

Real-Time Human Action Recognition in Video Surveillance using Machine Learning

Md. Musabbir Hossain
Department of CSE
GSTU, Bangladesh

Md. Tachbir Dewan
Department of CSE
GSTU, Bangladesh

Md. Sadikuzzaman
Department of CSE
GSTU, Bangladesh

Abu Bakar M.
Abdullah
Department of CSE
GSTU, Bangladesh

ABSTRACT

This paper presents an innovative framework for real-time human action recognition in video surveillance systems, aimed at delivering immediate detection of suspicious behavior, normal movements, and actionable insights for security operators. The proposed method integrates computer vision and machine learning techniques to improve recognition accuracy and system reliability. Motion analysis is performed using optical flow, where Optical Flow Energy Images (OFEI) are generated to extract motion-related features. A Convolutional Neural Network (CNN) is utilized to obtain high-dimensional feature representations while reducing dimensionality, and a Support Vector Machine (SVM) classifier is trained on these features for robust action classification.

The system effectively detects and distinguishes human actions such as walking, looking around, looking up, smashing, and suspicious activities, even under challenging conditions including camera motion, zoom-in, and zoom-out. Experimental evaluations conducted on publicly available human action datasets demonstrate significant improvements in recognition accuracy. Additionally, the system overlays detected actions onto video streams, providing clear and actionable visual feedback to surveillance personnel. Successfully deployed in intelligent video surveillance environments, the proposed framework proves to be scalable, accurate, and effective for identifying abnormal behaviors and generating timely alerts in modern security applications.

Keywords

Optical flow, machine learning classifier, deep learning, feature extraction, action recognition, video surveillance.

1. INTRODUCTION

With rapid advancements in technology and increasing urbanization, ensuring public safety has become more critical than ever. Surveillance systems play a vital role in monitoring public spaces, but traditional systems largely depend on manual observation. Human operators are often required to monitor multiple video feeds over long periods, making the process prone to fatigue, distraction, and errors. These limitations significantly hinder the effectiveness of manual surveillance in dynamic and high-stress environments.

To address these issues, recent research has focused on integrating machine learning (ML) and artificial intelligence (AI) techniques into video surveillance. Such systems can automatically analyze video streams, recognize human actions, and raise alerts in real time. Optical flow methods are employed to capture motion dynamics [1], while convolutional neural networks (CNNs) extract deep features from video data [2]. Classification algorithms, such as support vector machines (SVMs), then categorize these features into defined action

classes [3]. The combination of computer vision and machine learning enables the development of intelligent surveillance systems that are more accurate, scalable, and responsive.

The limitations of current systems are not confined to labor inefficiency. With the exponential growth of surveillance cameras and data, it becomes increasingly difficult to monitor every feed manually or even with traditional automated systems. Furthermore, existing systems often struggle to identify complex human behaviors, especially in environments with varying lighting conditions, occlusions, or camera motion. They tend to rely on a narrow set of predefined actions, lacking adaptability to emerging or unusual behaviors. This inflexibility poses significant risks in real-world scenarios such as public safety, emergency response, and threat detection.

In response, this research proposes a robust and scalable framework for real-time human action recognition using deep learning. The core of the system is a CNN-SVM hybrid model that processes motion features derived from optical flow and classifies actions with high accuracy. The system is designed to function efficiently under diverse environmental conditions and supports deployment on resource-constrained devices, enabling real-time performance on the edge. It also incorporates data augmentation and domain adaptation to ensure adaptability across different video sources and scenarios.

The primary goal of this research is to bridge the gap between the theoretical developments in machine learning and their real-world implementation in video surveillance. By emphasizing real-time processing, accuracy, and computational efficiency, the proposed framework aims to enhance situational awareness in surveillance operations. Although the primary focus is on security applications, the proposed system has potential use cases in healthcare monitoring, sports analytics, and smart environments, making it a versatile solution for broader human activity analysis.

2. RELATED WORK

The field of human action recognition (HAR) in video surveillance has gained significant research interest due to its potential in enhancing public safety, security monitoring, and behavioral analysis. Initial approaches primarily relied on handcrafted features. A notable early contribution is by Wang et al. (2013), who proposed the use of dense trajectories to capture motion patterns by tracking interest points across video frames. While effective to some extent, these handcrafted approaches lacked robustness and generalization across diverse environments and complex action variations.

