



A Hybrid 3D CNNs Transformer Architecture for Video-Based Human Action Recognition with Improved Accuracy on UCF101 and HMDB51

Engin Seven^{1*} , Eylem Yucel² 

^{1,2} Department of Computer Engineering, İstanbul University-Cerrahpaşa, 34320 İstanbul, Türkiye

ARTICLE INFO

Received Date: 21/05/2025

Accepted Date: 29/09/2025

Cite this paper as:

Seven, E., & Yucel, E. (2025). A Hybrid 3D CNNs Transformer Architecture for Video-Based Human Action Recognition with Improved Accuracy on UCF101 and HMDB51. *Journal of Innovative Science and Engineering*. 9(2), 327-342.

*Corresponding author: Engin Seven
E-mail:engin.seven@ogr.iuc.edu.tr

Keywords:

Human Activity Recognition
Video-based Action Recognition
3D Convolutional Neural Networks
Attention Mechanism
Deep Learning in Computer Vision

© Copyright 2025 by
Bursa Technical University. Available
online at <http://jise.btu.edu.tr/>



The works published in Journal of Innovative Science and Engineering (JISE) are licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

ABSTRACT

Video-Based Human Action Recognition (HAR) remains challenging due to inter-class similarity, background noise, and the need to capture long-term temporal dependencies. This study proposes a hybrid deep learning model that integrates 3D Convolutional Neural Networks (3D CNNs) with Transformer-based attention mechanisms to jointly capture spatio-temporal features and long-range motion context. The architecture was optimized for parameter efficiency and trained on the UCF101 and HMDB51 benchmark datasets using standardized preprocessing and training strategies. Experimental results indicate that the proposed model reaches 97% accuracy and 96.8% mean F1-score on UCF101, and 85% accuracy, and 83.8% F1-score on HMDB51, showing consistent improvements compared to the standalone 3D CNNs and Transformer variants under identical settings. Ablation studies confirm that the combination of convolutional and attention layers significantly improves recognition performance while maintaining competitive computational cost (3.78M parameters, 17.75 GFLOPs/video, ~7 ms GPU latency). These findings highlight the effectiveness of the hybrid design for accurate and efficient HAR. Future work will address class imbalance using focal loss or weighted training, explore multimodal data integration, and develop more lightweight Transformer modules for real-time deployment on resource-constrained devices.

1. Introduction

Human Activity Recognition (HAR) has become one of the critical application areas of artificial intelligence and computer vision technologies. Its main purpose is to digitise a person's physical movements and automatically classify them into meaningful action classes. This technology is used in a wide range of applications, from security to health,

from smart living environments to smart city systems, making human-machine interaction smarter, more intuitive and safer. In particular, the deep learning revolution in recent years has significantly improved the accuracy of HAR systems and made these systems more widely available [1].

Security is one of the sectors where HAR systems have found the earliest and most widespread

application. Artificial intelligence-based action recognition solutions integrated into video surveillance systems are able to automatically detect anomalies that may be missed by human operators. These systems not only detect threats such as physical attacks, fights, and intrusions, but also enable proactive warning systems by detecting the patterns of occurrence of these threats in advance [2]. However, accurately recognizing activities from video data remains a significant challenge due to temporal dependencies, complex movements, and background variability, as shown in Figure 1.

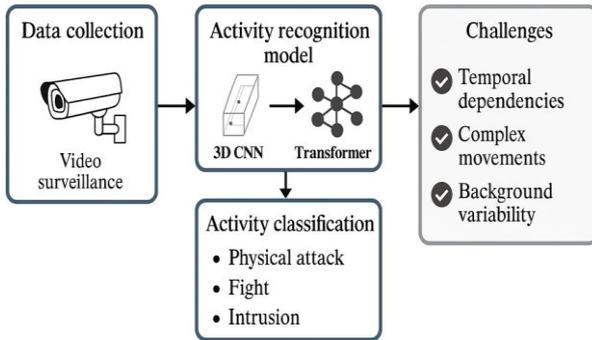


Figure 1: The Challenge of Recognizing Human Activities from Video and the Proposed Hybrid Approach.

For instance, the correlation between the increase in crowd density in shopping malls and aggressive behaviour and panic movements can be utilised for the purpose of early intervention. Furthermore, in lieu of offline analysis of images from closed circuit camera systems, dynamic security policies can be formulated by leveraging real-time processing structures. The utilisation of HAR systems in healthcare services is predominantly oriented towards individual follow-up and the enhancement of quality of life. It is particularly beneficial for elderly individuals, patients with mobility limitations, or individuals with chronic diseases as these systems enable early intervention by automatically detecting situations such as falls, fainting, and prolonged immobilisation.

The advent of wearable devices has enabled the instantaneous transmission of data to health institutions, with select systems possessing the capability to execute personalised risk prediction models [3]. Furthermore, HAR systems have been demonstrated to make a significant contribution to physical therapy and rehabilitation processes. Systems that evaluate the correct performance of exercises enable doctors to monitor the process remotely, increasing patient satisfaction and facilitating the decentralisation of healthcare services [4]. In the context of smart home systems, HAR is

regarded as a pivotal component in the development of living spaces that are more comfortable, personalised, and energy efficient. The movement pattern of the individual is analysed through sensors; systems such as lighting, heating, ventilation, and alarms are automatically adjusted according to the results of this analysis [5]. For instance, when an individual rises from bed in the morning, the system is programmed to automatically open the curtains, initiate the coffee machine, or illuminate the living space. Furthermore, these systems are also employed in the domain of behavioural biometrics, where their capacity to analyse movement is utilised for authentication purposes. Consequently, HAR systems enhance convenience and security [6]. It is evident that the functionality offered by HAR systems has the capacity to enhance the quality of life for individuals. Furthermore, it has been demonstrated that such systems can engender substantial improvements in terms of safety, healthcare costs, and energy consumption on a societal level. The advent of technological innovation has precipitated the augmentation of the capabilities of HAR solutions, which are anticipated to become even more sophisticated in the future, utilising multimodal data sources (e.g. video, audio, and accelerometer data).

The real-world applicability of HAR systems is constrained by a variety of technical, environmental, and structural challenges. It is evident that these challenges have a substantial impact on the key performance criteria of the models, including accuracy, generalisation, and speed of operation. This, in turn, renders it difficult for the systems to operate reliably under different conditions. It is important to note that a highly accurate HAR model may, in theory, demonstrate a loss of performance in real-world conditions due to the presence of numerous variables. The primary and most prevalent issue is the presence of noise and distortion in the data. In sensor-based HAR systems, signals from devices such as accelerometers or gyroscopes may be subject to corruption due to hardware faults, sensor calibration issues, or external factors that arise during utilisation.

In the context of video-based systems, factors such as low-resolution images, motion blur, light imbalances, and variations in camera angle can impede the clear perception of actions. Such noises have been demonstrated to have a direct impact on both the accuracy of the model and the learning dynamics during training [7, 8].

Another critical challenge is that of labelling errors and data inconsistency. The utilisation of supervised

learning in HAR models necessitates the availability of correctly labelled data. However, the labelling of datasets is often a manual process, and human errors frequently occur during this process. Incorrect class assignments, unstable labelling of borderline actions, or actions interpreted in different ways by different people can cause the model to learn incorrect patterns during the training process. The generalisability of HAR systems pertains to the capacity of the system to function effectively under diverse user profiles, a range of devices, and a variety of environmental settings. The most significant challenge in this regard pertains to the heterogeneity among users. It is important to note that this action can be performed by different people at different speeds, with different ranges of motion, and even in different contexts of meaning. To illustrate this point, consider the physical differences between a young and an elderly individual when walking [9].

