisation, the Adam optimiser outperformed stochastic gradient descent (SGD) by more than 5% in accuracy. K-fold cross-validation was applied to reduce overfitting, while anchor boxes of varying scales and aspect ratios improved the detection of objects of different sizes. The aspect ratio, which is used to mitigate the effect of object size and control edge-orientation issues by computing the ratio of width to height of the bounding boxes in YOLOv8, plays a crucial role in object detection.YOLOv8l, combined with DeepSORT, delivers state-of-the-art results, achieving a mean average precision of 96.1% in mixed datasets. The system assigns unique IDs to track individuals, enabling precise frame association during motion in crowded areas, which is crucial for maintaining accuracy in complex environments. This study presents a robust and

<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ipr2.70136?af=R>,

'To overcome the limitations of existing homography matrix prediction methods, we developed a two-stage prediction approach by optimizing the pixel position feature extraction structure. In the first stage, 4-keypoint coordinate deviation is leveraged to indirectly estimate the macroscopic values of the homography matrix, followed by direct fine-tuning of the microscopic 8-DOF numerical components in the second stage. We validated the effectiveness of this approach in three competitive sports scenarios: basketball, ice hockey, and handball.'

ABSTRACT

Homography estimation is a fundamental topic in computer vision, especially in scenarios that require perspective changes for intelligent analysis of sports fields, where it plays a crucial role. Existing methods predict the homography matrix either indirectly by evaluating the 4-keypoint coordinate deviation in paired images with the same visual content or directly by fine-tuning the 8 degrees of freedom numerical values that define the matrix. However, these approaches often fail to effectively incorporate coordinate positional information and overlook optimal application scenarios, leading to significant accuracy bottlenecks, particularly for paired images with differing visual content. To address these issues, we propose an approach that integrates both methods in a staged manner, leveraging their respective advantages. In the first stage, positional information is embedded to enhance convolutional computations, replacing serial concatenation in traditional feature fusion with parallel concatenation, while using 4-key-point coordinate deviation to predict the macroscopic homography matrix. In the second stage, positional information is further integrated into the input images to refine the direct 8 degrees of freedom numerical predictions, improving matrix fine-tuning accuracy. Comparative experiments with state-of-the-art methods demonstrate that our approach achieves superior performance, yielding a root mean square error as low as 1.25 and an average corner error as low as 14.1 in homography transformation of competitive sports image pairs.'

'Pan Zhang, Jiangtao Luo, Guoliang Xu, Xupeng Liang',

```
'7',
'2025',
'6']
```

['Vorticity Transport Equation-Based Shadow Removal Approach for Image Inpainting'],

<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ipr2.70114?af=R>,

(1) A new image inpainting algorithm using the VTE is developed; (2) FMM is introduced as an alternative to the iterative numerical solution of the VTE to improve the speed of image inpainting; and (3) the VTE equation is transformed into a weighted-average equation to be used in conjunction with FMM.

ABSTRACT

Shadows are common in many types of images, causing information loss or disturbance. Shadow removal can help improve the quality of the digital image. If there is no effective information available to restore the original image in the shaded area, the interpolation-based inpainting technique can be used to remove the shadow from the digital image. This image inpainting technique typically involves establishing and solving partial differential equations (PDEs), an iterative solving process that is very time-consuming. To solve the time-consuming problem, a method that introduces the fast marching method (FMM) into the vorticity transport equation (VTE) is demonstrated. VTE is a type of partial differential equation describing two-dimensional fluids. FMM is a numerical scheme for tracking the evolution of monotonically advancing interfaces via finite difference solution of the eikonal equation. The proposed method contains three main steps: (a) by investigating the relationship between VTE and the traditional PDE-based image inpainting method, a new image inpainting model using VTE is developed; (b) the area to be inpainted is divided into boundaries that shrink in layers from the outside inwards using FMM; and (c) the VTE image inpainting model is converted into a weighted average form to coordinate with FMM. The visual and quantitative evaluation of the experimental results of shadow removal shows that the proposed method outperforms PDE-based and state-of-the-art methods in terms of shadow-

