

Advancing Public Activity Recognition in Video Streams Using Hybrid Deep Learning Techniques: A Review

Onkar Tiwari^{1*}, Krishan Kumar²

^{1,2}Department of Computer Science and Engineering, Shri Krishna University, Chhatarpur, M.P. 471001, India

onkartiwarisku@gmail.com, 9893995166

Abstract: Public activity recognition in video streams has become a pivotal area of research due to its applications in surveillance, crowd monitoring, and smart city solutions. Recent advancements in deep learning, particularly hybrid architectures combining multiple learning paradigms, have significantly improved recognition accuracy and robustness in complex environments. This paper reviews the state-of-the-art hybrid deep learning approaches for public activity recognition in video streams, discusses their comparative merits, and outlines open research challenges. The study systematically analyzes methodologies, datasets, and performance metrics used across recent works, providing insights into future research directions.

Keywords: *Public Activity Recognition, Video Streams, Hybrid Deep Learning, Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), Surveillance.*

I. Introduction

The growing demand for automated surveillance and smart city monitoring has led to significant interest in public activity recognition in video streams. Traditional approaches struggled with occlusions, scale variations, and illumination changes, limiting their applicability in real-world scenarios [6], [8]. The rise of deep learning, especially hybrid architectures combining CNNs and RNNs [8], [16], has greatly improved robustness and recognition accuracy. Recently, transformer-based models have further advanced video understanding by capturing long-range dependencies [1].

A. Background and Motivation

The proliferation of video surveillance systems and the emergence of smart cities have dramatically increased the volume of video data requiring automated analysis [12]. Public activity recognition aims to interpret human behaviors and interactions from video streams, enabling applications like anomaly detection, crowd management, and security monitoring [11], [14]. Traditional computer vision methods struggled with variations in scale, illumination, and occlusions [6]. The advent of deep learning, particularly hybrid models that fuse spatial and temporal learning [3], [20], has revolutionized this domain, offering significant improvements in performance [5], [15].

B. Objectives and Scope

This review aims to advance understanding of hybrid deep learning applications for public activity recognition in video streams. The specific objectives are:

- To systematically review and synthesize hybrid deep learning techniques, including CNNs, RNNs, GCNs.
- To conduct a comparative analysis of state-of-the-art hybrid architectures by examining their datasets, methodologies, performance results, and limitations.
- To identify and discuss critical challenges including data scarcity and annotation bottlenecks.
- To propose future research directions, emphasizing self-supervised learning to mitigate annotation needs, lightweight and edge-deployable architectures.

The scope of this review is confined to studies that employ hybrid deep learning models integrating multiple neural paradigms—such as CNN+RNN, CNN+GCN, and transformer-graph hybrids—focused specifically on public or crowd activity recognition in video streams [3], [5], [15], [20].

II. Literature Review

A. Existing Approaches

Earlier methods relied on handcrafted trajectories and improved dense trajectories for activity recognition [6]. The two-stream CNN model proposed by Simonyan and Zisserman effectively integrated spatial and temporal streams for video analysis [8]. Temporal segment networks (TSN) provided a framework for capturing long-term temporal structures [16].

To reduce computation for edge deployment, lightweight models like MobileNets were proposed [2], while TSM introduced temporal shift modules to efficiently model temporal information [10]. Recent advances include transformer-based approaches, leveraging attention mechanisms to handle complex dependencies in video streams [1], [20].

Graph-based models such as spatial-temporal graph convolutional networks (ST-GCN) have been explored for skeleton-based activity recognition [15], showing promising results in modeling human pose dynamics.

Notable works include:

- CNN-LSTM hybrids for sequence modeling of activities.
- 3D CNNs that jointly model space-time features.
- Graph Convolutional Networks (GCNs) applied to human pose graphs for fine-grained activity recognition.

- Transformers adapted to video understanding, emphasizing long-range dependencies.

