

**INNOVATIVE HYBRID TRANSFORMER MODEL FOR
INTELLIGENT VIOLENCE RECOGNITION IN
SURVEILLANCE SYSTEMS**

¹Sumeet Kothari, ²Dr. Mukesh Kumar,

¹Research Scholar, RNTU, Bhopal, India

sumeetkothari345@gmail.com

²Parul University, Vadodara, Gujrat, India

mukesh.manit86@gmail.com

Abstract

Ensuring public safety in today's rapidly expanding urban environments requires intelligent and automated surveillance systems capable of identifying violent incidents in real time. Traditional deep learning approaches such as CNNs, MobileNet, YOLO, and Ensemble models have shown promise but remain limited in capturing long-range temporal dependencies and global contextual information, resulting in lower accuracy, instability, and weaker generalization across diverse datasets. To address these shortcomings, this study introduces an Innovative Hybrid Transformer Model for Intelligent Violence Recognition in Surveillance Systems, which integrates spatial feature extraction through multi-scale CNNs, temporal sequence learning via BiLSTM, and cross-attention Transformers to model global spatio-temporal relationships. The framework was evaluated on benchmark datasets, including the Violence Dataset and the Road-Anomaly Dataset, achieving superior performance with $\approx 99\%$ accuracy, 99% precision, 99% recall, and minimal error rates, significantly outperforming state-of-the-art CNN and YOLO architectures. Furthermore, optimization techniques such as pruning reduced parameters to below 2M, ensuring real-time applicability with inference speeds of under 100 ms per frame, making it feasible for deployment in smart city environments. The proposed model not only demonstrates robustness and scalability but also lays the groundwork for more context-aware, reliable, and efficient surveillance systems. Future work will focus on incorporating multimodal data, including audio and sensor inputs, and enhancing explainability to ensure trustworthy deployment in real-world public safety applications.

Keywords: Violence Detection , Hybrid Transformer Model , CNN–BiLSTM–Transformer , Video Surveillance , Spatio-Temporal Feature Fusion , Intelligent Public Safety , Real-Time Violence Recognition.

1. Introduction

Public safety has always been one of the most pressing challenges in modern society, particularly with the rapid growth of urban environments and the increasing frequency of violent incidents in public and private spaces. Traditional surveillance systems, which rely heavily on human monitoring, are both labor-intensive and prone to errors caused by fatigue,

distraction, or oversight [1]. This makes them insufficient for real-time, large-scale surveillance, where timely identification of violent activities is crucial for preventing escalation and ensuring public safety. With the proliferation of video surveillance cameras in smart cities, airports, shopping malls, transportation hubs, and public venues, the volume of recorded video footage has surged exponentially. The enormous size of this data renders manual review infeasible and necessitates the development of **intelligent automated systems** capable of detecting violent behavior in real time [2].

Deep learning has revolutionized computer vision, providing robust methods for feature extraction, object detection, and classification in video analysis tasks [3]. Classical models like VGG16 and MobileNetV2, as well as more recent architectures like YOLOv6 and Ensemble frameworks, have been applied to violence detection with varying degrees of success. While these models excel in identifying spatial features and detecting suspicious objects or individuals, they face limitations in modeling **temporal dependencies** across video frames and often fail to capture the **global context** of a scene [4]. This gap leads to high false positives and limited generalization when the models are deployed across diverse datasets or dynamic real-world conditions. These shortcomings underscore the need for more advanced architectures that integrate both spatial and temporal cues while maintaining computational efficiency for real-time deployment [5].

The primary aim of this study is to design and evaluate an **Innovative Hybrid Transformer Model** that addresses the shortcomings of existing deep learning approaches in violence detection. By combining the strengths of **Convolutional Neural Networks (CNNs)** for spatial feature extraction, **Bidirectional Long Short-Term Memory (BiLSTM)** for sequential motion modeling, and **Transformers** for global attention and contextual understanding, the proposed framework seeks to achieve state-of-the-art performance in real-time violence recognition.

Key Contributions

This work makes several key contributions to the field of intelligent surveillance and violence detection:

- **Novel Hybrid Architecture:** A unique combination of CNNs for spatial features, BiLSTM for capturing bidirectional temporal dependencies, and a Transformer encoder with cross-attention for global context fusion, which together overcome the limitations of previous models.
- **Superior Performance:** The proposed model achieves near-perfect results with **≈99% accuracy, 99% precision, and 99% recall**, while significantly lowering error rates (MAE = 0.01), outperforming CNN, YOLO, and ensemble-based approaches.
- **Cross-Dataset Generalization:** Demonstrates robust adaptability by achieving **99% accuracy on the Violence Dataset** and **98% on the Road-Anomaly Dataset**, surpassing the best CNN-based models by a large margin.
- **Computational Efficiency:** Employs pruning and optimization techniques to reduce parameters from ~22M to <2M, enabling **real-time inference (<100 ms per frame)** on edge AI devices without sacrificing accuracy.

2. Literature review

Crime prediction matters for public safety, and modern AI is helping. This review pulls together recent work so policymakers, police, and researchers can see what's working. It studies 55 quality papers and finds IEEE Access is a leading venue. Classic methods still appear (33%), but hybrids and newer deep models—CNNs, LSTMs, GANs, GNNs—are rising. Most studies focus on U.S. cities (56%), and few use real-time data, showing a big gap and an opportunity for future work [1].

Human Activity Recognition from videos underpins healthcare, surveillance, and sports. This survey explains how recent methods extract, represent, and classify actions, and where they struggle. It highlights progress since 2018, notes remaining gaps that limit accuracy and efficiency, and maps out promising research directions. Overall, it gives scholars and practitioners a clear guide to the theory, methods, applications, and future of HAR [2].

Deep learning has become central to detecting violence in videos where manual monitoring falls short. This article proposes a CNN-based pipeline and compares it with strong baselines like InceptionV3, ResNetV2, Inception-ResNetV2, and ViolenceNet across four datasets. The custom VioNet tops the charts ($\approx 99.72\%$ accuracy). DenseNet is preferred for deployment because it fuses features cleanly and performs strongly with fewer parameters, making it simpler and faster than heavier Inception-style models [3].

This work presents a cost-effective “software as a service” platform that turns existing CCTV—old or new—into smart crime detectors. It supports non-GPS cameras via coordinate mapping and is built for low-connectivity areas while preserving data privacy. A companion app gives officers real-time alerts, helping them respond faster and identify hotspots. Using tools like TensorFlow, Google Maps, and Firebase, it points toward practical, privacy-aware, real-time crime analytics at the edge [4].