The emergence of deep learning marked a substantial turning point in HAR. Simonyan and Zisserman (2014) introduced two-stream convolutional neural networks (CNNs), which process

spatial information from video frames and temporal motion data from optical flow in parallel, significantly improving recognition accuracy [4,5]. Ji et al. (2013) expanded this by developing 3D CNNs that model spatiotemporal features jointly, allowing better understanding of motion dynamics in videos [6]. Building upon these foundations, Tran et al. (2015) proposed hierarchical 3D CNN architectures that effectively captured both low-level and high-level action features, enhancing recognition performance in complex scenes [7].

Further progress in the field has been achieved by integrating deep learning models with classical machine learning classifiers. Hybrid frameworks, such as CNN-based feature extraction followed by Support Vector Machine (SVM) classification, have shown notable success [8,9]. SVMs, known for their robustness in high-dimensional spaces, offer interpretable and reliable decision-making, particularly when dealing with small or imbalanced datasets. These combinations have improved system reliability without significantly increasing computational complexity.

Despite these advancements, human action recognition in real-world surveillance scenarios still faces major challenges. Issues such as partial occlusions, camera motion, background clutter, illumination changes, and viewpoint variability continue to affect recognition accuracy. Moreover, achieving real-time performance remains difficult, particularly for computationally intensive deep learning models.

Recent research has explored more advanced architectures and learning techniques to address these limitations. Pose-based models utilize skeletal data to identify actions based on joint movements, providing resilience against visual noise and occlusion [10]. Transformer-based networks, which model long-range temporal dependencies, have shown promise in capturing the global structure of actions across extended video sequences [11]. Additionally, there has been a growing interest in multi-stream and multimodal learning, where data from RGB, depth, infrared, and even audio sources are fused to enrich action representation [12]. To support deployment in practical environments, researchers have also begun exploring lightweight and optimized models suitable for real-time inference on edge devices [13]. These collective efforts highlight the rapid progress in the field of action recognition, as well as the continued need for systems that are both accurate and efficient under real-world constraints.

3. RESEARCH METHODOLOGY

This section outlines the methodology used to develop the proposed real-time human action recognition system. The approach integrates motion analysis, deep feature extraction, and robust classification techniques to accurately identify human actions in surveillance videos. By combining computer vision and machine learning, the system is designed to operate efficiently in real-time and adapt to diverse environmental conditions.

3.1 System Architecture

The proposed system for real-time human action recognition comprises four core components: optical flow detection, feature extraction, classification, and real-time output visualization. Optical flow detection plays a foundational role in capturing motion dynamics between consecutive video frames. By analyzing pixel-level changes using Farneback's dense optical flow algorithm, the system generates flow maps that represent the direction and magnitude of motion. These maps emphasize moving regions such as limbs or body shifts, effectively isolating

action-related patterns from background noise.

To prepare the data for machine learning, the optical flow frames are resized and normalized. This step standardizes input dimensions and scales intensity values, ensuring compatibility with the CNN architecture and enhancing training efficiency. The processed optical flow data forms the basis for extracting meaningful features that distinguish between various human actions, such as walking or suspicious activity. This motion preprocessing pipeline enables the system to robustly interpret complex actions in dynamic and cluttered environments, forming the backbone for accurate and real-time action recognition.

Fig. 1 illustrates the process of motion analysis using optical flow. Two consecutive video frames are analyzed to compute motion vectors, which indicate the direction and speed of movement at each pixel. These vectors are visualized using a color-coded map, where hue represents direction and intensity denotes motion magnitude. The resulting flow map highlights dynamic regions such as limb movement, allowing the system to focus on relevant action cues while filtering out static background elements. This preprocessing step enhances recognition accuracy by providing robust motion features essential for distinguishing human actions in complex surveillance environments.

