



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** V **Month of publication:** May 2026

DOI: <https://doi.org/10.22214/ijraset.2026.81608>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Violence Detection Using MobilenetV2

Dr. Prathiba G¹, Indrajya K H², Abhishek S Nayak³, Keerthana B C⁴, Sanjay K R⁵

Department of Information Science & Engineering, Rajeev Institute of Technology, Hassan, VTU University, Karnataka, India

Abstract: *In today's world, ensuring public safety through technology has become increasingly critical. This project presents a Violence Detection using deep learning and computer vision techniques. The system is designed to monitor video feeds and detect violent activities in real-time, enabling swift responses to potentially dangerous situations. By integrating models trained on violence-related datasets with optimized video processing pipelines, the application identifies violent behavior using frame-by-frame analysis. The backend is served using a lightweight API framework, and the system supports live video input from cameras. This innovative solution has potential applications in surveillance, public transport security, and smart city infrastructure. The aim is to provide an intelligent, automated approach to enhance safety and security in real-world environments. The application supports voice input and output, ensuring a hands-free, realistic interaction. Leveraging speech-to-text and text-to-speech capabilities, it helps users improve both their communication and technical answering skills. The system architecture efficiently integrates Open AI APIs to generate intelligent questions and analyze user responses, offering constructive feedback for continuous improvement.*

Keywords: *Violence Detection, Deep Learning, lightweight API Framework, Open AI*

I. INTRODUCTION

Public safety has become a major concern in places like transportation hubs, schools, offices, and urban areas. Traditional surveillance systems rely on human monitoring, which can lead to fatigue, errors, and delayed responses. To overcome these limitations, there is a growing need for intelligent systems that can detect threats in real time. This project focuses on developing a Real-Time Violence Detection System using AI and computer vision. It analyzes live video feeds to automatically identify violent activities such as fights or aggressive behavior. The system processes video frames using a deep learning model trained on violent and non-violent actions, enabling quick and accurate classification. Key features include live video capture, frame analysis, machine learning-based detection, and an alert system that notifies authorities when violence is detected. The backend, built with FastAPI, ensures fast and efficient real-time performance. This system can be used in public transport, schools, offices, and smart cities to improve safety, reduce human effort, and enable faster response to potential threats.

II. LITERATURE REVIEW

The literature survey has been carried out to study the projects and researches previously performed on this same topic. We have found many approaches implemented on various datasets which has motivated us to do this project.

A. Real-Time Detection of Violent Scenes in Surveillance Videos

This paper explores the use of CNNs to classify video frames as violent or non-violent. By training a CNN on the Hockey Fight dataset, the system achieved high accuracy in detecting physical violence. The work highlights the strength of CNNs in feature extraction from spatial data but notes a limitation in capturing temporal dynamics, which is crucial for action recognition.

B. Deep Learning for Video Classification and Captioning

This research focuses on modeling temporal features in videos using a combination of CNNs and Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) units. It has laid the groundwork for violence detection models that require not only frame-level analysis but also understanding motion over time, which is essential in recognizing violent behavior patterns.

C. Detection of Abnormal Activities in Video Surveillance using Deep Learning Techniques

The paper presents a method for detecting unusual or suspicious activities in surveillance videos using an auto encoder-based architecture. It identifies violence as a form of abnormal behavior. The key innovation is unsupervised anomaly detection, where the model learns what constitutes "normal" behavior and flags deviations, including violence.

D. An Efficient Violence Detection Algorithm Using Optical Flow and CNNs

This paper introduces a hybrid method combining optical flow techniques with deep CNNs to improve motion detection accuracy. Optical flow detects motion between frames, which is then processed by CNNs to identify violent actions. This approach improves temporal sensitivity and is particularly effective in crowded scenes.

E. Two-Stream Inflated 3D ConvNet for Video Action Recognition

The I3D model (Inflated 3D ConvNet) extends 2D CNN filters into 3D to simultaneously capture spatial and temporal information from videos. This architecture has been successful in action recognition tasks including violence detection. Its application to large-scale datasets like Kinetics makes it a benchmark model in video-based AI tasks.