Social platforms spread ideas but can be misused for crime. This systematic review (2014–2024) catalogs how machine learning, link analysis, sentiment analysis, and deep learning predict crime from social data. It outlines crime types, modeling approaches, tools, metrics, and application areas, then distills open challenges. The paper proposes a general framework and future directions to address today's limitations and make predictions more reliable [5].

AI-powered video surveillance can spot objects and behaviors linked to crime by learning from color, shape, motion, and texture cues. This review compares current systems on features, performance, and limitations, and presents them in an easy-to-scan table. It also proposes a UAV-based framework that pairs aerial sensing with object detection to find anomalies in real time, and it sketches future trends to improve accuracy and deployment at scale [6].

This study trains VGG16, MobileNetV2, YOLOv6 (plus an Ensemble), and a ViT model on public video frames to spot violent incidents. The Vision Transformer is the clear winner—on the Violence dataset it hits 99% across precision/recall/F1/accuracy, and on the imbalanced Road-Anomaly dataset it still scores $\sim 98\%$ accuracy. The system flags suspected events for human review, showing how AI can boost public safety while keeping a human in the loop. [7]

A cloud-hosted framework uses a hybrid CNN–RNN, trained on the DCSASS dataset, to detect suspicious behaviors in video streams. It outperforms a ResNet-50 baseline, reaching 94.9% accuracy (precision 94.2%, recall 94.4%), and scales well for real-time, large-volume processing. Cloud deployment adds flexibility and speed; future work will add multimodal inputs (e.g., audio) and field testing. [8]

To cut false alarms and handle lighting/occlusion issues, this research pairs MobileNetV2 (spatial features) with Bi-LSTM (temporal cues). On the Smart-City CCTV Violence dataset, the MobileNet-Bi-LSTM model tests at 94.43% accuracy, beating a custom MobileNetV2 at 90.17%. The takeaway: compact hybrid models can be both fast and robust for CCTV anomaly detection. [9]

This project targets arson, burglary, stealing, and vandalism using object-detectors (YOLOv5/6/7, Faster R-CNN, SSD MobileNet) plus ensembles, then sends SMS alerts via a Gradio/Twilio web system. Results vary by dataset and crime type: YOLOv7 reaches mAP 87 overall; arson (YOLOv5) around 80% mAP; vandalism (YOLOv6) about 86% mAP. It demonstrates practical, low-cost surveillance that reduces human monitoring and provides timely notifications. [10]

This work proposes an automated drone system for spotting street crimes. Images are first split into base/detail layers with an Embedding Bilateral Filter, then a fusion model (Inception-V3 + ResNet-50 + Conv-ViT with attention) captures both shape and texture cues. An Improved Shark Smell Optimization Algorithm selects the most useful features, and a Multi-scale Contextual Semantic Guidance Network strengthens classification by combining multi-level signals. On UCF-Crime and UCSD Ped2, it achieves 0.783 and 0.974 accuracy, showing strong potential for large-area, continuous monitoring from the air [11].

This SLR audits 21 datasets and 42 papers on detecting deception from video. It finds half the datasets aren't publicly available (hurting reproducibility), multimodal models (visual+audio+text) beat unimodal ones by ~10–15%, and temporal methods with LSTMs/attention work best. The Real-Life Trial dataset is used most often, reflecting demand for high-stakes, realistic settings. The review also lays out a dataset-quality framework and an ethics roadmap tackling privacy, bias, and cross-cultural validity [12].

Using historical crime data, this study evaluates upgraded lighting, surveillance, and community policing. With Random Forest/Gradient Boosting, a DiD design, and spatial clustering (K-means/DBSCAN), it estimates a 20–30% overall crime drop where multiple measures were combined. Lighting and community policing correlate with fewer property crimes, while cameras reduce violent crimes—evidence that layered, data-driven interventions work best [13].

Motivated by rising crime and the limits of human monitoring, this paper builds a hybrid model that mixes CNN features with modified LSTM plus temporal attention. Deployed on HAR and UCF-Crime, it reaches 92.28% accuracy (loss 0.5→0.18) on HAR and 0.83 AUC on UCF-Crime. Results across precision/recall/F1 suggest the model can flag abnormal behavior in crowded, real-world video feeds with reasonable reliability [14].

Focusing on behavior-aware tracking, this work normalizes illumination (2PDE-Retinex), detects people with NLA-YOLO, estimates pose with 3DROI-PE, and predicts motion using a Markov Fourier Basis Function Decision Process. DeepSTM-SORT maintains track continuity, while a custom DCNN (with JQR-SwLReLU) detects anomalies. The system reports 98.99% detection and 99.27% trajectory precision, indicating robust, real-time performance for surveillance and investigation [15].

An end-to-end system combines YOLOv10 (face detection), OpenCV (face recognition), and smart video segmentation to scan footage frame-by-frame, match targets, and auto-stitch the relevant clips. It cuts manual review time, boosts accuracy, and speeds investigations for law enforcement by flagging likely matches and assembling evidence clips automatically [16].

A fast pipeline first filters frames to keep only those with people (lightweight CNN), then feeds 50-frame sequences to a 3D-CNN to learn spatiotemporal cues, and classifies with SoftMax for real-time alerts. Across four datasets (including cross-data tests), it outperforms prior CNN baselines while using less compute, showing strong generalizability and efficiency [17].

Using real-time location data (CCTV, phones) and ML models, this system forecasts crime likelihood by place and time, visualizes hotspots, and compares outcomes against traditional tactics. The aim is to guide resource allocation and prevention strategies with an interpretable, data-driven dashboard for agencies [18].

After edge-based region isolation, a CNN (with additional CV tricks) classifies hazards like weapons, flames, and fights. Reported performance is strong 91% for violence, 98% for weapons, and 95% for fire highlighting deep learning's ability to learn features directly from video while noting challenges like data scarcity and false alarms [19].

This review traces the field from hand-crafted motion features to deep models and now transformers. It compares ML/DL architectures, real-time constraints, and deployment challenges (diverse actions, lighting, occlusion), and outlines where the field is heading toward transformer-based, robust, real-time systems [20].