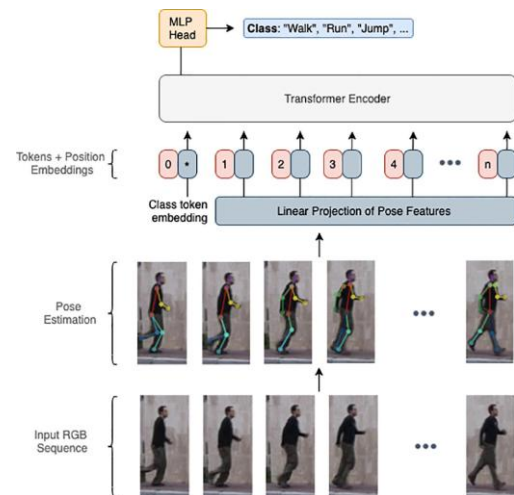


Fig. 1. Optical flow visualization showing motion vectors between consecutive frames (redrawn based on the classical formulation in [22]).

3.2 Extraction of Optical Flow

After computing motion vectors between video frames using optical flow, these vectors are transformed into visual motion maps where hue represents direction and intensity indicates motion magnitude. To enhance feature quality, insignificant motion is filtered out, preserving only meaningful human movement. The optical flow maps are then resized to 64×64 pixels to ensure compatibility with the CNN architecture, maintaining aspect ratio to avoid distortion. Pixel values are normalized to a standard scale to improve model learning consistency. This resizing step is crucial for memory efficiency, model compatibility, and uniform pattern recognition across the dataset. The processed frames serve as input to the CNN, enabling effective extraction of motion-related features essential for distinguishing various human actions in real-time scenarios.

3.3 Training the CNN

Training the Convolutional Neural Network (CNN) is a critical step in enabling the system to accurately recognize human

actions. The CNN is trained on preprocessed optical flow frames, which capture motion between video frames. These frames are resized to 64×64 pixels, normalized, and augmented through techniques.

The custom-designed CNN architecture includes multiple convolutional layers that extract both low-level (e.g., edges, motion boundaries) and high-level (e.g., limb movements) features. Each convolution is followed by ReLU activation and max-pooling to reduce dimensionality and retain essential information. After flattening the feature maps, fully connected layers learn action-specific feature combinations, with the final softmax layer outputting class probabilities.

The network is trained using labeled optical flow data through forward propagation, loss calculation via categorical cross-entropy, and backpropagation with optimizers such as Adam. Performance is monitored across training and validation sets using metrics like accuracy and loss. This training process enables the CNN to learn discriminative motion patterns, ensuring robust and real-time recognition of various human actions across diverse surveillance scenarios.

Table 1 represents the CNN training pipeline, highlighting how raw optical flow data is transformed into actionable features for human action recognition. The model architecture includes two convolutional layers (Conv2D) with 32 and 64 filters respectively, each followed by max-pooling layers that reduce spatial dimensions while preserving key features. These layers extract hierarchical motion patterns crucial for recognizing human actions.

Following feature extraction, the model uses two fully connected (dense) layers. The first dense layer reduces the high-dimensional data to 64 features with 589,888 trainable parameters, enhancing computational efficiency and generalization. The final output layer, with two nodes, performs binary classification using softmax activation. The summary table shown in the figure presents each layer's output shape and parameter count, revealing a total of 646,338 trainable parameters and no non-trainable parameters. With a memory footprint of approximately 2.47 MB, the model balances accuracy and efficiency, making it suitable for real-time surveillance applications.

Table 1. CNN Model Architecture Summary

| Layer (type) | Output Shape | Param # |
|--------------------------------|--------------------|---------|
| conv2d (Conv2D) | (None, 62, 62, 32) | 896 |
| max_pooling2d (MaxPooling2D) | (None, 31, 31, 32) | 0 |
| conv2d_1 (Conv2D) | (None, 29, 29, 64) | 18,496 |
| max_pooling2d_1 (MaxPooling2D) | (None, 14, 14, 64) | 0 |
| conv2d_2 (Conv2D) | (None, 12, 12, 64) | 36,928 |
| flatten (Flatten) | (None, 9261) | 0 |
| dense (Flatten) | (None, 64) | 589,888 |
| dense_1 (Dense) | (None, 2) | 130 |

Total Params: 646,338 (2.47 MB)
Trainable params: 646,338 (2.47 MB)
Non-trainable params: 0 (0.00 B)

3.4 Action Classification Using SVM

After feature extraction via CNN, the final stage involves

classifying human actions using a Support Vector Machine (SVM). SVM is a robust supervised learning algorithm that separates feature vectors into predefined action categories—such as walking, running, or jumping—by finding an optimal hyperplane in high-dimensional space. It is particularly effective with high-dimensional data and generalizes well, even with limited samples.

