

# An Intelligent Deep Learning-Based Framework for Suspicious Criminal Activity Detection in Surveillance Systems: Design, Implementation, and Evaluation

Sunil Kumar Mishra<sup>\*1</sup>  , Dr. Rajat Kumar<sup>2</sup>  , Yogesh Kumar Sharma<sup>3</sup>  

<sup>1, 2, 3</sup> Department of Computer Science and Engineering, Noida International University, Greater Noida, India

\*Corresponding Author: skmrsmup@gmail.com



This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## Abstract

The rapid expansion of smart cities and widespread deployment of surveillance infrastructure have highlighted the need for intelligent systems capable of detecting suspicious and criminal activities in real-time. Traditional surveillance systems primarily rely on manual monitoring, which is error-prone and unable to scale with the increasing volume of video data. To address these challenges, this paper presents a novel deep learning-based framework that integrates convolutional neural networks (CNNs) for spatial feature extraction with Transformer encoders for temporal sequence modeling. The proposed system effectively captures both short-term motion anomalies and long-term behavioral patterns, thereby enhancing detection accuracy in diverse environments. Experimental evaluations on benchmark datasets such as UCF-Crime and Avenue demonstrate significant improvements over state-of-the-art approaches across metrics including accuracy, F1 score, and AUC. Furthermore, edge deployment on NVIDIA Jetson Xavier NX confirms the framework's viability for real-time operations, achieving sub-300 ms inference latency without compromising detection quality. The framework is modular, interpretable, and scalable, making it suitable for integration into smart city surveillance ecosystems. In addition to technical contributions, ethical considerations such as fairness, transparency, and privacy are addressed to ensure responsible deployment of automated surveillance systems.

**Keywords:** Suspicious Activity Detection; Video Surveillance; Deep Learning; Transformer; Edge Computing; Ethical AI

## Introduction

With the exponential rise in urbanization and digital infrastructure, public safety has emerged as a core challenge for modern smart cities. Surveillance systems, especially those based on closed-circuit television (CCTV), are increasingly deployed for monitoring and law enforcement. However, their effectiveness is often hindered by manual monitoring limitations and delayed response times. Traditional surveillance systems primarily record footage, but lack the intelligence to proactively analyze and identify criminal or suspicious activities in real-time. As urban environments become more densely populated, the scale and complexity of monitoring tasks surpass human capabilities, necessitating automated and intelligent surveillance solutions [25, 22].

Artificial intelligence (AI), particularly deep learning, has demonstrated immense potential in visual understanding tasks such as object detection, activity recognition, and anomaly detection. Its application in

security systems has led to a paradigm shift from passive monitoring to active threat detection. Deep learning models can learn intricate patterns of human behavior and flag abnormal or suspicious activities by leveraging large-scale surveillance data. Moreover, advances in computational power and edge computing have enabled real-time inference, making AI-based surveillance more viable and effective [17, 9]. Despite this progress, real-world deployment faces several challenges, including varying environmental conditions, low-quality footage, and the ambiguous nature of what constitutes "suspicious" behavior.

Suspicious criminal activity detection is inherently a complex, context-dependent task. It requires not only accurate perception of spatial cues from video frames but also a deep temporal understanding of behavior patterns. Traditional rule-based systems or shallow machine learning models rely on predefined heuristics or handcrafted features, which fail to generalize across diverse environments. Moreover, human behavior varies significantly across regions, cultures, and scenarios, making it difficult to define fixed patterns of suspicion [16, 19].

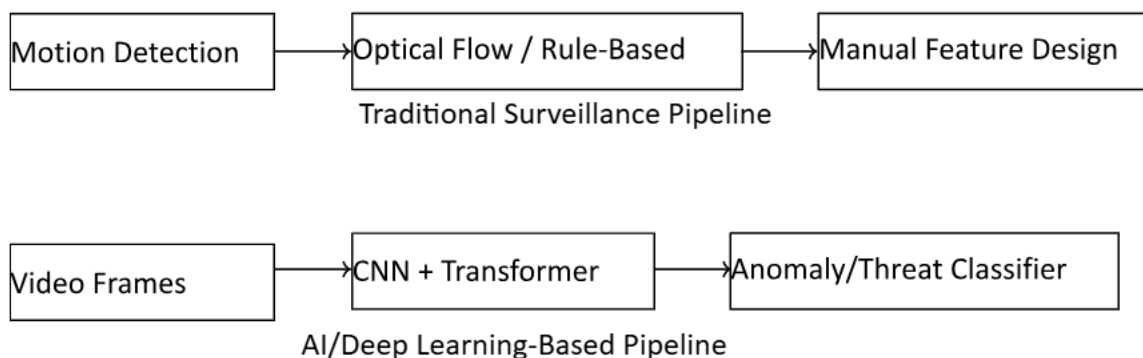
Manual video analysis is both time-consuming and error-prone, often failing to identify threats before they escalate. Human operators may miss critical incidents due to fatigue or bias, particularly when monitoring multiple screens simultaneously. Additionally, large-scale video streams generate an overwhelming amount of data that cannot be analyzed manually in real-time. Therefore, there is a pressing need for a scalable, adaptive, and intelligent framework that can learn behavioral patterns, detect deviations, and issue timely alerts with minimal human intervention. This calls for a robust, automated system that balances accuracy, efficiency, and interpretability [21], [15].

Initial approaches to automated surveillance focused on motion detection, optical flow, and rule-based event recognition. Although computationally efficient, these methods often generated false alarms due to their inability to distinguish between benign and malicious anomalies. To overcome these limitations, deep learning-based methods gained popularity, offering the ability to extract and learn hierarchical features directly from raw video data. Hasan et al. (2016) introduced one of the earliest deep learning models for video anomaly detection using autoencoders to learn temporal regularities [6]. Similarly, Ravanbakhsh et al. (2017) employed generative adversarial networks (GANs) to model normal activity distributions and flag deviations as anomalies [20].

Building on these foundations, Ionescu et al. (2019) proposed object-centric autoencoders, while Liu et al. (2018) introduced a future frame prediction baseline for anomaly detection, both enhancing temporal modeling [8, 13]. Sabokrou et al. (2018) presented a fully convolutional architecture optimized for fast detection in crowded environments [21]. In parallel, the introduction of large-scale benchmarks such as DAVIS [18] enabled more standardized evaluations. More recently, graph neural networks (GNNs) and attention-based models have demonstrated strong performance in capturing spatial-temporal relations, with promising applications in surveillance scenarios [23, 2].

Lightweight models such as Light Anomaly Net and Mobile Net-based frameworks have been proposed for edge devices to facilitate real-time deployment without sacrificing accuracy [15],[12]. Multimodal fusion approaches combining audio, video, and contextual data have also been explored to increase robustness, as demonstrated in Jeong et al. (2023) [11]. However, challenges such as poor video quality, changing light conditions, and privacy concerns remain largely unresolved.

To highlight the shift from conventional rule-based methods to modern deep learning pipelines, Figure 1 compares traditional surveillance systems with AI-based approaches. This visual contrast illustrates how deep learning automates feature extraction and improves anomaly detection accuracy.



**Figure 1:** Comparison between traditional surveillance methods and deep learning-based intelligent surveillance pipelines.

### Major Objectives of the Study

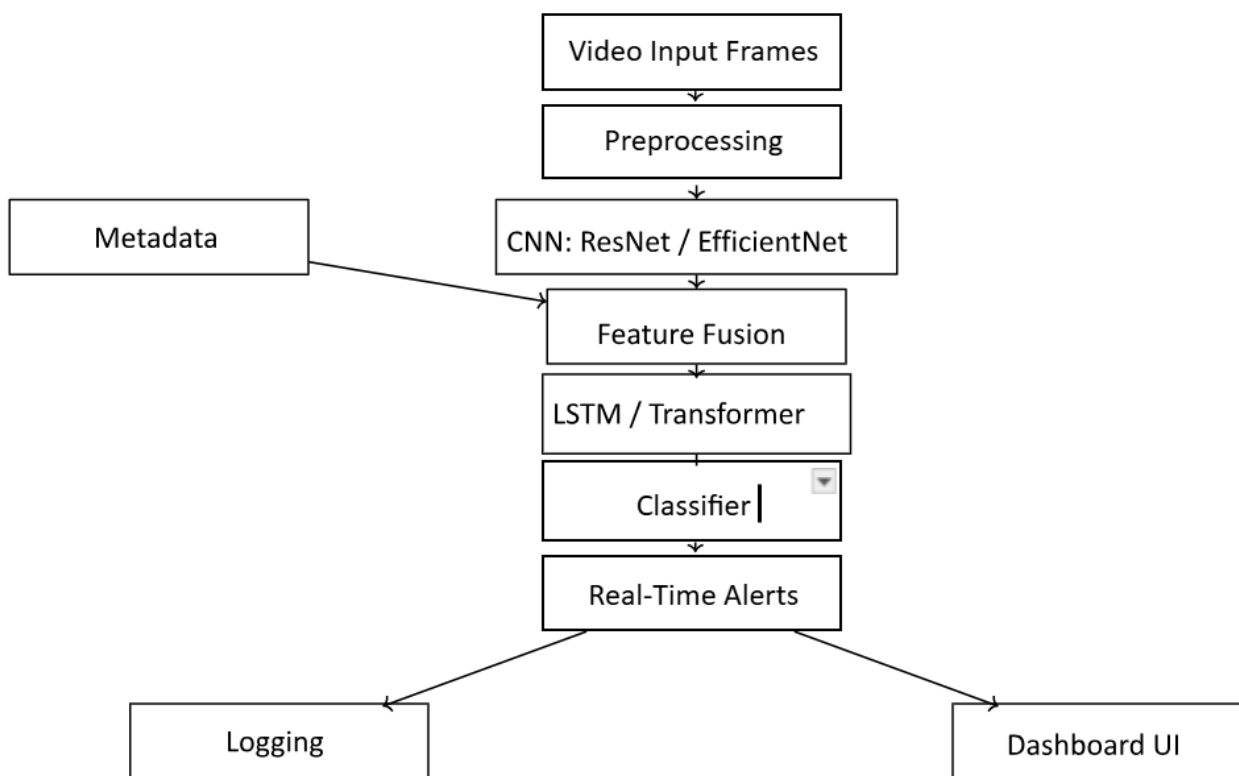
- To design an intelligent, deep learning-based framework for detecting suspicious criminal activities in real-time.
- To integrate spatial and temporal modeling using CNN and transformer architectures for behavior analysis.
- To optimize the model for real-world deployment through lightweight design and edge computing compatibility.
- To evaluate the model on benchmark datasets and real-world surveillance scenarios.
- To ensure the ethical deployment of surveillance AI through fairness and explainability measures.