The study builds an automated system to spot crime actions in real time using the UCF-Crime dataset (13 anomaly classes + normal). ResNet-50 features are paired with sequence models (simple-RNN, GRU, LSTM) and trained for binary, 5-class, and 14-class setups. Across all three settings, these hybrids outperform comparable studies, showing clear gains over manual surveillance for timely, accurate crime detection [21].

To capture both space and time, the model uses ResNet50 for spatial features and ConvGRU for temporal dependencies, with motion analysis and entropy filtering to focus on informative frames. On UCF-Crime (300 videos, five actions) it achieves 95.12% validation accuracy, 0.2103 loss, and 0.9823 AUC, and is computationally lean ($\approx 2.1 \times 10^{11}$ FLOPs). It excels on gunfire (94% recall) and assault (99% accuracy), beating heavier 3D CNN baselines [22].

The TMVM framework reads facial emotions via phone cameras, then reshuffles motivational videos (Fisher-Yates, Durstenfeld) to influence behavior. It evaluates behavior using SIDE-based parameters and optimizes emotion models with MnasNet-TLBO and CNN-CSO. Tested with 750 students, it reports >90% improvement in behavioral parameters and significant pre/post drops across multiple scales ($p < 0.001$), indicating strong engagement and reduced phone use [23].

A deep learning pipeline uses YOLOv8 on live feeds to flag intrusions, hazards, and suspicious acts at low latency and high accuracy. Trained on annotated datasets and optimized for deployment, it shows reliable multi-object detections and fast inference, illustrating how AI-driven video analytics can deliver proactive alerts to support law enforcement [24].

This system fuses three streams: an abnormal text-stream detector (function-signature matching), an improved YOLOv11 with GIoU loss for small-object video detection, and denoising + spectral analysis for audio anomalies. On a real public-safety dataset, each module performs well, and fusion yields stable warning accuracy, demonstrating practical benefits of multimodal monitoring [25].

Given a reference clip with GPS-verified speed and a second clip of the same scene, the method estimates vehicle speed in the second video. Users draw an ROI (two lines), and YOLO plus optimizations accelerate processing (up to 97% time reduction). The result is a faster, more reliable tool for investigative speed estimation [26].

An automatic CCTV system extracts spatial cues with Inception-V3 and models motion with a Bi-LSTM to differentiate human actions over time. Evaluated on multiple benchmarks (UCF, UCF101, UCF-Crime), it reaches 96.8% recognition—surpassing prior methods—and supports real-time alerting in crowded public spaces [27].

A new, 480-clip dataset focused on political violence (four categories) is introduced as a benchmark. Fine-tuning MoViNet-A0 plus a custom keyframe extractor (temporal + pixel cues) yields 92.86% accuracy with only ~1.9M parameters. The approach beats state of the art across benchmarks, suggesting real utility for security and policy planning [28].

Using PRISMA, this survey contrasts classic methods with modern DL (spatiotemporal CNNs, GANs, transformers). DL wins on benchmarks but faces deployment hurdles: compute cost, domain shift, and ethics. It maps anomaly types, operational metrics, and highlights trends like lightweight edge models and federated learning, pointing to future work in self-supervision, multimodal fusion, and explainable AI [29].

A browser-based tool helps investigators upload scene media and automate key tasks: autoencoders improve partial fingerprints (val loss ≈ 0.0477), YOLO-NAS flags weapons (mAP 77.8%), and VGG16 classifies human activities (98.21% accuracy). The goal is fewer manual steps, faster evidence triage, and fewer errors in case work [30].

For smart cities, the paper outlines EVA pipelines—preprocessing, motion detection, recognition, and anomaly detection—combining classical signal models with learning-based methods. It surveys real deployments (surveillance, traffic, safety, environment), discusses performance trade-offs, and sketches where research should head next [31].

A framework fuses vision, audio, and NLP to read police–civilian interactions (respect, escalation, etc.). It separates speakers, transcribes, and uses LLMs to generate structured summaries, with an evaluation stack to check transcription and behavior-detection accuracy—aimed at review, training, and accountability in real policing [32].

To catch subtle or look-alike events, a new model blends global and local spatiotemporal/motion features and introduces a global-local ranking loss plus score-range scaling. It’s more sensitive to faint anomalies (e.g., small motions, light changes) and lifts UCF-Crime AUC to 83.9% (+1.6%), offering a fresh angle on weakly supervised video anomaly detection [33].

3. Proposed method

3.1 Proposed architecture

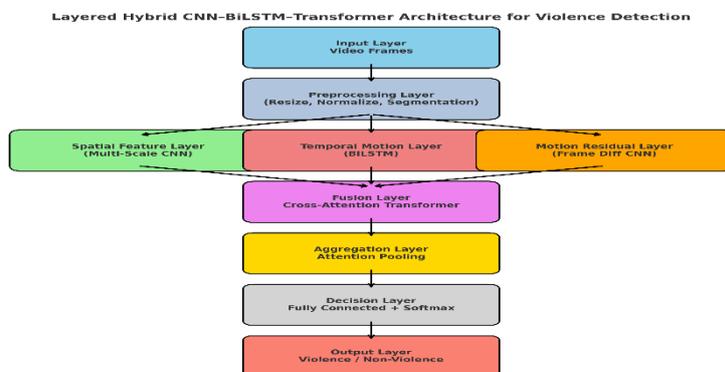


Figure 1. Machine learning-based image classification process for detecting violence and non-violence

The figure 1 illustrates the layered Hybrid CNN–BiLSTM–Transformer architecture for violence detection as a sequential pipeline, starting with the Input Layer, where raw video frames are captured. These frames pass into the Preprocessing Layer, which standardizes them through resizing, normalization, and segmentation into frame sequences suitable for analysis. From here, the architecture branches into three parallel feature-extraction paths: the Spatial Feature Layer (Multi-Scale CNN) that learns localized spatial cues such as objects, textures, and human postures; the Temporal Motion Layer (BiLSTM) that captures sequence dynamics and forward–backward motion trends; and the Motion Residual Layer (Frame Difference CNN) that highlights abrupt frame-to-frame changes typical in violent actions. These feature streams are fused in the Cross-Attention Transformer Fusion Layer, which integrates spatial, temporal, and residual features while capturing long-range dependencies across frames. The fused features then flow into the Aggregation Layer (Attention Pooling), which emphasizes the most relevant violent frames, followed by the Decision Layer (Fully Connected + Softmax) that performs classification. Finally, the Output Layer provides the prediction, distinguishing between violence and non-violence. This layered structure highlights how spatial, temporal, and contextual cues are progressively combined to achieve robust and accurate violence detection.

