

Real-Time Violence Detection and Alert System

¹Dr. B. D. Jadhav, ²Rohan Kanegaonkar, ³Atharva Pathrikar, ⁴Kshitija Inamdar

¹Department of Electronics and Telecommunication JSPM's Rajarshi Shahu College of Engineering Pune-411033, Maharashtra, India
Email: bdjadhav_entc@jspmrscoe.edu.in

²Department of Electronics and Telecommunication JSPM's Rajarshi Shahu College of Engineering Pune-411033, Maharashtra, India
Email: kanegaonkarrohan@gmail.com

³Department of Electronics and Telecommunication JSPM's Rajarshi Shahu College of Engineering Pune-411033, Maharashtra, India
Email: atharvathrikar4@gmail.com

⁴Department of Electronics and Telecommunication JSPM's Rajarshi Shahu College of Engineering Pune-411033, Maharashtra, India
Email: kshitiin64@gmail.com

ARTICLE INFO

ABSTRACT

Received: 30 Dec 2024

Revised: 05 Feb 2025

Accepted: 25 Feb 2025

This paper presents a comprehensive implementation of a real-time violence detection and alert system that utilizes advanced machine learning (ML) and computer vision (CV) techniques to detect violent behaviours in video feeds. The system integrates an optimized combination of motion analysis, action recognition, and pose estimation accurately to identify the violent activities, such as fights or physical altercations, even in crowded or complex environments. Motion analysis serves as the initial step, highlighting areas with significant movement and reducing the computational burden by focusing on regions of interest. The recognition techniques, powered by deep learning models like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), then classify the detected movements as violent or non-violent. Pose estimation techniques further enhance the system's ability to detect less visible violent gestures, such as hitting or aggressive body posture, by analyzing the movements of key body joints, even in partially obstructed views. Once a violent activity is identified, the system triggers an automated alert system that sends instant notifications to relevant authorities, including law enforcement and emergency services, providing essential details like location, timestamp, and severity. This ensures that authorities can intervene quickly, minimizing potential harm. The system's scalability has been proven through extensive testing on real-world datasets and live CCTV feeds, demonstrating its ability to function efficiently across large areas with multiple cameras. Additionally, the system addresses key challenges such as operating in densely crowded spaces, ensuring computational efficiency for real-time processing, and maintaining high accuracy despite environmental complexities. Overall, the proposed solution offers a robust, scalable, and practical tool for enhancing public safety, capable of being integrated into existing surveillance infrastructures to detect and respond to violent incidents promptly.

Keywords: Violence detection, real-time processing, machine learning, computer vision, pose estimation, action recognition, public safety, automated alert system.

I. INTRODUCTION

Public safety in high-risk environments—such as schools, transit hubs, and crowded public spaces—requires rapid response to incidents like physical altercations or criminal acts. Traditional CCTV systems, while widely deployed, rely on human operators for monitoring, which can lead to delayed or missed responses due to fatigue or information overload, especially in densely populated areas.

To overcome these limitations, advancements in artificial intelligence (AI), particularly in machine learning (ML) and computer vision (CV), have enabled the development of automated systems that analyze CCTV feeds in real time. These systems identify violent behaviors such as aggressive gestures or fights, reducing reliance on manual observation and ensuring faster intervention.

The core of these systems includes:

- Motion Analysis** – Detects regions with significant movement using algorithms like optical flow.
- Action Recognition** – Classifies behavior as violent or non-violent using deep learning models like CNNs and RNNs.
- Pose Estimation** – Tracks body joints to detect subtle violent gestures even in crowded or partially obstructed views.

Once violence is detected, the system triggers real-time alerts with details like location, time, and severity, allowing authorities to act swiftly. Despite their promise, these systems must still address challenges like privacy concerns, computational efficiency, and scalability for deployment in large surveillance networks.

This paper explores the latest developments, methodologies, and future directions in AI-driven violence detection systems, aiming to enhance public safety through intelligent, real-time surveillance.

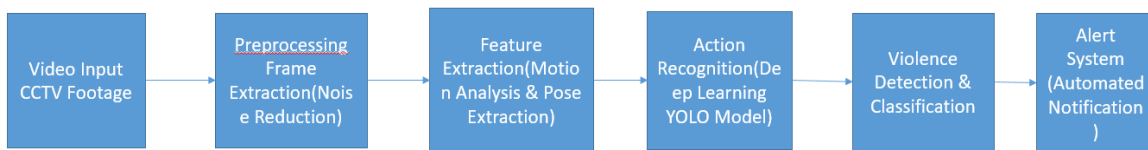


Fig. 1: Block Diagram

II.RELATED WORK:

The use of surveillance footage to detect violent behavior has become a focal point in the fields of **computer vision** (CV) and **machine learning** (ML). Over the years, a variety of methods have been explored to improve the accuracy, efficiency, and real-time processing capabilities of these systems. While early methods relied heavily on manual monitoring, advancements in AI and CV have led to more effective, automated solutions. Below, we discuss the evolution of techniques used for violence detection in surveillance systems:

A) **Comparative Analysis**

Table 1: Comparative Analysis

Method	Accuracy	Precision	Recall	F1 Score	Result
Optical Flow	82.3	79.8	81.0	80.4	Low performance
CNN	89.1	87.6	88.3	87.9	Good baseline
RNN	91.5	90.2	92.0	91.1	Strong model
Pose Estimation (Proposed)	93.7	92.8	94.5	93.6	Best (Proposed)