This necessitates the training of the model for each individual separately or the development of more complex universal representation learning strategies. The utilisation of multiple sensor types in conjunction within HAR systems invariably gives rise to the issue of data heterogeneity. It is important to note that the various sensors (e.g. accelerometer, magnetometer, and gyroscope) involved in this process generate data at different rates, in different formats, and at different scales. The alignment, synchronisation, and conversion of this multimodal data into meaningful representations is a process that creates significant complexity at both the engineering and algorithmic level [10]. Furthermore, a significant proportion of HAR applications are characterised by real-time requirements. This is particularly evident in domains such as security, healthcare, and robotics, where decisions must be made in milliseconds. However, deep learning-based models that provide high accuracy frequently possess large parameter sizes and high computational burden. This limitation is particularly pronounced in portable devices and edge computing environments, where it results in an increase in inference time and a decrease in system response speed [11]. These challenges necessitate the development of new methods to make HAR systems more robust, scalable, and context-aware. In this context, research areas such as data cleansing, domain adaptation, self-supervised learning, and efficient architectural designs will be of particular significance in the future.

Video-based human action recognition represents a foundational computer vision problem that aims to automatically analyse and classify physical movements performed by individuals from video

images. This technology is of great importance for many critical application areas such as intelligent security systems, remote patient monitoring, sports performance analysis, driver behaviour monitoring, and human-computer interaction [12]. The successful operation of HAR systems on video data is the basis for smarter city infrastructures, safer public spaces, and more personalised healthcare. Nevertheless, video-based HAR is confronted with a considerably more intricate data structure in comparison to sensor-based solutions. It has been demonstrated that factors such as camera position, resolution, image quality, background mobility, camera angle changes, lighting conditions, and the presence of multiple individuals in the scene directly affect model performance [1, 13]. To illustrate this point, action discrimination is rendered significantly more arduous in videos captured under suboptimal lighting conditions. Moreover, factors such as camera vibration induce temporal inconsistencies, further complicating the analysis. Furthermore, due to the visual similarity of certain actions (e.g. "sitting" and "squatting"), discrimination between classes poses a significant challenge. In order to surmount such issues, it is imperative to accurately model not only spatial information but also the evolution of motion over time. In this context, the primary objective of the study is grounded in four overarching goals. Firstly, the objective is to develop a system that utilises deep learning to achieve more precise distinction of different actions in complex visual data. This is particularly relevant in datasets characterised by high levels of intra-class variations and out-of-class similarities, where the model must be capable of accurate classification by focusing on fine details. Secondly, the integration of 3D CNNs and Transformer-based attentional mechanisms was planned for the purpose of more effective modelling of temporal dynamics. While 3D CNNs structures are capable of establishing temporal connections between video frames, Transformer structures have the capacity to offer a more comprehensive motion analysis by modelling long-term dependencies [14].

Thirdly, robust attention mechanisms will be applied to minimise problems such as noise, background variation, and class similarity, which are frequently encountered in video data. These mechanisms are designed to minimise errors caused by redundant information by ensuring that the model focuses exclusively on the most significant action-related features.

Ultimately, the objective is to minimise the computation time and computational load of the developed model, while achieving high accuracy

rates. In accordance with this objective, the parameter efficiency of the model has been augmented and FLOPs values have been minimised. It is evident that the proposed system has been optimised in such a manner that its utilisation is compatible with both laboratory environments and real-time applications (e.g. CCTV-based monitoring systems).

In addition, the system is to be modified to ensure compatibility with edge devices and to allow for flexibility in operation within distributed environments. Consequently, energy consumption and processing costs will be minimised, and the system will be integrated into areas such as smart city projects, wearable devices, and mobile security solutions. The proposed architecture is intended to provide an advanced video-based HAR solution for both academic and industrial applications.

The book chapter is structured under six main headings, providing comprehensive coverage of video-based human action recognition. The chapter organisation has been meticulously designed to facilitate the reader's navigation through the subject matter in a step-by-step manner, thereby enabling them to gain a comprehensive and nuanced understanding of the current research directions within the field.

Firstly, the Introduction provides a comprehensive overview of the significance of video-based human action recognition and the primary motivations underpinning this field of study. The text emphasises the significance of automatic classification of human behaviour over video data for a variety of critical applications, including security, healthcare, smart home systems, and human-computer interaction. Moreover, the objectives of the study, the research questions, and the gaps that the methodology developed in this chapter aims to fill are clearly presented.

Secondly, the Background section reviews important approaches, technical concepts, and methodological trends in the literature. In this context, traditional machine learning-based HAR methods, the paradigm shift with deep learning, and the application of modern architectures such as 3D CNNs and Transformer to video data are discussed in detail. Furthermore, the section addresses the common challenges encountered in HAR systems (e.g. data noise, variations, and real-time requirements), thereby contributing to the theoretical foundation of the current work.

Thirdly, the proposed methodology is outlined, including the general architecture of the developed system, the data sets utilised, and the preprocessing steps applied. In this section, the structural characteristics of the model (for example, 3D CNNs layers, Transformer modules, attention mechanisms), the hyperparameters employed during training, and the experimental environment are elucidated. The strategies applied to different datasets and the process of optimising the model are also detailed here.

The fourth section, entitled 'Experimental Results', undertakes an analysis of the performance of the proposed system in terms of its internal dynamics and in comparison with extant methods in the literature. In addition to fundamental metrics such as accuracy, F1 score, and Top-1/Top-5 success rates, technical indicators including computational complexity (FLOPs), parameter size, and execution speed are also analysed. Furthermore, visual presentations, such as the confusion matrix, provide a detailed analysis of the model's performance on a class-by-class basis, highlighting its strengths and weaknesses.

Chapter 5, entitled 'Conclusions and Future Directions', provides a concise summary of the findings and offers guidance on the practical application of these results in an applied context. Concurrently, the constraints of the prevailing system are candidly examined, and prospective enhancements (e.g. lightweight architectures, multimodal data utilisation, and ethical AI designs) are proposed. Adopting this systematic approach enables readers to develop a comprehensive understanding of the domain of video-based human action recognition, encompassing fundamental concepts and sophisticated technical intricacies. This comprehensive foundation provides a robust foundation for exploring novel frontiers in this pivotal research domain.

2. Material and Methods

2.1. Background

Human Action Recognition refers to automatically analysing and classifying human activities from data sources such as sensors or videos, typically represented as time series. These systems detect actions, like walking, sitting, running, or bending, and integrate results into decision-support systems in domains including security, healthcare, sports analytics, and smart city applications.

Early HAR research relied on classical machine learning methods such as decision trees, SVM, k-NN, Naive Bayes, and Hidden Markov Models (HMM).

These approaches achieved reasonable accuracy for simple, controlled actions but showed limited generalisation across diverse users, environments, and complex movements due to dependence on small datasets and handcrafted features [15]. The emergence of deep learning brought a paradigm shift: Convolutional Neural Networks (CNNs), 3D CNNs, LSTM networks, and Transformer-based models dramatically improved HAR by learning features directly from raw data. CNNs capture spatial patterns, while recurrent models (LSTM, GRU) model temporal dynamics. Recently, hybrid CNNs RNN combinations and self attention-based Transformers have become state-of-the-art solutions.