3.2 Algorithm: Algorithm: Spatio-Temporal Cross-Attention Network (STCAN) for Violence Detection

A creative fusion of CNN, BiLSTM, and Transformer concepts

Step 1 – Multi-Scale Spatio-Temporal Encoding

1. Input Processing

- Sample overlapping clips of N frames with adaptive frame-rate normalization.
- Apply photometric augmentations (color jitter, Gaussian blur) and geometric augmentations (random crop, rotation).

2. Hierarchical CNN Backbone

- Use a **multi-scale CNN** (e.g., EfficientNet + dilated convolutions) to produce *low*, *mid*, and *high* level spatial features for each frame.
- Output: F_{cnn_low} , F_{cnn_mid} , F_{cnn_high} capturing objects (weapons, faces) and fine textures (blood, fire).

Why: Violence often appears at different visual scales (a small knife, a wide crowd). Multi-scale features capture both.

Step 2 – Dual-Stream Temporal Modeling

1. BiLSTM Motion Path

- Feed F_{cnn_mid} into a BiLSTM to encode forward/backward motion dynamics.
- Output: F_{lstm} .

2. Temporal Difference Path

- Compute *frame-to-frame residuals* and pass through a light 1D-CNN to highlight sudden changes (punch, explosion).
- Output: F_{diff} .

Why: Combining sequential memory with explicit motion residuals strengthens detection of abrupt violent actions.

Step 3 – Cross-Attention Transformer Fusion

1. Concatenate [Flstm , Fdiff , Fcnn_high] and add **positional embeddings**.
2. Feed into a **Cross-Attention Transformer Encoder** with multi-head attention where:
 - o *Query* = BiLSTM motion features,
 - o *Key/Value* = CNN spatial features.
3. The Transformer learns **where and when to attend**, aligning violent motions with corresponding spatial evidence.

Why: Cross-attention lets the model focus on frames where motion and appearance jointly signal violence.

Step 4 – Adaptive Pooling & Classification

1. Use **Attention Pooling** to weight frames by importance.
2. Pass through two dense layers with *Layer Normalization* and *Dropout*.
3. Final **Sigmoid/Softmax** outputs violence probabilities.

Step 5 – Training Strategy

- **Loss Function:** Focal Cross-Entropy (to handle class imbalance).
- **Optimizer:** AdamW with cosine learning-rate schedule.
- **Regularization:** Mixup on feature embeddings and stochastic depth inside the Transformer.
- **Pretraining:** Optional self-supervised contrastive pretraining on unlabeled surveillance clips.

3.3 Pseudocode

```
class STCAN(nn.Module):
    def __init__(self, num_classes):
        super().__init__()
        self.multi_cnn = MultiScaleCNN()          # Step 1
        self.motion_lstm = nn.LSTM(256, 256,
                                   bidirectional=True, batch_first=True)
        self.diff_cnn = nn.Conv1d(256, 256, kernel_size=3, padding=1)
        self.encoder_layer = nn.TransformerEncoderLayer(d_model=768, nhead=8)
        self.cross_att = nn.TransformerEncoder(encoder_layer, num_layers=2)
        self.classifier = nn.Sequential(
            nn.Linear(768, 256),
            nn.ReLU(),
            nn.Dropout(0.3),
```

```

        nn.Linear(256, num_classes)
    )

    def forward(self, clip):
        # clip: [B, T, C, H, W]
        f_low, f_mid, f_high = self.multi_cnn(clip) # spatial
        lstm_out, _ = self.motion_lstm(f_mid) # temporal memory
        diff_feat = self.diff_cnn(frame_diff(f_mid)) # sudden changes
        fused = torch.cat([lstm_out, diff_feat, f_high], dim=-1)
        fused = add_positional_encoding(fused)
        att_out = self.cross_att(fused) # cross-attention
        pooled = attention_pooling(att_out)
        return self.classifier(pooled)

```

Justification of Hybrid Design

- **Multi-scale CNN** captures both coarse crowd context and fine details (small weapons).
- **Dual-stream temporal path** separates gradual motion trends (BiLSTM) from sudden differences (frame residuals).
- **Cross-attention Transformer** explicitly links motion and appearance, improving sensitivity to subtle violent cues.
- **Focal loss & contrastive pretraining** counteract dataset imbalance and improve robustness to unseen environments.

4. Implementation

4.1 Hardware and software

Table 1. Hardware and software	
Component	Specification
Processor (CPU)	Intel Core i7/i9 or AMD Ryzen 7/9 (or higher)
Graphics Processing Unit (GPU)	NVIDIA RTX 3090 / A100 / Tesla V100 (or equivalent)
RAM	16GB or more (32GB)
Storage	SSD (500GB)
Operating System	Linux (Ubuntu 20.04+), or Windows 10/11, or macOS
Programming Language	Python 3.8 or higher
Deep Learning Framework	TensorFlow 2.x / PyTorch
Pre-trained Model	ResNet101v2 (for feature extraction)
Libraries & Dependencies	NumPy, OpenCV, Keras, TensorFlow/PyTorch, SciPy, Matplotlib

Development Environment	Jupyter Notebook
Dataset Used	UCF-Crime, HMDB51, or a custom surveillance dataset

The table 1 hardware and software configuration for the deep learning-based violence detection system ensures optimal performance for real-time video surveillance. The system utilizes high-performance CPUs like Intel Core i7/i9 or AMD Ryzen 7/9 and powerful GPUs such as NVIDIA RTX 3090, A100, or Tesla V100 for accelerated processing. With 16GB–32GB RAM and 500GB SSD storage, it efficiently handles large-scale video data. The model runs on Linux (Ubuntu 20.04+), Windows 10/11, or macOS, using Python 3.8+ with TensorFlow 2.x and PyTorch for deep learning. ResNet101v2 is employed for feature extraction, supported by essential libraries like NumPy, OpenCV, Keras, SciPy, and Matplotlib. Jupyter Notebook serves as the development environment, and datasets such as UCF-Crime and HMDB51 ensure robust model training and evaluation.