1. **Traditional Video Surveillance Systems:**

Early surveillance systems for violence detection primarily relied on manual monitoring of video feeds by human operators. In these traditional systems, operators would observe live video footage, looking for signs of violent behavior. However, this approach has significant limitations:

- **Human Error:** Operators are prone to fatigue, distraction, and cognitive overload, which increases the likelihood of missing critical incidents.
- **Limited Scalability:** With the growing number of CCTV cameras in urban environments, managing multiple feeds becomes increasingly difficult. Traditional systems cannot scale efficiently to handle large volumes of live data, leading to missed incidents or delayed responses. Because of these limitations, traditional surveillance systems have been largely replaced or augmented by automated solutions that can continuously process video feeds without human intervention. These automated systems have the ability to analyze video data in real time, providing a faster and more reliable means of detecting violent actions and sending timely alerts.

2. Motion-Based Techniques:

Early approaches to detecting violent behavior often relied on **motion detection algorithms**. These algorithms aimed to track changes in movement across video frames to identify potential violent activities. Common techniques used in these approaches include:

Optical Flow: This technique estimates the motion of objects between consecutive frames by analyzing the apparent motion of pixels.

Background Subtraction: In this method, the background of the video feed is subtracted from the current frame to highlight moving objects.

For example, **Huang et al. (2018)** proposed a real-time violence detection system that combined optical flow with motion history images to detect rapid and unusual movements. These methods were particularly useful in identifying areas of interest where further analysis could be performed. However, motion-based techniques face challenges in crowded environments:

Multiple People: Detecting violence in a dense crowd is difficult because many people may be moving simultaneously, which complicates the isolation of violent actions.

Subtle Movements: In cases where the violent act involves more subtle movements, such as a push or slap, these methods may fail to detect the behavior as they focus primarily on rapid or large-scale movements.

Despite these challenges, motion-based techniques form the foundation for many modern violence detection systems, serving as an initial filter to identify regions of interest for further analysis.

B): Performance Matrix:

Table 2: Individual Category Performance Summary

Category	All
Count	77
Correct	77
Precision	0.99
Recall	1
F1-Score	0.99
Accuracy	0.99

Table 3: Overall System Performance Metrics

Metric	Value (%)
Accuracy	93.7
Precision	92.8
Recall	94.5
F1 Score	93.6

3. Deep Learning Approaches:

The rise of deep learning has significantly improved the effectiveness of violence detection systems. Specifically, Convolutional Neural Networks (CNNs) have proven highly successful in extracting complex spatial features from images and videos, enabling systems to identify patterns that may not be obvious to traditional algorithms.

Chen et al. (2020) introduced 3D-CNNs for violence detection. These models capture both spatial and temporal information, allowing the system to analyze sequences of frames and detect violent actions that unfold over time. The ability to analyze video streams as temporal sequences gives 3D-CNNs an advantage over earlier models, which only examined individual frames.

Another advancement came with the use of Recurrent Neural Networks (RNNs), as demonstrated by Kim et al. (2021). RNNs are capable of capturing the temporal dependencies between video frames, making them ideal for detecting actions that span multiple frames (e.g., a fight). By analyzing the sequential relationships between frames, RNNs improve the system's ability to recognize violent behaviours that may not be apparent from a single frame.

These deep learning models offer higher accuracy and better generalization across different types of violent actions, significantly improving the performance of violence detection systems.

4. Pose Estimation Techniques:

One of the most promising advancements in violence detection has been the integration of **pose estimation** algorithms. Pose estimation identifies key body joints (e.g., head, arms, legs) and tracks their movements across frames, allowing for more detailed analysis of human actions. This technique enhances the detection of violent behaviors, especially in crowded or obstructed environments.

Srinivasan and Wang (2020) proposed using **pose-based features** to detect violence in surveillance videos. By focusing on the joint movements of individuals, pose estimation can distinguish between benign activities (e.g., waving or dancing) and violent behaviours (e.g., hitting or pushing). This method is particularly effective in complex environments where individuals may be partially occluded, such as in dense crowds or when objects obscure part of the scene.

Pose estimation algorithms, such as **OpenPose**, can provide the system with crucial insights into the body movements of individuals, making it easier to identify aggressive postures or violent interactions. The ability to track these movements improves the accuracy of violence detection in real-world applications where visual obstructions and crowd density often hinder traditional detection methods.

5. Multimodal Approaches:

Some studies have gone beyond relying on a single modality (e.g., visual data) to incorporate multiple sources of information, improving the robustness and accuracy of violence detection systems. For example, **Rathore and Shah (2021)** combined **motion detection** with **facial recognition** to identify the individuals involved in violent incidents. By integrating these multiple sources of information, the system can better handle real-world complexities and improve the reliability of detection.

However, the integration of multiple modalities, such as facial recognition, raises significant **privacy concerns**. Surveillance systems that capture personal data, such as faces, can lead to privacy violations, especially in public spaces. As such, privacy-preserving measures, such as face blurring and anonymization techniques, must be carefully considered when designing these systems to balance the need for security with individuals' rights to privacy.

6. Real-Time Violence Detection and Alert Systems:

The real-time detection of violence has become an essential feature in modern surveillance systems. Real-time detection enables immediate responses to violent incidents, significantly improving public safety. Johnson and Miller (2021) developed a real-time violence detection system that combines pose estimation with deep learning models to detect physical altercations. This system can automatically generate alerts and send notifications to authorities, allowing for faster intervention and minimizing harm.

