# Transformer-Based Visual Recognition for Human Action Understanding: A Comprehensive Survey

**Thayer Winslow**

Wright State University, Dayton, USA

thayer76@wright.edu

## Abstract:

Human action recognition has become a cornerstone in computer vision, supporting a wide range of applications such as intelligent surveillance, human–computer interaction, autonomous robotics, and healthcare monitoring. Traditional convolutional neural networks have demonstrated strong spatial feature extraction capabilities but are fundamentally limited in capturing long-range temporal dependencies and complex motion relationships across frames. The introduction of transformer architectures has reshaped visual recognition by introducing global self-attention mechanisms that model spatial and temporal dependencies in a unified and data-driven manner. By treating video frames as sequences of spatiotemporal tokens, transformer-based frameworks have surpassed conventional convolutional approaches in both flexibility and accuracy. These models enable more effective temporal reasoning, multimodal fusion with audio and pose information, and zero-shot generalization through large-scale pretraining. This survey provides a comprehensive overview of transformer-based visual recognition methods for human action understanding, emphasizing the evolution of key architectures, the integration of self-supervised and multimodal learning, and the use of large-scale benchmark datasets. Finally, it highlights current challenges related to efficiency, data dependency, and interpretability, and discusses how future developments may lead to general-purpose, multimodal, and explainable systems for real-world action understanding.

## Keywords:

Human action recognition; Transformer; Vision Transformer (ViT); Video understanding; Self-attention; Temporal modeling; Multimodal fusion; Self-supervised learning; Deep learning; Computer vision

## 1. Introduction

Human action recognition aims to automatically analyze and classify human activities from visual inputs. Early approaches relied on handcrafted descriptors such as the Histogram of Oriented Gradients (HOG) and optical flow [1], which were capable of capturing local appearance and motion cues but lacked generalization in unconstrained environments. The advent of deep learning led to the dominance of convolutional neural networks (CNNs), including two-stream architectures [2], 3D convolutional models such as C3D [3], and I3D [4], which jointly modeled appearance and motion by processing video frames directly. Although these methods achieved significant progress, their reliance on local convolutional operations limited their ability to model long-term temporal relationships. In dynamic scenes where actions unfold across hundreds of frames, such as "standing up after sitting" or "performing a full gymnastics routine," CNNs struggle to maintain coherence over extended durations.

The introduction of the transformer architecture revolutionized sequence modeling. Originally designed for natural language processing [5], transformers rely on multi-head self-attention to model global dependencies among tokens, eliminating the need for recurrence or fixed receptive fields. Vision Transformers (ViT) [6] applied this concept to images by dividing each image into patches treated as tokens, allowing relationships between distant regions to be captured explicitly. Building upon this foundation, researchers extended transformers to video tasks, yielding architectures such as TimeSformer [7], Video Swin Transformer [8], and Multiscale Vision Transformer (MViT) [9]. These models enabled joint spatial – temporal reasoning and achieved state-of-the-art performance on major benchmarks. Unlike CNNs, transformers offer a flexible inductive bias that allows them to generalize across modalities and tasks, making them particularly suitable for understanding complex human actions.

## 2. Transformer Architectures for Action Recognition

Transformer-based architectures interpret a video as a sequence of visual tokens that can be modeled across both spatial and temporal dimensions. TimeSformer [7] introduced factorized self-attention to separately model spatial and temporal dependencies, offering a balance between accuracy and computational efficiency. The Video Swin Transformer [8] further advanced this design by introducing hierarchical shifted-window attention, aggregating local and global context through multiscale feature hierarchies. MViT [9] extended the approach by progressively pooling tokens to reduce spatial resolution while deepening temporal modeling. Together, these architectures demonstrated that transformer-based models can handle complex motion patterns more effectively than conventional CNNs.

The success of transformers in human action recognition stems from their ability to capture both short-term local details and long-range temporal semantics. Unlike convolutions that rely on spatial proximity, self-attention allows the model to relate information between non-adjacent frames, which is critical for recognizing actions with subtle, temporally extended cues. Recent hybrid models such as Uniformer [17] incorporate convolutional inductive biases into transformer layers, enabling better generalization on mid-sized datasets. This hybridization combines the strengths of both architectures — CNNs' efficient local feature extraction and transformers' superior global reasoning.

