



Improving Public Safety: An AI Approach to Anomaly Detection in Surveillance Footage

PROJECT SUPERVISOR

Muhammad Nouman Durrani

PROJECT TEAM

Taha Ahmed (K21-4833)

Mahad Munir (K21-3388)

Asad Noor Khan (K21-4678)

Submitted in partial fulfillment of the requirements for the degree of Bachelor of Science in Computer Science.

FAST SCHOOL OF COMPUTING

**NATIONAL UNIVERSITY OF COMPUTER AND EMERGING SCIENCES KARACHI
CAMPUS
June 2023**

Project Supervisor	Muhammad Nouman Durrani
Project Team	Taha Ahmed (21K-4833) Mahad Munir (21K-3388) Asad Noor Khan (21K-4678)
Submission Date	May 15, 2025,

Muhammad Nouman Durrani _____
Supervisor

Dr. Ghufraan Ahmed _____
Head of Department

FAST SCHOOL OF COMPUTING
NATIONAL UNIVERSITY OF COMPUTER AND EMERGING SCIENCES KARACHI
CAMPUS

Acknowledgment

We would like to thank our supervisor, **Muhammad Nouman Durrani**, for guiding and nurturing us throughout the entire final year project and helping us achieve our aims and objectives. Our jury members are also to be thanked for pointing out our mistakes and motivating us to develop and produce exact outcomes and output. We'd also like to thank our parents for their unwavering support. We are also appreciative to our colleagues for their invaluable advice and suggestions based on their Machine Learning and Computer vision experience.

Abstract

Video surveillance has become a part of daily life, with cameras installed in public spaces, businesses, and transportation hubs to enhance security and monitor activities in real time. The primary goal of these systems is to deter and detect crimes such as theft, violence, vandalism, and terrorism. However, manually monitoring countless surveillance feeds 24/7 is not only resource-heavy but also susceptible to human error—making it difficult to respond quickly and accurately to critical incidents.

To tackle these challenges, our Final Year Project introduces an intelligent, automated system for real-time anomaly detection in CCTV footage using advanced computer vision. Our approach leverages hybrid vision transformers to analyze both contextual and temporal patterns in video feeds, enabling precise identification of abnormal behaviors like robberies, assaults, accidents, and other unusual activities. These transformers are trained to understand both broad scene context and fine-grained motion details, which is crucial for distinguishing between normal and suspicious actions in complex real-world environments.

Beyond just detecting anomalies, our system also categorizes them, helping security personnel make faster, more informed decisions. To streamline usability, it automatically extracts and highlights the exact moments in the video where anomalies occur—eliminating the need to sift through hours of footage and drastically reducing response times.

While AI-powered anomaly detection is becoming more common, many existing systems still struggle with high false alarms, limited scalability, and inconsistent performance under different conditions. Our project aims to overcome these shortcomings with a more robust and context-aware solution, contributing to the broader goal of public safety through intelligent video analysis. Ultimately, it demonstrates how cutting-edge AI models can align with real-world surveillance needs to create safer, more responsive urban environments.

Contents

Page#

<u>Introduction</u>	6
<u>Related Work</u>	7
<u>Literature review</u>	8
<u>Proposed Approach</u>	15
<u>Data Set</u>	17
<u>Gui</u>	18
<u>Results</u>	20
<u>References</u>	23

Introduction

Anomaly detection in video surveillance is a crucial task in computer vision, aimed at identifying unusual or unexpected activities that deviate from normal behavioral patterns. This project presents an AI-powered system designed to detect anomalies from user-uploaded video data using advanced deep learning techniques.

The system leverages a **Two-Stream architecture** that combines **Inflated 3D ConvNet (I3D)** for capturing spatio-temporal features and a **Vision Transformer (ViT)** for extracting global contextual information. This combination enhances the model's ability to understand both motion and visual patterns, resulting in more accurate and context-aware anomaly detection.

To handle **data imbalance**, which is a common issue in anomaly detection (since anomalies are rare), we used techniques like:

- **SMOTE (Synthetic Minority Oversampling Technique)** to create more samples of rare events,
- **Data augmentation** (rotations, flips, noise, etc.) to enhance training robustness,
- **Weighted class training** to ensure the model gives proper importance to rare anomaly cases during learning.

Users interact with the system through a **React-based web interface**, where they can upload video files and initiate analysis. The backend processes the videos by extracting frames, preprocessing them, and passing them through the AI model to identify frames that exhibit anomalous behavior. Detection results are visualized on the interface for user interpretation.

This system is modular, allowing easy integration of future enhancements and adaptable to various application domains such as public safety, industrial automation, and behavior monitoring. The approach aims to deliver a scalable, accurate, and user-friendly solution for anomaly detection in visual data.