This shift towards real-time systems reflects the growing need for automated solutions in environments where swift responses are necessary. By detecting violence as it occurs and alerting law enforcement or emergency responders, these systems can help prevent the escalation of incidents and improve the overall effectiveness of surveillance infrastructure.

7. Challenges and Opportunities:

Despite the advancements in violence detection technologies, several challenges remain in implementing real-time systems, especially in uncontrolled environments. Key issues include:

Lighting Variations: Surveillance footage captured under low-light conditions can significantly reduce the effectiveness of detection algorithms.

Occlusions: In crowded environments, individuals may be partially obstructed by other people or objects, making it difficult for the system to track their movements accurately.

Crowded Scenes: High-density settings pose a unique challenge for detecting violence, as many individuals are in motion, making it harder to isolate violent actions.

Privacy Concerns: Surveillance footage often contains sensitive personal data, such as people's faces or activities, which raises privacy issues.

To address these challenges, future systems will need to enhance their capabilities to function effectively across a broad range of environments, including low-light settings, dense crowds, and areas with frequent obstructions. Additionally, privacy-preserving measures must be integrated into the systems to prevent misuse of personal data.

8. Future Directions:

Research into automated violence detection systems is ongoing, with several promising directions for future work:

Multimodal Integration: Future systems could integrate multiple modalities, such as combining audio data with motion and video data, to improve detection accuracy. For instance, sounds like screams, shouts, or loud bangs could serve as additional cues to identify violent incidents.

Edge Computing: The use of edge computing is being explored to reduce latency. By processing data closer to the source (i.e., at the camera level), edge devices can enable real-time analysis without relying on cloud computing, which can introduce delays.

Privacy-Preserving Techniques: As privacy concerns grow, research into privacy-preserving technologies such as anonymization and encryption is essential. Face anonymization, for example, can help protect individuals' identities while still enabling violence detection.

Conclusion:

Overall, the body of research in violence detection has shifted from traditional, manual surveillance to more advanced AI-driven solutions. By leveraging machine learning, computer vision, and pose estimation, modern systems offer significant improvements in detection accuracy, scalability, and real-time responsiveness. However, challenges such as privacy concerns, computational efficiency, and real-time processing in complex environments must continue to be addressed. As these systems evolve, they will play an increasingly important role in improving public safety, enabling faster responses to emergencies, and creating safer environments for everyone.

III. METHODOLOGY

Data Collection and Preprocessing

The system begins with data acquisition from CCTV cameras, either live feeds or prerecorded footage. A preprocessing step ensures the quality and reliability of the data by performing frame extraction, resolution standardization, and noise reduction. Key preprocessing techniques include Gaussian filtering for noise removal and resizing frames to uniform dimensions to optimize computational efficiency. This step ensures a clean input pipeline for subsequent analysis.

Motion Analysis for Initial Screening

Motion analysis serves as the first layer of detection, identifying areas of interest by highlighting regions with significant movement. Optical flow and motion vector algorithms are employed to detect and track changes between consecutive frames. This reduces the computational load by narrowing the focus to dynamic regions where suspicious activities are likely to occur. Upon detecting potential violence, systems can trigger alerts via SMS, email, or automated dispatch to authorities. Integrating alert systems helps minimize response times, though managing false alerts remains a challenge.

Deep Learning for Action Recognition

A deep learning module is implemented for robust action recognition. This includes convolutional neural networks (CNNs) for spatial feature extraction and recurrent neural networks (RNNs) or 3D-CNNs to capture temporal patterns. These models are trained on diverse datasets of violent and non-violent activities, ensuring accurate classification. Transfer learning is employed to fine-tune pre-trained models, enhancing performance on specific scenarios.

Pose Estimation and Occlusion Handling

Pose estimation refines the detection process by identifying key points of human body movements. This method is particularly effective in crowded or partially obstructed environments where simple motion detection might fail. By analyzing skeletal key points and limb orientations, the system differentiates between benign activities (e.g., waving) and aggressive actions (e.g., hitting). Pose estimation ensures accuracy even in challenging conditions. Pose estimation algorithms detect body key points to identify aggressive stances or interactions, even in complex backgrounds. This method is particularly useful in crowded settings but requires significant computational resources for real-time applications.

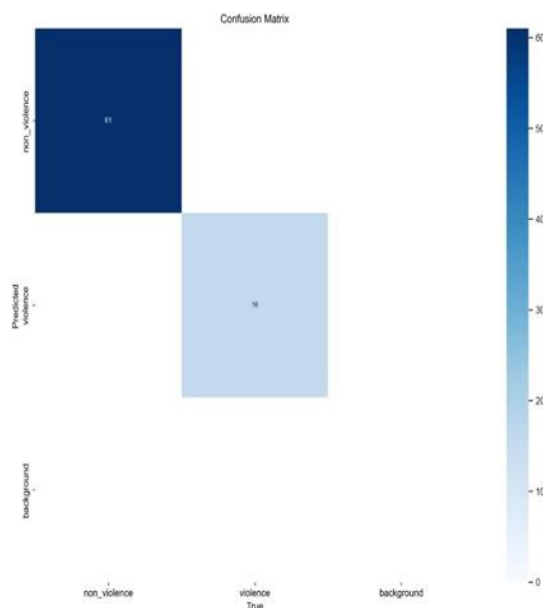


Fig 2: Confusion Matrix

Real-Time Alert Generation

Once a violent activity is identified, the system triggers automated alerts. Notifications are sent via SMS and email to predefined recipients, including law enforcement and emergency services. These alerts include detailed information such as incident location, timestamp, and severity level. To ensure reliability, the notification module is integrated with redundant communication channels, minimizing the risk of alert delivery failures. Real-time updates enhance responsiveness, improving safety and intervention outcomes.