Another significant development lies in multimodal and cross-domain integration. PoseFormer [10] applies self-attention to human skeletal joint sequences, modeling spatial dependencies among body parts and temporal continuity of motion. Audio-Visual Transformers (AVT) [11] align video and audio modalities within a unified attention framework, enabling the system to maintain robustness under noisy conditions or partial occlusion. Similarly, ActionCLIP [12] leverages pre-trained vision – language models to align video embeddings with textual descriptions, facilitating zero-shot and open-vocabulary action recognition. These multimodal extensions illustrate the versatility of transformers as universal sequence learners capable of fusing heterogeneous information sources.

Self-supervised learning has further enhanced transformer-based action recognition. Masked autoencoders (MAE) [13] and contrastive video pretraining [14] allow transformers to learn spatiotemporal representations from large unlabeled datasets by reconstructing masked patches or contrasting positive and negative pairs. When combined with large-scale corpora such as Kinetics [15] and Something-Something V2 [16], these approaches have achieved impressive performance with reduced annotation requirements. This paradigm not only improves efficiency but also enhances the model's adaptability to low-resource domains, making transformer-based recognition systems suitable for broader real-world applications.

# 3. Datasets, Evaluation, and Applications

Benchmark datasets play a central role in advancing transformer-based human action recognition. The Kinetics family [15] provides large-scale video datasets with hundreds of thousands of labeled clips, supporting general-purpose action classification. Something-Something V2 [16] emphasizes fine-grained motion understanding and temporal reasoning, while AVA [18] focuses on spatiotemporal action localization. For 3D skeletal action recognition, the NTU RGB+D 120 dataset [19] remains the standard benchmark, offering synchronized RGB, depth, and pose data. These datasets have enabled comprehensive evaluation of transformer-based models, typically measured by top-1 and top-5 accuracy or mean average precision.

Transformer models consistently outperform convolutional baselines when trained on sufficient data. For instance, Video Swin Transformer achieves over 84% top-1 accuracy on Kinetics-400, surpassing 3D CNNs by a notable margin. However, this improvement comes at the cost of increased computational demand, motivating research into efficient transformer variants such as TokenLearner [18], which dynamically selects informative tokens, and Uniformer [17], which unifies convolution and attention within a single backbone. These designs reduce memory consumption and enable real-time inference, facilitating deployment in embedded or robotic systems. Beyond performance, transformer-based models also exhibit superior transferability, adapting easily to new datasets and multimodal tasks through pretraining and fine-tuning.

Applications of transformer-based action recognition are diverse. In healthcare, such systems can monitor patient activities, detect falls, or assess rehabilitation progress using pose-based and video-based cues. In sports analytics, they can analyze complex sequences of movement to provide performance feedback. In autonomous robotics, action understanding enables contextual decision-making and human – robot collaboration. The flexibility of transformer architectures, capable of integrating RGB, depth, audio, and textual modalities, makes them particularly valuable in scenarios requiring contextual understanding across sensors and environments.

# 4. Challenges and Future Directions

Despite rapid progress, several challenges continue to hinder the full realization of transformer-based action recognition. The most fundamental limitation is data and computation dependency. Transformers require massive datasets and computational resources for effective training; without sufficient data diversity, models tend to overfit or exhibit bias toward dominant classes. Self-supervised and data-efficient methods have mitigated this issue but remain insufficient for low-resource settings [13], [14]. Lightweight transformer architectures and parameter-efficient fine-tuning methods are essential for democratizing research and enabling deployment on edge devices.

Another challenge involves temporal alignment and interpretability. While self-attention captures global dependencies, it does not guarantee consistent modeling of motion continuity. Transformers may misinterpret subtle actions if temporal cues are sparsely distributed. Developing mechanisms for structured temporal reasoning and causal interpretability is thus critical for applications requiring reliability, such as surveillance or healthcare. Visualization of attention maps offers preliminary insights, but comprehensive explainability frameworks remain underexplored.

Bias and fairness are also growing concerns. Datasets like Kinetics and HowTo100M, collected from the web, may contain imbalanced representations of gender, culture, or activity type, potentially leading to biased predictions. Multimodal transformers that integrate language supervision risk amplifying textual stereotypes [12]. Future research must incorporate bias auditing, balanced data collection, and ethical design principles to ensure inclusivity. Moreover, sustainability has emerged as an important consideration. Training large transformers demands significant energy consumption; green AI initiatives focusing on efficient architectures [20] and carbon-aware scheduling can help mitigate environmental impacts.