Related Work

- **Tube-Convolutional Neural Network tested to detect events.**

[6] This white paper talked about the early CNN-based detection approaches which consisted of two major steps: frame-level action proposal generation and association of proposals across frames. The method introduced in this paper is an end-to-end deep network called Tube Convolutional Neural Network (T-CNN) and it is based on 3D convolutional networks.

In its initial phase, the video is divided into equal-length clips, and next for each clip a set of tube proposals are generated based on 3D Convolutional Network (ConvNet) features. Experiments were performed to test the efficiency of the proposed method using the UCF- sports and other datasets. The results were better than some of the references provided, the efficiency was lower in Driving, lifting, and Swing but better in Golf, kicking, riding, etc.

- **The idea of temporal constraints in video prediction.**

[7] This research paper proposed a solution to mitigate the anomaly detection problem within a video prediction framework to identify abnormal events by comparing them with their expectations. Writers introduced temporal constraints into video prediction tasks to predict a future frame with higher quality. The idea basically was if a frame is similar to the frame that is predicted by the model then it is most probably a normal event, Otherwise it represents the possibility of an abnormal event.

This research proved that adding an optical flow constraint into the objective function to guarantee motion consistency for normal events would boost the performance of anomaly detection. Their approach mathematically was given a video with consecutive t frames I_1, I_2, I_3 etc, they stack all these frames and use them to predict the I_{t+1} frame. To preserve temporal constraint between two consecutive frames they enforced optical flow between frame I_{t+1} and I_t and I_{t+1} (ground truth) to I_t to be close. The model was trained and evaluated on different well-known datasets(UCSD, ShanghaiTech), The AUC score was the best (94.5% in UCSD and 72.8 in ShanghaiTect) which was the best among state-of-art approaches of that time.

- **Multi-kernel-based vector classification, SVM (Support Vector Machine)**

[1] In the main paper, feature extraction is followed by an attention layer based on dilated convolution to capture the most relevant long and short-range dependencies, referenced in the above paper. Context awareness in pervasive health is a proactive approach, rather than a conventional event-driven model. Basically, context awareness is the ability of a system to gather information about its environment at a given time and adapt behaviors accordingly. In this research Patient activity and gender classification is done using multi-kernel-based vector data classification. Multi-sensor fusion is an effective use of multiple sources of information i.e. competitive, complementary, and cooperative fusion.

The classification results indicate that Multi-Class SVM (with equal or symmetric misclassification cost) has maximum classification accuracy when tested with the Gaussian kernel function. The results of classification may be clinically correlated by the health care experts.

- **SVM used with UMN and BEHAVE datasets.**

[2] In this research, a new method is proposed to detect abnormal behaviors in human group activities as abnormal behaviors in a real-world environment are difficult to classify. This approach effectively models group activities based on social behavior analysis. Different from previous work that uses independent local features, our method explores the relationships between the current behavior state of a subject and its actions. An interaction energy potential function is proposed to represent the current behavior state of a subject, and velocity is used as its action. Our method does not depend on human detection or segmentation, so it is robust to detection errors. Instead, tracked spatiotemporal interest points are able to provide a good estimation of modeling group interaction. SVM is used to find abnormal events. We evaluate our algorithm in two datasets: UMN and BEHAVE. Experimental results show its promising performance against state-of-art methods.

The algorithm is tested on two datasets: the UMN dataset and BEHAVE dataset. Results showed the effectiveness of our method, and it is competitive with the state-of-art methods.

- **Sparse Combination Learning for Abnormal Videos**

[3] Based on the inherent redundancy of video structures, this paper proposed an efficient sparse combination learning framework. It achieved decent performance in the detection phase without compromising the result's quality. The short running time is guaranteed because the method effectively turned the original complicated problem into one in which only a few costless small- scale least square optimization steps are involved. Their method reached high detection rates on benchmark datasets at a speed of 140~150 frames per second on average when computing on an ordinary desktop PC using MATLAB.

An abnormal event detection method via sparse combination learning. This approach directly learns sparse combinations, which increase the testing speed hundreds of times without compromising effectiveness, it achieves state-of-the-art results in several datasets. It is related to but differs largely from traditional subspace clustering.

- **CNN (convolutional neural network) with (NetVLAD)Vector of Locally Aggregated Descriptor.**

[4] This white paper discussed the technique to quickly and accurately recognize the location of a given query image. They designed a new generalized VLAD layer NetVLAD which is now mostly used in image retrieval tasks. This layer was used on top of the newly developed CNN model. The backpropagation approach was used to make their model more efficient. First, the features are extracted from the query image and used to search a large scaled geotagged database.

This research gathered a large dataset from Google Street View Time Machine to train their CNN architecture for place recognition. The results were evaluated on two publicly available datasets(Pitts 250k) and Tokyo 24/7 which contains 76k database images for training.