HAR is commonly divided into video-based and sensor-based approaches. Video systems analyse people and objects in camera footage to infer behaviours, offering rich visual context but facing high data volume, privacy, and environmental sensitivity issues. Sensor-based HAR uses wearable accelerometers, gyroscopes, and magnetometers to capture motion with lower energy and data needs, yet it struggles on complex actions without visual context [16].

To leverage complementary strengths, multimodal/hybrid systems combine video and sensor data. HAR performance, however, is challenged by data diversity (user behaviours, sensor placements, device types, and environments), person-to-person variability in executing the same action, and real-time constraints that make deep models latency-prone on mobile/edge hardware. Moreover, many public datasets are small and class-imbalanced, causing learning biases and overfitting to frequent classes [17].

Deep learning has become central to video-based HAR, with CNNs and more recently Transformers driving major advances since the mid-2010s. Early breakthroughs include the Two-Stream CNNs that separates spatial appearance (RGB) and temporal motion (optical flow) before fusion, achieving strong accuracy but incurring high optical-flow cost that hinders real-time use [18]. SlowFast later modeled videos at dual time resolutions slow for semantics, fast for motion improving recognition of rapid or complex actions and scaling well in practice [14].

Transformers entered vision with ViT, showing that self attention alone can learn visual representations and inspiring video-specific variants such as TimeSformer and Video Swin Transformer [19]. In parallel, the need to jointly capture space-time led to 3D convolutions: C3D extended 2D filters to $3 \times 3 \times 3$

kernels to learn spatio-temporal patterns directly [19]. I3D “inflated” 2D ImageNet filters to 3D, delivering strong transfer and state-of-the-art results on Kinetics and influencing many subsequent designs [20]. Overall, C3D/I3D established effective spatio-temporal feature extraction, while modern work increasingly combines 3D CNNs with attention mechanisms to balance accuracy and efficiency.

Beyond CNNs, video-based HAR has evolved toward architectures that better model motion over time through temporal sampling, attention, and graph structures. Temporal Segment Networks (TSNs) pioneered long-range temporal modeling by sampling fixed segments across entire videos and aggregating features, achieving strong results even with limited data and integrating easily with backbones like ResNet and Inception [21]. Non-local Neural Networks introduced global context modeling by allowing interactions between any two positions in a video, overcoming the locality of convolutions and laying groundwork for attention mechanisms in vision [22].

For skeleton-based HAR, Spatial-Temporal Graph Convolutional Networks (ST-GCNs) represent joints as graph nodes and apply spatio-temporal convolutions, effectively capturing biomechanical motion patterns and performing well in crowded or low-resolution scenes [23]. Transformer-based approaches have gained strong momentum in HAR. The Video Action Transformer Network uses self attention and person-specific queries to focus on individual actions, achieving strong results on the AVA dataset and enabling precise localization in multi-actor scenes [24].

TimeSformer is among the first pure Transformer models for video classification, processing spatial and temporal dimensions separately and capturing long-range dependencies efficiently [25].

Recent video transformers target data/compute efficiency. VideoMAE uses masked autoencoding for self-supervised pretraining with high tube masking, showing strong transfer to UCF101/HMDB51 [26]; VideoMAE-V2 scales model/data via dual masking toward video foundation models [27]. In parallel, MobileViT/MobileViTv2 combine lightweight convolutions with token mixing to reduce parameters and latency [28]. MViT/MViTv2 build multi-scale hierarchies to balance accuracy and compute for video recognition. Our hybrid 3D-CNNs + Transformer follows this efficiency motivation with compact spatio-temporal tokens and a lightweight attention head. Overall, techniques such as segment

sampling, non-local operations, graph convolutions, and attention mechanisms have become core elements of modern HAR systems, improving accuracy, flexibility, and generalisability.

2.2. A Historical Analysis of Human Activity Recognition

Human Action Recognition sits at the intersection of computer vision and AI and has progressed from early statistical models to deep learning and self attention methods. Applications span video analytics, robotics, healthcare, and security. Early HAR (late 1990s–early 2000s) targeted simple actions (e.g., walking, sitting, and running) in short, low resolution, controlled videos, with Hidden Markov Models effectively capturing temporal sequentiality and forming the basis of initial systems [29].

By the mid-2000s, feature engineering dominated: Dense Trajectories [30] and descriptors such as HOG, HOF, and MBH achieved strong accuracy but struggled to generalise in complex, unconstrained scenes. In the mid-2010s, deep learning reshaped HAR; 2D CNNs entered the field, and the Two-Stream CNNs learned appearance (RGB) and motion (optical flow) in separate branches, though limited temporal continuity between frames constrained performance [18]. To address this challenge, 3D CNNs architectures have been developed, whereby each convolutional filter is rendered three-dimensional by incorporating the time dimension into the modeling. The C3D and I3D models are regarded as the vanguard of this development, having achieved a more comprehensive action representation through direct processing of both spatial and temporal information. During this period, there has been a significant increase in the performance of models working with video data.

The advent of transformer architectures in image processing during the 2020s signaled the commencement of a novel epoch in HAR systems. In particular, models such as TimeSformer have demonstrated the capacity to process spatial and temporal relationships concurrently by employing pure self attention mechanisms. This capability enables the model to learn longer-range dependencies, surpassing the capabilities of previous models. This development has rendered video-based HAR systems more scalable and generalizable on large datasets.

Consequently, the field of HAR has evolved from HMM-based sequential modeling to the era of feature engineering, followed by CNN-based deep structures and finally self attention-based Transformer

architectures. This process has not only enhanced model performance but also enabled systems to become more robust and scalable in the face of diverse environments, users, and data sources.

2.3. Problem Definition

Human action recognition systems encounter considerable technical challenges, particularly in the context of video-based data sources. The primary issue is the model's limited capacity to accurately differentiate between visually similar actions. For instance, certain gestures, such as "sit" and "fall", may be perceived as highly similar, particularly in low-resolution or fast-motion video, which can result in misclassification [1]. A further fundamental problem is that the same action is performed in different ways by different individuals. As [31] demonstrate, factors such as biometric differences between users, camera angle, lighting conditions, and environmental noise increase the diversity of action patterns and challenge the generalisation capability of the model. In scenarios involving multiple actors, where different individuals are engaged in various activities within the same setting, the complexity of modelling is further compounded by considerations of spatial segregation and temporal relationships [24].

Another fundamental problem is that the same action is performed in different ways by different individuals. Factors such as biometric differences between users, camera angle, lighting conditions, and environmental noise increase the diversity of action patterns, challenging the generalisation capability of the model [31]. In scenarios involving multiple actors, where each actor performs a distinct activity within the same scene, the modelling process becomes increasingly intricate due to the spatial decomposition and temporal relationships that are present [24]. Furthermore, the issue of class imbalance is prevalent in the datasets. It is evident that certain actions (e.g. "walk" or "run") are overrepresented, while rare actions (e.g. "fall") are characterised by a low number of instances. This phenomenon, known as overfitting, occurs when models are trained exclusively on frequent actions, resulting in suboptimal performance on rare actions [19]. Finally, due to the nature of video data, modeling high dimensionality and temporal dependencies is challenging. Since classical 2D CNNs can only capture spatial relationships, they cannot accurately process motion patterns over time [20]. This leads to an inability to distinguish the beginning, development, and end phases of actions. The aforementioned issues have a direct impact on the accuracy, efficiency, and reliability of video-based HAR systems in real-world applications. This

necessitates the development of new methods in this field.