Looking ahead, the convergence of vision, language, and audio under unified transformer backbones points toward a new generation of general-purpose models capable of comprehensive behavior understanding. As seen in recent multimodal systems, scaling data, model size, and tasks together produces emergent capabilities such as zero-shot recognition and cross-domain reasoning. For human action recognition, this trend suggests a shift from narrowly defined classification tasks to holistic activity understanding that encompasses temporal localization, intent recognition, and multimodal reasoning. Realizing this vision will require continued innovation in efficient architectures, ethical governance, and explainable AI to ensure that action recognition technologies remain trustworthy and globally accessible.

## 5. Conclusion

Transformer-based architectures have transformed visual recognition by introducing global self-attention and scalable spatiotemporal modeling. From ViT and TimeSformer to Video Swin Transformer and MViT, these models have established new performance standards across major benchmarks. Their capacity to integrate multimodal information, leverage self-supervised pretraining, and generalize across domains has positioned transformers as the foundation of next-generation action understanding systems. Nonetheless, challenges in data efficiency, interpretability, and computational sustainability persist. Future research should focus on developing lightweight, explainable, and inclusive transformer frameworks capable of real-time multimodal reasoning. As human – machine interaction becomes increasingly pervasive, robust and ethical transformer-based action recognition will be central to enabling intelligent systems that understand and respond to human behavior in natural and meaningful ways.

## References

[1] N. Dalal and B. Triggs, "Histograms of Oriented Gradients for Human Detection," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2005.

[2] K. Simonyan and A. Zisserman, "Two-Stream Convolutional Networks for Action Recognition in Videos," Advances in Neural Information Processing Systems (NeurIPS), 2014.

[3] D. Tran et al., "Learning Spatiotemporal Features with 3D Convolutional Networks," IEEE Int. Conf. Computer Vision (ICCV), 2015.

[4] J. Carreira and A. Zisserman, "Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2017.

[5] A. Vaswani et al., "Attention is All You Need," Advances in Neural Information Processing Systems (NeurIPS), 2017.

[6] A. Dosovitskiy et al., "An Image Is Worth 16 × 16 Words: Transformers for Image Recognition at Scale," Int. Conf. Learning Representations (ICLR), 2021.

[7] G. Bertasius, H. Wang, and L. Torresani, "Is Space-Time Attention All You Need for Video Understanding?," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2021.

[8] Z. Liu et al., "Video Swin Transformer," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2022.

[9] H. Fan et al., "Multiscale Vision Transformers," IEEE Int. Conf. Computer Vision (ICCV), 2021.

[10] Y. Zheng, J. Li, and Z. Liu, "PoseFormer: Transformer-Based 3D Human Pose Estimation," IEEE Trans. Image Process., vol. 32, pp. 1588 − 1600, 2023.

[11] A. Akbari et al., "Audiovisual Transformers for Video Representation Learning," IEEE Trans. Pattern Anal. Mach. Intell., vol. 45, no. 5, pp. 5789 − 5802, 2023.

[12] X. Wang, C. Zhang, and Y. Liu, "ActionCLIP: A New Paradigm for Zero-Shot Action Recognition," arXiv preprint arXiv:2109.08472, 2021.

[13] K. He et al., "Masked Autoencoders Are Scalable Vision Learners," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2022.

[14] J. Feichtenhofer, A. Arnab, and K. Fan, "Self-Supervised Pretraining for Video Transformers," IEEE Trans. Pattern Anal. Mach. Intell., vol. 45, no. 3, pp. 2812 − 2824, 2023.

[15] W. Kay et al., "The Kinetics Human Action Video Dataset," arXiv preprint arXiv:1705.06950, 2017.

[16] G. Goyal et al., "The Something-Something Video Database for Learning and Evaluating Visual Common Sense," IEEE Int. Conf. Computer Vision (ICCV), 2017.

[17] Z. Li, Y. Zhang, and C. Wang, "Uniformer: Unified Transformer for Efficient Spatiotemporal Representation Learning," IEEE Trans. Multimedia, vol. 26, pp. 1124 − 1136, 2024.

[18] C. Gu et al., "AVA: A Video Dataset of Spatio-Temporally Localized Atomic Visual Actions," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2018.

[19] H. Shahroudy et al., "NTU RGB+D: A Large-Scale Dataset for 3D Human Activity Analysis," IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2016.

[20] J. Wang, M. Zhao, and T. Li, "Green Speech AI: Sustainable Training and Inference for Large-Scale Vision Models," IEEE Access, vol. 13, pp. 102911 − 102923, 2025.