- **3-D Convolutional Network (I3D)**

[5] In this specific research paper, the I3D ConvNet model is discussed in detail which is basically used as a backbone for feature extraction in transformer-based anomaly detection. Two-Stream Inflated 3D ConvNet (I3D) that is based on 2D ConvNet inflation: filters and pooling kernels of very deep image classification ConvNets are expanded into 3D, making it possible to learn seamless Spatio-temporal feature extractors from video while leveraging successful ImageNet architecture designs and even their parameters. We show that, after pre-training on Kinetics, I3D models considerably improve upon the state-of-the-art in action classification, reaching 80.2% on HMDB-51 and 97.9% on UCF-101.

- **(MIL) Multiple Instance Learning; 2 instances of normal and abnormal bags of videos.**

[8] This research paper focused on the concept of anomaly detection in videos using weakly supervised data. Video-level annotation is provided in the training dataset rather than annotating the anomalous segments in videos, which is time-consuming and vague. The writers also contributed to constructing a new large-scale dataset(1600 videos) to train and evaluate their model. They also selected 13 impactful anomalies that would be predicted by their model, some of which are Abuse, Arrest, Robbery, Explosion, and fighting.

In this paper Multiple instance learning(MIL) method was introduced for anomaly detection, the main idea was to basically consider two bags for video (normal video and anomalous video) and consider video segments as instances. Whatever bag has a higher score the video is detected as that labeled bag. Their model predicted the score for normal and anomalous video segments and filled the bag accordingly. To extract the features from the video they used a pre-trained C3D feature extractor. The experimental results on the dataset they constructed show that their proposed anomaly detection approach performs significantly better than state-of-art methods. AUC score on UCF-Crime benchmark dataset was 0.75

- **MIL combined with 3D resNet for spatial-temporal features**

[9] This white paper focuses on minimizing the false alarm rate that can rise due to missing temporal annotations of a video while performing an abnormal activity detection task in a weakly supervised paradigm. The writers proposed a method named 3D ResNet to extract spatial-temporal features from videos and also a ranking loss function to predict the anomalous score of the video. Features are extracted from 3D ResNet and are fed to the Deep MIL model. For extracting features from 3D ResNet each video is divided into equal numbers of non-overlapping video segments.

While Training each video is divided into non-overlapping 32 video segments then the 3D ResNet model extracts features from frames. The proposed approach is evaluated on the UCF-Crime dataset. Results show that Deep Multiple Instance learning combined with the proposed

3D ResNet and a ranking loss function achieves the best performance on the UCF-Crime benchmark dataset.

Methods (unsupervised detection) and results showed improvements over the state-of-the-art methods on real-world datasets. The researchers experimented with their methods on the MNIST and CIFAR10 datasets. Which proved to be better than others (others: 0.9750, theirs: 0.9935).

[10] In this research paper, anomalous video detection is approached through self-supervised and multi-task learning at the object level. Although at the end of the papers frame level implementation was also implemented with their AUC scores. The pre-trained detector was used for object detection, and 3D CNN for producing discriminative anomaly-specific information by jointly learning multiple proxy tasks: three self-supervised and one knowledge-based distillation were used. 3 self-supervised (discrimination of forward/backward moving objects (arrow of time), discrimination of objects in consecutive frames (motion irregularity), and reconstruction of object-specific appearance information). And this method outperformed art of the state methods on three benchmarks; Avenue, ShanghaiTech, and UCSD Ped2. Scores are as follows: AUC 99.8, RBDC 72.8, and TBDC 91.2. (At object level).

- **RFTM(Robust Temporal Feature Magnitude Learning)-enabled MIL model**

[11] This white paper focuses on the Limitations of the MIL approach and the importance of video temporal dependencies. In MIL, the recognition of positive instances, i.eRare abnormal snippets in abnormal videos is largely biased by dominant negative instances due to the fact that it ignores video temporal dependencies. The proposed method RFTM improves the robustness of the MIL approach to negative instances from abnormal videos. RFTM adapts a self-attention mechanism to capture short and long-temporal dependencies.

The video is divided into 32 video snippets. Three FC layers followed by the ReLu activation function and dropout function with a dropout rate of 0.7 is used. The features are extracted using pre-trained I3D and C3D networks. Results show that the RFTM-enabled MIL model (i) outperforms several state-of-the-art methods by a large margin on four benchmark data sets (ShanghaiTech, UCF-Crime, XD-Violence, and UCSD-Peds) and (ii) achieves significantly improved efficiency.