B. Comparative Studies

Hybrid CNN-RNN models have been compared against 3D CNNs, showing benefits for long sequences [16]. Graph and attention-based models have demonstrated improved performance in dense scenarios but still face challenges like generalizing across environments [15], [17].

Table 1. Comparative Study of Reviewed Paper

Paper Title	Dataset(s) Used	Method Used	Key Contribution	Results	Limitations
<i>"Temporal Segment Networks for Action Recognition in Videos"[16]</i>	UCF101, HMDB51, Kinetics	Temporal Segment Networks (TSN)	Sparse temporal sampling to capture long-term structure	~94.2% on UCF101 (3 splits avg), ~69.4% on HMDB51	Misses very fine-grained temporal transitions; static sampling may overlook fast local actions.
<i>"Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition"[15]</i>	NTU RGB+D	ST-GCN (Graph Conv. on skeleton data)	First to use GCN for modeling spatial-temporal relations in human joints	~81.5% cross-subject accuracy on NTU RGB+D	Needs high-quality pose estimation; limited for RGB-only data or heavy occlusion.
<i>"SlowFast Networks for Video Recognition"[3]</i>	Kinetics-400, AVA	SlowFast Networks (dual pathways)	Simultaneously models slow semantic content & fast motion	~79.8% top-1 on Kinetics-400, SOTA on AVA detection	Computationally heavy; requires dual network paths increasing memory & compute.
<i>"OpenPose: Realtime Multi-Person 2D Pose Estimation Using Part Affinity"</i>	COCO, MPII	OpenPose (Pose estimation)	Real-time 2D multi-person pose estimation via part affinity fields	~61.8 mAP on COCO keypoints	Not for action recognition directly; performance drops under occlusion or

<i>Fields"[21]</i>					dense crowds.
<i>"Two-Stream Convolutional Networks for Action Recognition in Videos"[8]</i>	UCF101, HMDB51	Two-Stream CNN (RGB + Optical flow)	Explicitly fuses spatial (RGB) & temporal (optical flow) streams	~88.0% on UCF101, ~59.4% on HMDB51	Needs precomputed optical flow; increases storage & compute, struggles with long videos.

Several comparative studies have benchmarked these methods on datasets such as UCF101, HMDB51, and large-scale crowd activity datasets [7], [16]. Results generally indicate that hybrid models outperform single-stream networks by effectively leveraging both spatial and temporal information [3], [5], [15], [20].

Studies comparing CNN-LSTM and 3D CNN architectures have found that while 3D CNNs excel at capturing short-term motion dynamics, CNN-LSTM hybrids are more effective for modeling long temporal sequences, providing better performance in scenarios requiring extended temporal context [3], [8], [16]. Similarly, graph-based models have demonstrated superior capability in handling complex interactions such as crowd behaviors; however, they often depend on accurate pose estimation as a prerequisite, which can limit their applicability in environments with significant occlusions or dense crowds [15], [21].

III. Methodology of Review

A. Selection Criteria

This review focuses on research articles published between 2017 and 2025 that specifically:

- Employ hybrid deep learning models for activity recognition.
- Target public or crowd scenarios in video streams.
- Report quantitative performance metrics on recognized benchmark datasets.

Exclusions were made for studies focusing solely on private indoor activities (e.g., smart home monitoring) or using purely traditional machine learning techniques.

B. Data Sources and Search Strategy

Relevant publications were retrieved from digital libraries including IEEE Xplore, ACM Digital Library, SpringerLink, and ScienceDirect. Searches combined keywords such as:

- "public activity recognition,"
- "video streams,"
- "hybrid deep learning,"
- "CNN LSTM,"
- "graph neural networks activity."

Reference chaining from key papers was also used to capture seminal and recent contributions.

C. Inclusion and Exclusion Criteria

Included studies were required to:

- Present an original hybrid deep learning architecture or a novel application of existing architectures to public activity recognition.
- Provide empirical evaluations on public datasets.
- Discuss challenges, limitations, or propose future directions.