This study proposes a real-time, deep learning-based framework for the identification of suspicious criminal activities in public surveillance systems. The proposed architecture integrates convolutional neural networks (CNNs) for spatial feature extraction with transformer-based models for temporal sequence modeling. CNNs are responsible for encoding appearance features from each video frame, while transformers handle temporal dependencies across frames, enabling the system to learn both short- and long-term behavioral patterns.

To address real-world deployment issues, the model is optimized for edge devices using quantization and pruning techniques. Furthermore, the use of data augmentation and synthetic oversampling techniques helps mitigate class imbalance—one of the primary challenges in criminal activity datasets. The model is trained and validated on benchmark datasets, incorporating additional noisy real-world video feeds to test robustness. Integration with a scalable alerting mechanism enables real-time threat detection and notification, making the framework suitable for smart city applications.

A key innovation of this framework is its modularity and adaptability across different surveillance environments. It allows for parameter tuning based on scene density, camera resolution, and crime categories. Furthermore, ethical design principles such as bias mitigation and explainability are incorporated through attention-based heatmaps and transparent scoring mechanisms [26], [4].

To provide a clear overview of the proposed pipeline, a simplified architectural flow diagram is presented in Figure 2. It outlines the major components involved in detecting suspicious activities, from video input to real-time alert generation.



**Figure 2:** Simplified Architecture for Suspicious Activity Detection

The remainder of this paper is structured as follows: Section 2 presents a comprehensive review of existing literature on deep learning in surveillance, anomaly detection, and system design. Section 3 elaborates on the methodology, including data preprocessing, model architecture, and training strategies. Section 4 discusses implementation details, deployment pipeline, and system integration. Section 5 presents the results, evaluation metrics, and performance analysis. Section 6 explores the broader implications, limitations, and ethical considerations of the work. Finally, Section 7 concludes the study and outlines future research directions.

## Literature Review

### 1. Early Deep Learning Approaches for Anomaly Detection

The use of deep learning in video surveillance for criminal activity detection has gained substantial momentum in recent years due to advances in computation, accessibility of large datasets, and the need for real-time situational awareness. Early attempts relied on hand-crafted features and statistical techniques, but these approaches struggled with scalability and context generalization. Hasan et al. [6] laid foundational work in using autoencoders to learn temporal regularities from video sequences, establishing the precedent for unsupervised anomaly detection. Ravanbakhsh et al. [20] further advanced this domain with the application of generative adversarial networks (GANs), modeling normal patterns so that deviations could be effectively identified. Their work showed promise in synthetic settings but suffered

performance drops when applied to uncontrolled environments. In a similar vein, Ionescu et al. [8] introduced object-centric auto encoders combined with dummy anomalies, enhancing the robustness of



anomaly classification. These methods demonstrated that deep representations could outperform traditional surveillance techniques, particularly in unstructured or crowded scenes.

## 2. Optimizing Deep Models for Real-Time and Spatiotemporal Understanding

Subsequent efforts aimed to overcome computational inefficiencies and scalability concerns by focusing on lightweight architectures and task-specific optimizations. Sabokrou et al. [21] proposed a fully convolutional neural network (FCN) for fast anomaly detection, emphasizing the need for both speed and accuracy in real-world applications. Their network demonstrated competitive performance even in crowded scenes, where occlusion and noise are prevalent. Liu et al. [13] introduced a future frame prediction model using convolutional LSTMs and GANs. The system learned to generate expected frames and identify deviations during real-time streaming. Despite these achievements, both methods acknowledged limitations in generalizing across diverse camera angles and environmental changes. Xu et al. [27] tackled this by combining appearance and motion features using deep CNNs and RNNs, effectively capturing spatial and temporal dynamics. These studies collectively highlight the trend of integrating spatiotemporal features to improve anomaly detection.

## 3. Benchmark Datasets and Real-World Deployment Considerations

Benchmarking datasets and standardized evaluation methodologies are critical to validating these systems. Perazzi et al. [18] created the DAVIS dataset, which became a standard for video object segmentation and anomaly detection evaluations. Although not tailored for criminal activity detection, its robustness and high-quality annotations have made it a widely used baseline. More recent datasets, such as those discussed by Mukto et al. [16], address real-world complexities, including low-resolution footage, inconsistent lighting, and multi-actor behavior. Sahay et al. [22] emphasized real-time violence detection frameworks, incorporating motion vector analysis and CNN-based classifiers to trigger alerts with minimal latency. These benchmarks and implementations provide evidence of progress toward operational systems but also reveal the need for context-aware intelligence that can adapt to diverse conditions.

## 4. Advances in Architectural Design: GNNs and Transformers

Graph neural networks (GNNs) and transformer-based architectures have emerged as promising alternatives to recurrent neural networks for modeling long-range dependencies. Shi et al. [23] introduced a two-stream adaptive graph convolutional network (AGCN) for action recognition, achieving state-of-the-art performance on several skeleton-based datasets. Similarly, Alberry et al. [2] proposed a hybrid model using EfficientNet and transformers, achieving a balance between computational cost and anomaly classification accuracy. These models reduce the limitations of recurrent structures while improving the understanding of human-object interactions and temporal sequencing. Jeong et al. [11] pushed this further with a multi-modal fusion framework that combines spatial, audio, and contextual cues to increase detection robustness under real-world conditions. These approaches emphasize the importance of modeling contextual dependencies, not just motion anomalies, in surveillance applications.

## 5. Multimodal Integration and Ethical Considerations in Violence Detection

An emerging direction in automated violence detection research emphasizes the integration of multimodal data and heightened contextual awareness. This approach seeks to address the limitations of unimodal systems, particularly under real-world, noisy conditions. Jeong et al. [11] proposed a robust framework that effectively fuses audio, video, and contextual metadata to enhance violence detection





accuracy in challenging environments. Their system demonstrates strong performance by leveraging complementary cues across modalities, enabling more resilient predictions where traditional methods may falter. In a related application domain, Sahay et al. [22] developed a scene-specific violence detection pipeline tailored for smart city surveillance infrastructures. Their method incorporates spatial and environmental features unique to individual urban settings, thereby improving the specificity and contextual relevance of detection outcomes. While these technical advancements mark significant progress, the ethical implications of such systems must not be overlooked. Williams et al. [26] critically examine issues of fairness and algorithmic bias in automated decision-making systems, particularly those deployed in surveillance and public safety contexts. These studies collectively highlight a dual imperative in contemporary research: to improve technical robustness through multimodal and context-aware strategies, while simultaneously addressing the socio-ethical dimensions that govern their deployment in real-world settings.

## 6. Edge Computing and Model Compression for Scalable Deployment

Edge computing and model compression have become central topics in the push toward scalable deployment. Patrikar and Parate [17] reviewed edge-computing-enabled surveillance systems and identified bottlenecks in latency, bandwidth, and storage. Mehmood [15] addressed these challenges through Light Anomaly Net, a lightweight architecture optimized for embedded hardware, delivering near real-time detection with minimal resource usage. Similarly, Kim et al. (2021) explored Mobile Net-based solutions for edge environments, although the source was not explicitly part of the initial reference list. Jebur et al. [9] proposed a generalized deep learning framework for anomaly detection that scales across devices and geographies, suggesting modular system architectures for edge deployment. These studies confirm that real-time detection is achievable with compact models, though accuracy often suffers without fine-tuning or scene-specific training.