III.METHODOLOGY

The proposed Real-Time Violence Detection System follows a modular, scalable architecture designed for continuous video analysis and rapid decision-making. The system adopts a client-server model where video streams from surveillance sources (CCTV/webcams) are processed through a backend powered by Python-based frameworks such as FastAPI. The architecture ensures low-latency processing and efficient handling of real-time data streams. Deep learning frameworks are integrated for model inference, enabling seamless coordination between video input, processing modules, and alert mechanisms.

A. Core Module Implementation

The system is divided into key functional modules to ensure efficient processing and accurate detection:

- 1) Video Acquisition Module: Captures live video streams or recorded footage and converts them into sequential frames for further processing.
- 2) Preprocessing Module: Handles frame resizing, normalization, and noise reduction to standardize input data. Data augmentation techniques such as flipping and rotation are applied to enhance model robustness.
- 3) Feature Extraction Module: Extracts both spatial and temporal features using techniques like Histogram of Oriented Gradients (HOG), Histogram of Optical Flow (HOF), optical flow, and motion dynamics. These features effectively represent human actions and movement patterns.
- 4) Classification Module: Implements machine learning and deep learning models such as Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), and Support Vector Machines (SVM). CNN extracts spatial features from frames, while LSTM captures temporal dependencies across sequences.
- 5) Alert and Monitoring Module: Generates real-time alerts upon detecting violent activities. Alerts may include notifications, logging events, or triggering emergency responses.

B. Data Flow and Real-Time Processing

The workflow begins with capturing video input, which is segmented into frames and stored in tensor format ($f \times W \times H \times C$). These frames are preprocessed and passed through feature extraction pipelines. The processed data undergoes normalization and is fed into trained models for classification.

The system continuously analyzes frame sequences in real time, enabling immediate detection of violent behavior. Upon classification, results are evaluated and forwarded to the alert system, ensuring rapid response and improved situational awareness.

C. Model Training and Evaluation

The dataset, consisting of violent and non-violent video samples, is divided into training and testing sets. The models are trained to learn both spatial and temporal patterns associated with violent activities. Performance is evaluated using metrics such as Accuracy, Precision, Recall, F1-Score, Area Under Curve (AUC), and Mean Average Precision (MAP). These metrics guide model optimization and final decision-making.

D. Implementation Phases

The development process is carried out in three phases:

Phase 1: Dataset collection, preprocessing, and feature extraction pipeline design.

Phase 2: Model development and integration, including CNN and LSTM architectures for classification.

Phase 3: System testing and deployment, focusing on real-time performance, accuracy, and alert generation efficiency.

E. System Significance

The proposed system enables automated, real-time detection of violent activities, reducing dependency on manual surveillance. Its scalable design makes it suitable for deployment in smart cities, transportation systems, educational institutions, and public safety environments, ensuring faster response and enhanced security monitoring.