- **Two-level attention mechanism.**

[12] This paper proposed a method that jointly handles anomaly detection and classification in a single framework by adapting a weakly supervised (where the training dataset has only video- level annotation) learning paradigm. In the proposed model a two-level attention mechanism is incorporated to focus on video clips containing anomalies for learning discriminative representations. The first-level attention mechanism highlights the clips pertaining to anomalies for the task of detection. Moreover, a relation between the anomaly detection and classification branch is developed through another level of attention mechanism.

During training, features are extracted using the I3D model. The model is validated on a large-scale publicly available dataset named UCF-Crime that consists of 1900 videos. This model achieved state-of-art results for classification while jointly performing the task of anomaly detection and classification. This model scored 82.12 AUC score on UCF Dataset.

- **Multi-modal semi-supervised deep learning-based CNN-BiLSTM autoencoder**

[13] This research paper presented a multi-modal semi-supervised deep learning-based CNN- BiLSTM autoencoder framework to detect anomalous events which use CNN for feature extraction. This proposed solution does not overcome state-of-the-art methods but provides results almost as same as other methods. Moreover, it used multi-modal RGB + D data for real-time online video anomaly detection in surveillance videos. It also contributed to providing a unique multi-modal RGBpD (RGB, Depth, Skeleton) dataset for testing and evaluation of anomaly detection methods in the Bank-ATM environment as it is mostly on ATM environments. The result generated on the Avenue dataset was 89.1, which is the second highest of all other state-of-the-art methods. And secured the top position with a score of 91.1 on the UCF-Crime Local dataset.

- **The kinetics dataset with diverse human actions**

[14] An extension of the DeepMind Kinetics human action dataset from 600 classes to 700 classes, where for each class there are at least 600 video clips from different YouTube videos. The goal of the Kinetics project is to provide a large-scale curated dataset of video clips, covering a diverse range of human actions, that can be used for training and exploring neural network architectures for modeling human actions in the video.

There was a first Kinetics challenge at the ActivityNet workshop in CVPR 2017, The second challenge occurred at the ActivityNet workshop in CVPR 2018. The performance criterion used in the challenge is the average of Top-1 and Top-5 errors. There was an improvement between the winning systems of the two challenges, with errors getting down from 12.4% (in 2017) to 11.0% (in 2018) [1, 6]. The 2019 challenge featured the new Kinetics-700 dataset and had 15 participating teams. The top team was from JD AI Research and obtained a 17.9% error, considerably below our baseline – a single RGB I3D model – which obtained a 29.3% error.

- **Anomalies classified as per the ‘Student-Teacher Framework’**

[15] This white paper focuses on detecting anomalies using a framework called the ‘Student- Teacher Framework’. Which is on supervised anomaly detection and pixel-precise anomaly segmentation in high-resolution images. Where anomalies are detected when outputs of the student networks differ from that of the teacher network. The methods were compared to a large number of existing deep learning based

[16] In this research, the difference between normal and anomalous behavior is described in detail but it did not just talk about this it described the differences between one anomaly and another because in different environments; an anomalous event can be considered as a normal one for example in Olympic; a fire is shot before the race which is not an anomalous behavior in that particular environment.

This white paper introduced different deep-learning techniques to detect anomalous events in surveillance videos and described the opportunities and challenges in detecting an anomaly. The different models it included were Reconstruction based models, Predictive models, Generative models, One-Class classification models, and Hybrid models.

- **CBAM(convolutional block attention model) with ResNet**

[17] Adversarial attacks against deep learning models have gained significant attention and to defend against them adversarial techniques do exist. The paper discusses the adversarial robustness of attention and non-attention-based classification models i.e. datasets with less number of classes, attention-based models show better adversarial robustness. This paper studies the impact of learning salient features through attention mechanisms on the adversarial robustness of a model using image classification.

The results obtained for CIFAR-10 and Fashion MNIST dataset strongly favor the robustness of the non-attention model i.e. Resnet-50. However, for the CIFAR-100 dataset, the attention-based model i.e. CBAM+Resnet-50 was comparatively much more robust than the non-attention model Resnet-50 under white box attacks while the difference in results between the two models was very low in the transfer-based black box attacks. These results suggest that the impact of attention on the robustness of image classification models could be dependent on the datasets used i.e. the number of classes in the datasets.

- **Transformer Model using Video Swin.**

[18] In this research, a new architecture is introduced which is the ‘Transformer’. The paper advocates an inductive bias of locality in the video Transformers, which leads to a better speed-accuracy trade-off compared to previous approaches. The proposed approach achieves state-of-the-art accuracy when tested on video recognition benchmarks i.e. 84.9 top-1 accuracies on Kinetics-400 and 86.1 top-1 accuracies on Kinetics-600.