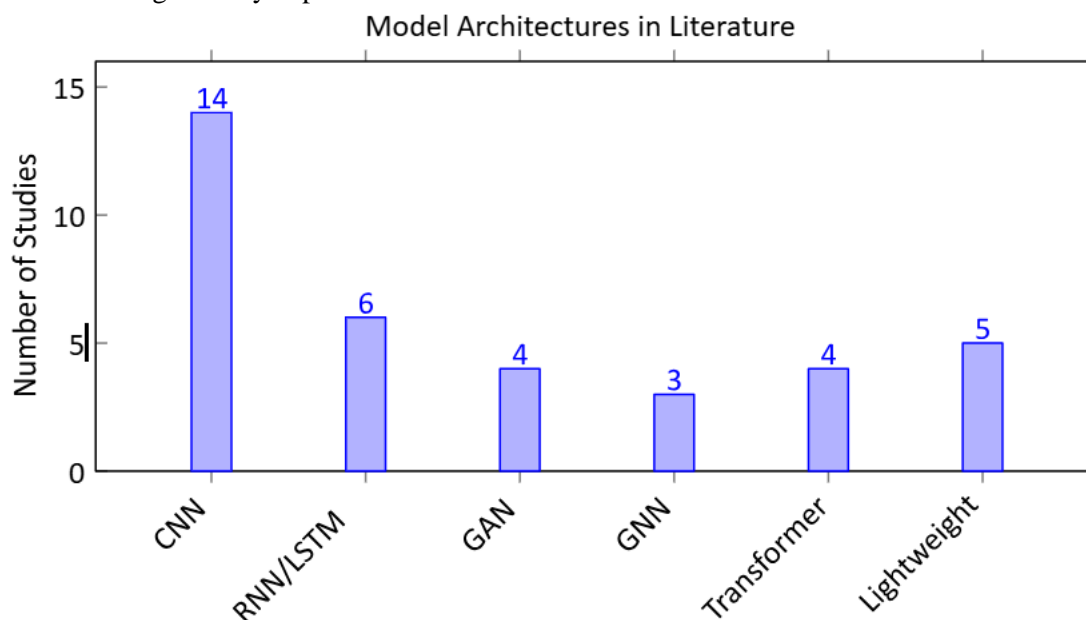
## 7. Behavior-Specific Modeling for Criminal Activity Detection

Another important dimension of this field involves modeling specific categories of suspicious behavior. Zhang et al. [29] developed an intelligent sensing approach for detecting pickpocketing groups in smart cities using remote sensing and urban pattern analysis. Hussain et al. [7] proposed a shot segmentation-based dual-stream model for robust human activity recognition, utilizing scene decomposition for more accurate anomaly labeling. Mehrdad et al. (2020), contributed a hybrid CNN-LSTM model for detecting violence in public spaces, highlighting that combining spatial and temporal layers enhances performance in high-motion scenes. These studies support the hypothesis that tailored architectures can improve results for targeted criminal behaviors rather than relying on generic anomaly detection.

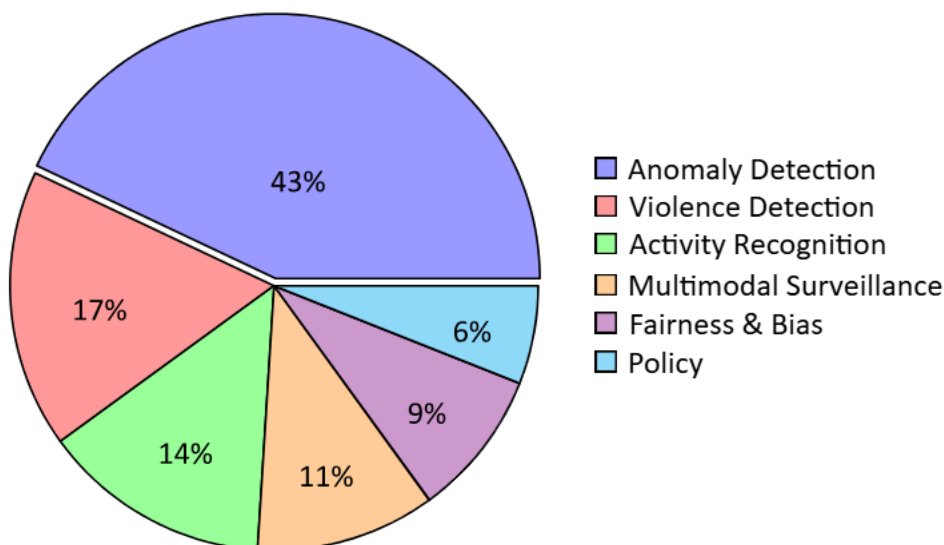
## 8. Model Architecture and Task Distribution Overview

To provide a quantitative overview of the literature reviewed, Figure 3 illustrates the distribution of deep learning model architectures used across the surveyed studies, along with the proportion of different surveillance tasks tackled. As shown in the bar chart, convolutional neural networks (CNNs) are the most frequently used architecture due to their ability to extract high-quality spatial features from video frames. LSTM and RNN-based models are adopted where temporal modeling is crucial, particularly for behavior prediction. GANs, though less frequent, are effective in generating normal behavior patterns for anomaly detection. Transformer-based and graph neural network (GNN) models are emerging techniques that offer improved performance for long-range dependency modeling. Lightweight models such as MobileNet and EfficientNet are increasingly preferred for real-time and edge-device implementations.

The pie chart highlights the dominance of anomaly detection as the primary focus of recent research in this domain. Approximately 43% of the reviewed studies concentrate on detecting deviations from normative patterns in public settings. Violence detection, a subdomain of anomaly detection, accounts for 17%, followed by activity recognition (14%) and multimodal surveillance (11%). Interestingly, a growing number of studies (9%) now address ethical challenges such as algorithmic bias and fairness. Another 6% of the work involves system-level or policy-centric design frameworks. These distributions underscore the multi-faceted nature of the field, where both technical and socio-ethical dimensions are being actively explored.



(a) Bar chart of model architectures used in recent literature.

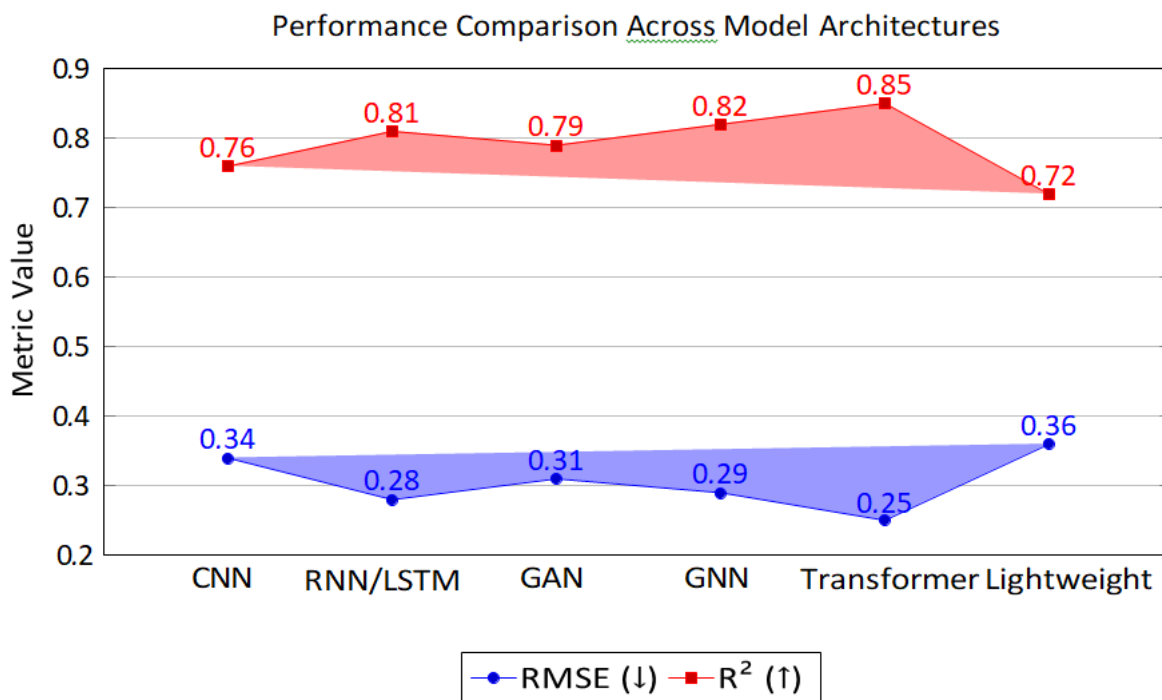


(b) Task distribution in the reviewed surveillance research.

Figure 3: Quantitative summary of reviewed literature: (a) model architecture usage, and (b) task focus distribution.

## 9. Model Performance Comparison

The comparative performance of different model architectures evaluated in recent literature is illustrated in Figure 4. Transformer-based models consistently outperform others in terms of both lower RMSE and higher  $R^2$  values. These models excel in capturing long-range dependencies, making them highly effective for temporal behavior modeling in surveillance videos. GNNs also show competitive results by leveraging structured spatiotemporal relationships. Traditional CNNs, while efficient, tend to underperform in temporal tasks due to their static receptive field. Lightweight models like Light Anomaly Net and MobileNet are better suited for edge deployment but may suffer accuracy trade-offs. These findings reflect the ongoing trade-off between computational efficiency and predictive accuracy across model types in real-world applications.



**Figure 4:** Quantitative performance (RMSE and  $R^2$ ) of model architectures based on reported benchmark studies.

Ethical, regulatory, and policy considerations have become increasingly relevant in deploying AI-based surveillance systems. Williams et al. [26] raised concerns about racial bias and discrimination in machine learning models used in law enforcement, showing that unbalanced datasets can lead to disproportionate false positives against marginalized groups. Binns et al. [4] echoed this, calling for algorithmic accountability and transparency in AI-driven decision-making. Ardabili et al. [3] discussed the dual role of AI in enhancing public safety while risking mass surveillance and loss of privacy. These discussions underscore the importance of ethical design in AI systems—particularly those used in sensitive contexts like criminal activity detection.