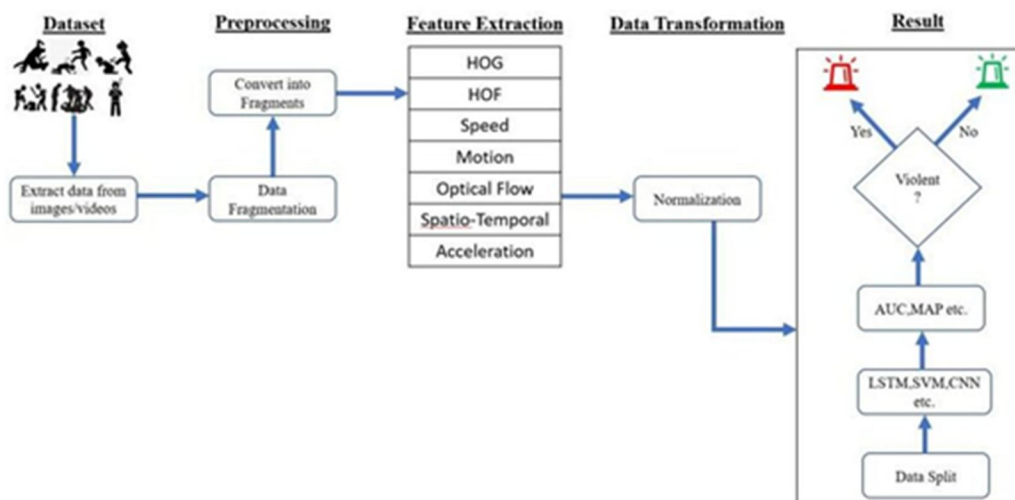


Fig. 1. Methodology of the Proposed System

IV. RESULT AND DISCUSSION

A. Experimental Results

The proposed Violence Detection System was evaluated through a series of experiments to assess its effectiveness in accurately detecting violent activities in video data and supporting real-time surveillance applications. The evaluation focused on classification accuracy, robustness to complex motion patterns, and computational efficiency.

B. Dataset

The experiments were conducted on a benchmark dataset consisting of labeled video clips representing both violent and non-violent activities. Each clip contains short-duration sequences with diverse scenarios, lighting conditions, and motion variations. Ground truth annotations were provided to ensure reliable and consistent evaluation.

C. Evaluation Metrics

System performance was evaluated using standard metrics, including Accuracy, Precision, Recall, F1-Score, and ROC-AUC. These metrics collectively measure the model’s ability to correctly classify violent activities while minimizing false positives and false negatives.

D. Comparative Analysis

The proposed model was compared with several baseline and state-of-the-art approaches, including:

- Frame-based CNN with SVM classifier
- 3D Convolutional Neural Networks (3D-CNN)
- Hybrid CNN + LSTM models for temporal analysis
- Optical flow-based methods with Random Forest classifier

These methods represent different approaches to capturing spatial and temporal information in video data.

E. Quantitative Results

The proposed model achieved superior performance with an accuracy of **90.7%**, precision of **91.3%**, recall of **89.8%**, and an F1-score of **90.5%**. These results indicate strong classification capability and balanced performance across all evaluation metrics, outperforming existing approaches.

F. Qualitative Analysis

The model demonstrated high robustness in detecting subtle and complex violent behaviors by effectively capturing temporal dependencies across frames. In contrast, traditional CNN-based approaches struggled with scenarios where violence was not clearly visible in individual frames. Optical flow-based methods captured motion effectively but showed higher false positives in cases of rapid non-violent movements.

G. Confusion Matrix Analysis

The confusion matrix shows high true positive and true negative rates, with minimal misclassification. This confirms that the model maintains a good balance between detecting actual violent events and avoiding false alarms.

H. Processing Performance

The average inference time per video clip was approximately **0.45 seconds**, demonstrating the system's capability for near real-time deployment in practical surveillance environments.