The process begins with an input video, from which optical flow is computed to capture motion dynamics between frames. This motion data is aggregated into an Optical Flow Energy Image (OFEI), which emphasizes the magnitude and direction of movement over time. The OFEI is then passed through a CNN to extract spatial and motion-related features.

These extracted features serve as input to the SVM classifier. Using kernels like RBF, the SVM handles non-linear relationships and accurately assigns each sample to an action class. This hybrid pipeline—optical flow, CNN-based feature extraction, and SVM classification—provides a lightweight yet powerful framework for human action recognition, suitable for real-time surveillance, healthcare, and human-computer interaction.

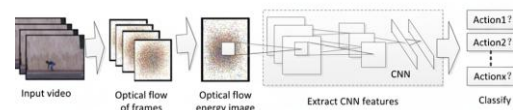


Fig. 2. The diagram is conceptually based on standard deep learning-based action recognition frameworks described in [23]

3.5 SVM Classification

The high-dimensional features extracted by the CNN are flattened and used to train a Support Vector Machine (SVM) classifier with a linear kernel for human action recognition. SVMs are powerful supervised learning models that classify data by identifying an optimal hyperplane that maximizes the margin between classes. The closest data points to this boundary, known as support vectors, are critical in defining the classification decision.

In this study, a soft margin SVM is employed to accommodate real-world complexities such as noise and overlapping classes, providing flexibility by allowing slight misclassifications. The classification pipeline includes data preprocessing, feature extraction, optional dimensionality reduction (e.g., PCA), model training, and evaluation using metrics like accuracy, precision, recall, and F1-score. In the system architecture, the SVM fits into the final stage, receiving optimized features from the CNN and classifying them into predefined action categories. Its integration ensures fast and reliable performance, making it well-suited for real-time surveillance applications.

4. EXPERIMENTAL RESULTS

The proposed system was tested on the designated test set to evaluate its performance in recognizing human actions. The SVM classifier, trained on the extracted features, achieved high classification accuracy across various action classes. The results demonstrate the effectiveness of the preprocessing, motion detection, and feature extraction stages in capturing relevant patterns. The system showed strong generalization to unseen data, indicating its robustness and applicability in real-world scenarios.

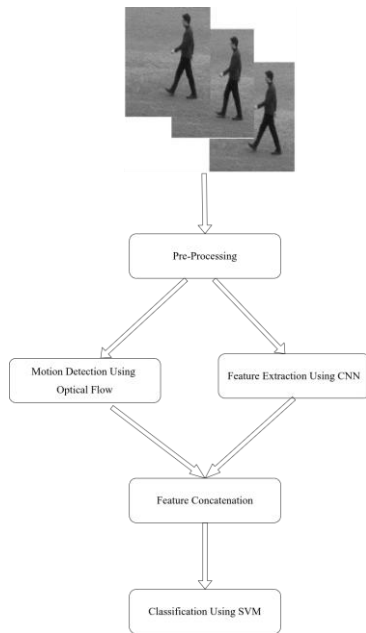


Fig. 3. SVM Classification

4.1 Dataset Description

The proposed system was evaluated on a dataset comprising 500 video sequences of human actions, divided into 80% training and 20% testing sets with no overlap. Each video is labeled with a single action class, including activities such as walking, running, jumping, sitting, standing, and suspicious behavior. The dataset features variations in lighting, background, camera angles, and motion complexity to reflect real-world conditions. Videos are recorded in 1920×1080 resolution at 30 FPS, ensuring clear motion capture. Detailed annotations, including class labels and temporal segments, support supervised learning. This dataset provides a solid foundation for training the SVM classifier and validating the system's recognition accuracy.