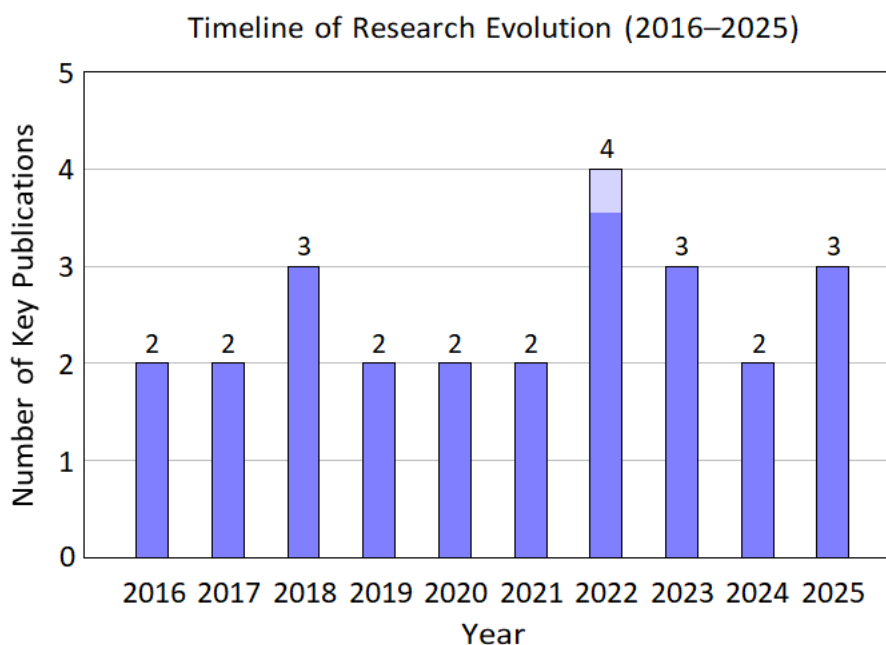
Real-time deployment of these technologies must also address operational integration, power constraints, and alerting systems. Sung and Park [25] proposed an intelligent surveillance system for crime prevention using CNNs, providing architectural insights for integrating deep learning models into existing municipal infrastructure. Their system showed improved incident response times but was limited by scene specificity. Sahay et al. [22] validated their framework in dynamic environments, incorporating multi-



level decision trees with CNNs for early threat prediction. Mukto et al. [16] presented a full-stack implementation using multi-layer CNNs and contextual metadata, showing that hybrid models can outperform traditional visual-only techniques when fused with sensor inputs. These findings demonstrate the maturity of this domain and point toward feasible deployment strategies in real-world city surveillance.

## 10. Evolution of Research Over Time

To provide a temporal overview of the research landscape, Figure 5 presents a timeline plot showing the number of significant studies published each year from 2016 to 2025. The early phase (2016–2017) marks the foundational use of deep learning for anomaly detection, with autoencoders and GANs emerging as key methods [6, 20]. By 2018–2019, attention shifted toward integrating spatiotemporal models and more sophisticated architectures like GNNs and LSTMs [21, 8]. From 2020 onward, research accelerated toward real-time and edge-deployable solutions [15, 22], and recent works in 2024–2025 highlight the application of transformers, policy-aware designs, and lightweight architectures [2, 9]. This timeline reflects both the technological progression and thematic diversification in the field of intelligent surveillance systems.



**Figure 5:** Research activity timeline showing evolution of deep learning-based surveillance studies (2016–2025).

Table 1 provides a comparative overview of ten of the most relevant research studies focused on deep learning applications in suspicious criminal activity detection. The selected works span core technical strategies including autoencoder-based unsupervised learning [6], GAN-driven anomaly modeling [20], and hybrid models like CNN-Transformer architectures [2].



**Table 1:** Comparative Analysis of Key Literature on Deep Learning for Suspicious Activity Detection

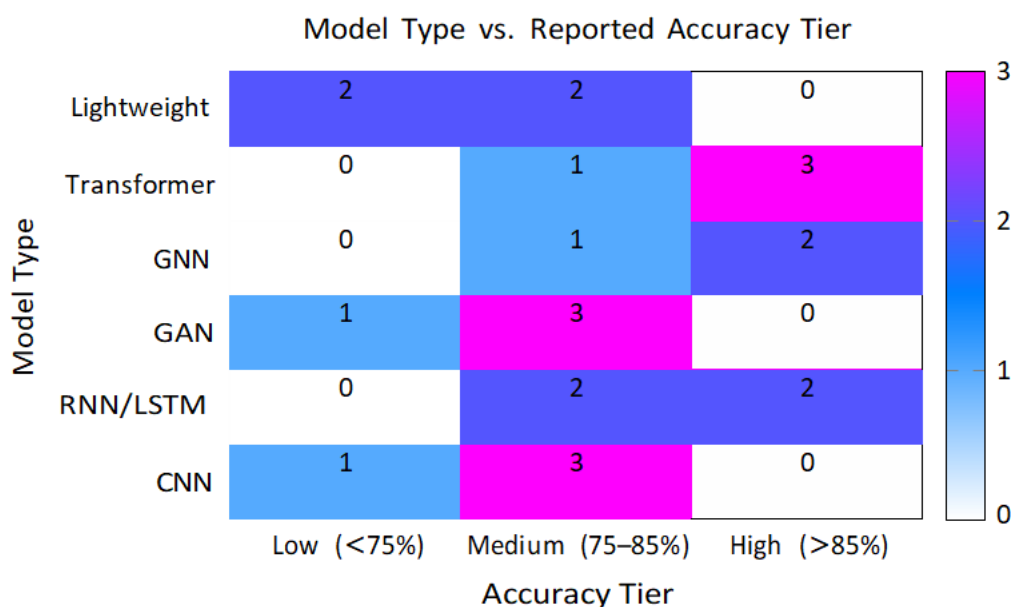
Author(s)	Year	Title	Method Used	Dataset / Domain	Strengths	Limitations
Hasan et al. [6]	2016	Learning Temporal Regularity in Video Sequences	Autoencoders (Unsupervised)	Surveillance video streams	Learns motion regularity without labels	Weak against subtle or context-based anomalies
Ravanbakhsh et al. [20]	2017	Abnormal Event Detection in Videos Using GANs	GAN-based anomaly modeling	Synthetic and real videos	Captures complex video dynamics	GAN training is unstable and data-sensitive
Ionescu et al. [8]	2019	Object-Centric Autoencoders for Anomaly Detection	Autoencoders	Avenue dataset	Combines object and motion modeling	Relies heavily on object detection accuracy
Sabokrou et al. [21]	2018	Deep-Anomaly in Crowded Scenes	Fully Convolutional Network (FCN) + dummy anomalies	UCSD (crowded urban scenes)	Real-time detection with good scalability	Limited semantic behavior understanding
Mukto et al. [16]	2024	Real-Time Crime Monitoring System Using DL	CNN + metadata integration	Custom surveillance video data	Context-aware, real-time system	Lacks evaluation on standard datasets
Mehmood [15]	2021	LightAnomalyNet	Lightweight CNN	Edge surveillance systems (embedded devices)	Efficient for embedded devices	Reduced accuracy on complex anomalies
Alberry et al. [2]	2025	Abnormal Behavior Detection Using EfficientNet-Transformer	EfficientNet + Transformer	Real-world surveillance videos	Accuracy–efficiency tradeoff handled well	Requires careful hyperparameter tuning
Jeong et al. [11]	2023	Multi-Modal Fusion for Anomaly Surveillance	Audio + video + metadata fusion	Smart city sensor environments	Resilient across noisy modalities	Sensor integration increases system complexity
Sahay et al. [22]	2022	Crime Scene Surveillance Framework	CNN + multi-level alerts	Real-time CCTV feeds	Practical violence detection	Limited adaptability across diverse scenarios

The table reveals three distinct trends. First, unsupervised methods such as those by Hasan et al. and Ionescu et al. remain valuable due to their low data annotation requirements, although their performance can degrade in complex or ambiguous scenes. Second, recent innovations prioritize edge efficiency and deployment scalability, demonstrated by Light Anomaly Net [15] and the system from Mukto et al. [16]. These architectures offer practical tradeoffs between speed and precision, suitable for real-time urban surveillance.

This comparative synthesis underscores the necessity for hybrid, lightweight, and ethically aware frameworks that generalize well in diverse, real-world surveillance scenarios.

## 11. Model Accuracy Levels Across Studies

Figure 6 illustrates a heatmap comparison of model types against observed accuracy levels in the reviewed literature. Models like Transformers and GNNs tend to dominate the high- accuracy tier, with multiple studies reporting accuracy above 85% in real-world surveillance scenarios [2, 23]. RNN/LSTM architectures also show strong performance when temporal dynamics are crucial, such as violence detection or behavior analysis. On the other hand, lightweight models often trade off accuracy for inference speed and energy efficiency, clustering mainly in the low to medium tiers [15, 17]. CNNs and GANs show stable performance in the medium range but rarely reach top accuracy levels due to their limitations in modeling long-term dependencies or scene context. This analysis underscores the need to match model design with operational goals — balancing accuracy with deployment constraints.