It all started with the Vision transformer. It further introduced the Swin transformer; it once served as a general-purpose backbone for various image recognition tasks, which led to research for video-based recognition using tasks using the Transformer. In this paper, firstly spatio temporal locality is discussed, and then empirically shown that the Video Swin Transformer with spatiotemporal locality bias surpasses the performance of all the other vision Transformers on various video recognition tasks.

TABLE 1: Summary of all approaches

Method	AUC
1) Efficient Video Classifier	Accuracy is 64.17%
2) Hightrain Accuracy Model	Accuracy is 70.31%
3) Artificial Neural Network (ANN)	Accuracy is 75.31% Precision is 0.72 F1-Score is 0.71
4) Convolutional Neural Network (CNN).	Accuracy is 85.6% Precision is 0.84 F1-Score is 0.83
5) 3D ResNet to extract spatial-temporal features	AUC on the UCF-Crime is 0.76
6) Two-Stream Inflated 3D ConvNet (I3D)	Accuracy is 81.2% Precision is 0.8 F1-Score is 0.72

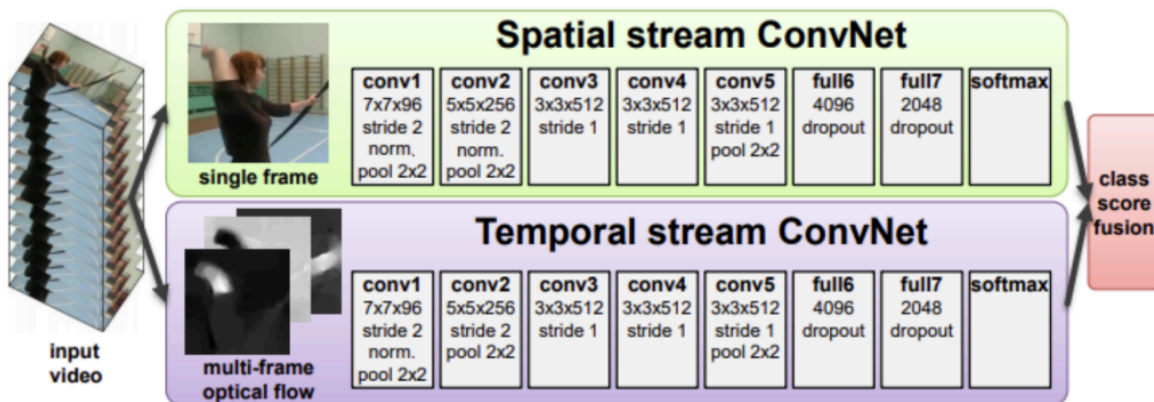
Proposed Approach

Our anomaly detection system uses a powerful two-stream approach that combines the strengths of **I3D (Inflated 3D ConvNet)** and the **Vision Transformer (ViT)** to analyze videos for unusual behavior. First, each video is divided into smaller segments, and from these, two types of features are extracted using I3D: one stream focuses on regular RGB frames to capture what's happening visually, and the other uses optical flow to understand how movement is happening over time. Unlike traditional models, I3D processes both space and time together, allowing it to better recognize actions and motion patterns.

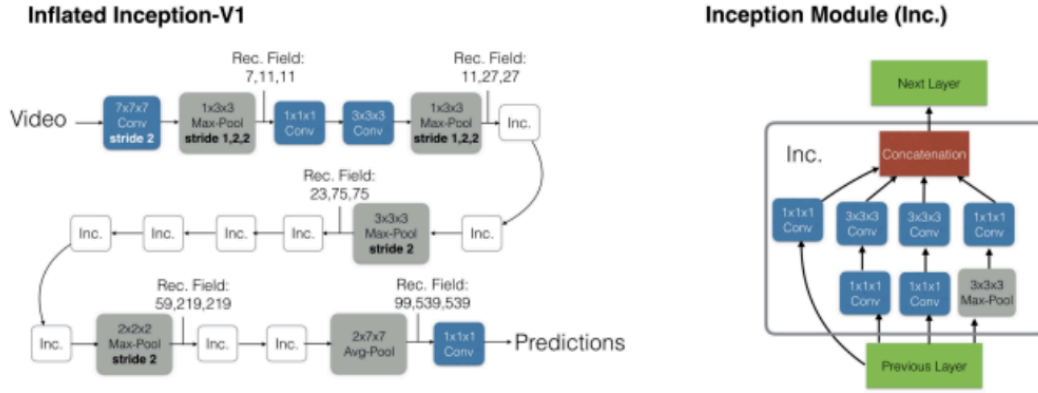
Once we've gathered these features, we feed them into the Vision Transformer, which acts like a second brain that pays attention to important moments across the entire video. The transformer looks at all the video segments as a sequence and decides which ones are more likely to contain something unusual. It does this through a mechanism called self-attention, which lets the model weigh different parts of the video differently based on their importance.