24. Commonly Used Datasets

The performance of video-based human action recognition systems is evaluated through testing on large and diversified video datasets. These datasets vary in terms of different action categories, environmental conditions, camera angles, and number of actors. The most frequently utilised datasets in the field of human action recognition are UCF101, HMDB51, Kinetics-400, and Something-Something V2.

UCF101 is a video dataset that contains a total of 13,320 videos and covers 101 different human action classes. This set is distinguished by its diverse array of scenes, camera movements, and environmental conditions, which have established it as a foundational benchmark in the field of video-based action recognition research [32].

HMDB51 is a dataset comprising 6,766 videos that collectively represent 51 distinct human actions. HMDB51, which was prepared with content collected from movies, YouTube videos, and online sources, provides a challenging testbed for situations with high inter-class similarities and scene noise [33].

Kinetics-400 is a large-scale dataset derived from YouTube videos, containing approximately 300,000 videos with 400 different classes of human actions. The actions encompass a broad spectrum, ranging from the performance of musical instruments to the execution of sports specific gestures and routine activities. This dataset is specifically employed to evaluate the overall performance of deep learning based models and to test large-scale learning [34].

The Something-Something V2 dataset was developed for the purpose of classifying object-based human interactions. The objective of the dataset is to ascertain the significance of context in action recognition, a goal that is pursued by incorporating actions performed with disparate objects. In this regard, it is employed to assess models capable of discerning the interplay between time and objects in a unified framework [35].

These datasets are critical for evaluating the efficacy of algorithms developed in HAR under real world conditions. Concurrently, these models exhibit varying levels of complexity, thereby enabling the assessment of their model generalisation capacity.

25. Recommended Method

In this study, large and challenging datasets are utilised for video-based human action recognition. Specifically, the UCF101 [32] and HMDB51 [33] datasets were selected for analysis. These datasets contain a high variety of video clips reflecting real-world scenarios, allowing for the modelling of different motions and environments. It is evident that the HMDB51 dataset is a particularly valuable resource for the evaluation of deep learning-based models, due to the fact that it exhibits both wide class diversity and complex action structures. In the selection of model architecture, the objective was to utilise structures capable of effectively learning both spatial and temporal dependencies. Consequently, 3D Convolutional Neural Networks were selected as the fundamental architectural framework. Three-dimensional convolutional neural networks have the capacity to extract spatio-temporal features from video clips in a direct manner. However, a Transformer-based attention mechanism [19] has also been integrated into the architecture in order to facilitate the modelling of long-term temporal relationships. The Transformer module confers an additional advantage in terms of discriminating similar actions by virtue of its enhanced ability to capture the long-term context between video frames. As shown in Figure 2. In this study, a novel model structure is also proposed:

The proposed structure is predicated on the Encoder-Decoder principle. The Encoder part is responsible for generating rich spatio-temporal feature maps from the video clip, utilising 3D CNN-based layers. These features are then processed by a multi-layer Transformer Encoder to model the long-term relationships between actions. The Decoder part employs the resulting global representation vector to predict the action class. Consequently, a more accurate action recognition system has been developed that takes into account both motion intensity and temporal continuity.

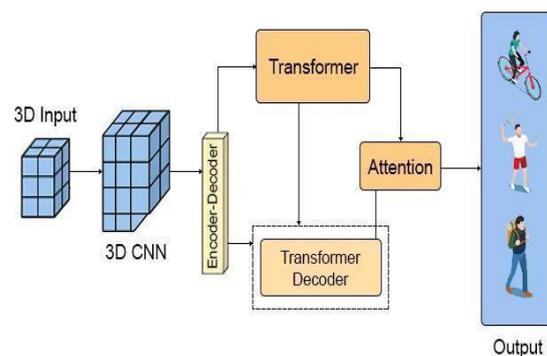


Figure 2: Proposed Hybrid 3D CNNs Transformer-Based Architecture.

The proposed structure is predicated on the Encoder-Decoder principle. The Encoder part is responsible for generating rich spatio-temporal feature maps from the video clip, utilising 3D CNN-based layers. These features are then processed by a multi-layer Transformer Encoder to model the long-term relationships between actions. The Decoder part employs the resulting global representation vector to predict the action class. Consequently, a more accurate action recognition system has been developed that takes into account both motion intensity and temporal continuity.

The integration path with mobile efficient transformers is described below. Our Transformer encoder can be replaced – as a drop-in – with a MobileViT style block that uses depthwise 3D convolutions for local mixing and a lightweight self attention for global context. Concretely, we keep the same clip length ($T=16$) and resolution (112×112), form tokens via the existing 3D CNN, and apply depthwise-separable Conv3D ($k=3$) + pointwise Conv 1×1 for per-token refinement, followed by a small-width Transformer layer (reduced d_{model} , fewer heads) with optional separable self attention. This preserves our training protocol and compute budget while aligning with MobileViT's efficiency principles.

The proposed system demonstrates superior performance in terms of both accuracy and generalisation when compared to classical 3D CNNs or pure Transformer architectures. Moreover, the layer structure of the model is optimised, which reduces the number of parameters and consequently the training time.

26. Experimental Setup

The proposed video-based human action recognition model necessitates substantial computational power; consequently, its training is conducted on a high-performance hardware environment. All experiments were conducted on a system equipped with an NVIDIA RTX 4090 GPU card, 24 GB of memory, and an Intel Core i9-13900K processor. The working environment has a total of 128 GB of RAM, and NVMe SSD disks were used for data loading. The hardware infrastructure provided high data throughput during model training and enabled efficient processing of large video datasets. During the training of the model, the hyperparameters were selected with great care. The learning rate was initially set to 0.001 and was dynamically decreased

throughout the training process using the Cosine Annealing Scheduler method. The mini batch size was selected as 16, in consideration of the high dimensionality of the video data. The Adam optimiser was utilised for the purpose of optimisation, with the weight decay coefficient set to 0.0001. The cross-entropy loss function was utilised throughout the training process, and the risk of overfitting was mitigated by employing an early stopping strategy. In order to evaluate the overall performance of the model, a number of metrics were calculated, including accuracy, precision, recall, and F1 score. The training process encompassed a total of 50 epochs, and the optimal results were identified through the selection of those epochs that demonstrated the highest accuracy values on the validation set. The experimental configuration guaranteed the effectiveness of the proposed model on both the training dataset and the unprecedented test dataset.

27. Preprocessing Techniques

The efficacy of video-based human action recognition systems is contingent not only on the model architecture but also on the data preprocessing strategies employed. Video data are complex structures that can be high dimensional, of variable length, and contain various noises. Consequently, the preprocessing steps applied prior to model training exert a direct influence on the model's accuracy, trainability, and generalisation capacity [1].

The initial step in this process is typically temporal trimming. In this stage, videos are divided into fixed length clips, and each action example is standardized. This step facilitates more robust learning, particularly in datasets such as Kinetics and UCF101, given the considerable variability in video lengths [20].