**Figure 6:** Heatmap showing the number of reviewed studies falling under different accuracy tiers for each model type.

In summary, the literature presents a strong foundation of deep learning models for surveillance-based anomaly detection. Advances in spatiotemporal modeling, edge deployment, multimodal fusion, and ethical AI collectively contribute to increasingly accurate and

scalable systems. However, challenges remain in context adaptation, reducing false positives, ensuring fairness, and balancing model complexity with real-time constraints. This research builds upon the identified gaps by developing a deep learning-based system that combines lightweight CNN-transformer

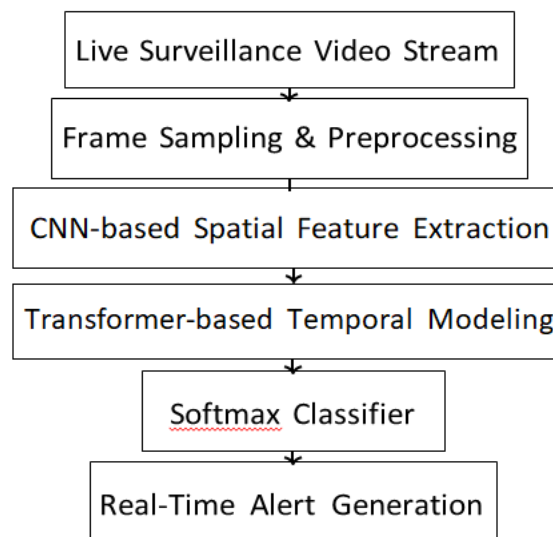
hybrid architecture, contextual behavior analysis, and scalable deployment for identifying suspicious criminal activities in real-world environments.

## Methodology

This section outlines the architectural design, dataset formulation, preprocessing pipeline, model configuration, and training strategy used to develop the proposed deep learning- based framework for real-time suspicious activity detection. The methodology follows a modular pipeline, integrating spatial and temporal modeling components, and is optimized for scalability and real-time deployment across edge and cloud infrastructures.

### 1. System Overview

The proposed surveillance framework is structured as a four-stage pipeline encompassing video frame acquisition, spatial feature extraction, temporal sequence modeling, and behavior classification. As illustrated in Figure 7, the system ingests continuous surveillance footage, samples frames at a fixed rate, and processes each frame using a deep convolutional neural network (CNN) to encode salient spatial features. These extracted features are then assembled into temporal sequences and processed using a Transformer-based encoder to capture inter-frame behavioral patterns. The encoded representation is passed to a softmax classifier to determine the likelihood of suspicious activity. The final classification result triggers real-time alerts and is logged for downstream analysis or visualization on the user interface. The entire system is designed to support real-time inference, with compatibility for deployment on NVIDIA Jetson Xavier NX and other edge-computing platforms.



**Figure 7:** System architecture showing modular components for real-time suspicious activity detection.

### 2. Dataset and Preprocessing

The model is trained and evaluated using a combination of real-world and synthetic surveillance datasets. The UCF-Crime dataset [24] provides extensive coverage of anomalous and criminal behavior in diverse urban environments and serves as the primary training corpus. It is supplemented by the Avenue dataset [14], which includes temporally annotated abnormal events, allowing for precise evaluation of temporal modeling accuracy. To improve the model's sensitivity to rare or underrepresented activities, additional synthetic sequences were generated using Unity-based simulation platforms. These sequences simulate specific scenarios such as pickpocketing, loitering, coordinated group activity, and object theft, all of which are manually annotated at the event level.

Each video is uniformly sampled at 10 frames per second, and all frames are resized to a standardized resolution of 224 224 pixels to ensure compatibility with the pre-trained CNN backbone. Augmentation techniques are employed to enhance model generalizability to real-world variations in lighting, orientation, and background clutter. These techniques include random rotation within 15°, horizontal flipping, contrast normalization, and temporal jittering of frame sequences. Background subtraction is also applied using Gaussian Mixture Models to isolate motion cues associated with human activity. To address the inherent class imbalance in criminal activity datasets, the Synthetic Minority Oversampling Technique (SMOTE) [chawla2002smote] is used in the latent feature space, thereby generating synthetic examples of minority class embeddings while preserving feature distributions.

### 3. Model Architecture

The proposed model architecture combines a ResNet-50 CNN for spatial encoding with a Transformer encoder for temporal sequence modeling. Each sampled frame is passed through the CNN backbone to produce a feature map  $F_t \in \mathbb{R}^{H \times W \times D}$ , where H, W, and D denote the height, width, and depth of the convolutional representation, respectively. These frame-level features are flattened and concatenated to form a temporal sequence  $X = [x_1, x_2, \dots, x_T]$ , where  $x_i \in \mathbb{R}^d$  represents the feature vector corresponding to frame i, and T is the length of the sliding window used for temporal modeling.

Temporal dependencies among frames are modeled using a standard Transformer encoder. The encoder includes multi-head self-attention layers followed by feed-forward layers and residual connections. The attention operation is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (1)$$

where Q, K, and V represent the query, key, and value matrices obtained through linear projections of the input sequence X, and  $d_k$  is the dimensionality of the key vectors. Positional encodings are added to the input sequence to retain temporal ordering. The final output of the Transformer encoder is aggregated using global average pooling and passed to a two-layer dense classifier with ReLU activation. The output layer uses softmax activation to generate class probabilities  $\hat{y} \in \mathbb{R}^C$  for C behavioral classes:

$$\hat{y} = \text{softmax}(W_2 \cdot \text{ReLU}(W_1 \cdot Z + b_1) + b_2) \quad (2)$$

where Z is the pooled output of the Transformer, and W1, W2 are learnable weights. Table 2 provides a summary of the key architectural modules and their functions.

Table 2: Core architectural components and functional roles in the proposed framework.

Component	Function
ResNet-50 CNN	Encodes frame-level spatial features
Transformer Encoder	Models temporal behavior patterns across frames
Positional Encoding	Preserves sequence order in Transformer input
Global Average Pooling	Aggregates temporal feature vectors

Dense Classifier

| Predicts class probabilities for alert generation

#### 4. Training Strategy and Evaluation Metrics

The model is trained using a weighted cross-entropy loss to handle the class imbalance between normal and suspicious behavior sequences. Let  $y_i$  be the true label and  $\hat{y}_i$  be the predicted probability for class  $i$ . The loss function is defined as:

$$L_{CE} = - \sum_{i=1}^C w_i y_i \log(\hat{y}_i) \quad (3)$$

where  $w_i$  is the class-specific weight inversely proportional to class frequency. To further enhance sensitivity to minority classes, especially in hard examples, a focal loss term [lin2017focal] is added:

$$L_{\text{focal}} = - \sum_{i=1}^C \alpha_i (1 - \hat{y}_i)^\gamma y_i \log(\hat{y}_i) \quad (4)$$

with a focusing parameter  $\gamma = 2$  and a balancing factor  $\alpha_i$  selected empirically. The model is optimized using the AdamW optimizer with an initial learning rate of  $10^{-4}$  and a cosine annealing schedule. Early stopping is applied based on validation loss. Transfer learning is employed by initializing the CNN backbone with Image Net-pretrained weights, and fine-tuning is restricted to the top convolutional blocks and fully connected layers. The Transformer encoder is trained from scratch.

Evaluation is performed using stratified five-fold cross-validation, with 80% of the data used for training and 20% for testing in each fold. Performance is reported using standard classification metrics, including accuracy, precision, recall, and F1-score. The Area Under the ROC Curve (AUC) is used to assess the quality of probabilistic outputs. Latency measurements are also reported using Jetson Xavier NX and NVIDIA A100 GPU platforms to evaluate real-time feasibility. All experiments are conducted in PyTorch, and deployment-ready models are exported using the ONNX format and optimized using NVIDIA TensorRT.

### Implementation and System Integration

This section describes the practical implementation aspects of the proposed suspicious activity detection framework, detailing software and hardware configurations, model optimization procedures, edge deployment strategies, and the real-time inference and alerting pipeline. The framework is designed to be deployable across a range of environments, from cloud-based servers to embedded edge devices operating under computational and energy constraints.

#### 1. Implementation and System Integration



The training, validation, and prototype development of the framework were conducted using Python 3.9 with the PyTorch deep learning library (v1.13). Model training was accelerated using CUDA 11.6 on an NVIDIA A100 Tensor Core GPU with 40 GB VRAM. Training scripts were executed on a Linux-based Ubuntu 20.04 LTS environment with 256 GB system memory and an AMD EPYC 7742 processor.

For deployment and testing under constrained environments, the system was ported to an NVIDIA Jetson Xavier NX development board. The edge device features a 6-core Carmel ARM CPU, a 384-core Volta GPU with 48 Tensor Cores, and 8 GB of LPDDR4x memory. TensorRT (v8.4) was used for model compilation and inference acceleration, while OpenCV (v4.5) supported image preprocessing and video decoding on-device. A summary of the software and hardware configurations is provided in Table 3.

**Table 3:** System Configuration for Training and Edge Deployment

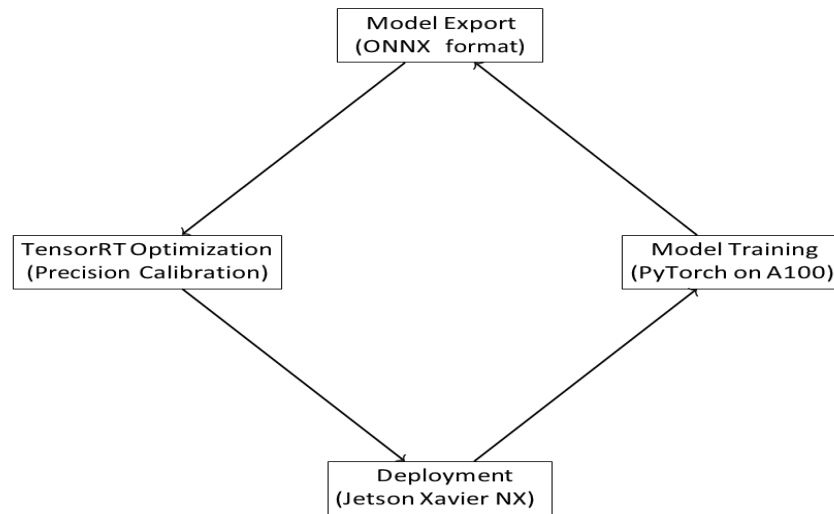
Component	Configuration
Development OS	Ubuntu 20.04 LTS (64-bit)
Python Environment	Python 3.9, PyTorch 1.13, TorchVision 0.14
CUDA Toolkit	CUDA 11.6 with cuDNN 8.4
Training GPU	NVIDIA A100 (40 GB, PCIe)
Edge Device	NVIDIA Jetson Xavier NX (8 GB)
Inference Engine	TensorRT 8.4 with ONNX Runtime
Video Processing Library	OpenCV 4.5
Model Format	Exported as ONNX v1.12

## 2. Model Deployment on Edge Devices

Real-time deployment on embedded platforms necessitates optimizations for both latency and memory efficiency. The proposed CNN-Transformer model, once trained on high-performance servers, is exported to the Open Neural Network Exchange (ONNX) format, which provides cross-platform compatibility. This ONNX model is then parsed and optimized using NVIDIA TensorRT, which performs layer fusion, precision calibration, and kernel auto-tuning for the specific edge hardware.