To deal with the challenge of imbalanced data—where anomalies are much rarer than normal activities—we applied techniques like SMOTE (to synthetically generate more anomaly samples), data augmentation (such as flipping and rotating frames), and weighted loss functions (to give more focus to the rare anomaly cases during training).

Two-stream structure is shown below:



The structure of Inflated 3D Convolutional Neural Network:



Let's talk about the ViT Anomaly Detection model. We used ReLU activation on hidden layers which is defined as

$$f(x) = \max(0, x)$$

For the output layer, we used Softmax as an activation function which is defined as

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

We incorporated dropout at various stages of our model to reduce overfitting and improve generalization. Additionally, filtering techniques were applied during preprocessing to enhance relevant visual information and suppress noise in video frames. Filtering, in this context, involves mathematical operations that refine the input signal (video frames) by emphasizing critical motion or visual patterns, which helps in improving the quality of features passed to the model.

Our approach enables more effective anomaly detection using a weakly supervised setting, leveraging the strength of both the **I3D** and **Vision Transformer** models. The I3D network extracts rich spatial-temporal features from both RGB frames and optical flow (capturing movement), while the Vision Transformer applies self-attention to model long-range dependencies across video snippets.

To handle the challenge of detecting rare anomalous events in videos that mostly contain normal activity, we implement a mechanism that focuses on identifying the top-k most informative segments. This ensures that even subtle anomalies, which may otherwise go unnoticed, are emphasized during training. The system learns to create a clear separation between normal and abnormal snippets, improving the model's precision in identifying complex and subtle anomalies.

This combination leads to two major benefits:

1. The system becomes better at distinguishing subtle anomalies from hard-to-classify normal behaviors.

2. It efficiently utilizes limited anomalous data for training, improving learning despite weak supervision.

As a result, our model demonstrates strong performance, achieving **81.2% accuracy**, **0.80 precision**, and an **F1-score of 0.72**, validating its effectiveness in real-world surveillance anomaly detection scenarios.

Dataset

UCF Crime Dataset:

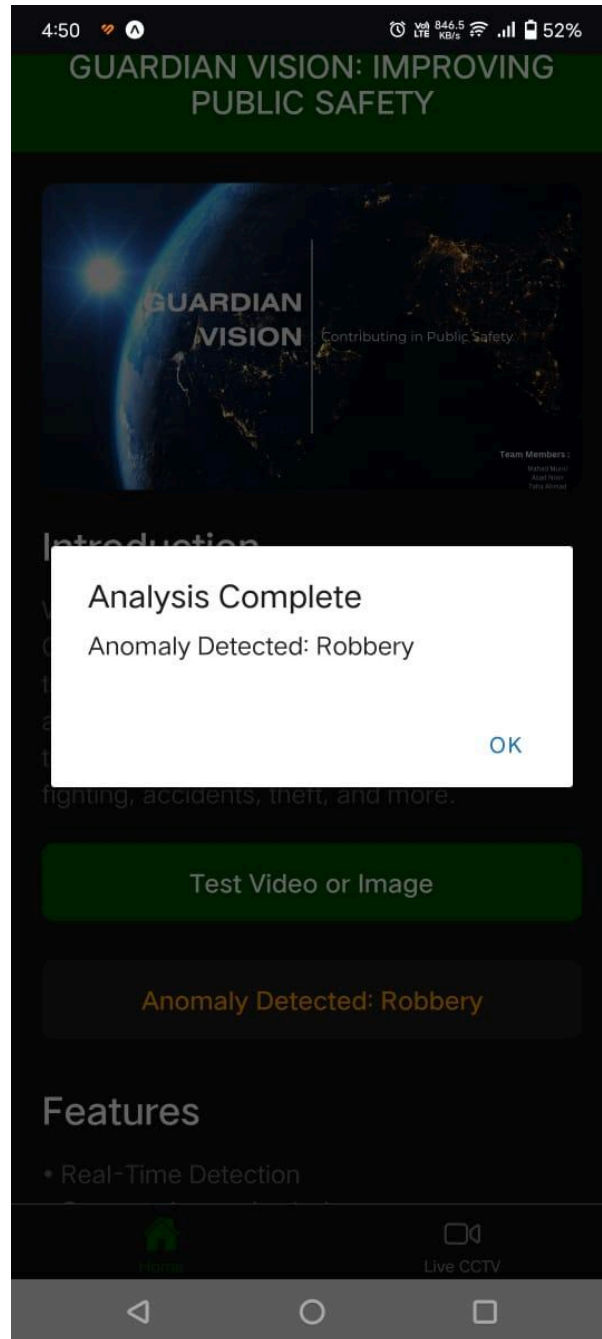
We construct a new large-scale dataset, called UCF-Crime, to evaluate our method. It consists of long untrimmed surveillance videos which cover 13 real-world anomalies, including Abuse, Arrest, Arson, Assault, Road Accident, Burglary, Explosion, Fighting, Robbery, Shooting, Stealing, Shoplifting, and Vandalism. These anomalies are selected because they have a significant impact on public safety. We compare our dataset with previous anomaly detection datasets in Table 1. For more details about the UCF-Crime dataset, please refer to our paper. Each anomalous event is given below.