"We reconceptualize visual tracking as a multivariate time-series forecasting (MTSF) problem. In addition, we propose a principled approach that models the dynamic nature of target motion using a regime-switching framework. This method employs an underlying Markov jump process (MJP) to govern transitions between multiple latent motion patterns, each characterized by its own stochastic differential equation (SDE).
In this paper, we reconceptualize visual tracking as a multivariate time-series forecasting (MTSF) problem. Specifically, the goal of visual tracking—predicting the target's state over time, including its (x, y) center coordinates and scale—can be naturally framed as forecasting future states from a sequence of past observations. Viewed through this lens, visual tracking aligns with the challenges of MTSF, where the objective is to capture complex temporal dependencies among multiple variables. However, applying MTSF to visual tracking introduces new difficulties due to the inherently intricate nature of object motion, which often involves abrupt and nonlinear variations in direction, velocity, and behavioral patterns. To address these complexities, we propose a principled approach that models the dynamic nature of target motion using a regime-switching framework. This method employs an underlying Markov jump process (MJP) to govern transitions between multiple latent motion patterns, each characterized by its own stochastic differential equation (SDE). By doing so, our model adapts to diverse temporal dynamics in a data-driven manner, enabling robust and precise prediction of future target states. Experimental results demonstrate that our method outperforms state-of-the-art visual tracking approaches, particularly in scenarios where target objects exhibit diverse and dynamic motion patterns over time."

'Seonghak Lee, Jisoo Park, Radu Timofte, Junseok Kwon',

'8',

'2025',

'11'],

['QMBOC-HFM: Enhanced QMBOC for Large-Scale LEO Navigation Augmentation Systems',

'<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ell2.70398?af=R>',

'A novel QMBOC signal with improved code tracking performance, multipath suppression, and anti-jamming capabilities.
Large-scale low earth orbit (LEO) constellations can enhance traditional global navigation satellite systems (GNSS) in many ways, but the rapid increase in the number of LEO satellites also poses challenges to the anti-interference performance of navigation signals. To improve the performance of the B1C signal in the future BeiDou global navigation satellite system, this paper proposes a modulation method: it enhances the existing quadrature-multiplexed binary offset carrier (QMBOC) by using periodic binary hyperbolic frequency modulation (HFM) signals as subcarriers. The enhanced QMBOC signal (QMBOC-HFM) demonstrates superior anti-interference and positioning accuracy compared to the original QMBOC signal. The proposed QMBOC-HFM provides a potential signal design for future LEO navigation-augmented global satellite navigation systems.'

'Shugan Zhang, Xinming Huang, Bofang Chen',

'8',

'2025',

'9'],

['Maximum Circumnavigation Radius for a Class of Bearings-Only Target Tracking Problem',

'<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ell2.70390?af=R>',

'Under a virtual intersecting localization and tracking algorithm, the maximum circumnavigation radius is proposed for the bearings-only observer by analysing the relationship between the target position estimation error, observation error and the velocity of the observer and the target. The proposed method can provide an engineering application guidance for a class of bearings-only circumnavigation tracking problems.
For the problem of bearings-only target tracking using circumnavigation method, to enhance the safety and stealth, the observer needs to maintain the maximum possible distance from the target. But an excessively large circumnavigation radius may lead to significant tracking error or even no solution in target localization. Under a virtual intersecting localization and tracking algorithm, the maximum circumnavigation radius is proposed for the bearings-only observer by analysing

'2025',
 '2'],
 ['Extended target tracking using neural network and Gaussian process',
 '<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/ell2.70151?af=R>',
 "In extended target tracking, Gaussian Process (GP) is utilized to model unknown contour functions based on the model-predicted target center and contour measurements. However, model prediction relies on accurate prior knowledge. When the model-predicted target center is inaccurate, it will affect the modelling of the measurement model. To address issue, this letter introduces a hybrid-driven approach that combines extended Kalman filter using GP with neural network; proposes an extended target tracking algorithm using neural network and GP. The algorithm predicts the target center according to the neural network and the target's kinematic model, and takes the prediction center and the contour measurements at the current moment as the input of the neural network, which in turn provides real-time estimates for the predicted center compensation. The simulation results show that the algorithm has a significant improvement in tracking performance and better accuracy in estimating the center position and extent state of the target.",
 'Hao Wang, Liping Song',
 '8',
 '2025',
 '1']]