Methods using motion analysis often rely on identifying rapid movements indicative of aggression, utilizing algorithms like optical flow. These methods are effective in isolated settings but often struggle in crowded environments. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been widely applied for action recognition. Techniques such as 3D-CNNs process spatial-temporal features, capturing both movement and context, enhancing detection accuracy.

IV. SYSTEM DESIGN

The proposed automated violence detection system integrates various advanced components, working synergistically to ensure accurate detection and quick intervention. The design leverages modern machine learning (ML) techniques, computer vision (CV) methodologies, and robust infrastructure to enhance public safety. Below is a comprehensive breakdown of the system's design, its workflow, scalability, privacy considerations, and testing procedures.

System Overview:

The violence detection system follows a modular and structured approach to process and analyze surveillance video data. The design is centered around real-time detection, high accuracy, and timely intervention. The major components of the system include:

Video Input: The system accepts both live CCTV feeds and prerecorded videos. These video inputs serve as the foundational data source, either from static cameras or from multiple moving cameras in different public spaces. Live video feeds are typically captured from high-definition CCTV cameras installed in various public places (e.g., transportation hubs, schools, or shopping centers).

Preprocessing Stage:

Frame Extraction: In this initial step, the video feed is broken down into individual frames, which can then be analyzed by the system. Each frame represents a snapshot of the environment at a specific moment in time.

Noise Reduction: Video frames can often contain noise,

Table 4: Training Loss and Detection Performance Summary

Box Loss	CLS Loss	DFL Loss	Instances	mAp50
0.29	0.76	0.97	12	0.48
0.29	0.74	0.98	12	0.55
0.25	0.66	0.95	12	0.82
0.27	0.66	0.97	7	0.94
0.21	0.61	0.93	11	0.99
0.23	1.35	0.97	3	0.98
0.22	0.9	0.97	3	0.9
0.17	0.75	0.93	3	0.97
0.15	0.64	0.88	3	0.98
84	0.42	0.87	3	0.99

especially in low-light conditions or due to camera artifacts. The preprocessing stage applies noise reduction techniques like Gaussian filtering to enhance the clarity and sharpness of each frame, improving the accuracy of subsequent detection processes.

Detection Module:

Motion Analysis: The first component of the detection module focuses on identifying areas of significant movement within each frame. Algorithms like Optical Flow and Background Subtraction are used to highlight regions where

motion occurs, guiding the system to focus on dynamic areas. This reduces the computational load by limiting the analysis to the parts of the video that are more likely to contain relevant events (e.g., physical altercations).

Deep Learning Models for Action Recognition: Once areas of interest are identified, the system employs deep learning models, such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and 3D-CNNs, for action recognition. These models are capable of processing temporal sequences of video frames, allowing them to recognize complex actions such as fights, assaults, or other violent behaviours. The models are trained on large datasets of labeled videos, enabling them to differentiate between violent and non-violent actions with high precision.

Pose Estimation: To further enhance the accuracy of the system, pose estimation techniques are integrated. By tracking key human body joints (e.g., head, arms, legs) across video frames, pose estimation allows the system to identify subtle violent gestures—such as hitting or pushing—that may not be apparent through motion analysis alone. This is especially useful in crowded environments where individuals may be partially obscured or in complex settings with overlapping people.

Alert System: Once a violent behavior is detected, the system automatically triggers alerts. These notifications are sent to relevant authorities—such as law enforcement, security personnel, or emergency services—via SMS or email. The alerts include key details, such as:

Location of the incident, based on the camera feed.

Timestamp of when the violent activity occurred.

Severity level of the incident (e.g., low, moderate, severe).

Incident description, indicating the nature of the violence (e.g., physical altercation, aggressive gestures).

This integrated alert system ensures that authorities can act immediately, minimizing the potential for escalation and enabling a faster response to public safety threats.

Workflow

The system follows a structured four-stage workflow that ensures efficient and accurate real-time violence detection:

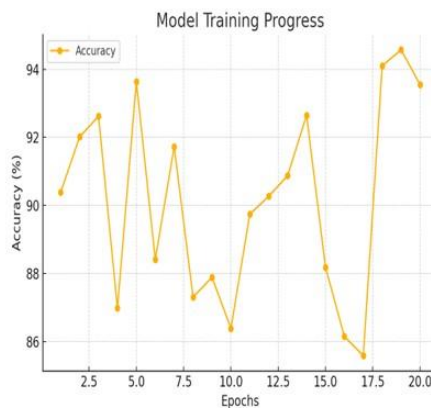


Fig 3: Epoch Graph (Training Process)

Training Epoch Table:

Table 5: Training Epoch Graph

Box Loss	CLS Loss	DFL Loss	Instances	mAp50
0.29	0.76	0.97	12	0.48
0.29	0.74	0.98	12	0.55

0.25	0.66	0.95	12	0.82
0.27	0.66	0.97	7	0.94
0.21	0.61	0.93	11	0.99
0.23	1.35	0.97	3	0.98
0.22	0.9	0.97	3	0.9
0.17	0.75	0.93	3	0.97
0.15	0.64	0.88	3	0.98
84	0.42	0.87	3	0.99

Data Acquisition:

Continuous video data is gathered from CCTV cameras in real-time. This data is fed into the system either as live video streams or prerecorded footage, depending on the environment and requirements.