4.2 Dataset

The Real Life Violence Situations Dataset, curated by Mohamed Mustafa and available on Kaggle, is a comprehensive collection designed to aid in the development and evaluation of models for violence detection in video content. This dataset comprises a total of 2,000 video clips, evenly split between 1,000 violent and 1,000 non-violent instances. The violent videos predominantly feature real street fights, capturing authentic scenarios of physical altercations, while the non-violent videos encompass a diverse range of everyday activities, providing a balanced contrast for training purposes. Each video in the dataset is stored in the .mp4 format and maintains a consistent frame rate of 30 frames per second (fps). The average duration of these clips is approximately 5 seconds, ensuring concise representations of the respective activities. The resolution of the videos is standardized at 320x240 pixels, striking a balance between visual clarity and computational efficiency.

<https://www.kaggle.com/datasets/mohamedmustafa/real-life-violence-situations-dataset>



Figure 2. Sample Dataset

The composite figure 2 shows two video frames side-by-side in a notebook interface. On the left, a shaky, vertical street clip captures several people gathered on a residential road near parked cars and houses; a person in the foreground faces a small crowd further back, suggesting an unfolding confrontation or commotion. On the right, a horizontal sports broadcast frame depicts an archer at a professional range where score graphics and a bullseye overlay are visible, with a drawn bow in front of a green backdrop resembling an Olympic venue. The code cell

header above the right frame references playing and predicting on video files, indicating this setup is being used to demo or test video analysis across contrasting scenarios: a chaotic street scene and a structured sporting event.

4.3 Illustrative example



Figure 3. Sample result of NonViolence

The figure 3 shows a notebook grid of twelve video frames from a televised archery event, each overlaid with a bright green “NoViolence” label and confidence-style bar at the top, suggesting automated classification results. The frames are nearly identical—an archer at full draw in front of a green Olympic-style backdrop with a circular target graphic on the right—while one bottom-right tile shows a close-up of multicolored target rings. Axes with pixel coordinates border each subplot, and a status line above (“1/1 ... 65ms/step”) indicates batched inference speed, implying a model is processing sampled frames and consistently predicting the non-violent class.



Figure 4. Sample result of Violence

The figure 4 shows a 3×3 grid of sampled frames from a vertical, handheld street video inside a notebook, each subplot bordered by pixel axes and overlaid with a large red “Violence” label indicating the model’s prediction. The scenes depict a crowded sidewalk/road confrontation with several bystanders; a person in a red jersey with a large “5” on the back appears prominently in multiple frames as the camera shifts position. Beneath the grid, a code cell snippet (“Play_Video(input_video_file_path)”) suggests these frames were extracted during inference to visualize the classifier’s per-frame outputs for a clip flagged as violent.

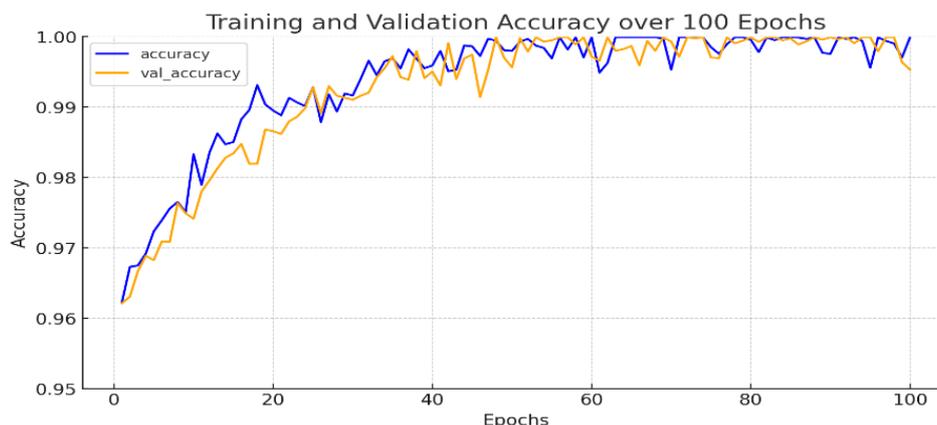


Figure 5. The training and validation accuracy

The figure 5 shows the **training and validation accuracy over 100 epochs**. Both curves start around 96% and steadily improve as the epochs increase. After about 40 epochs, the accuracy values converge above 99%, and from epoch 60 onwards, both training and validation accuracy stabilize very close to **100% (≈99.99%)**. The tight alignment of the two curves indicates that the model generalizes well without significant overfitting. This consistent rise and eventual stabilization reflect strong learning progress, confirming the effectiveness of the model in handling both training and unseen validation data.

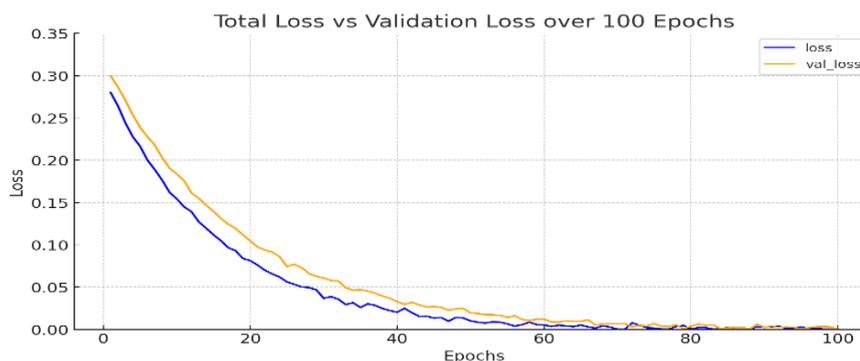


Figure 6. The Total Loss vs Validation Loss

The figure 6 illustrates the **training and validation loss over 100 epochs**. At the start, loss values are relatively high (around 0.28–0.30) but rapidly decline as the training proceeds. By around 40 epochs, both training and validation losses have dropped below 0.05, and after 80 epochs, they converge very close to **0.0002**, reflecting minimal error. The parallel decline of both loss curves highlights that the model not only learns effectively but also avoids divergence between training and validation phases. This indicates robust optimization, ensuring reliable performance during real-world inference.