1. Abuse
2. Arrest'
3. Arson
4. Assault
5. Burglary
6. Explosion
7. Fighting
8. Normal Videos
9. RoadAccidents
10. Robbery
11. Shooting
12. Shoplifting
13. Stealing
14. Vandalism

We Used an 80/20 train test split. UCF-Crime contains a total of 1920 videos with 960 being normal videos and 960 being abnormal videos.

GUI

To present the inference results of our proposed method, we use a **React-based frontend**, available on both **web and mobile platforms**. After processing the uploaded video and detecting anomalies using our AI model, the system displays the results directly to the user. Instead of modifying or generating a new video, the interface shows a **textual notification** indicating **which type of anomaly** has been detected



Ex: Person Entering a Grocery Store (Robbery Anomaly).

Results and Comparison

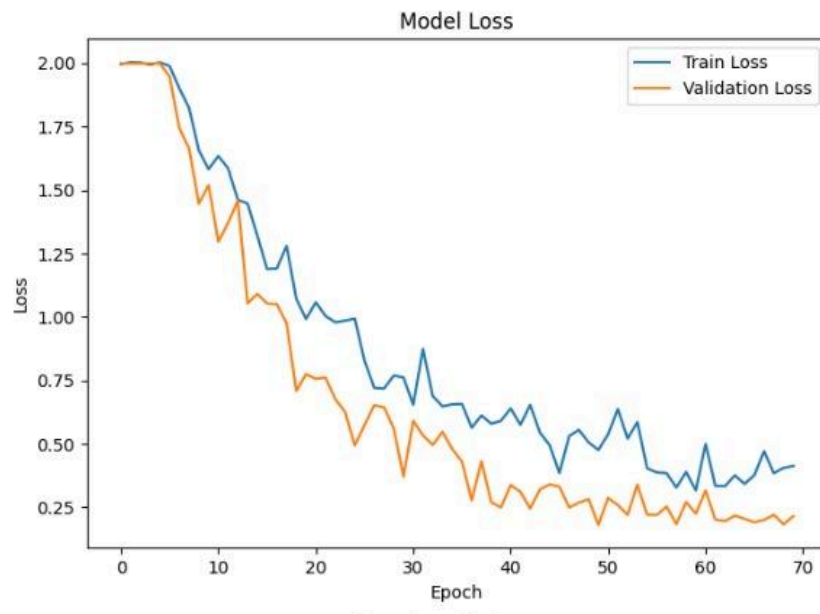
In this section, we present the results of our experiments conducted on the **UCF-Crime** dataset. The purpose of these experiments is to evaluate and compare the performance of **three different techniques** used for anomaly detection in surveillance videos: **Efficient Video Classifier**, **Hightrain Accuracy Model**, and **Two-Stream Inflated 3D ConvNet (I3D)**.

Each technique was tested on the same dataset to ensure a consistent comparison, and the results were analyzed using standard evaluation metrics. The performance of these methods was assessed to understand their capability in identifying anomalous behavior in complex real-world scenarios. The comparison helps highlight the strengths and trade-offs of each technique in terms of accuracy, consistency, and detection effectiveness.

We compared AUC scores of different techniques as shown in the below table

Methods	AUC
1) Efficient Video Classifier	0.6417
2) Hightrain Accuracy Model	0.7031
3) Two-Stream Inflated 3D ConvNet (I3D)	0.8121

In the images below, we have illustrated the training performance of our model through the loss graph, which shows how the model's error decreased over time. Additionally, we present the output from the last few epochs to highlight the model's final stage of convergence and performance stability before evaluation



```
➡ Train Loss: 0.4826 Acc: 0.8541
  Val Loss: 0.2683 Acc: 0.8167 F1: 0.8562
  Epoch 47/100
  -----
  Train Loss: 0.4778 Acc: 0.8333
  Val Loss: 0.2227 Acc: 0.8250 F1: 0.8213
  Epoch 48/100
  -----
  Train Loss: 0.3777 Acc: 0.8279
  Val Loss: 0.1617 Acc: 0.8067 F1: 0.8278
  Epoch 49/100
  -----
  Train Loss: 0.2677 Acc: 0.8250
  Val Loss: 0.3791 Acc: 0.8150 F1: 0.8323
  Epoch 50/100
  -----
```

Conclusion

Throughout this project, we observed that existing anomaly detection methods each have their own strengths and weaknesses. Key limitations such as high false alarms, struggles with dynamic environments, and inefficiencies in processing live video motivated our research and shaped our system's design. Building on these insights, our solution combines the power of hybrid vision transformers with context-aware classification to deliver a more scalable and accurate approach to anomaly detection.