Papers were excluded if they:

- Lacked experimental validation.
- Focused solely on anomaly detection without general activity classification.
- Were non-peer-reviewed or workshop abstracts with insufficient technical detail.

V. Challenges and Open Issues

The lack of large annotated public activity datasets makes it difficult to train data-hungry hybrid deep models [11], [13]. Moreover, computational demands of complex architectures like SlowFast networks and transformers pose challenges for real-time applications [1], [3]. Ensuring privacy in widespread surveillance remains a key societal concern [4], [12].

Domain shift issues also cause models to underperform when tested in unseen environments, despite attempts with domain-adversarial training and CORAL alignment [9], [19].

Despite remarkable progress, several challenges and open issues remain in advancing public activity recognition in video streams using hybrid deep learning techniques.

Challenges in advancing public activity recognition using hybrid deep learning techniques remain significant. A key issue is data scarcity and annotation complexity, as large-scale, diverse datasets with annotated public activities are limited, and manual labeling of lengthy video sequences is both time-consuming and susceptible to inconsistencies, ultimately hindering

model generalization [1], [2]. Additionally, scalability and real-time constraints pose serious hurdles; hybrid models, particularly those integrating CNNs with RNNs or GCNs, demand considerable computational resources, making their deployment in real-time applications such as city-wide surveillance systems challenging [3]. Another major concern is handling occlusions and dense crowds, where frequent visual obstructions and high crowd densities degrade the performance of models that depend heavily on accurate pose estimation or trajectory tracking [4]. Furthermore, achieving generalization to unseen environments remains difficult, as models trained on specific datasets often underperform when faced with different camera angles, lighting conditions, or cultural contexts [5]. Finally, the proliferation of automated surveillance raises important privacy and ethical considerations, necessitating robust frameworks to ensure data protection and responsible AI deployment [6].

VI. Future Directions

Future work can leverage self-supervised learning to reduce annotation needs [13], while variational knowledge distillation and lightweight CNN architectures enable deployment on edge devices [2], [18]. Combining spatial-temporal graphs with transformer mechanisms may yield more robust models for crowded scenes [1], [15]. Further, privacy-preserving frameworks such as differential privacy and federated learning are critical for ethical deployments [4], [12].

Several promising research directions can address these challenges and further advance the field. Self-supervised and semi-supervised learning approaches that leverage unlabeled video data could substantially reduce reliance on large annotated datasets, thereby easing data bottlenecks [7]. In parallel, efforts toward lightweight and edge-AI models, such as through model pruning or knowledge distillation, are essential for enabling the deployment of compact hybrid architectures on edge devices, ensuring low latency and real-time processing capabilities [8]. Additionally, domain adaptation and transfer learning techniques hold significant potential for improving the generalization of models across diverse environments, enhancing their robustness to variations in camera viewpoints, lighting, and cultural contexts [9]. Integrating spatio-temporal graphs with transformer architectures represents another compelling direction, as it could more effectively capture long-range dependencies and intricate interactions in crowded scenes [10]. Finally, advancing privacy-preserving methods, including federated learning and differential privacy, will be critical to developing responsible activity recognition systems that uphold individual rights and ethical standards in pervasive surveillance settings [11].

VII. Conclusion

This review highlighted how hybrid deep learning architectures, from CNN-RNN combinations to graph and transformer-based approaches, have transformed public activity recognition in video

streams. Despite these advances, challenges like data scarcity, computational constraints, generalization, and privacy remain open. Future research must aim to develop scalable, generalizable, and ethically responsible solutions to enable widespread, real-time adoption of public activity recognition systems.

This paper reviewed the evolution and current state of hybrid deep learning techniques for public activity recognition in video streams. Hybrid models integrating CNNs, RNNs, GCNs, and transformers have significantly advanced recognition capabilities, outperforming traditional methods by jointly modeling spatial and temporal dynamics. However, challenges such as data annotation burdens, scalability, and ethical concerns persist. Future work should explore self-supervised learning, efficient architectures, and privacy-preserving frameworks to build robust, adaptable, and socially responsible activity recognition systems.