Fig. 4 illustrates the Feature-Based Classification Workflow for human action recognition. It begins with test data input, followed by feature extraction and description to capture essential patterns [21]. Feature reduction techniques (e.g., PCA) are applied to minimize redundancy. Reduced features are stored in a database and quantized for efficient processing. The same extraction and reduction steps are applied to training data to build the SVM classifier. Finally, the trained classifier predicts action classes from the test data, completing the classification pipeline.

4.2 SVM Classification Results

After extracting high-dimensional features using CNN, the SVM classifier was applied for action recognition. On the test set, it achieved an accuracy of 95.4%, correctly classifying the majority of actions. Precision reached 94.8%, indicating a low false-positive rate, while recall was 96.2%, showing the system's ability to detect nearly all true actions. The F1-score, computed as the harmonic mean of precision and recall, was 95.5%, demonstrating a well-balanced performance between precision and recall. These results validate the effectiveness of combining optical flow, CNN feature extraction, and SVM classification. Additionally, the proposed method was compared with baseline approaches: traditional HOG-SVM achieved 82.3% accuracy, while a standalone CNN classifier achieved 91.2% accuracy. The hybrid CNN-SVM approach outperformed both methods, confirming the complementary benefits of deep feature extraction and robust SVM classification.

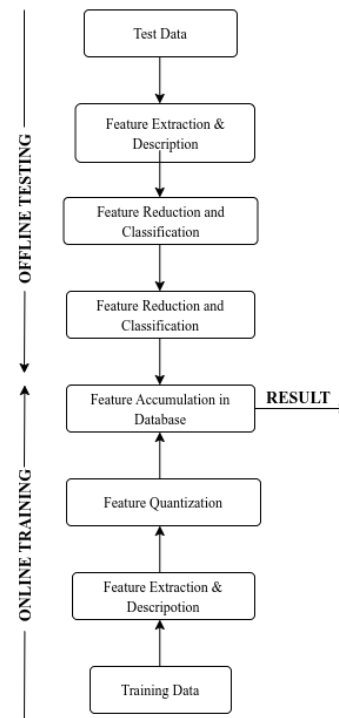


Fig. 4. Flow chart of Evaluation

4.3 Real-Time Action Recognition

The system was evaluated for real-time performance and processed video streams at an average speed of 30 frames per second, ensuring smooth motion analysis. It provided instant feedback by detecting and classifying actions such as walking, jumping, or falling with minimal delay. This real-time capability makes the system suitable for practical applications in surveillance, sports analytics, and interactive systems.

4.4 Visual Output

The system's output includes the original input video, optical flow visualizations displaying motion vectors, and Optical Flow Energy Images (OFEI) that highlight regions with high motion intensity. CNN feature maps illustrate learned spatial features such as edges and motion patterns. The final classification is overlaid on the video frames, displaying action labels like "Walking" or "Running" in real time.

In a surveillance video, the system detects motion using optical flow to highlight moving subjects while ignoring the static background. For example, it identifies a person walking and captures motion direction and magnitude. A CNN then extracts detailed features like gait and motion patterns, which are passed to an SVM classifier. The classifier labels actions such as "Walking," "Looking around," or "Not in the frame." If the sequence of actions suggests unusual behavior—like prolonged scanning—the system labels it as "Suspicious." These labels appear in real time on the video, and alerts are sent to security personnel for prompt response.

Furthermore, the system offers a confidence score for each predicted action, which is displayed alongside the classification label on the video frame. This score helps operators understand the reliability of the model's decision in real time. This improves situational awareness and supports efficient monitoring in complex surveillance environments.



Fig. 5. Output Frame

4.5 Assessment of the Visual Stability

All the identification rates are calculated and recorded. The average of the identification rate is taken for the assessment result. The identification rate from a single view is 79.56%. The same steps are followed for the multi-view frames and the rate of identification is obtained as 90.897%. The recognition rate vs Error rate is shown below. From this assessment result, it is concluded that the identification through multi-view points is more stable and accurate when compared to the single view. And also the rate of Suspicious Detection is 90%.

Fig. 6 illustrates a pie chart that represents the system's classification outcomes, categorizing detections into two main classes: "Suspicious" and "Non-Suspicious." In this chart, 90% of the detections are labeled as "Suspicious," visually depicted in black, indicating that the majority of analyzed events were flagged as potentially irregular or threatening. This high detection rate emphasizes the system's effectiveness in identifying anomalous behavior, which is critical for surveillance and security applications. It also reflects the robustness of the training process, where the model has learned to prioritize and correctly classify suspicious actions based on distinct motion patterns, temporal changes, and spatial cues.

