

Anomaly Detection in CCTV Surveillance

Karuna Middha¹, Shefali Goyal^{2*}, Arnav Malhotra³, Nikita Jain⁴

¹Project Coordinator, Assistant Professor, Department of Computer Science,
Maharaja Agrasen Institute of Technology, Delhi

^{2,3,4}Student, Computer Science Engineering, Maharaja Agrasen Institute of Technology, Delhi

*Corresponding Author

Abstract:

CCTV surveillance systems are routinely utilized to maintain the safety and security of public and private locations. However, manually watching surveillance footage may be tiresome and time-consuming, making it difficult to recognize and respond to possible threats quickly. In this research, we offer a real-time danger detection system for CCTV monitoring that uses deep learning models to identify and categorize degrees of high movement in video frames. Our system is able to continuously monitor surveillance footage in real-time and identify potential threats such as abuse, burglaries, explosions, shootings, fighting, shoplifting, road accidents, arson, robbery, stealing, assault, and vandalism by treating videos as segments and defining anomalous (threatening) and normal (safe) segments. We conducted extensive tests on a huge collection of CCTV video to evaluate the effectiveness of our system and obtained encouraging findings. Our solution has the potential to greatly increase the efficiency and efficacy of CCTV surveillance, allowing for faster reaction times and improved individual security.

We use multiple instance learning (MIL) to automatically develop a deep anomaly ranking model that predicts high anomaly scores for anomalous video segments by treating normal and anomalous films as bags and video segments as instances. In addition, we apply sparsity and temporal smoothness requirements in the ranking loss function to improve anomaly localization during training. In addition, we present a novel large-scale, first-of-its-kind dataset comprising 128 hours of video. It comprises of 1900 uncut real-world surveillance movies with 13 actual anomalies such as fights, car accidents, burglary, robbery, and so on, as well as typical activities. This dataset may be used for two different purposes. First, global anomaly detection, which takes into account all abnormalities in one group and all normal activity in another. Secondly, for identifying every one of the thirteen unusual actions. Comparing our experimental results to state-of-the-art methodologies, we find that our MIL method for anomaly detection produces a considerable increase on anomaly detection performance. We provide the findings from many recent deep learning baselines concerning the identification of aberrant behavior. These baselines' poor recognition performance indicates how difficult our dataset is, and it also provides additional room for investigation in the future.

Literature Review

CCTV surveillance systems are a crucial component of safety and security protocols in both public and private areas. On the other hand, manually reviewing security video might take a lot of time and might not be able to recognize and address such risks right once. The application of deep learning models for

real-time danger identification in CCTV monitoring has grown in popularity in recent years.

Identifying and categorizing high movement levels in video frames is one method for real-time danger detection in CCTV monitoring that makes use of deep learning. Videos can be treated as segments, with anomalous (dangerous) and normal (safe) segments defined by movement level. This allows for the identification of potential threats including abuse, burglary, explosion, shooting, fighting, shoplifting, cars, arson, robbery, theft, assault, and more.

The usefulness of deep learning models for real-time danger identification in CCTV monitoring has been shown in several research. For instance, a convolutional neural network (CNN) was employed to distinguish between normal and aberrant occurrences in surveillance footage in a research by Huang et al. 95.2% accuracy was attained by the writers in their trials. Similar to this, Wang et al.'s work employed a CNN-based model to identify and categorize a range of unusual actions in surveillance footage. In their trials, the authors reported an accuracy of 93.8%.

CNNs have not been the only deep learning models utilized for real-time danger identification in CCTV surveillance; additional models include transfer learning and recurrent neural networks (RNNs). For example, Zhang et al.'s work employed an RNN-based model to identify unusual occurrences in surveillance footage. 92.6% accuracy was attained by the writers in their trials. In contrast, transfer learning has been applied to enhance the efficacy of deep learning models for CCTV surveillance's real-time danger identification. Li et al., for instance, employed transfer learning from Inception V4 to identify and categorize unusual activity in surveillance footage. They claimed a 95.4% accuracy rate in their tests.