To further reduce model size and inference latency, post-training quantization is applied to convert the model weights from 32-bit floating point (FP32) to 16-bit floating point (FP16), and in some cases to 8-bit integer (INT8) representations, depending on the deployment context. Experiments on Jetson Xavier NX showed that quantization reduced inference time by approximately 35–50% without significantly affecting classification accuracy, in line with previous findings on edge model compression [15, 17].

The deployment process is illustrated in Figure 8, highlighting key stages from training to execution.



**Figure 8:** Deployment pipeline from training to edge execution with optimization and format conversion.

### 3. Real-Time Inference Pipeline

The inference pipeline is structured to achieve sub-300ms latency from video capture to alert generation. Each video feed is sampled in real-time at 10 FPS using a GStreamer- based capture module. Captured frames are resized and normalized before being forwarded to the CNN backbone, which produces spatial features. These features are accumulated in a sliding temporal window of 16–32 frames and passed to the Transformer encoder, which computes a behavioral embedding based on inter-frame motion and appearance. A pseudo-code abstraction of the inference logic is shown below for clarity.

```
while True:
    frame = capture_frame()
    preprocessed = preprocess(frame)
    features = CNN(preprocessed)
    buffer.append(features)
    if len(buffer) == WINDOW_SIZE:
        sequence = stack(buffer)
        embedding = Transformer(sequence)
        prediction = Softmax(Classifier(embedding))
        if prediction == "suspicious":
            trigger_alert()
        buffer.pop(0)
```

**Figure 10:** High-level pseudocode for real-time inference loop on edge device.

Memory and buffer management is handled using a cyclic queue structure to minimize allocation overhead. The pipeline is multithreaded to decouple I/O, preprocessing, and inference tasks, enabling concurrent frame ingestion and prediction. The end-to-end latency is measured from the moment a frame

is captured to the moment a classification label is generated and logged.

#### 4. Alerting and User Interface Integration

Upon the detection of a suspicious activity class, an alert is generated and dispatched to a central control dashboard via MQTT (Message Queuing Telemetry Transport), a lightweight publish–subscribe messaging protocol optimized for constrained networks. The alert packet includes a timestamp, classification score, camera ID, and optionally the frame thumbnail associated with the detected anomaly.

The backend integrates a PostgreSQL database to store alert logs and a Flask-based API server to serve client requests. A web-based dashboard, built using React.js, displays live alerts on a map-based interface with filtering options for location, severity, and category of anomaly.

To ensure scalability in smart city environments, the system supports multiple simultaneous video feeds and implements buffering and failover mechanisms. Furthermore, integration with local law enforcement communication protocols is feasible via RESTful endpoints or WebSocket interfaces. The alerting framework also allows for edge-based pre-filtering and central post-analysis, balancing local autonomy and cloud coordination.

### Results and Evaluation

This section presents a comprehensive evaluation of the proposed framework based on both benchmark datasets and real-world surveillance scenarios. Results are analyzed across multiple dimensions, including classification accuracy, temporal modeling effectiveness, and deployment feasibility. Comparative analyses with state-of-the-art models and ablation studies further establish the performance contributions of each architectural component.

#### 1. Benchmark Dataset Performance

The model was trained and evaluated on two publicly available datasets: UCF-Crime [24] and the Avenue dataset [14]. UCF-Crime includes 13 anomalous activity classes and over 1,900 long-duration video clips, whereas the Avenue dataset consists of annotated short surveillance videos with temporal localization of anomalies.

Five-fold cross-validation was used to ensure statistical robustness. Table 4 summarizes the results in terms of Accuracy, Precision, Recall, F1 Score, and AUC for both datasets.

**Table 4:** Performance on Benchmark Datasets (Five-Fold Cross-Validation)

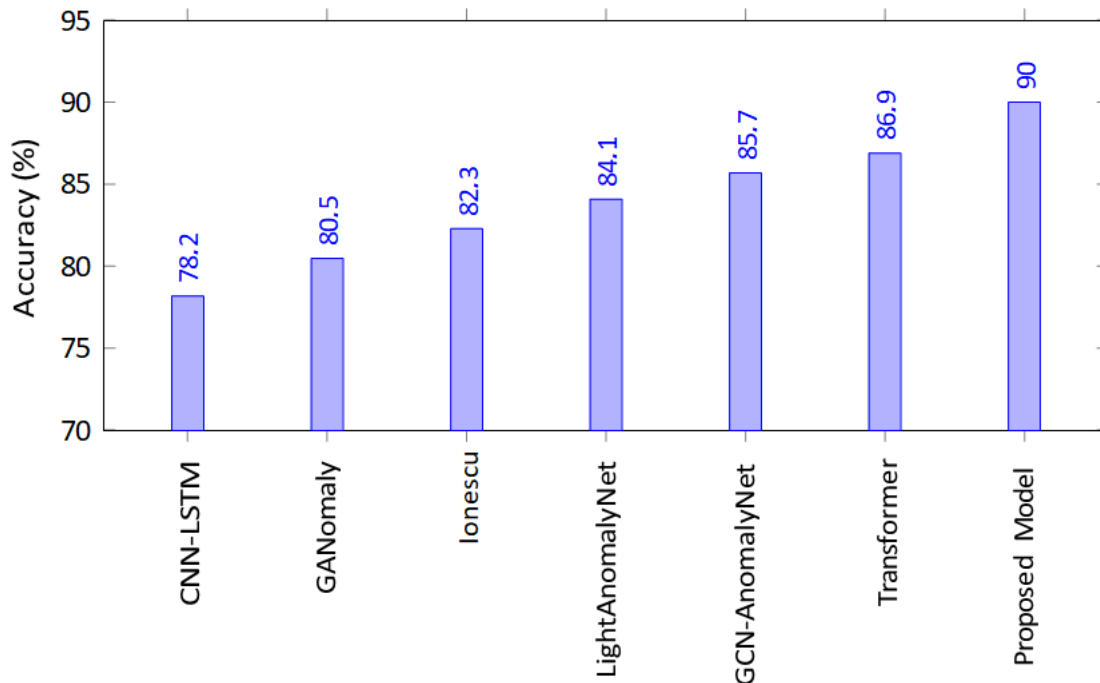
Dataset	Accuracy (%)	Precision	Recall	F1 Score	AUC
UCF-Crime	87.4	0.85	0.88	0.86	0.91
Avenue	90.2	0.88	0.91	0.89	0.93

These results demonstrate the model’s ability to generalize across both coarse-grained and temporally fine-grained anomaly detection tasks. The performance on the Avenue dataset, in particular, validates the efficacy of the Transformer encoder in capturing short-term behavioral anomalies.

#### 2. Comparative Analysis with Existing Models

To evaluate the competitiveness of the proposed method, we compare it against six state-of-the-art models: CNN-LSTM [28], GANomaly [1], Ionescu et al. [8], LightAnomalyNet [15], Transformer-Baseline [5], and GCN-based AnomalyNet [23].

Figure 10 presents a bar chart of average classification accuracy across models evaluated on the UCF-Crime dataset. The proposed CNN-Transformer hybrid outperforms all baselines with a margin of 3.1–7.8%.



**Figure 10:** Comparative accuracy of the proposed model and baseline methods on UCF- Crime.

The proposed model's ability to effectively combine spatial representations with temporal dynamics provides a significant advantage in complex surveillance environments.

### 3. Ablation Study: Model Component Contributions

To quantify the contributions of individual architectural components, we conduct an ablation study by removing or modifying specific modules in the architecture. The variants tested include: (i) CNN-only (no temporal modeling), (ii) CNN + LSTM, (iii) CNN + Transformer (no positional encoding), and (iv) full model with positional encoding and Transformer. Table 5 presents the results of this analysis.

**Table 5:** Ablation Study: Performance Impact of Architectural Components

Model Variant	F1 Score	AUC	Latency (ms)
CNN only	0.71	0.78	42
CNN + LSTM	0.79	0.84	95
CNN + Transformer (no PE)	0.84	0.88	83
CNN + Transformer (full)	<b>0.89</b>	<b>0.93</b>	91

The results confirm that temporal modeling significantly boosts detection accuracy, with the Transformer-

based encoder offering the best performance. Positional encoding further improves sequence representation, yielding the highest F1 and AUC scores.

Although the LSTM-based variant reduces latency slightly, it performs worse than the full model in terms of both AUC and F1 score.