References

- [1] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *Proc. Int. Conf. Learning Representations (ICLR)*, 2021.
- [2] A. G. Howard et al., "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [3] B. Feichtenhofer et al., "SlowFast Networks for Video Recognition," *Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV)*, pp. 6202–6211, 2019.
- [4] C. Dwork and A. Roth, "The Algorithmic Foundations of Differential Privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, nos. 3–4, pp. 211–407, 2014.
- [5] H. Chen et al., "Rethinking Temporal Fusion for Video-based Human Activity Recognition," *IEEE Trans. Circuits and Systems for Video Technology*, vol. 32, no. 6, pp. 3904–3918, Jun. 2022.
- [6] H. Wang and C. Schmid, "Action Recognition with Improved Trajectories," *Proc. IEEE Int. Conf. Computer Vision (ICCV)*, pp. 3551–3558, 2013.
- [7] J. Carreira and A. Zisserman, "Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 6299–6308, 2017.
- [8] K. Simonyan and A. Zisserman, "Two-Stream Convolutional Networks for Action Recognition in Videos," *Proc. Advances in Neural Information Processing Systems (NIPS)*, vol. 27, pp. 568–576, 2014.
- [9] L. Sun and K. Saenko, "Deep CORAL: Correlation Alignment for Deep Domain Adaptation," *Proc. European Conf. Computer Vision Workshops (ECCVW)*, pp. 443–450, 2016.
- [10] M. Lin et al., "TSM: Temporal Shift Module for Efficient Video Understanding," *Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV)*, pp. 7082–7092, 2019.

- [11] M. Sabokrou, M. Fathy, M. Hoseini, and R. Klette, "Deep-anomaly: Fully convolutional neural network for fast anomaly detection in crowded scenes," *Computer Vision and Image Understanding*, vol. 172, pp. 88–97, 2018.
- [12] M. Z. Adnan, N. A. Nizam, and M. N. M. Nasir, "A Review on Privacy and Ethical Challenges in Video Surveillance," *IEEE Access*, vol. 9, pp. 125485–125499, 2021.
- [13] S. Misra et al., "Self-Supervised Learning for Video-based Person Re-Identification," *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 7154–7163, 2020.
- [14] S. Sudhakaran and O. Lanz, "Learning to Recognize Abnormality with Deep Spatiotemporal Networks," *Proc. IEEE Int. Conf. Advanced Video and Signal-Based Surveillance (AVSS)*, pp. 1–6, 2018.
- [15] S. Yan, Y. Xiong, and D. Lin, "Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition," *Proc. AAAI Conf. Artificial Intelligence*, vol. 32, no. 1, 2018.
- [16] X. Wang et al., "Temporal Segment Networks for Action Recognition in Videos," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 41, no. 11, pp. 2740–2755, Nov. 2019.
- [17] X. Zhang et al., "Multi-Scale Attention with Dense Encoder for Weakly Supervised Temporal Action Localization," *IEEE Trans. Image Processing*, vol. 31, pp. 5112–5124, 2022.
- [18] Y. Choi et al., "Variational Knowledge Distillation for Lightweight Human Activity Recognition," *IEEE Internet of Things Journal*, vol. 9, no. 1, pp. 279–290, Jan. 2022.
- [19] Y. Ganin et al., "Domain-Adversarial Training of Neural Networks," *J. Machine Learning Research*, vol. 17, no. 59, pp. 1–35, 2016.
- [20] Y. Sun, M. Lin, and X. Tang, "Spatial Temporal Attention Based Method for Human Action Recognition in Videos," *Proc. IEEE Int. Conf. Image Processing (ICIP)*, pp. 3218–3222, 2019.
- [21] Z. Cao et al., "OpenPose: Realtime Multi-Person 2D Pose Estimation Using Part Affinity Fields," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 43, no. 1, pp. 172–186, Jan. 2021.