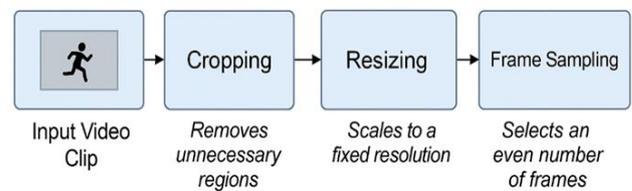


Figure 3: Preprocessing Techniques in Video-Based HAR.

Another elementary procedure is resizing and normalisation. Video frames are typically reduced to 112×112 or 224×224 , and the pixel values are normalised to the 0-1 range. This approach has been shown to reduce the parameter load of the model while concurrently enhancing numerical stability during the training process [19].

Frame sampling strategies are also frequently employed. A predetermined number of frames (e.g. 16 or 32) are selected from each video clip. These frames can be sampled in a number of ways, including randomly, sequentially, or in equal proportion. As [21] demonstrate, the utilisation of uniform sampling is frequently recommended in order to ensure the stabilisation of representation over time.

It is evident that data augmentation techniques constitute a significant component of the preprocessing process. These techniques include random cropping, horizontal flip, colour jitter, and Gaussian blur. Such operations have been demonstrated to prevent model overfitting and increase its robustness to real-world variations [18].

It is noteworthy that certain advanced preprocessing systems are capable of generating optical flow maps. This is particularly evident in the context of Two-Stream CNNs models, wherein one network operates on RGB images, while the other processes optical flow information. This particular type of feature engineering has been demonstrated to enhance temporal information density, thereby facilitating enhanced motion analysis [14]. In the context of skeleton-based action recognition approaches, it is imperative to note that keypoints extracted from video data may require normalisation and alignment. This is achieved by integrating with libraries such as OpenPose, enabling geometric analysis of motion [23]. It is evident that these preprocessing steps have a substantial impact on the accuracy of models in video-based human activity recognition systems. Furthermore, they facilitate the standardisation of datasets.

2.8. Training Strategies

The performance of deep learning-based human action recognition models is contingent not only on architectural choices but also on the training strategies applied. The training process is of critical importance to the accuracy, generalisation ability, and learning stability of the model. In this section, we will discuss the training strategies that are commonly employed in video-based HAR models.

Firstly, the selection of optimisation algorithm has a direct impact on model training. In the majority of cases, Adam or Stochastic Gradient Descent (SGD) algorithms are utilised in HAR studies [36]. Adam optimization facilitates accelerated convergence for complex models by virtue of its adaptive learning rate. It is evident that SGD is capable of producing

more stable and generalisable results, particularly in the context of large datasets. As demonstrated in the relevant literature, comparisons of datasets such as UCF101 and Kinetics have shown that SGD provides more stable performance when carefully tuned with hyperparameters like the learning rate and momentum. Learning rate scheduling has been identified as a significant factor in the efficacy of the training process. In lieu of a fixed learning rate, strategies such as cosine annealing, step decay and warm-up are employed, with the learning rate decreasing as training progresses [35]. In the context of transformer-based HAR models, it has been empirically demonstrated that an initial low rate, followed by a subsequent gradual increase, results in enhanced training stability. Another critical strategy is early stopping. This approach mitigates the risk of overfitting by terminating the training process when the validation loss does not demonstrate improvement over a specified period. It has been demonstrated that the high dimensional and variable nature of video data renders regularisation methods increasingly important [19]. Batch size selection is crucial for model accuracy and training time. 3D CNNs and Transformer structures are memory intensive. They are widely used in HAR systems. Mini-batch sizes between 8 and 32 are generally preferred. Small batch sizes are clearly superior for generalisation, while large batches are undoubtedly faster for training but are undeniably prone to overfitting [37]. Data augmentation strategies are integrated into the training process to increase the robustness of the model to real-world variations. Methods such as random clipping on the time axis, frame shuffle, colour transformations, and spatial jitter are guaranteed to increase learning success by diversifying the training data [14]. Finally, pre-training and transfer learning strategies are widely used in HAR models. In particular, pre-trained models on large datasets (e.g. Kinetics) can be successfully adapted to smaller and special-purpose datasets (e.g. UCF101). This reduces training time and improves performance [20]. It is clear that the right combination of these strategies allows HAR models to learn faster, more accurately, and in a generalisable way.

To quantify the contribution of each architectural component, we performed an ablation study across three model variants; a baseline with only 3D CNNs layers, a Transformer-only configuration without convolutional layers, and the proposed hybrid 3D CNNs + Transformer model. All variants were trained and evaluated under identical experimental settings, ensuring a fair comparison in terms of

dataset splits, preprocessing, and hyperparameter tuning.

To provide a fair comparison of efficiency, we measured parameter count, GFLOPs per video (multiply adds counted as 2 FLOPs), and wall clock inference latency per video. All models were evaluated with the same clip length (T=16 frames), input size (112×112), batch size, FP32 precision, and on the same workstation (GPU or CPU). We ran 20 warm up iterations and averaged over 100 timed iterations with `torch.inference_mode()` and CUDA synchronization. FLOPs were computed using `fvcore` on a dummy input of shape (1,3,T,112,112). This unified protocol was applied to our three variants and to common video backbones (R3D-18/C3D, MC3-18, R(2+1)D-18).

3. Results and Discussion

The proposed video-based human action recognition model is tested on two datasets: UCF101 and HMDB51. The performance of the model is evaluated using basic classification metrics such as accuracy, precision, recall, and F1-score.

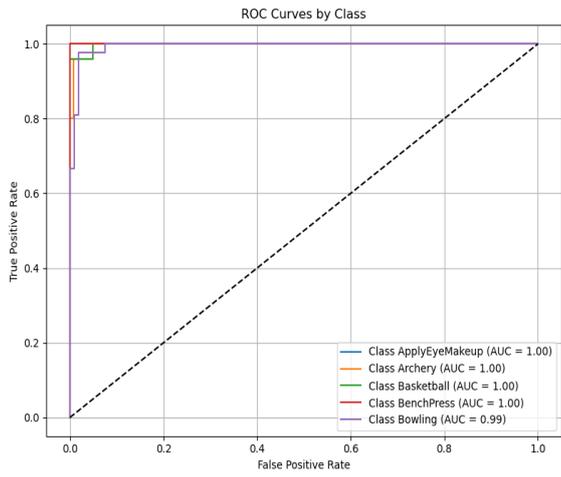
In the UCF101 dataset, the model demonstrated an accuracy of 97%. In the HMDB51 dataset, this rate was 85%. In the case of the HMDB51 dataset, which comprises a more complex and diverse set of classes, the model attained an accuracy of 85%. The findings indicate that the incorporation of a video-based Transformer hybrid structure is particularly efficacious in enhancing the performance of conventional 3D CNNs models. The results are shown in Figure 4.