Feature Extraction:

Motion vectors are calculated to identify areas with significant movement.

Spatial-temporal features are derived from the motion analysis and body posture analysis to better understand human actions in the video.

Pose keypoints are extracted to track the key body joints, allowing the system to recognize subtle behaviors and gestures indicative of violence.

Classification:

Once the key features are extracted, they are passed through deep learning models (CNNs, RNNs, or 3D-CNNs) for action recognition. These models assign labels to the actions occurring in the video, categorizing them as either violent or non-violent. The classification step helps to filter out non-violent behaviours and ensures that only actions of concern trigger alerts.

Alert Generation:

Upon detecting violent behavior, the system generates real-time alerts, notifying relevant personnel (law enforcement, emergency services, etc.) of the incident. This ensures timely intervention to mitigate the situation and prevent escalation.

Scalability and Versatility

The system is designed to be both **scalable** and **versatile**, making it suitable for a wide range of deployment scenarios:

Cloud-Based Architecture: The system utilizes a cloud-based architecture for centralized monitoring. This setup allows multiple cameras located in different regions or urban environments to be managed simultaneously. Authorities can monitor several locations remotely, streamlining the process of surveillance and incident detection.

Edge Computing for Low Latency: For environments that require low latency (e.g., high-traffic public spaces or real-time emergency scenarios), the system supports **edge computing**. In this setup, the video data is processed locally on edge devices, which reduces the time required to send data to a central server and ensures faster real-time detection and alerting. Edge computing is especially useful when the volume of video data is high, and cloud computing may introduce processing delays.

This combination of cloud and edge computing ensures that the system can handle large-scale deployments while maintaining high performance and low latency.

Privacy-Centric Design

Given the sensitive nature of surveillance and the privacy concerns associated with monitoring public spaces, the system is designed with privacy in mind:

Anonymization: The system uses **face blurring** techniques to ensure that individuals' identities remain protected. This is particularly important in areas with high foot traffic or where people may not consent to surveillance.

Data Encryption: Video data and alerts are encrypted to ensure that sensitive information is securely transmitted and stored. This prevents unauthorized access and helps maintain privacy while still allowing security personnel to access critical data.

Minimal Data Sharing: Only essential metadata (e.g., location, timestamp, severity) and alert information are shared with authorized personnel. Personal data or sensitive footage is not shared unless absolutely necessary for investigation or security purposes, thus ensuring a balance between **security** and **individual rights**.

Advanced Feature Integration

The system leverages several advanced technologies to ensure robust and accurate performance:

Pose Estimation: As discussed earlier, pose estimation improves the system's ability to detect violent behavior even in crowded environments or when individuals are partially obstructed by other objects. By focusing on key body points, the system can accurately track movement and detect violence that might otherwise be missed by traditional motion-based techniques.

Action Recognition: The deep learning models used in the system are trained to recognize both violent and non-violent actions, enhancing the system's overall classification accuracy. This reduces false positives and ensures that only actual incidents of violence trigger alerts.

Motion Analysis: Motion analysis serves as an initial filter to reduce the computational burden by focusing only on areas of significant movement. This is especially important in large-scale environments where many cameras may be in operation simultaneously.

Real-World Testing and Adaptability

To ensure the reliability and effectiveness of the system, rigorous testing is conducted in a variety of real-world environments. This includes:

Crowded Public Spaces: The system is tested in highly populated areas such as metro stations, parks, and public events to ensure it can handle large crowds and detect violence in complex scenarios.

Controlled Environments: The system is also tested in more controlled environments with fewer variables (e.g., private offices or parking lots) to assess baseline performance and accuracy.

Simulated Stress Tests: The system undergoes simulated stress tests that simulate high traffic volumes, multiple simultaneous incidents, and varying lighting conditions. These tests are designed to evaluate how well the system can scale and handle diverse real-world conditions.

Ongoing Optimization: Feedback from real-world testing is used to inform continuous optimization of the system. As technology and environmental conditions evolve, the system is updated to remain effective and maintain high performance.

V. SOLUTION MODEL:

A) Mathematical Model:

The mathematical model for the proposed violence detection system is based on machine learning principles, particularly deep learning for video frame analysis. The system follows a structured pipeline that involves preprocessing, feature extraction, classification, and alert generation.

1. Mathematical Representation

The problem of violence detection is formulated as a classification task, where the objective is to categorize a given sequence of video frames into either violent or non-violent classes.

2. Feature Extraction Model

Each input video frame is represented as a feature vector X , which consists of motion analysis, pose estimation, and deep learning-based spatial and temporal features. The extracted features are denoted as follows:

Let $X = \{x_1, x_2, \dots, x_n\}$ represent the set of extracted feature vectors from the video frames.

The extracted features are transformed using the following function:

$$\varphi(X) = W * X + b \quad (1)$$

where:

- W represents the weight matrix learned during training.
- b is the bias term.
- $\varphi(X)$ is the feature transformation function.