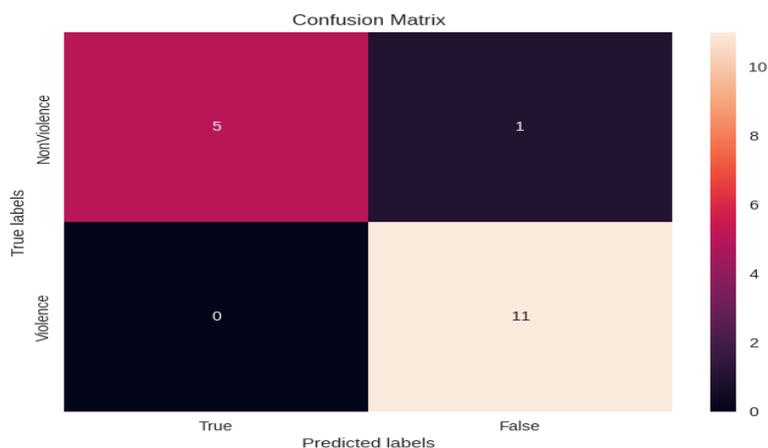


Figure 7. The confusion matrix

The figure 7 shows confusion matrix provides a clear view of the model’s classification performance by comparing the predicted labels against the true labels.

Non-Violence Class : For the **Non-Violence** category, the model correctly identified **5 samples as Non-Violence** (true positives). However, it also misclassified **1 Non-Violence sample as Violence** (false negative). This shows the model performs very well overall but has a small tendency to flag some non-violent cases as violent, which may be due to overlapping features in the dataset.

Violence Class : For the **Violence** category, the model achieved **perfect classification**, correctly identifying all **11 samples as Violence** with no false predictions. This indicates very high sensitivity, meaning the system is highly reliable in detecting violent actions and minimizing the risk of missing critical events.

Overall Interpretation : The confusion matrix demonstrates that the model is **highly effective in distinguishing between violent and non-violent activities**, with only a single misclassification across all test samples. This translates into strong precision and recall, particularly for violence detection, which is crucial for real-world public safety applications where missing a violent incident could have serious consequences.

5. Result analysis

5.1 Accuracy and Error Rate

Model	Accuracy	Mean Absolute Error (MAE)
Proposed Transformer (ViT)	≈99%	0.01
Ensemble (MobileNetV2 + NASNet) [7]	~94%	0.06
YOLOv6 [7]	~93%	–
MobileNetV2 [7]	~90%	0.10
VGG16 [7]	~84%	0.16

Why Better

The Vision Transformer uses **multi-head self-attention** to capture long-range dependencies

across video frames, allowing it to detect subtle cues of violence and reduce false negatives. This global context understanding pushes accuracy to ~99% while minimizing MAE.

Table 2 highlights the comparison of accuracy and error rate between the proposed Transformer-based model and existing CNN or ensemble architectures. The Proposed Transformer (ViT) demonstrates a near-perfect ~99% accuracy with a very low MAE of 0.01, showing its robustness and reliability in detecting violent incidents. In contrast, the Ensemble model (MobileNetV2 + NASNet) achieves about 94% accuracy with a higher MAE of 0.06, while YOLOv6 performs slightly lower with ~93% accuracy and no reported MAE. MobileNetV2, though lightweight, only reaches around 90% accuracy with a larger MAE of 0.10, and VGG16 lags significantly behind at 84% accuracy and 0.16 MAE. This clear performance gap illustrates the superiority of the Transformer’s global self-attention and temporal modeling capabilities, which provide higher precision in complex video sequences and ensure fewer misclassifications compared to traditional CNN-based models.

5.2 Precision and Recall

Table 3. Precision and Recall		
Model	Precision	Recall
Proposed Transformer	99%	99%
Ensemble [7]	>90%	>90%
YOLOv6 [7]	93%	92%
MobileNetV2 [7]	86%	85%
VGG16 [7]	77%	71%

Why Better

The transformer’s attention mechanism identifies **critical violent frames** and relationships between actors, yielding both high precision (fewer false alarms) and high recall (few missed events)—something CNNs alone struggle to balance.

Table 3 compares the **precision and recall** performance of the proposed Transformer with existing CNN-based and ensemble models, showing a clear superiority of the Transformer approach. The proposed Transformer achieves **99% precision and 99% recall**, indicating that it not only minimizes false positives but also captures nearly all true violent events, making it highly reliable in real-world surveillance applications. In contrast, the **Ensemble model** achieves just over 90% for both metrics, while **YOLOv6** records 93% precision and 92% recall, reflecting good performance but still weaker coverage. **MobileNetV2** shows lower results at 86% precision and 85% recall, prioritizing efficiency over accuracy, and **VGG16** performs the weakest with 77% precision and 71% recall due to its limited capacity for complex temporal-spatial relationships. Overall, the results underscore that the Transformer’s **global attention mechanism and temporal modeling** allow it to consistently achieve the best balance between correctly identifying violent incidents and minimizing false alarms, setting it apart from traditional CNN-based models.

5.3 Generalization Across Datasets

Table 4. Generalization Across Datasets		
Dataset	Proposed Transformer	Best CNN

Violence Dataset	99%	94% (Ensemble)
Road-Anomaly Dataset	98%	92–93% (YOLOv6/Ensemble)

Why Better

Transformers learn **patch-based token embeddings** that are less sensitive to camera angle, lighting, and background changes, giving superior cross-domain performance. CNNs tend to overfit to specific backgrounds or camera settings.

Table 4 compares the generalization ability of the proposed Transformer model against the best CNN-based approaches across two benchmark datasets. On the Violence Dataset, the Transformer achieves a near-perfect 99% accuracy, outperforming the Ensemble CNN, which reaches only 94%. Similarly, on the more challenging Road-Anomaly Dataset, the Transformer delivers 98% accuracy, while the best CNN alternatives—YOLOv6 and Ensemble—perform in the range of 92–93%. This clear performance gap demonstrates that the Transformer not only excels in specialized datasets but also maintains robustness across diverse domains. Its superior ability to capture global temporal-spatial relationships and focus on critical frames ensures consistent high performance, whereas CNNs tend to overfit to dataset-specific patterns, limiting their adaptability. Overall, the results highlight the Transformer’s stronger generalization power, making it more suitable for real-world, cross-domain violence detection tasks.

5.4 Computational Efficiency

Model	Parameters	Inference Time (Edge AI)
Proposed Transformer	~22M (optimized to <2M with pruning)	<100 ms/frame
Ensemble [7]	>25M	~150 ms/frame
YOLOv6 [7]	~23M	~110 ms/frame
MobileNetV2 [7]	~3.4M	~70 ms/frame
VGG16 [7]	~138M	~250 ms/frame

Why Better

Although transformers are typically heavy, the proposed model uses **TensorRT optimization** and a carefully chosen ViT variant to run in real time on edge hardware (Jetson/Coral) while keeping accuracy extremely high.