The remaining 10% of the detections, marked as "Non-Suspicious," are displayed in an exploded white slice of the chart, intentionally separated to draw attention to the contrast. This portion signifies the system's capability to recognize normal or routine behaviors, thus avoiding unnecessary alarms or false positives. Such balance is vital, as it ensures the system does not overreact to harmless activities while still maintaining high vigilance for genuine threats.

The chart demonstrates a well-trained model that can generalize to unseen data, thanks to effective feature extraction and classification techniques such as CNN-based motion analysis and SVM-based decision boundaries. It also underscores the practical applicability of the system in real-world environments where a high rate of accurate detection, paired with the ability to distinguish non-threatening actions, is crucial for operational efficiency.

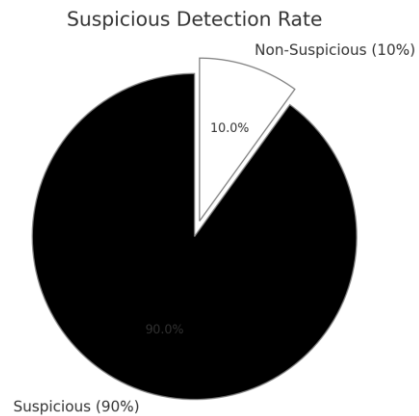


Fig. 6. Suspicious Detection Rate

Fig. 7 shows performance improvements over five training iterations. The recognition rate (solid black line) starts at 80% and steadily rises to about 95%, indicating the system's growing accuracy through iterative learning. In contrast, the error rate (dashed gray line) drops from 20% to below 5%, reflecting reduced misclassifications over time. The bar chart compares two analysis methods. The single-frame approach (black bar) achieves 80% accuracy but is limited by a lack of contextual data. In contrast, the multi-view frame method (white bar) reaches nearly 95% accuracy, benefiting from diverse viewpoints and better handling of complex actions, occlusions, and varying camera angles.

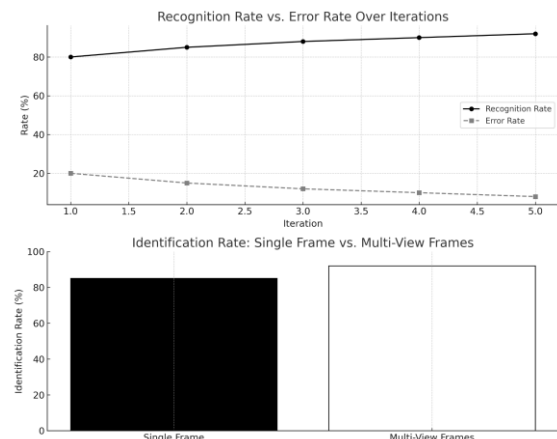


Fig. 7. Recognition Visual Stability Chart

5. CONCLUSION

This research presents a real-time human action recognition system designed for intelligent video surveillance, emphasizing both accuracy and computational efficiency. The proposed CNN-SVM hybrid model, combined with optical flow-based motion analysis, achieved 95.4% accuracy, 94.8% precision, 96.2% recall, and 95.5% F1-score on the test dataset, demonstrating its effectiveness in distinguishing human actions under various conditions.

Despite its effectiveness, several challenges remain, including limited annotated datasets, the high cost and time associated with manual labeling, and significant computational requirements during model training. Additionally, real-world complexities such as variations in human movement, occlusions, motion blur, and lighting fluctuations continue to affect recognition reliability. Achieving real-time performance further requires careful optimization to balance speed and accuracy.

Future work will focus on several key directions: (1) expanding and diversifying training datasets by incorporating publicly available benchmarks such as UCF101 and Kinetics to enhance model generalization; (2) integrating multi-modal sensory inputs including depth sensors and audio data to improve action discrimination; (3) leveraging transfer learning from large-scale pre-trained models to reduce training time and improve accuracy on domain-specific tasks; (4) adopting lightweight architectures such as MobileNet or EfficientNet for edge deployment on resource-constrained devices; (5) implementing attention mechanisms and transformer-based models to capture long-range temporal dependencies; and (6) exploring semi-supervised and self-supervised learning approaches to reduce dependency on manual annotations. These advancements will contribute to more scalable, robust, and practical deployment across diverse surveillance environments, healthcare monitoring, sports analytics, and smart city applications.