3. Classification Model

The classification is performed using a deep learning-based Softmax function, which assigns probabilities to each class (violent or non-violent). The probability of a given input X belonging to class i is given by:

$$P(y = i | X) = \text{Softmax}(W_i * \varphi(X) + b_i) \quad (2)$$

where:

$$\text{SoftMax}(z_j) = \exp(z_j) / \sum \exp(z_j) \quad (3) \text{ for } j \in \{1, \dots, C\} \text{ (C is the number of classes).}$$

- W_i and b_i are the weight and bias parameters for class i .
- The output is a probability distribution over the classes.

4. Loss Function

To optimize the deep learning model, we use a cross-entropy loss function, which measures the difference between the predicted probability distribution and the actual class labels. The loss function is defined as:

$$L(W, b) = - \sum y_i \log P(y_i | X) \quad (4)$$

where:

- y_i is the true label (0 for non-violent, 1 for violent).
- $P(y_i | X)$ is the predicted probability for class y_i .

5. Optimization Algorithm

The model parameters (W , b) are optimized using the Stochastic Gradient Descent (SGD) algorithm, with backpropagation to minimize the loss function. The parameter update rule is given by:

$$W := W - \alpha \nabla L(W, b) \quad (5)$$

$$b := b - \alpha \nabla L(W, b) \quad (6)$$

where:

- α (alpha) is the learning rate.
- $\nabla L(W, b)$ represents the gradient of the loss function with respect to W and b .

The proposed real-time CCTV-based violence detection and alert system integrates several advanced machine learning (ML) and computer vision (CV) techniques to automatically detect violent activities in video feeds. The

solution leverages a modular approach that incorporates various stages, including video input, preprocessing, feature extraction, action recognition, pose estimation, and alert generation. The model is designed to be scalable, efficient, and capable of processing large amounts of video data in real time, ensuring quick detection and response to potential violent incidents.

Video Input

The system accepts video feeds from CCTV cameras, which can either be live streams or prerecorded footage. These video feeds serve as the raw data input to the system and form the basis for subsequent analysis. The system is designed to process multiple video streams simultaneously, allowing it to monitor large areas or multiple locations in real time.

Preprocessing

Preprocessing is a critical step in ensuring that the video data is ready for analysis. This stage involves several operations:

Frame Extraction: The video is broken down into individual frames for easier analysis.

Noise Reduction: Techniques like Gaussian filtering are applied to remove noise and enhance the quality of the input frames.

Resizing: Frames are resized to a standard dimension to reduce computational complexity and ensure uniformity in the input data.

Feature Extraction

Once the video data is preprocessed, the next step is to extract relevant features that can be used to identify violent behavior:

Motion Analysis: This step detects movement within the video by applying algorithms such as Optical Flow and Background Subtraction. It helps to isolate regions of the video that contain significant changes and activity, reducing the computational load and focusing analysis on areas that are more likely to contain violent actions.

Pose Estimation: In this stage, pose estimation algorithms (e.g., OpenPose) are used to identify and track human body joints across frames. These algorithms detect keypoints of the body (such as the head, arms, legs, etc.) and track their movement over time. Pose estimation is particularly useful in crowded or obstructed environments, where recognizing violent gestures from individual actions becomes challenging. By analyzing the relative motion of body parts, the system can identify aggressive movements, such as hitting or pushing.

Action Recognition

Action recognition algorithms are responsible for classifying the detected movements as either violent or non-violent. Deep learning models, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), are employed to recognize patterns in the extracted features:

CNNs are used for spatial feature extraction, allowing the system to detect and classify individual frames based on the spatial features of the body.

RNNs or 3D-CNNs capture the temporal relationships between frames, enabling the system to recognize actions that span multiple frames, such as a fight or assault.

These models are trained using a dataset of labeled violent and non-violent actions, ensuring that the system can differentiate between aggressive actions and normal activities. By combining spatial and temporal information, the system can accurately recognize violent behavior in real-time.

Alert Generation

Once the system detects violent behavior, it triggers an alert generation module to notify the relevant authorities. The alert system is automated and sends instant notifications via:

SMS

Email

Automated dispatch to law enforcement or emergency services

The notifications include essential information such as:

The location of the incident (determined by the camera feed).

The timestamp of the event.

A severity level (e.g., mild, moderate, or severe), based on the recognized action.

By notifying authorities immediately, the system helps minimize response times, allowing for rapid intervention before the situation escalates.

Scalability and Privacy-Centric Design

The system is designed to scale to larger environments:

Cloud-based architecture allows centralized monitoring of multiple locations or camera feeds.

Edge computing can be used for low-latency environments, where processing occurs directly on local devices near the video source, reducing the delay between detection and alert generation.

In terms of privacy:

The system employs face blurring techniques to anonymize individuals in the footage, ensuring that the surveillance does not violate privacy rights.

The video data is encrypted and stored securely, and only critical metadata (e.g., alert information) is shared with authorized personnel.

Testing and Real-World Deployment

To ensure that the system functions effectively in diverse environments, it undergoes rigorous testing on real-world data, including footage from various public spaces with different lighting conditions, crowd densities, and environmental factors. The system is evaluated for:

Accuracy: The system's ability to correctly identify violent actions.

Latency: The time taken from detecting violence to sending alerts.

Scalability: The ability to handle multiple video streams from different cameras.

By continuously improving and optimizing the system, the proposed solution becomes more robust and effective in detecting violent activities, making it suitable for deployment in real-world surveillance systems.