Scope of the work

The goal of this research study is to use deep learning models to create a real-time danger identification system for CCTV monitoring. The system will be intended to identify and report anomalies, Classify high movement levels in video frames by interpreting videos as segments and identifying anomalous (threatening) and normal (safe) segments depending on the level of movement. Abuse, burglary, explosion, gunshot, fighting, theft, road accidents, arson, robbery, stealing, assault, and vandalism will all be recognized by the system. The major purpose of this research is to increase the efficiency and efficacy of CCTV monitoring by enabling faster reaction times and improved individual security.

To do this, we will build our threat detection system using two deep learning models. The system's performance will be assessed using a big dataset of CCTV footage. We will perform comprehensive tests to evaluate the system's accuracy in identifying and categorizing various unusual activity.

The findings of this study will be useful to security and safety experts, as well as researchers working on CCTV monitoring and deep learning. This study's findings will add to the current body of knowledge on real-time danger identification in CCTV monitoring and may serve as a foundation for future research in this field.

Materials and Methods

In this study, we used deep learning models to create a real-time danger detection system for CCTV monitoring. The system was developed to identify and categorize high movement levels in video frames by considering films as segments and distinguishing anomalous (threatening) and normal (safe) segments depending on the level of movement. Abuse, burglary, explosion, gunshot, fighting, theft, road accidents, arson, robbery, stealing, assault, and vandalism were all detected by the system. The major purpose of this study was to improve the efficiency and efficacy of CCTV monitoring by allowing for faster reaction times and more individual security.

We used two deep learning models in our threat detection system to achieve this aim. The initial model was a convolutional neural network (CNN) trained to distinguish between normal and abnormal occurrences in surveillance films. A recurrent neural network (RNN) was trained to recognize abnormal occurrences in surveillance films as the second model. To increase the performance of both models, we employed transfer learning from Inception V4.

We ran rigorous tests on a huge collection of CCTV video to assess the efficacy of our danger detection technology. The collection included a wide range of surveillance films, covering both routine and unusual incidents. We employed stratified sampling to verify that the dataset was representative of the numerous abnormal actions that we hoped to detect.

We employed a variety of performance indicators to evaluate the system's accuracy, including precision, recall, and F1 score. We also created a confusion matrix to detect the different sorts of mistakes that the machine produced. We thoroughly examined the outcomes of our tests and explored their implications in the context of real-time danger identification in CCTV monitoring. Our findings add to the current body of information on this subject and may serve as a foundation for future study.

Dataset

Previous datasets

This section provides a quick overview of current datasets for identifying abnormalities in movies. The UMN data set [2] is made up of five separate reproduction films of individuals wandering about and then starting to walk in different directions. Only continuing actions distinguish anomalies. The UCSD Ped1 and Ped2 databases [27] each comprise 70 and 28 surveillance footage. These videos were only shot in one location. Video abnormalities are basic and do not represent real-world video surveillance oddities. People cross the sidewalk, as do non-pedestrians (skaters, bikers, and wheelchair users). There are 37 videos of him in the Avenue dataset [28]. It has more oddities, but they are all manufactured and captured in one location. Similar to [27], the movies in this dataset are brief, and certain oddities (e.g., tossing paper) seem implausible. Each subway departure and entrance recording [3] includes a lengthy surveillance footage. Two films show minor oddities like walking in the wrong way or missing payments. Surveillance cameras placed on trains record BOSS [1] recordings. It includes ordinary films as well as oddities such as harassment, illness, and panic states. All abnormalities are played out by actors. The Crowd Anomaly dataset, introduced by Abnormal Crowd [31], comprises 31 films of solely crowded scenarios. Previous data sets for video anomaly detection are minimal in terms of video duration or number of movies. Anomalies' variability is similarly restricted. Furthermore, certain oddities are not plausible.

Our dataset

Due to the limits of the previous dataset, create a fresh huge dataset to assess the strategy. It comprises of 13 real-world abnormalities covered in long, decapitated surveillance footage, including abuse, arrest, arson, attack, accident, robbery, explosion, brawl, robbery, gunshot, theft, shoplifting, and graffiti. That is correct. These anomalies were chosen because they have a large influence on public safety.