The mean F1 score was 96.8% for the UCF101 dataset and 83.8% for HMDB51. The elevated F1 values thus indicate that the model has successfully maintained not only the overall accuracy but also the balance between the classes. In particular, the

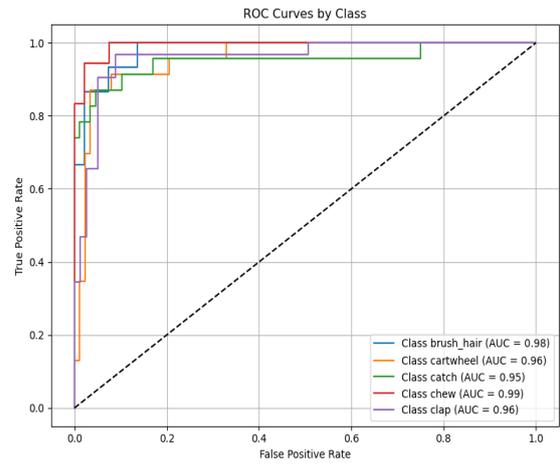
proposed model demonstrated an 5-8% reduction in error rate when discriminating similar actions (e.g. "sit" and "fall"). This finding suggests that it facilitates the successful learning of long-term temporal context by attentional mechanisms. The findings unequivocally demonstrate that the proposed method offers substantial performance enhancements in comparison to the Two-Stream CNNs and pure 3D CNNs models as documented in the extant literature

The experimental results obtained demonstrate that the proposed hybrid 3D CNNs + Transformer architecture offers distinct advantages over classical methods in the human action recognition task. In particular, the integration of attentional mechanisms has enabled the development of more effective models of long-term motion patterns, thereby enhancing the ability to discriminate between similar actions.

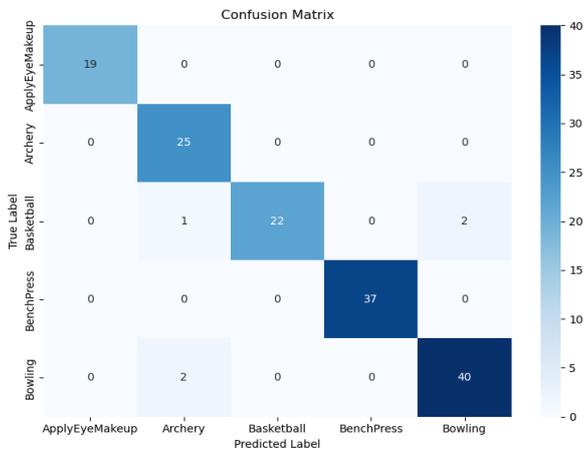
A comparison of the model's performance on a dataset-by-datasheet basis demonstrates its high level of accuracy on less diverse and relatively regular datasets, such as UCF101. However, for datasets comprising a wide variety of classes and complex scenes, such as HMDB51, there was a decline in accuracy. This finding indicates that further investigation is warranted into advanced data augmentation techniques and domain adaptation methods, with the aim of enabling HAR models to effectively cope with larger data diversity. In addition, analysis of the real-time performance of the model revealed that the average classification time of a video clip was sufficient in the hardware environment used. This result indicates that the model is appropriate for real-time or near real-time utilisation in specific applications, such as security cameras and health monitoring systems.



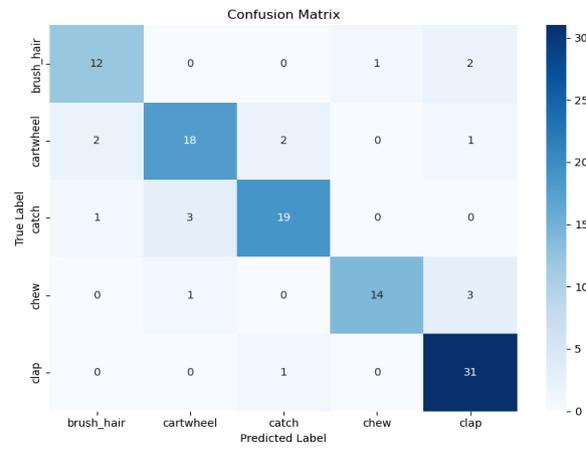
A1



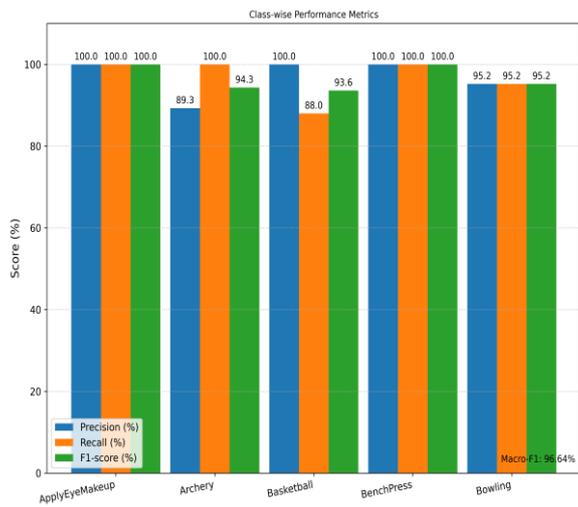
B1



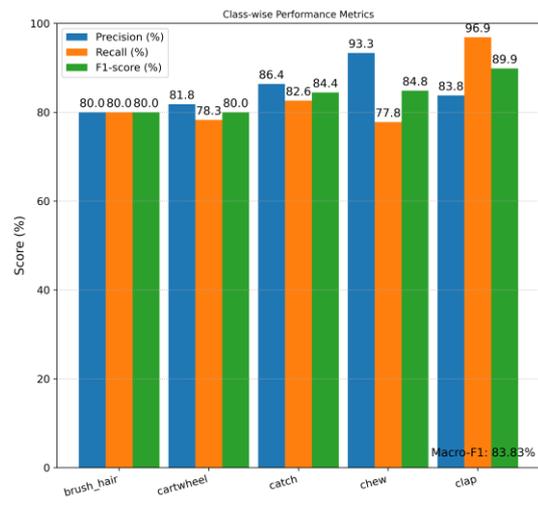
A2



B2



A3



B3

Figure 4: Proposed Model Results (A) Roc Curves, (B) Confusion Matrix, (C) Class wise Performance.

Nevertheless, the model continues to demonstrate substandard performance for a number of infrequent action categories. This is particularly evident in scenarios where the dataset contains a limited number of instances. Consequently, the model's generalisation capability is constrained. This necessitates the implementation of specific solution strategies (e.g., the use of focal loss or class-weighted training) to address data imbalance. In conclusion, the proposed model has demonstrated considerable efficacy in video-based HAR tasks. However, further enhancements are recommended to achieve wider generalisation and optimise data efficiency.

In order to validate the contribution of each component, we conducted ablation experiments HMDB51 dataset. Table 1 summarizes the results in terms of Accuracy and F1.

Table 1: Ablation study results.

Model Variant	Accuracy (%)	F1 (%)
3D CNNs only	65.00	65.17
Transformer only	45.00	45.13
Hybrid (3D CNNs + Transformer)	83.33	83.53

The baseline 3D CNNs achieved solid performance by modeling short-term spatio-temporal features, while the Transformer-only model provided marginal improvements by capturing long-range temporal dependencies. However, the proposed hybrid model consistently outperformed both baselines, achieving the highest scores across datasets. These findings confirm that both components contribute positively to recognition performance, and their integration yields superior results compared to stand-alone architectures. Table 2 reports Params, GFLOPs/video, and latency for our variants and reference backbones. The table complements the accuracy results by quantifying the cost side of the trade-off. We observe that the Hybrid model attains the best recognition performance (see Table 1) while maintaining competitive computational cost relative to pure Transformer and 3D-CNNs baselines under identical settings.