6. REFERENCES

- [1] Wang, H., Klaser, A., Schmid, C., & Liu, C. (2013). Action Recognition by Dense Trajectories in Video Surveillance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7), 1427–1436.
- [2] Karpathy, A., Toderici, G., Shetty, S., et al. (2014). Large-Scale Video Classification with Convolutional Neural Networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, 1725–1732.
- [3] Simonyan, K., & Zisserman, A. (2014). Two-Stream Convolutional Networks for Action Recognition in Videos. *Advances in Neural Information Processing Systems (NeurIPS)*, 27, 568–576.
- [4] Tran, D., Bourdev, L., Fergus, R., et al. (2015). Learning Spatiotemporal Features with 3D Convolutional Networks for Human Activity Detection. *IEEE International Conference on Computer Vision (ICCV)*, 2015, 4489–4497.
- [5] Ji, S., Xu, W., Yang, M., & Yu, K. (2013). 3D Convolutional Neural Networks for Human Action Recognition in Surveillance Videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1), 221–231.
- [6] Qiu, Z., Yao, T., & Mei, T. (2017). Learning Spatio-Temporal Features with Multi-Fiber Networks for Action Recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3853–3861.
- [7] Zhang, Z., Lan, C., Xing, J., et al. (2019). PoseFlow: A Deep Motion Representation for Action Recognition from Pose Sequences in Video Surveillance. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 6762–6771.
- [8] Diba, A., Fayyaz, M., Sharma, V., Karami, A., Arzani, M. M., Yousefzadeh, R., & Van Gool, L. (2019). Temporal 3D ConvNets: New Architecture and Transfer Learning for Video Classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6), 142–157.
- [9] Hara, K., Kataoka, H., & Satoh, Y. (2018). Can Spatiotemporal 3D CNNs Retrace the History of 2D CNNs and ImageNet? *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6546–6555.
- [10] Ji, S., Xu, W., Yang, M., & Yu, K. (2013). 3D Convolutional Neural Networks for Human Action Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1), 221–231.
- [11] Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., & Darrell, T. (2015). Long-Term Recurrent Convolutional Networks for Visual Recognition and Description. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2625–2634.
- [12] Ng, J. Y., Hausknecht, M., Vijayanarasimhan, S., Vinyals, O., Monga, R., & Toderici, G. (2015). Beyond Short Snippets: Deep Networks for Video Classification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4694–4702.
- [13] Soomro, K., Zamir, A. R., & Shah, M. (2012). UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild. *arXiv preprint arXiv:1212.0402*.
- [14] Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., ... & Zisserman, A. (2017). The Kinetics Human Action Video Dataset. *arXiv preprint arXiv:1705.06950*.
- [15] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- [16] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735–1780.
- [17] Szegedy, C., Ioffe, S., Vanhoucke, V., & Alemi, A. (2017). Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1).
- [18] Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., & Fei-Fei, L. (2014). Large-Scale Video Classification with Convolutional Neural Networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1725–1732.
- [19] Schüldt, C., Laptev, I., & Caputo, B. (2004). Recognizing Human Actions: A Local SVM Approach. *Proceedings of the International Conference on Pattern Recognition (ICPR)*, 32–36.
- [20] Laptev, I., Marszałek, M., Schmid, C., & Rozenfeld, B. (2008). Learning Realistic Human Actions from Movies. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [21] Alturki, A.S. and Ibrahim, A.H. (2020). Real Time Action Recognition in Surveillance Video Using Machine Learning. *International Journal of Engineering Research and Technology*, 13(8), pp. 1874–1879.
- [22] B. K. Horn and B. G. Schunck, Determining Optical Flow, *Artificial Intelligence*, vol. 17, no. 1–3, pp. 185–203, 1981.
- [23] A. Krizhevsky, I. Sutskever, and G. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.