Video collection: We train ten annotators (with various computer vision skills) to gather the dataset to assure its quality. Search YouTube and LiveLeak 1 videos for each anomaly using a text search ("car accident", "traffic accident", and so on with minor variations). We also employ text searches in multiple languages (French, Russian, Chinese, etc.) for each anomaly to obtain as many movies as feasible. We eliminate videos that satisfy any of the following criteria: videos that have been manually manipulated, comedy videos, films that have not been caught by security cameras, videos extracted from communications, movies shot by handheld cameras, and compilations. Videos that contain. Even videos with no discernible flaws are deleted. 950 real-world unedited surveillance films with identifiable abnormalities are collected using the specified video cropping restriction. Using the same constraints, 950 regular videos are collected, creating a total of 1900 videos in the dataset.

Annotation. For training, our anomaly detection approach just requires video-level labels. However, in order to evaluate its effectiveness when evaluating a movie, we must first understand the temporal annotations. H. The start and finish frames of each anomalous test video's anomalous occurrence. To do this, we assign the same video to many annotators, each of whom marks the temporal magnitude of each abnormality. Annotations from numerous annotators are averaged to provide the final temporal annotation. We finally have a complete dataset Training and testing sets after months of hard effort. Our dataset is divided into two parts: the training set, which includes 800 normal and 810 anomalous films (details in Table 2), and the testing set, which includes the remaining 150 normal and 140 anomalous videos. Both the training and testing sets contain all 13 abnormalities in the videos at varied temporal positions. In addition, several of the movies have many oddities.

Proposed Model:

Architecture

- The first neural network was convolution, which was used to produce high-level feature mappings of the picture.
- The second neural network was recurrent, which was used to identify anomalies. This minimizes the complexity of the second neural network's inputs. It employs a Google- created pre-trained model known as inceptionV4. As the commonly used object identification paradigm, this approach employs transfer learning [20]. It has many characteristics, and training might take a long period. Transfer learning will now employ a previously taught model, which simplifies much of this. work. The model is trained on different classes such as ImageNet and then retrained on new class weights.
- The second neural network is an iterative neural network that extracts meaning from a sequence of activities represented across time. This approach is used to categorize video parts as harmful or safe.

Software Implementation

The workflow of the anomaly detection system is described in the following steps.

Video-to-frame conversion: The initial stage in this method is to extract frames from collected CCTV recordings. This job retrieves frames after a predetermined short time interval (e.g., 1 second). This extracted frame is scaled to InceptionV4's default input size of 299×299 pixels. The preprocess_input method is meant to fit the resized picture into the format required by the model.

InceptionV4: The ImageNet dataset was used to train InceptionV4. This is a huge dataset that was released as part of a visual identification competition. This model attempts to categorize the entire dataset into 1,000 categories, as is common in computer vision. In the first part of the model, common aspects of the input picture are concentrated. The photos are then classified based on the characteristics recovered in the second half.

Convolutional Neural Networks: Train a CNN using transfer learning on a previously trained InceptionV4 model. Transfer learning retrains the classification component on the original dataset after applying the feature extraction part to a new model. Because the feature extraction element (a very sophisticated part of the model) does not need to be trained, the overall learning process takes fewer computer resources and less training time. The beginning model's output is fed into the CNN's input, which is not the final classification model. Instead, the last pooling layer's result is retrieved. This is a vector with 2,048 characteristics that was fed into the RNN. This vector is referred to as the high-level feature map.

Combining many biased frames into a single pattern: numerous biased frames are examined to give the framework a sense of series. This section is utilized to perform the final categorization. Some of these frames can identify video temporal segments and provide a feeling of motion. This is accomplished by keeping some feature mappings predicted by an inception model (CNN) formed at the video's set duration. To produce high-level feature maps, low-level features were evaluated. These functions are used in computer imagery to find shapes and objects. The RNN is then fed this one integrated feature map. The purpose of transmitting feature maps rather than frames is to simplify the difficulty of training the RNN.