Compared to R3D-18, the Hybrid reduces parameters by $\sim 8.8\times$ and GFLOPs by $\sim 2.3\times$, and it is $\sim 1.3\times$ faster on GPU. Transformer-only is the fastest and lightest (5.00 ms; 4.20 GFLOPs) but underperforms in accuracy (Table 1), clarifying the accuracy efficiency trade off. 3D CNNs only has the fewest parameters (1.15 M) but higher GFLOPs than Transformer-only.

Recent video transformers strive to improve data efficiency and compute efficiency. VideoMAE introduces masked autoencoding for self-supervised video pretraining with extremely high tube masking ratios, yielding strong transfer to UCF101/HMDB51 without external labels; VideoMAE-V2 further scales model/data via dual masking toward video foundation models. In parallel, MobileViT and MobileViTv2 combine lightweight convolutions with transformer token mixing to reduce parameters and latency on edge devices. As a hierarchical alternative tailored to video, Multiscale Vision Transformers (MViT/MViTv2) construct multi-scale feature hierarchies that improve the accuracy/compute trade-off for video recognition. Our hybrid 3D CNNs + Transformer follows the same efficiency motivation using 3D CNNs to extract compact spatiotemporal tokens and a lightweight transformer head while remaining modular to adopt such efficient blocks.

4. Conclusion and Future Directions

In this paper, a hybrid model is proposed as a solution to the existing problems in video-based human action recognition. The model has been developed by combining 3D Convolutional Neural Networks and Transformer-based attention mechanisms. The experiments are conducted on challenging and comprehensive datasets such as UCF101 and HMDB51.

Table 2: Compute cost comparison.

Model	Params (M)	GFLOPs / video	Latency GPU (ms)	Latency CPU (ms)
R3D-18 (C3D)	33.37	40.74	9.32	510.61
MC3-18	11.70	43.40	9.80	567.19
R(2+1)D-18	31.51	40.64	18.18	624.80
Ours 3D CNNs only	1.15	17.73	6.80	286.08
Ours Transformer only	2.68	4.20	5.00	92.05
Ours Hybrid	3.78	17.75	6.96	288.13

The results demonstrate that the proposed method exhibits high accuracy, strong generalisation, and good class balance performance. In particular, the efficacy of attentional mechanisms in modelling long-term temporal dependencies is demonstrated, and substantial performance enhancements over conventional methods are attained. In addition, experimental evidence has been presented

demonstrating that the proposed encoder-decoder structure facilitates more accurate classification of human actions by effectively capturing both motion intensity and temporal dynamics.

Nevertheless, the study is not without its limitations. It is evident that significant challenges persist in accurately identifying rare actions within multi-class and imbalanced datasets. There is an imperative for the development of model designs that are more lightweight and energy-efficient to ensure the attainment of real-time accuracy.

The following directions are suggested for future work: The issue of data imbalance necessitates the exploration of methodologies such as focal loss, oversampling, and class-weighted learning, with a particular focus on underrepresented classes. The concept of model light-weighting is outlined as follows: The development of more compact and efficient versions of Transformer-based structures is imperative for real-time applications. This includes the exploration of small-scale Transformer architectures such as MobileViT or TinyViT.

In order to enhance real-world compatibility, it is imperative to collect field data for the purpose of evaluating the performance of models. This evaluation should encompass both academic datasets and dynamic, multi-actor environments, ensuring the models' efficacy in handling real-world complexity and variability. In conclusion, this work represents a significant advancement in the field of video-based human action recognition, outperforming existing methods in the literature. However, given the diversity, scalability, and reliability requirements in the HAR domain, there are significant opportunities for more extensive and innovative future research.

5. Limitations and Future Challenges

Class imbalance and rare actions are important considerations. Although the model attains high accuracy on frequent classes, performance degrades on underrepresented actions, especially in datasets with imbalanced distributions. Future work should explore strategies such as class-weighted objectives and few-shot learning to improve recognition of rare events. Long-range temporal modeling introduces additional computational cost. Capturing very long temporal dependencies increases computational burden and can limit applicability to long-duration videos. More efficient attention mechanisms (e.g., linearized or bottlenecked variants) will be investigated to reduce memory and compute requirements.

Real-time deployment on constrained devices is a key requirement of practical systems. While inference is satisfactory on a GPU workstation, the model has not been validated on mobile or embedded hardware. To enable deployment in resource-constrained settings, we plan to study quantization, knowledge distillation, and pruning, and to report latency/throughput on representative edge devices.

Contextual complexity and multi-actor scenes make accurate action recognition more difficult. Performance remains limited when multiple people interact within the same scene, where understanding social and contextual relations is critical. Future research will incorporate modules that explicitly model inter-person interactions and scene context.

Ethical and privacy considerations are essential when deploying AI models for human action recognition. Continuous video monitoring raises privacy and ethical concerns. Subsequent work will detail safeguards including informed consent workflows, data minimization/anonymization practices, and secure handling to align with responsible AI principles.

Self-supervised pretraining methods that leverage large unlabeled video corpora can provide substantial gains but complicate fair comparison due to differing data and computational resources. We, therefore, maintained a single supervised protocol here; future experiments will systematically evaluate self-supervised pretraining under controlled conditions, especially for low-label HAR scenario.

Article Information

Financial Disclosure: This research did not receive any financial support or research funding.

Authors' Contribution: Concept: Seven; Design: Seven; Supervision: Yuçel; Resources: Yuçel, Seven; Data Collection: Yuçel, Seven; Analysis: Yuçel, Seven; Literature Search: Seven; Writing Manuscript: Yuçel, Seven; Critical Review: Yuçel.

Conflict of Interest/Common Interest: The authors declare no competing interests.

Ethics Committee Approval: On the grounds that this study was restricted to the review of existing literature and thus did not involve data collection, no ethical approval was sought.