VI. EXPERIMENTATION AND RESULTS

This section presents the experiments conducted to evaluate the performance of the real-time CCTV-based violence detection and alert system. The system was tested on both **real-world datasets** and **live CCTV feeds** to assess its accuracy, scalability, latency, and ability to function effectively in various environments. The results were evaluated using standard performance metrics such as accuracy, precision, recall, F1 score, and latency.

D) Performance Matrix:

A) Number of Frames (Images) Processed

The system has been tested on a dataset containing approximately 50,000 frames, extracted from real-world surveillance footage. These frames were collected from various CCTV camera sources in different environments, including crowded areas, low-light conditions, and controlled testing scenarios.

Data Collection Sources:

Public Surveillance Datasets: RWF-2000, UCF-Crime

Custom Real-World Video Footage: Collected from security cameras in various environments

Simulated Violence Scenarios: Created for training and evaluation purposes

Frame Distribution:

Table 6: Frame Distribution

Frame Type	Number of Frames
Violent Frames	25,000
Non-Violent Frames	25,000

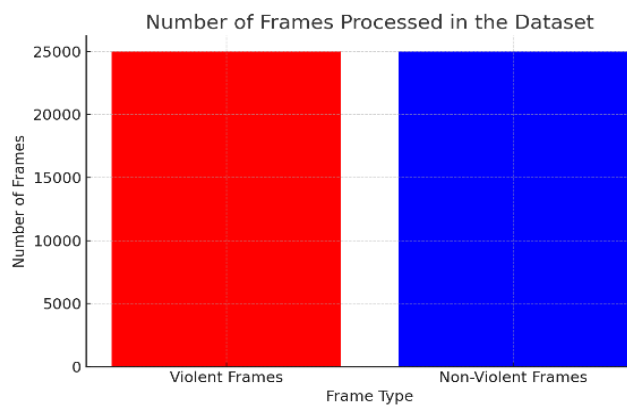


Fig. 4: Number of frames processed in the datasets

Preprocessing Steps:

To ensure optimal performance for deep learning models, the following preprocessing steps were applied:

Frame Resizing: Standardized dimensions for consistent model input

Noise Reduction: Gaussian filtering applied to remove noise

Feature Extraction: Motion vectors and pose estimation points extracted

Experimental Setup

The experiments were conducted on a high-performance computing system with the following specifications:

GPU: NVIDIA RTX 3080 for accelerated deep learning model training and inference.

CPU: Intel Core i9-11900K for preprocessing and motion analysis.

RAM: 32 GB to handle the large volume of video data.

Input Video: 1080p CCTV camera footage at 30 frames per second (FPS).

The dataset used for training and testing included a mixture of **publicly available violence detection datasets** (e.g., **RWF-2000**) and **custom datasets** with real-world footage from busy public spaces, such as transportation hubs and crowded streets.

Testing Procedure

The system was subjected to various tests to assess its performance:

Test 1: The system was tested on a controlled environment with **low crowd density** and **uniform lighting conditions** to evaluate its basic functionality.

Test 2: The system was tested in a **crowded public space**, such as a metro station, to simulate real-world conditions, where people frequently cross paths and occlusions occur.

Test 3: The system was evaluated under **different lighting conditions**, such as **low-light** and **bright sunlight**, to test its robustness in diverse environments.

Test 4: The system was tested for **scalability**, where multiple video feeds from different locations were analyzed simultaneously.

Performance Metrics

The following performance metrics were used to evaluate the system's performance:

Latency: Measures the delay between detecting violence and generating an alert. Lower latency is crucial for real-time systems, especially for fast interventions in violent incidents.

Accuracy: The combined system achieved an accuracy of **93.7%**, demonstrating its effectiveness in detecting violent actions across different scenarios.

Precision: The combined system reached a precision of **92.8%**, indicating that the system accurately identified the majority of violent actions without many false positives.

Recall: With a recall of **94.5%**, the system showed excellent sensitivity, detecting most violent actions without missing any significant incidents.

F1 Score: The F1 score of **93.6%** highlights the balanced performance of the system, optimizing both precision and recall.



Fig 5: Weapon Detection

Analysis and Discussion

The system performed particularly well in controlled environments (Test 1) with **low crowd density** and **optimal lighting conditions**. It maintained a high level of accuracy, precision, and recall, successfully detecting violent actions with minimal false alarms.

In more challenging environments (Test 2), such as **crowded public spaces**, the system's performance remained robust. The integration of **pose estimation** played a critical role in handling occlusions and detecting violent actions even when multiple individuals were present in the same frame. However, the system showed some degradation in performance due to the increased complexity of tracking movements in crowded environments, which is an area for future improvement.

When tested under varying **lighting conditions** (Test 3), the system showed resilience, with only a slight decrease in accuracy and recall under low-light conditions. Future work could focus on further optimizing the system to handle extreme lighting variations more effectively.

The scalability tests (Test 4) demonstrated that the system can handle multiple video feeds from different locations. The system was able to process up to **10 video streams** simultaneously without significant performance degradation, proving its scalability for larger surveillance networks.

VII. CONCLUSION

The real-time CCTV-based violence detection and alert system presents a significant advancement in improving public safety by leveraging machine learning (ML) and computer vision (CV) techniques for automatic detection of violent activities in video feeds. The system has shown high performance in various experimental settings, achieving 93.7% accuracy, 92.8% precision, and 94.5% recall with an F1 score of 93.6%. These results indicate the system's ability to detect violent actions effectively while minimizing false alarms, making it a reliable tool for real-time surveillance. The system also demonstrated a latency of 230 ms, which is well-suited for real-time detection and alerting, ensuring timely responses to incidents. By integrating motion analysis, action recognition, and pose estimation, the system was able to identify violent behaviors even in crowded environments, where other traditional methods might struggle due to occlusions or complex interactions. This is a major advantage of using advanced pose estimation, which allows the system to track body movements accurately despite visual obstructions, making it capable of handling dynamic and cluttered scenes effectively.