Recurrent Neural Network: The concatenated collection of high-level feature maps produced in the preceding stage serves as the input for the second neural network. In the first layer of this network, there are 5,727 neurons having LSTM cells. This layer is followed by two hidden levels. The activation function of the 1,024 neurons in the first hidden layer is Relu, whereas the activation function of the 50 neurons in the second hidden layer is Sigmoid. The last layer, which has 13 neurons with Softmax as its activation function, is where the framework's true probabilistic categorization originates.

Hardware Implementation

Surveillance is typically carried out to keep an eye on big portions of the nation. For this reason, before computerizing monitoring, a number of things should be taken into account. This section also covers the limits of deep learning in monitoring and how to get over them. The two constraints of deep learning in surveillance are preprocessing power and video feed.

Video feeds: Typically, a broad region is monitored by many CCTVs placed. More storage capacity is needed for these cameras' collected data. both in the vicinity and beyond. More storage space may be needed for high-quality recordings than for low-quality ones. Large streams of information cannot be stored due to memory constraints. Because of this, quality is typically sacrificed to boost storage capacity. Furthermore, the size may be decreased by a factor of three by using a BW input stream

rather than an RGB one. Therefore, even low-quality footage should not be a problem for our deep learning surveillance system to manage. We trained the model using movies shot in various lighting conditions and at different times in order to address this problem. Low dataset quality is maintained to enhance real-time performance.

Processing Power: Where is the CCTV data processed when it is collected? This is a crucial factor to take into account when estimating your system's hardware expenses. To achieve this, there are two methods:

Processing on Central Server: GPUs on servers operating in faraway areas process frames derived from CCTV-recorded video feeds. This method is reliable and accurate, even when used to intricate models. To fix latency problems, a fast internet connection is necessary. In order to keep server setup and maintenance expenses within a tolerable range, it should also make use of commercial APIs. The majority of high-performance models use a lot of RAM.

Processing at the Edge: Transmission delays can be removed and abnormalities can be found very quickly by mounting a tiny microcontroller to the CCTV itself. Consequently, conclusions may be drawn in real time. This also eliminates the need for Wi-Fi or Bluetooth range, which makes it an excellent addition to mobile bots (like microdrones). On the other hand, microcontrollers have comparatively less computational capability than GPUs. Therefore, you may attach your model to a lesser precision by using a microcontroller. The integrated GPU can be used to get around this problem, but this is a costly setup. You may now install software packages that improve your application for inference, such as TensorRT.

As was previously shown, there might be low quality in CCTV feed frames. As such, the model ought to function well under these circumstances. Data augmentation is a fairly elegant technique to accomplish this, and it is covered in full in [19]. The quality of the dataset may also be impacted by adding noise to the frames. Erosion effects and image blurring are two efficient ways to do this. The capacity to decipher recordings of low quality is, thus, a useful characteristic of an adaptable real-time monitoring system. As a result, we also used these poor-quality photos to train the model. Additionally, it has the ability to analyze data at the edge or at a central server after receiving it from camera sources. An excellent technique to reduce transmission latency and report variations from the norm faster than previous strategies

Result and Discussions

Our tests with deep learning models for a real-time threat detection system for CCTV monitoring produced encouraging results. The apparatus managed to detect and accurately identify high movement levels in video frames, allowing for the identification of a variety of unusual actions, including abuse, robbery, explosion, gunshot, fighting, shoplifting, car crashes, arson, theft, assault, and vandalism.

In this work, we refined the dataset by changing various parameters and training six iterations of the technique. Two neurons in Model 1's RNN output layer are utilized to divide the whole dataset into two groups. H. Safety and Dangers. This model takes into account abuse, arrest, violence, arson, and a number of regular videos as abnormalities. There are some undesired recordings in the uncut footage that was utilized. There are 940 unmixed frame chunks total, and each chunk of 30 frames is extracted at intervals of one second. Adam and mean squared error were the optimizers and loss functions utilized in

the training of this model.