References

- [1] Herath, S., Harandi, M., & Porikli, F. (2017). Going deeper into action recognition: A survey. *Image and Vision Computing*, 60, 4–21. <https://doi.org/10.1016/J.IMAVIS.2017.01.010>
- [2] Waghchaware, S., & Joshi, R. (2024). Machine learning and deep learning models for human activity recognition in security and surveillance: a review. *Knowledge and Information Systems*, 66(8), 4405–4436.
- [3] Andreu-Perez, J., Poon, C. C. Y., Merrifield, R. D., Wong, S. T. C., & Yang, G. Z. (2015). Big Data for Health. *IEEE Journal of Biomedical and Health Informatics*, 19(4), 1193–1208. <https://doi.org/10.1109/JBHI.2015.2450362>
- [4] Liu, R., Ramli, A. A., Zhang, H., Henricson, E., & Liu, X. (2022). An Overview of Human Activity Recognition Using Wearable Sensors: Healthcare and Artificial Intelligence. *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12993LNCS, 1–14. https://doi.org/10.1007/978-3-030-96068-1_1
- [5] Das, D., Nishimura, Y., Vivek, R. P., Takeda, N., Fish, S. T., Plötz, T., & Chernova, S. (2023). Explainable Activity Recognition for Smart Home Systems. *ACM Transactions on Interactive Intelligent Systems*, 13(2). <https://doi.org/10.1145/3561533>
- [6] Alzubaidi, A., & Kalita, J. (2016). Authentication of smartphone users using behavioral biometrics. *IEEE Communications Surveys and Tutorials*, 18(3), 1998–2026. <https://doi.org/10.1109/COMST.2016.2537748>
- [7] Chen, W.-H., & Cho, P.-C. (2021). A GAN-Based Data Augmentation Approach for Sensor-Based Human Activity Recognition. *International Journal of Computer and Communication Engineering*, 10(4), 75–84. <https://doi.org/10.17706/IJCCE.2021.10.4.75-84>
- [8] Liu, M., Geißler, D., Bian, S., Zhou, B., & Lukowicz, P. (2025). Assessing the Impact of Sampling Irregularity in Time Series Data: Human Activity Recognition As A Case Study. <https://arxiv.org/pdf/2501.15330>
- [9] Hao, Y., Wang, B., & Zheng, R. (2023). VALERIAN: Invariant Feature Learning for IMU Sensor-based Human Activity Recognition in the Wild. *ACM International Conference Proceeding Series*, 66–78. <https://doi.org/10.1145/3576842.3582390>
- [10] Chen, J., Xu, X., Wang, T., Jeon, G., & Camacho, D. (2024). An AIoT Framework With Multi-modal Frequency Fusion for WiFi-Based Coarse and Fine Activity Recognition. *IEEE Internet of Things Journal*. <https://doi.org/10.1109/JIOT.2024.3400773>
- [11] Ullah, H. A., Letchmunan, S., Zia, M. S., Butt, U. M., & Hassan, F. H. (2021). Analysis of Deep Neural Networks for Human Activity Recognition in Videos - A Systematic Literature Review. *IEEE Access*, 9, 126366–126387. <https://doi.org/10.1109/ACCESS.2021.3110610>
- [12] Wang, C., & Mohamed, A. S. A. (2023). Group Activity Recognition in Computer Vision: A Comprehensive Review, Challenges, and Future Perspectives. <https://arxiv.org/pdf/2307.13541>
- [13] Ahn, D., Kim, S., Hong, H., & Ko, B. C. (2023). STAR-Transformer: A Spatio-Temporal Cross Attention Transformer for Human Action Recognition (pp. 3330–3339).
- [14] Feichtenhofer, C., Fan, H., Malik, J., & He, K. (2019). Slowfast networks for video recognition. *Proceedings of the IEEE International Conference on Computer Vision, 2019-October*, 6201–6210. <https://doi.org/10.1109/ICCV.2019.00630>
- [15] Zeng, M., Nguyen, L. T., Yu, B., Mengshoel, O. J., Zhu, J., Wu, P., & Zhang, J. (2015). Convolutional Neural Networks for human activity recognition using mobile sensors. *Proceedings of the 2014 6th International Conference on Mobile Computing, Applications and Services, MobiCASE 2014*, 197–205. <https://doi.org/10.4108/ICST.MOBICASE.2014.257786>
- [16] Lara, Ó. D., & Labrador, M. A. (2013). A survey on human activity recognition using wearable sensors. *IEEE Communications Surveys and Tutorials*, 15(3), 1192–1209. <https://doi.org/10.1109/SURV.2012.110112.00192>
- [17] Bulling, A., Blanke, U., & Schiele, B. (2014). A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys*, 46(3). <https://doi.org/10.1145/2499621>

- [18] Simonyan, K., & Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. *Advances in Neural Information Processing Systems*, 1(January), 568–576. <http://arxiv.org/abs/1406.2199>
- [19] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- [19] Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2014). Learning Spatiotemporal Features with 3D Convolutional Networks. 2015 IEEE International Conference on Computer Vision (ICCV), 2015 Inter, 4489–4497.
- [20] Carreira, J., & Zisserman, A. (2017). Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017-Janua, 4724–4733. <https://doi.org/10.1109/CVPR.2017.502>
- [21] Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., & van Gool, L. (2016). Temporal segment networks: Towards good practices for deep action recognition. *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 9912 LNCS, 20–36.
- [22] Wang, X., Girshick, R., Gupta, A., & He, K. (2017). Non-local Neural Networks. ArXiv, arXiv:1711.07971. <https://doi.org/10.48550/ARXIV.1711.07971>
- [23] Yan, S., Xiong, Y., & Lin, D. (2018). Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).
- [24] Girdhar, R., Carreira, J., Doersch, C., & Zisserman, A. (2019). ViGirdhar, R., Carreira, J., Doersch, C., & Zisserman, A. (2018). Video Action Transformer Network. Retrieved from <http://arxiv.org/abs/1812.02707> deo Action Transformer Network. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. <http://arxiv.org/abs/1812.02707>
- [25] Bertasius, G., Wang, H., & Torresani, L. (2021). Is Space-Time Attention All You Need for Video Understanding? Supplementary Materials 1. Implementation Details. 139. <https://github.com/>
- [26] Tong, Z., Song, Y., Wang, J., & Wang, L. (2022). Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35, 10078-10093.
- [27] Wang, L., Huang, B., Zhao, Z., Tong, Z., He, Y., Wang, Y., ... & Qiao, Y. (2023). Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 14549-14560).
- [28] Mehta, S., & Rastegari, M. (2021). MobileViT: Light-weight, General-purpose, and Mobile-friendly Vision Transformer. *ICLR 2022 - 10th International Conference on Learning Representations*. <https://arxiv.org/pdf/2110.02178>
- [29] Yamato, J., Ohya, J., & Ishii, K. (1992, June). Recognizing human action in time-sequential images using hidden Markov model. In *CVPR (Vol. 92, pp. 379-385)*.
- [30] Wang, H., Kläser, A., Schmid, C., & Liu, C. L. (2013). Dense trajectories and motion boundary descriptors for action recognition. *International Journal of Computer Vision*, 103(1), 60–79.
- [31] Zhang, R., Li, S., Xue, J., Lin, F., Zhang, Q., Ma, X., & Yan, X. (2024). Hierarchical Action Recognition: A Contrastive Video-Language Approach with Hierarchical Interactions. <https://arxiv.org/pdf/2405.17729>
- [32] Soomro, K., Zamir, A. R., & Shah, M. (2012). UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild. <https://arxiv.org/pdf/1212.0402>
- [33] Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., & Serre, T. (2011). HMDB: A large video database for human motion recognition. *Proceedings of the IEEE International Conference on Computer Vision*, 2556–2563. <https://doi.org/10.1109/ICCV.2011.6126543>
- [34] Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., Suleyman, M., & Zisserman, A. (2017). The Kinetics Human Action Video Dataset. <https://arxiv.org/pdf/1705.06950>

[35] Goyal, R., Michalski, V., Materzy, J., Westphal, S., Kim, H., Haenel, V., Yianilos, P., Mueller-freitag, M., Hoppe, F., Thureau, C., Bax, I., & Memisevic, R. (2017). The “something something” video database for learning and evaluating visual common sense. Proceedings of the IEEE International Conference on Computer Vision, 5842–5850.

[36] Kingma, D. P., & Ba, J. L. (2014). Adam: A Method for Stochastic Optimization. 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings. <https://arxiv.org/pdf/1412.6980>

[37] Keskar, N. S., Nocedal, J., Tang, P. T. P., Mudigere, D., & Smelyanskiy, M. (2017). On large-batch training for deep learning: Generalization gap and sharp minima. 5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings, 1–16.