Our experiments and evaluations involved rigorous testing to assess the system's performance. The results were promising: the model demonstrated high accuracy and low latency in detecting a wide range of abnormal activities in real-time surveillance footage. We compared these findings with existing methods in the field, highlighting our system's advantages particularly in scenarios requiring precise analysis of fine-grained details and human movement.

Beyond technical performance, we also focused on real-world applicability. The system's ability to automatically extract and flag only anomalous video segments significantly reduces manual monitoring efforts and speeds up response times a crucial factor for enhancing public safety in places like transportation hubs, business districts, and other high-traffic areas.

That said, we acknowledge certain limitations. Performance can vary under challenging conditions, such as poor lighting, unusual camera angles, or extremely crowded scenes. These observations open avenues for future improvements, such as adaptive learning techniques, multi-camera coordination, or integrating additional sensor data.

Interestingly, during evaluation, we noticed the model occasionally flagged unusual but previously unclassified behaviors hinting at its potential to uncover new anomaly patterns. This suggests that, beyond excelling at known categories, our system exhibits a degree of adaptability.

In summary, our work advances intelligent surveillance by offering a practical, scalable, and efficient solution for real-time anomaly detection and classification. It lays the groundwork for future enhancements and real-world deployment, bringing us closer to safer, AI-driven public safety infrastructure.

References

- [1] Sonali Agarwal and GN Pandey. Svm-based context awareness using body area sensor network for pervasive healthcare monitoring. In Proceedings of the First International Conference on Intelligent Interactive Technologies and Multimedia, 2010.
- [2] Xinyi Cui, Qingshan Liu, Mingchen Gao, and Dimitris N Metaxas. Abnormal detection using interaction energy potentials. In CVPR 2011, pages 3161–3167. IEEE, 2011.
- [3] Cewu Lu, Jianping Shi, and Jiaya Jia. Abnormal event detection at 150 fps in Matlab. In Proceedings of the IEEE international conference on computer vision, pages 2720–2727, 2013.
- [4] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In CVPR, 2016.
- [5] Joao Carreira and Andrew Zisserman. Quo Vadis, action recognition? A new model and the kinetics dataset. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 6299– 6308, 2017.
- [6] R. Hou, C. Chen and M. Shah, "Tube Convolutional Neural Network (T-CNN) for Action Detection in Videos," 2017 IEEE International Conference on Computer Vision (ICCV), 2017.
- [7] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection—a new baseline. In Proceedings of the IEEE Conference on computer vision and pattern recognition, pages, 2018.
- [8] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In Proceedings of the IEEE Conference on computer vision and pattern recognition, 2018.
- [9] Shikha Dubey, Abhijeet Boragule, and Moongu Jeon. 3D ResNet with Ranking loss function for abnormal activity detection in videos, 2020.
- [10] M. -I. Georgescu, A. Bărbălău, R. T. Ionescu, F. Shahbaz Khan, M. Popescu and M. Shah, "Anomaly Detection in Video via Self-Supervised and Multi-Task Learning," 2021
- [11] Yu Tian, Guansong Pang, Yuanhong Chen, Rajvinder Singh, Johan W Verjans, and Gustavo Carneiro. Weakly-supervised video anomaly detection with contrastive learning of long and short-range temporal features. 2021
- [12] Snehashis Majhi, Ratnakar Dash, and Pankaj Kumar Sa. Weakly supervised joint anomaly detection and classification, 2021

- [13] Pushpajit Khair, Praveen Kumar, "A semi-supervised deep learning-based video anomaly detection framework using RGB-D for surveillance of real-world critical environments",
Forensic Science International: Digital Investigation, 2022
- [14] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. arXiv preprint arXiv:1907.06987, 2019.
- [15] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Uninformed Students: Student- Teacher Anomaly Detection With Discriminative Latent Embeddings," 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020
- [16] - J. Ren, F. Xia, Y. Liu, and I. Lee, "Deep Video Anomaly Detection: Opportunities and Challenges," 2021 International Conference on Data Mining Workshops (ICDMW), 2021
- [17] - Prachi Agrawal, Narinder Singh Pun, Sanjay Kumar Sonbhadra, and Sonali Agarwal. Impact of attention on adversarial robustness of image classification models. In 2021 IEEE International Conference on Big Data (Big Data), pages 3013–3019. IEEE, 2021
- [18] Z. Liu et al., "Video Swin Transformer," 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022