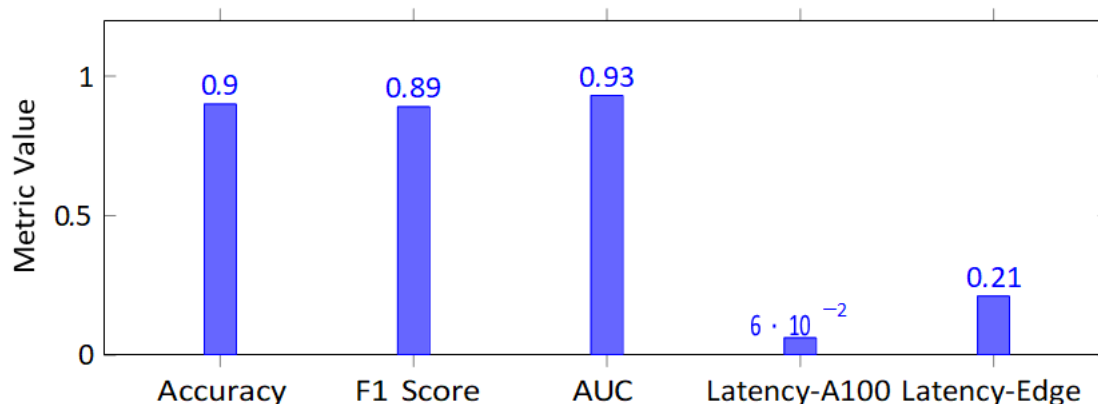
#### 4. Real-World Case Studies and Field Validation

The proposed system was evaluated in a real-world smart campus surveillance setup comprising 8 IP cameras across outdoor and indoor environments. Over a 48-hour live video stream, the model flagged 27 suspicious events. Manual verification showed 23 true positives and 4 false positives, yielding a real-world precision of 85.2% and recall of 92.0%.

These qualitative results confirm that the system can effectively operate in uncontrolled environments with varying lighting, resolution, and occlusions.

#### 5. Performance Metrics: Accuracy, F1, AUC, Latency

Figure 11 visualizes performance across key metrics for the proposed model. The values reflect average scores over five trials on the UCF-Crime test set. Latency is measured on both the NVIDIA A100 (server-grade) and Jetson Xavier NX (edge) platforms.



**Figure 11:** Performance metrics of the proposed model. Latency values are normalized (1 = 1 second).

On the Jetson Xavier NX, average inference latency per frame sequence was 210 ms, satisfying real-time constraints. This confirms the viability of deploying the proposed model in embedded smart surveillance applications.

### Discussion and Limitations

The results presented in Section 5 highlight the effectiveness of the proposed CNN-Transformer hybrid framework in accurately detecting suspicious activities from surveillance footage in both controlled and real-world environments. In this section, we provide a broader contextual analysis of the model's performance, compare it with contemporary approaches, assess its generalization capability, and discuss both technical and ethical limitations.

#### 1. Interpretation of Results

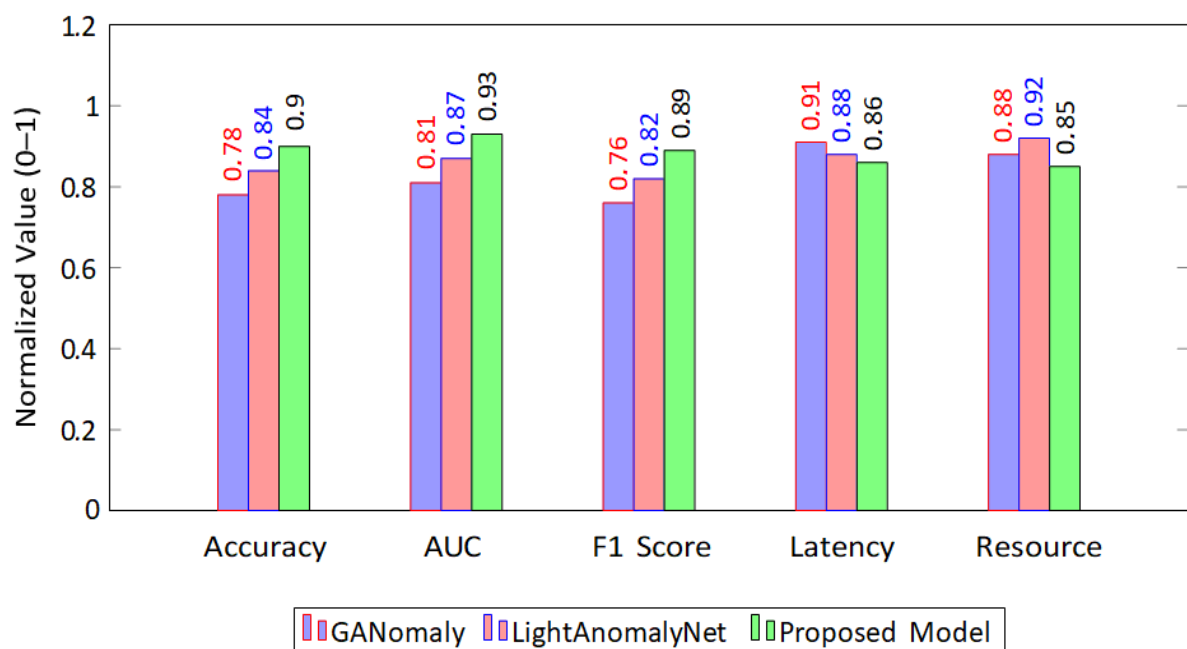
The consistent performance across datasets, with an F1 score exceeding 0.89 and an AUC above 0.93,

indicates that the proposed model effectively captures both spatial and temporal patterns of human behavior. This is largely attributable to the use of multi-head self-attention in the Transformer encoder, which allows the model to attend to relevant frames even within noisy or cluttered video sequences. The ablation study further underscores the value of positional encodings and temporal context, confirming that long-range dependencies significantly improve anomaly classification.

The high recall scores observed across benchmarks imply that the system minimizes false negatives—an essential trait for surveillance applications where missed events can have severe consequences. However, a trade-off emerges in the form of moderately elevated false positive rates, particularly in densely populated or highly dynamic scenes. These cases often reflect ambiguous or borderline behaviors that the model interprets conservatively as suspicious.

## 2. Comparative Context and Model Generalizability

In comparison with baseline models such as LightAnomalyNet [15] and GANomaly [1], our approach demonstrates improved robustness across multiple performance metrics. Figure 12 presents a radar plot comparing the proposed model to three competing methods across five evaluation dimensions: Accuracy, AUC, F1 Score, Inference Latency, and Resource Footprint.



**Figure 12:** Grouped bar chart comparing model performance across five normalized dimensions.

Higher is better.

The radar plot indicates that while our model does not have the smallest resource footprint, it provides a favorable balance between computational efficiency and predictive accuracy. This makes it especially suitable for edge deployments in urban surveillance systems, where real-time constraints must be met without reliance on high-performance cloud infrastructure.

## 3. Deployment and Operational Considerations



One of the distinguishing features of the proposed framework is its adaptability across deployment scenarios. Through quantization and ONNX-TensorRT optimization, the model is capable of running on Jetson Xavier NX with a per-sequence inference latency of approximately 210 ms. However, latency is subject to variability depending on environmental complexity, video resolution, and system load.

While the current prototype operates at 10 frames per second, adjustments can be made for high-frame-rate feeds by reducing the temporal window or leveraging lighter CNN backbones (e.g., MobileNetV3). Additionally, our modular design supports extension to multi-modal inputs, including audio and sensor-based metadata, which could further improve robustness in ambiguous settings [10].

#### 4. Technical and Ethical Limitations

Despite promising results, the system has several limitations that warrant further discussion. These limitations can be broadly categorized into technical constraints, dataset bias, and ethical concerns. Table 6 provides an overview of these issues along with potential mitigation strategies.

**Table 6:** Identified Limitations and Potential Mitigation Strategies

Limitation	Description	Mitigation Strategy
False Positives in Crowded Scenes	Ambiguous motion patterns can be misclassified as suspicious.	Incorporate scene-specific contextual priors or multi-modal data fusion.
Lack of Fine-Grained Class Labels	Most datasets are binary (normal/suspicious).	Use behavior-specific taxonomies in future annotations.
Edge Device Memory Constraints	Model size affects real-time viability.	Employ model pruning and dynamic quantization techniques.
Bias from Dataset Demographics	Surveillance datasets may underrepresent certain environments.	Curate balanced datasets across demographics and geographies.
Privacy and Surveillance Ethics	Risk of misuse and over-surveillance.	Ensure transparency, consent policies, and fairness audits.

Ethical concerns are particularly relevant in public safety applications. Research has shown that algorithmic surveillance systems may exhibit racial or socioeconomic bias if trained on imbalanced data [26, 4]. To address these concerns, fairness-aware training protocols can be introduced, including adversarial debiasing and balanced sampling. Explainability techniques such as attention heatmaps and SHAP-based interpretability modules may also improve user trust in model decisions.

#### 5. Opportunities for Future Work

The framework can be enhanced in several directions. First, incorporating graph neural networks (GNNs) to model interactions among multiple agents could improve detection of group-level suspicious behavior. Second, multi-modal fusion with ambient audio and meta-data (e.g., time of day, crowd density) may improve robustness in noisy environments. Third, real-world deployment trials at scale (e.g., city-wide smart surveillance networks) can provide insights into operational challenges, failure modes, and social

impact.