[[['Curvi-Tracker: Curvilinear structure segmentation refinement by iterative tracking',
 '<https://www.sciencedirect.com/science/article/pii/S0031320325014608>',
 '',
 'Zhan Heng, Maurice Pagnucco, Erik Meijering, Yang Song',
 '2',
 '2026',
 '5'],
 ['SNNTracker: Online High-Speed Multi-Object Tracking With Spike Camera',
 '<http://ieeexplore.ieee.org/document/11165142>',
 "Multi-object tracking (MOT) is crucial for applications such as autonomous driving and robotics, yet traditional image-based methods struggle in high-speed scenarios due to motion blur and temporal gaps caused by low frame rates. Spike cameras, with their ability to continuously record spatiotemporal signals, overcome these limitations. However, existing spike-based methods often rely on intermediate image reconstruction or discrete clustering, limiting real-time performance and temporal continuity. To address this, we propose SNNTracker, the first fully spiking neural network (SNN)-based MOT algorithm tailored for spike cameras. SNNTracker integrates a dynamic neural field (DNF)-based attention mechanism for target detection and a winner-take-all (WTA)-based tracking module with online spike-timing-dependent plasticity (STDP) for adaptive learning of object trajectories. By directly processing spike streams without reconstruction, SNNTracker reduces latency, computational overhead, and dependency on image quality, making it ideal for ultra-high-speed environments. It maintains robust, continuous tracking even under occlusions, severe lighting variations, or temporary object disappearance, by leveraging SNN-estimated motion predictions and long-term online clustering. We construct three types of spike-camera MOT datasets covering dense and sparse annotations across diverse real-world scenarios, including camera ego-motion, deformable and ultra-fast motion (up to 2600 RPM), occlusion, indoor/outdoor lighting changes, and low-visibility tracking. Extensive experiments demonstrate that SNNTracker consistently outperforms state-of-the-art MOT methods—both ANN- and SNN-based—achieving MOTA scores above 96% and up to 100% in many sequences. Our results highlight the advantages of spike-driven SNNs for low-latency, high-speed, and label-free multi-object tracking, advancing neuromorphic vision for real-time perception.",
 'Yajing Zheng, Chengen Li, Jiyuan Zhang, Zhaofei Yu, Tiejun Huang',
 '3',
 '2026',
 '12']],

['SUIT: Spatial-Spectral Union-Intersection Interaction Network for Hyperspectral Object Tracking',
<http://ieeexplore.ieee.org/document/11267013>,
 'Hyperspectral videos (HSVs), with their inherent spatial-spectral-temporal structure, offer distinct advantages in challenging tracking scenarios such as cluttered backgrounds and small objects. However, existing methods primarily focus on spatial interactions between the template and search regions, often overlooking spectral interactions, leading to suboptimal performance. To address this issue, this paper investigates spectral interactions from both the architectural and training perspectives. At the architectural level, we first establish band-wise long-range spatial relationships between the template and search regions using Transformers. We then model spectral interactions using the inclusion-exclusion principle from set theory, treating them as the union of spatial interactions across all bands. This enables the effective integration of both shared and band-specific spatial cues. At the training level, we introduce a spectral loss to enforce material distribution alignment between the template and predicted regions, enhancing robustness to shape deformation and appearance variations. Extensive experiments demonstrate that our tracker achieves state-of-the-art tracking performance. The source code, trained models and results will be publicly available via <https://github.com/bearshng/suit> to support reproducibility',
 'Fengchao Xiong, Zhenxing Wu, Jun Zhou, Sen Jia, Yuntao Qian',
 '4',
 '2025',
 '12']],

['Unlocking the power of multi-modal fusion in 3D object tracking',
<https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cvi2.12335?af=R>,
 '3D Single Object Tracking plays a vital role in autonomous driving and robotics, yet traditional approaches have predominantly focused on using pure LiDAR-based point cloud data, often neglecting the benefits of integrating image modalities. To address this gap, we propose a novel Multi-modal Image-LiDAR Tracker (MILT) designed to overcome the limitations of single-modality methods by effectively combining RGB and point cloud data. Our key contribution is a dual-branch architecture that separately extracts geometric features from LiDAR and texture features from images. These features are then fused in a BEV perspective to achieve a comprehensive representation of the tracked object. A significant innovation in our approach is the Image-to-LiDAR Adapter module, which transfers the rich feature representation capabilities of the image modality to the 3D tracking task, and the BEV-Fusion module, which facilitates the interactive fusion of geometry and texture features. By validating MILT on public datasets, we demonstrate substantial performance improvements over traditional methods, effectively showcasing the advantages of our multi-modal fusion strategy. This work advances the state-of-the-art in SOT by integrating complementary information from RGB and LiDAR modalities, resulting in enhanced tracking accuracy and robustness.',
 'Yue Hu',
 '6',
 '2025',
 '2']]