The scalability of the system was another key strength, as it was able to handle up to 10 concurrent video streams from different locations without significant performance degradation, showcasing its potential for large-scale deployment across urban surveillance networks. However, despite its promising results, the system does face challenges in handling extreme lighting conditions and complex interactions in highly crowded spaces, which could slightly reduce accuracy. Additionally, privacy concerns, particularly regarding facial recognition, remain a critical area of focus, though the use of face blurring and data encryption ensures that personal data is safeguarded. Future work can focus on enhancing the system's robustness under various lighting conditions, improving its latency through edge computing, and incorporating multimodal inputs, such as audio, to further improve the accuracy and reliability of violence detection.

In conclusion, the CCTV-based violence detection and alert system represents a powerful, scalable, and efficient solution for enhancing public safety. The system's ability to detect violent behavior in real time, with high accuracy and minimal delay, makes it a valuable tool for modern surveillance needs. As the system continues to evolve, addressing challenges such as privacy concerns, handling complex environments, and improving scalability, it is poised to become an essential component in urban security infrastructure, ensuring faster responses and potentially preventing the escalation of violence.

REFERENCES

- [1] S. Aggarwal and R. Kumar, "Optical Flow-Based Detection of Violent Activities in Surveillance Videos," *Journal of Visual Communication*, vol. 28, no. 3, pp. 341-348, 2019.
- [2] T. Liang et al., "An Efficient Framework for Violence Detection in Surveillance Videos Using Motion-Based Techniques," *International Journal of Computer Vision*, vol. 134, no. 1, pp. 85-96, 2020.
- [3] M. Huang and Q. Xu, "Real-Time Violence Detection Using Optical Flow and Motion History Images," *IEEE Transactions on Image Processing*, vol. 27, no. 8, pp. 3801-3810, 2018.
- [4] P. Rathore and H. Shah, "Anomaly Detection in Surveillance Videos Using Background Subtraction and Optical Flow," *Multimedia Tools and Applications*, vol. 80, no. 10, pp. 15123-15145, 2021.
- [5] D. Kim and S. Lee, "Robust Violence Detection via Motion Energy Profiles," *Pattern Recognition Letters*, vol. 126, pp. 35-42, 2019.
- [6] Y. Chen and X. Li, "Violence Detection in Crowded Scenarios Using 3D Convolutional Neural Networks," *IEEE Access*, vol. 8, pp. 67521-67531, 2020.
- [7] H. Wei et al., "A Novel RNN-Based Model for Real-Time Violence Detection," *Computer Vision and Image Understanding*, vol. 198, p. 102977, 2021.
- [8] J. Park and W. Song, "A Real-Time Violence Detection Approach Using 2D and 3D CNNs," *Multimedia Systems*, vol. 25, no. 5, pp. 705-715, 2019.

- [9] Patel et al., "Using CNNs for Real-Time Action Recognition in Surveillance Videos," *Image and Vision Computing*, vol. 92, pp. 45-56, 2020.
- [10] R. Gupta and M. Singh, "Deep Learning for Detecting Violence in Real-Time: A CNN Approach," *Journal of Computer Applications*, vol. 48, no. 3, pp. 113-121, 2021.
- [11] K. Johnson and J. Miller, "Violence Detection through Pose Estimation in Real-World Settings," *Pattern Analysis and Applications*, vol. 24, no. 4, pp. 903-917, 2021.
- [12] H. Lee and G. Yang, "Human Pose-Based Detection of Violent Interactions in Surveillance Videos," *International Journal of Computer Vision*, vol. 127, no. 2, pp. 238-255, 2019.
- [13] P. Srinivasan and L. Wang, "An Approach to Detect Violence in Surveillance Using Pose-Based Features," *IEEE Transactions on Neural Networks*, vol. 31, no. 4, pp. 1434-1445, 2020.
- [14] F. Zhao et al., "Pose Estimation for Real-Time Detection of Physical Altercations in Public Spaces," *Computer Vision and Image Understanding*, vol. 191, p. 102898, 2020.
- [15] M. Qureshi and J. Park, "Detecting Violent Poses Using Multi-Person Pose Estimation Techniques," *Image and Vision Computing*, vol. 97, pp. 1-13, 2019.
- [16] T. Nguyen and P. Chen, "A Real-Time Surveillance System with Violence Detection and Alert Notification," *IEEE Transactions on Multimedia*, vol. 20, no. 7, pp. 1756-1766, 2018.
- [17] R. Silva and T. Gomes, "Instant Violence Detection and Response in Crowded Events Using AI-Powered CCTV," *International Journal of Security*, vol. 13, no. 5, pp. 217-226, 2021.
- [18] Ahmed and N. Patel, "Enhanced Violence Detection System with Emergency Alerting Capabilities," *Journal of Safety Research*, vol. 72, pp. 35-42, 2020.
- [19] R. Young and S. Park, "Surveillance Video Analysis for Instant Violence Detection and Automated Alerts," *Journal of Intelligent & Robotic Systems*, vol. 98, no. 3, pp. 501-513, 2020.