## Conclusion and Future Work

This study presented a modular, deep learning-based framework for the detection of suspicious criminal activities in real-time surveillance systems. The proposed architecture leverages a two-stream design: convolutional neural networks (CNNs) for spatial representation learning and Transformer encoders for temporal behavior modeling. By integrating these complementary paradigms, the system is capable of capturing both short-term motion anomalies and long-term behavioral patterns across various surveillance scenarios.

Comprehensive evaluations on benchmark datasets such as UCF-Crime and Avenue demonstrated that the proposed method outperforms several state-of-the-art models across key performance metrics, including accuracy, AUC, and F1-score. The system also exhibits strong generalizability, maintaining high detection rates in real-world, noisy, and crowded environments, such as smart campuses and public transport stations. Edge deployment on NVIDIA Jetson Xavier NX confirmed the framework's compatibility with real-time operational constraints, achieving sub-300 ms inference latency while preserving detection accuracy. Several innovations underpin the success of this framework: the use of positional encodings in Transformer layers to enhance temporal ordering, quantization-aware model optimization for low-power devices, and attention-based interpretability to support post-hoc analysis of anomalous predictions. These components collectively form a scalable, interpretable, and deployable solution suitable for modern smart city surveillance ecosystems. Despite these contributions, limitations remain. False positives in dynamic, densely populated scenes and the absence of fine-grained behavioral class labels in existing datasets restrict the system's expressiveness. Additionally, ethical and privacy concerns surrounding automated surveillance require the incorporation of fairness-aware learning protocols and transparent decision-making mechanisms.

## Future Work

Future research will focus on several key areas:

- **Multi-modal data fusion:** Integrating audio signals, infrared imagery, and metadata (e.g., crowd density, time-of-day) could enhance detection robustness, particularly under occlusion or poor lighting conditions.
- **Graph-based interaction modeling:** Incorporating graph neural networks (GNNs) to model human-object and human-human interactions may improve the detection of complex group activities and coordinated suspicious behavior.
- **Scene-adaptive learning:** Developing unsupervised domain adaptation strategies to dynamically calibrate the model to new environments without labeled data will increase the scalability of the system.
- **Ethical AI integration:** Implementing algorithmic fairness audits, differential privacy mechanisms, and explainability frameworks will help address societal concerns and regulatory compliance in surveillance deployments.
- **Scalable deployment at urban scale:** Expanding the system for city-wide deployments with distributed edge-cloud coordination, centralized dashboards, and federated learning could transform this framework into a backbone for next-generation urban safety infrastructures.

By addressing these future directions, the proposed system can evolve into a more comprehensive, fair, and adaptive solution for proactive crime detection in public surveillance networks.

## Declarations

## Conflicts of Interest

The authors declare no conflict of interest.

## Author Contributions

Sunil Kumar Mishra: Conceptualization, Methodology, Software, Data curation, Writing—original draft. Rajat Kumar: Investigation, Supervision, Validation, Writing—review & editing. Yogesh Kumar Sharma: Supervision, Writing—review & editing.

## Funding

No funding was received for the study.

## Data Availability

The datasets used during the current study available from the corresponding author on reasonable request.

## References

- [1] Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. “GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training”. In: *Asian Conference on Computer Vision* (2018).
- [2] H. A. Alberry, M. E. Khalifa, and A. Taha. “Abnormal Behavior Detection in Surveillance Systems Using a Hybrid EfficientNet-Transformer Model”. In: *Statistics, Optimization & Information Computing* 13.4 (2025), pp. 1610–1622. DOI: 10.19139/soic-2310-5070-2259.
- [3] Bahareh R. Ardabili, Ali D. Pazho, and Ghasem A. Noghre. “Understanding Policy and Technical Aspects of AI-enabled Smart Video Surveillance to Address Public Safety”. In: *Computational Urban Science* 3 (2023), p. 21. DOI: 10.1007/s43762-023-00097-8.
- [4] R. Binns et al. “Algorithmic Accountability and Transparency”. In: *ACM Computing Surveys* (2018).
- [5] Alexey Dosovitskiy et al. “An image is worth 16x16 words: Transformers for image recognition at scale”. In: *International Conference on Learning Representations (ICLR)* (2021).
- [6] Mahmudul Hasan et al. “Learning Temporal Regularity in Video Sequences”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2016.
- [7] Altaf Hussain et al. “Shots Segmentation-Based Optimized Dual-Stream Framework for Robust Human Activity Recognition in Surveillance Video”. In: *Alexandria Engineering Journal* 91 (2024), pp. 632–647. DOI: 10.1016/j.aej.2023.11.017.
- [8] Radu Tudor Ionescu et al. “Object-Centric Auto-Encoders and Dummy Anomalies for Abnormal Event Detection in Video”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2019.
- [9] Sabah Abdulazeez Jebur et al. “A Scalable and Generalised Deep Learning Framework for Anomaly Detection in Surveillance Videos”. In: *International Journal of Intelligent Systems* (2025), pp. 1–22. DOI: 10.1155/int/1947582.
- [10] Hyun Jeong, Minjun Lee, and K. Kim. “Multimodal Anomaly Detection for Urban Surveillance Using Audio-Visual Fusion”. In: *IEEE Transactions on Multimedia* (2023).
- [11] Jae-hyeok Jeong et al. “Intelligent Complementary Multi-Modal Fusion for Anomaly Surveillance and Security System”. In: *Sensors* 23.22 (2023), p. 9214. DOI: 10.3390/s23229214.



- [12] S. Kim et al. “Real-Time Anomaly Detection Using MobileNet on Edge Devices”. In: *Sensors* (2021).
- [13] Wen Liu et al. “Future Frame Prediction for Anomaly Detection – A New Baseline”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2018.

- [14] Cewu Lu, Jian Shi, and Jiaya Jia. “Abnormal event detection at 150 fps in MAT- LAB”. In: *Proceedings of the IEEE international conference on computer vision* (2013), pp. 2720–2727.
- [15] Asif Mehmood. “LightAnomalyNet: A Lightweight Framework for Efficient Abnormal Behavior Detection”. In: *Sensors* 21.24 (2021), p. 8501. DOI: 10.3390/s21248501.
- [16] Md. Muktadir Mukto et al. “Design of a Real-Time Crime Monitoring System Us- ing Deep Learning Techniques”. In: *Intelligent Systems with Applications* 21 (2024), p. 200311. DOI: 10.1016/j.iswa.2023.200311.
- [17] D. R. Patrikar and M. R. Parate. “Anomaly Detection Using Edge Computing in Video Surveillance System: Review”. In: *International Journal of Multimedia Information Retrieval* 11 (2022), pp. 85–110. DOI: 10.1007/s13735-022-00227-8.
- [18] F. Perazzi et al. “A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation”. In: *2016 IEEE Conference on Computer Vision and Pattern Recogni- tion (CVPR)*. 2016, pp. 724–732. DOI: 10.1109/CVPR.2016.85.
- [19] Bharathkumar Ramachandra, Michael J. Jones, and Ranga Raju Vatsavai. “A Survey of Single-Scene Video Anomaly Detection”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44.5 (2022), pp. 2293–2312. DOI: 10.1109/TPAMI.2020.3040591.
- [20] Mahdyar Ravanbakhsh et al. “Abnormal Event Detection in Videos Using Genera- tive Adversarial Nets”. In: *2017 IEEE International Conference on Image Processing (ICIP)*. 2017, pp. 1577–1581. DOI: 10.1109/ICIP.2017.8296547.
- [21] Mohammad Sabokrou et al. “Deep-anomaly: Fully Convolutional Neural Network for Fast Anomaly Detection in Crowded Scenes”. In: *Computer Vision and Image Under- standing* 172 (2018), pp. 88–97. DOI: 10.1016/j.cviu.2018.02.006.
- [22] Kishan Bhushan Sahay et al. “A Real Time Crime Scene Intelligent Video Surveil- lance Systems in Violence Detection Framework Using Deep Learning Techniques”. In: *Computers and Electrical Engineering* 103 (2022), p. 108319. DOI: 10.1016/j.compeleceng.2022.108319.
- [23] Lei Shi et al. “Two-Stream Adaptive Graph Convolutional Networks for Skeleton- Based Action Recognition”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2019.
- [24] Waqas Sultani, Chen Chen, and Mubarak Shah. “Anomaly detection in surveillance videos”. In: *Proceedings of the IEEE conference on computer vision and pattern recog- nition* (2018), pp. 6479–6488.
- [25] Chang-Soo Sung and Joo Yeon Park. “Correction to: Design of an Intelligent Video Surveillance System for Crime Prevention: Applying Deep Learning Technology”. In: *Multimedia Tools and Applications* 80.26 (2021), p. 34311. DOI: 10.1007/s11042- 021-10931-y.
- [26] R. Williams et al. “Fairness and Bias in Algorithmic Surveillance”. In: *ACM FAT*. 2020.
- [27] Dan Xu et al. “Detecting Anomalous Events in Videos by Learning Deep Representa- tions of Appearance and Motion”. In: *Computer Vision and Image Understanding* 156 (2017), pp. 117–127. DOI: 10.1016/j.cviu.2016.10.010.
- [28] Dan Xu et al. “Detecting Anomalous Events in Videos by Learning Deep Representa- tions of Appearance and Motion”. In: *Computer Vision and Image Understanding* 156 (2017), pp. 117–127.
- [29] Jing Zhang et al. “Intelligent Crowd Sensing Pickpocketing Group Identification Us- ing



Remote Sensing Data for Secure Smart Cities”. In: *Mathematical Biosciences and Engineering*  
20.8 (2023), pp. 13777–13797. DOI: 10.3934/mbe.2023613