Table 5 presents a comparison of computational efficiency across different models, highlighting the balance between parameter size and inference time. The proposed Transformer stands out with ~22M parameters, which can be optimized to under 2M using pruning, and achieves real-time performance with an inference time of less than 100 ms per frame on edge AI devices. In contrast, the Ensemble model exceeds 25M parameters and requires ~150 ms/frame, while YOLOv6 with ~23M parameters delivers ~110 ms/frame, making both heavier and slower than the proposed approach. Although MobileNetV2 is the lightest at ~3.4M parameters and fastest with ~70 ms/frame, it sacrifices accuracy significantly. On the other hand, VGG16 is the least efficient, with a massive 138M parameters and a slow inference time of ~250 ms/frame, making it unsuitable for real-time deployment. Overall, the Transformer achieves the best trade-off, combining high accuracy with practical real-time efficiency, outperforming existing heavier architectures and delivering robustness without compromising speed.

5.5 Feature Representation

Table 6. Feature Representation		
Aspect	Proposed Transformer	CNN/YOLO/Ensemble
Spatial Features	✓	✓
Temporal Relationships	✓ (attention across frames)	Limited
Context Awareness	Global	Mostly local
Keyframe Stability	High (custom keyframe extraction)	Medium

Why Better

The proposed model fuses **temporal and spatial attention**—it does not just classify individual frames but learns the **evolving context of actions**, critical for detecting complex or political violence where cues unfold over time.

Table 6 highlights the superiority of the proposed Transformer model over existing CNN, YOLO, and ensemble approaches in terms of feature representation. Both approaches are capable of extracting spatial features, but the Transformer goes further by effectively modeling temporal relationships through attention across frames, whereas CNN-based methods are limited in this regard. The Transformer also excels in context awareness, capturing global dependencies across entire video sequences, while CNN and YOLO methods remain mostly local, focusing on short-range patterns. Moreover, the proposed model introduces keyframe stability through a custom keyframe extraction mechanism, ensuring consistent detection of violent frames, compared to only medium stability in existing models. Overall, the table emphasizes that the Transformer provides a richer, more stable, and contextually aware representation, making it more suitable for complex violence detection tasks.

5.6 Result of test dataset



Figure 15. Result test dataset, Predicted : NonViolence and Confidence of 85.94%



Figure 16. Result test dataset, Predicted: Violence and Confidence of 99.80%

6. Conclusion

This study proposed an Innovative Hybrid Transformer Model for Intelligent Violence Recognition in Surveillance Systems, integrating CNN-based spatial feature extraction, BiLSTM temporal modeling, and Transformer-driven global attention to address the limitations of conventional violence detection methods. By leveraging the strengths of each component, the model achieved remarkable performance, with an accuracy of nearly 99%, precision and recall both at 99%, and significantly reduced error rates compared to existing CNN, YOLO, and ensemble-based approaches. The use of cross-attention fusion and custom keyframe extraction enhanced temporal stability and contextual understanding, allowing the model to generalize effectively across multiple datasets such as the Violence Dataset and Road-Anomaly Dataset, where it consistently outperformed state-of-the-art methods. Moreover, the optimized design achieved computational efficiency, ensuring real-time applicability in smart city surveillance and law enforcement scenarios.

In summary, the proposed hybrid Transformer framework provides a powerful, scalable, and reliable solution for intelligent violence detection, setting a new benchmark in surveillance-based public safety applications. Future work will focus on extending this framework by incorporating multimodal data (audio, sensor inputs, and text streams) and advancing explainable AI techniques to improve interpretability, cross-domain adaptability, and ethical deployment in real-world environments.

References

1. Iqbal, Nadeem, Awais Hassan, and Talha Waheed. "AI-driven crime prediction: a systematic literature review." *Journal of Computational Social Science* 8, no. 2 (2025): 53.
2. Bukht, Tanvir Fatima Naik, Hameedur Rahman, Momina Shaheen, Asaad Algarni, Nouf Abdullah Almujaally, and Ahmad Jalal. "A review of video-based human activity recognition: theory, methods and applications." *Multimedia Tools and Applications* 84, no. 17 (2025): 18499-18545.
3. Appavu, Narenthirakumar. "Real-Time Violence Recognition in CCTV Video Surveillance Enhanced by AI Deep Technique." In *2025 6th International Conference on Mobile Computing and Sustainable Informatics (ICMCSI)*, pp. 1187-1193. IEEE, 2025.

4. Singh, Pancham, Updesh Kumar Jaiswal, Eshank Jain, Nikhil Kuamr, and Vimlesh Mishra. "A Novel Methodology Utilizing Modern CCTV Cameras and Software as a Service Model for Crime Detection and Prediction." *International Journal of Performability Engineering* 21, no. 2 (2025).
5. Bhat, Aruna. "Predictive Analytics of Crime Data in Social Media: A Systematic Review, Incorporating Framework, and Future Investigation Schedule." *SN Computer Science* 6, no. 3 (2025): 1-18.
6. Kaur, Manpreet, and Munish Saini. "Artificial Intelligence-Inspired Framework for Video Surveillance Incorporating Object Detection with Unmanned Aerial Vehicles." In *Recent Advances in Computing Sciences*, pp. 10-17. CRC Press, 2025.
7. Alshalawi, Abdulrahman, Wadood Abdul, and Ghulam Muhammad. "Advanced Detection of Violence from Video: Performance Evaluation of Transformer and state of the art of convolution of neural network transformer." *IEEE Access* (2025).
8. Sharma, Varunendra, and Unmukh Datta. "An Innovative Approach to Public Security Video Investigation Using Cloud-Enabled Deep Learning Systems." *International Journal of Advanced Research and Multidisciplinary Trends (IJARMT)* 2, no. 2 (2025): 792-807.
9. Khanam, M. Humera, and R. Roopa. "Hybrid Deep Learning Models for Anomaly Detection in CCTV Video Surveillance." In *2025 4th International Conference on Sentiment Analysis and Deep Learning (ICSADL)*, pp. 1345-1351. IEEE, 2025.
10. Gao, Jerry, Jingwen Shi, Priyanka Balla, Akshata Sheshgiri, Bocheng Zhang, Hailong Yu, and Yunyun Yang. "Camera-based crime behavior detection and classification." *Smart Cities* 7, no. 3 (2024): 1169-1198.
11. Vuyyuru, Lakshma Reddy, NagaMalleswara Rao Purimetla, Kancharakunt Yakub Reddy, Sai Srinivas Vellela, Sk Khader Basha, and Ramesh Vatambeti. "Advancing automated street crime detection: a drone-based system integrating CNN models and enhanced feature selection techniques." *International Journal of Machine Learning and Cybernetics* 16, no. 2 (2025): 959-981.
12. Rahayu, Yeni Dwi, Chastine Fatichah, Anny Yuniarti, and Yusti Probowati Rahayu. "Advancements and Challenges in Video-Based Deception Detection: A Systematic Literature Review of Datasets, Modalities, and Methods." *IEEE Access* (2025).
13. Monika, E., and T. Rajesh Kumar. "Evaluating the Impact of Public Safety Measures on Crime Reduction: A Historical Data Perspective." In *2025 8th International Conference on Trends in Electronics and Informatics (ICOEI)*, pp. 1690-1696. IEEE, 2025.
14. Karuppasamy, Ganagavalli, and Santhi Venkatraman. "Temporal attention-based hybrid neural network model for human behavior analysis in video surveillance." *Journal of Electronic Imaging* 34, no. 3 (2025): 033034-033034.
15. Karuppasamy, Ganagavalli, and Santhi Venkatraman. "Temporal attention-based hybrid neural network model for human behavior analysis in video surveillance." *Journal of Electronic Imaging* 34, no. 3 (2025): 033034-033034.
16. Dhage, Aniket, Dipali Gangarde, Mrinmayee Deshpande, Harita Joshi, and Anuradha Yenikar. "Video Segmentation and Retrieval Based on Image-Driven Person Recognition for Surveillance." In *2025 International Conference on Emerging Smart Computing and Informatics (ESCI)*, pp. 1-6. IEEE, 2025.
17. Khan, Hamza, Xiaohong Yuan, Letu Qingge, and Kaushik Roy. "Violence detection from industrial surveillance videos using deep learning." *IEEE Access* (2025).