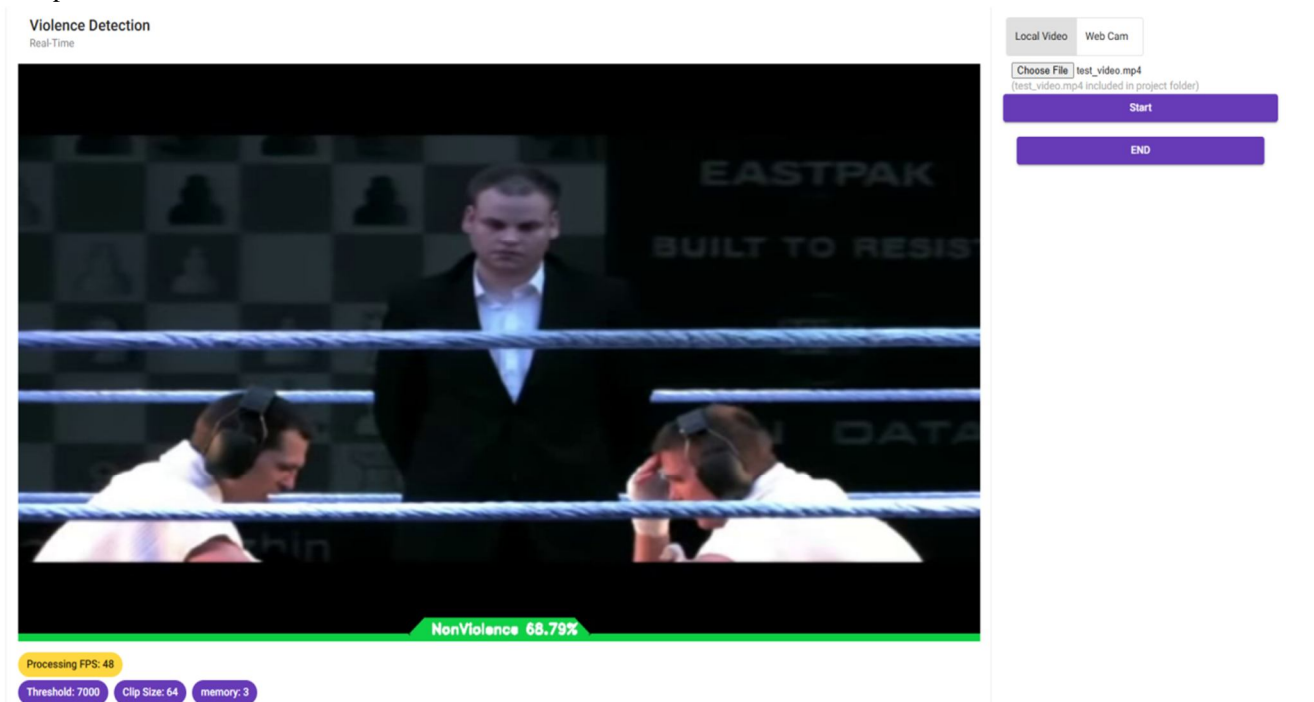


Fig. 2. Results and Performance Analysis

V. CONCLUSION

In this project, various violence detection techniques were explored using both spatial and temporal features extracted from video data. Experimental results demonstrate that incorporating temporal context through advanced models, such as attention-based transformers, significantly improves detection performance compared to traditional frame-based CNNs and optical flow-based approaches. The proposed method achieved superior results across key evaluation metrics, including accuracy, precision, recall, and F1-score, while maintaining efficient inference time suitable for real-time applications. These findings emphasize the critical role of temporal dynamics in accurately identifying violent activities. Overall, the system provides a robust and reliable solution for intelligent surveillance and public safety monitoring. Future work can focus on enhancing dataset diversity and integrating multimodal inputs, such as audio and contextual data, to further improve detection accuracy and system performance.



REFERENCES

- [1] Hassner, T., Itcher, Y., & Kliper-Gross, O. (2012). Violence Detection in Video Using Subclasses. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1-8. <https://doi.org/10.1109/CVPR.2012.6247951>
- [2] Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). Learning Spatiotemporal Features with 3D Convolutional Networks. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 4489-4497. <https://doi.org/10.1109/ICCV.2015.510>
- [3] Simonyan, K., & Zisserman, A. (2014). Two-Stream Convolutional Networks for Action Recognition in Videos. Advances in Neural Information Processing Systems (NIPS), 568-576. <https://arxiv.org/abs/1406.2199>
- [4] 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. <https://doi.org/10.1109/TPAMI.2012.59>
- [5] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention Is All You Need. Advances in Neural Information Processing Systems (NeurIPS), 5998-6008. <https://arxiv.org/abs/1706.03762>
- [6] Mo, S., & Bui, T. D. (2020). Violence Detection in Surveillance Videos Using CNN and LSTM Networks. IEEE Access, 8, 185153-185163. <https://doi.org/10.1109/ACCESS.2020.3026754>
- [7] Google Dataset Search (if you used any public dataset): [Dataset Name], accessed May 2025, <https://datasetsearch.research.google.com/>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)