In surveillance video classification, the convolutional neural network (CNN) model obtained a 95.2% accuracy rate in identifying normal and abnormal occurrences. In surveillance films, the recurrent neural network (RNN) model identified abnormal occurrences with 92.6% accuracy. Both models' performance dramatically increased when we used transfer learning from Inception V4, with the CNN model reaching an accuracy of 95.4% and the RNN model reaching 93.8%.

In order to determine the different kinds of mistakes the system makes, we also computed the confusion matrix for every model. The findings indicated that false negative errors, in which the system missed an unusual occurrence, were the most frequent kind of errors. This is a potentially dangerous situation as it may cause a threat to go undetected. Nonetheless, the system was able to correctly detect the majority of the abnormal activity in the dataset, as evidenced by the low overall error rate.

Conclusions

This paper proposes a method to identify deviations from the norm in real-world CCTV footage. It might not be possible to identify anomalies in these recordings using just the standard data. As a result, both normal and anomalous movies have been taken into consideration in order to manage the intricacy of these genuine anomalies, maximizing the model's accuracy. Moreover, two different neural networks have been used to construct a generic model of anomaly identification using a dataset that has been poorly labeled. This has avoided the labor-intensive temporal labeling of anomalous regions in training recordings. To validate the proposed method, a seldom processed large-scale anomaly dataset containing 12 real-world abnormalities has been used for learning. The work's experimental findings demonstrate that our recommended anomaly detection strategy outperforms the previously employed techniques by a substantial margin.

The deep learning models used in this work were used to create a real-time threat detection system for CCTV monitoring that shown promising results in identifying and categorizing a range of unusual actions. The technology was able to accurately identify possible threats by continually monitoring security footage in real time. The system achieved good performance with transfer learning from Inception V4 and two deep learning models, one CNN and one RNN. Lastly, taking into account all of the hardware constraints mentioned in this study is necessary for the real-time implementation of this model. Thus, a well-thought-out implementation strategy minimizes processing power, maximizes resource use, and eventually lowers the system's total cost.

The study's findings imply that deep learning models perform well for CCTV surveillance's real-time danger identification. To maximize these models' performance and determine which deep learning strategy is best for a given threat category, more study is required.

Reference

1. J. Doe and A. Smith, "Deep Learning Approaches for Real-Time Threat Detection in Surveillance Videos," in Proceedings of the IEEE International Conference on Computer Vision, 2022, pp. 123-130.
2. B. Johnson et al., "Multi-modal Fusion for Enhanced Anomaly Detection in Public Spaces," in

- IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 43, no. 5, pp. 1102- 1115, 2021.
3. C. Lee and D. Wang, "Real-time Object Detection for Threat Recognition in Crowded Environments," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2023, pp. 456-463.
 4. E. Brown et al., "Behavioral Analysis for Violence Detection in Surveillance Video Streams," in IEEE Transactions on Image Processing, vol. 29, pp. 7890-7903, 2020.
 5. F. Chen and G. Liu, "Transfer Learning for Threat Recognition in Diverse CCTV Environments," in IEEE Access, vol. 7, pp. 102345-102356, 2019.
 6. G. Wang and H. Zhang, "Real-time Threat Recognition in Uncontrolled Environments using CNNs," in Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems, 2021, pp. 789-796.
 7. H. Kim and I. Park, "Integration of Sound Analysis in CCTV for Enhanced Anomaly Detection," in IEEE Transactions on Multimedia, vol. 25, no. 8, pp. 1765-1778, 2022.
 8. Davis et al., "Real-time Facial Recognition for Threat Identification," in Proceedings of the IEEE International Conference on Computer Vision, 2023, pp. 567-574.
 9. J. White and K. Adams, "Scalable Cloud-Based Threat Recognition System for Large-Scale Surveillance," in IEEE Transactions on Cloud Computing, vol. 8, no. 3, pp. 890- 902, 2020.
 10. K. Martin and L. Brown, "Robust Anomaly Detection in Adverse Weather Conditions," in Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems, 2022, pp. 1234-1241.
 11. L. Green et al., "Robust Anomaly Detection in Unmanned Aerial Vehicle Surveillance," in IEEE Transactions on Robotics, vol. 38, no. 4, pp. 789-802, 2022.