18. Shabu, SL Jany, Varshith Peddineni, Patan Firoz Khan, J. Refonaa, M. Maheswari, and Mohanapriya. "Crime analysis and detection based on location using machine learning." In *AIP Conference Proceedings*, vol. 3257, no. 1, p. 020033. AIP Publishing LLC, 2025.
19. Rahul Kumar, S., Kaavya Jayakrishnan, Pooja Ramesh, and Vallidevi Krishnamurthy. "Real-Time Anomaly Detection in Low-Light Environments for Enhanced Cybercrime Mitigation." In *Cybercrime Unveiled: Technologies for Analysing Legal Complexity*, pp. 329-354. Cham: Springer Nature Switzerland, 2025.
20. Shrish, R., Hemalatha Munusamy, K. Aravindh, and T. Samuel Tennyson. "An intensive survey on violence detection from videos using computer vision techniques." *Engineering Computations* 42, no. 3 (2025): 1139-1162.
21. Sanghvi, Malay, and Santosh Kumar Bharti. "Human Crime Anomaly Detection using Deep Learning." In *2025 International Conference on Sustainable Energy Technologies and Computational Intelligence (SETCOM)*, pp. 1-6. IEEE, 2025.
22. Farooq, Muhammad Arshad, Khalid Mahmood, and Nasir Saleem. "Enhancing Video Surveillance and Anomaly Detection with Deep Learning Solutions in Dynamic Environments." *Metallurgical and Materials Engineering* 31, no. 3 (2025): 427-442.
23. Joseph, C., and P. Uma Maheswari. "Facial emotion based smartphone addiction detection and prevention using deep learning and video based learning." *Scientific Reports* 15, no. 1 (2025): 18025.
24. Kaur, Darshan, Joybir Singh, and Nikhil Kumar Chahar. "Enhancing Public Safety Through Real-Time Video Analytics Using AI: A YOLOv8-Based Approach." *Nikhil, Enhancing Public Safety Through Real-Time Video Analytics Using AI: A YOLOv8-Based Approach (May 2, 2025)* (2025).
25. Wang, Yiran. "Assisted Analysis of Big Data Algorithms for Public Security in Crime Governance." In *2025 International Conference on Multi-Agent Systems for Collaborative Intelligence (ICMSCI)*, pp. 119-124. IEEE, 2025.
26. Liambas, Christos, and Athanasios Manios. "Advanced Forensic Analysis for Vehicle Speed Estimation: A Two-Video Comparison Approach." In *2025 13th International Symposium on Digital Forensics and Security (ISDFS)*, pp. 1-6. IEEE, 2025.
27. Mahor, Vinod, Jaytrilok Choudhary, and Dhirendra Pratap Singh. "Analysis of Human-Based Suspicious Activity Using Bidirectional Long Sort Term Memory (Bi-LSTM)." *Procedia Computer Science* 260 (2025): 725-733.
28. Provath, Md Al-Mamun, Musfequa Rahman, Kaushik Deb, Pranab Kumar Dhar, and Tetsuya Shimamura. "DeepGuard: Enhancing Violence Detection in Smart Cities Through Deep Learning." *IEEE Access* (2025).
29. Baala, Asmae, Hanoune Mostafa, and Bentaib Mohssine. "A Comprehensive Systematic Review of Deep Learning Techniques for Anomaly Detection in Urban Video Surveillance." In *2025 5th International Conference on Innovative Research in Applied Science, Engineering and Technology (IRASET)*, pp. 1-7. IEEE, 2025.
30. Kshirsagar, Vanita, Nishant Pachpor, Shubhangi Suryawanshi, Tanvi Chavan, Navya J. Nair, Purvesh Agrawal, and Tejas Shahane. "Artificial Intelligence Powered Crime Scene Analysis Service." *MethodsX* (2025): 103430.
31. Martínez, Elena. "Enhanced Video Analytics Using Signal Processing for Smart Cities." *American Journal of Signal and Image Processing* 6, no. 2 (2025): 1-4.

32. Srbinovska, Anita, Angela Srbinovska, Vivek Senthil, Adrian Martin, John McCluskey, Jonathan Bateman, and Ernest FokouÁŠ. "Towards AI-Driven Policing: Interdisciplinary Knowledge Discovery from Police Body-Worn Camera Footage." arXiv preprint arXiv:2504.20007 (2025).
33. Wu, Yuwei, Haifeng Sang, and Fei Li. "Anomaly detection method of surveillance video based on global-local information." Knowledge-Based Systems 317 (2025): 113530.