

Deep Learning-Based Intelligent Video Surveillance System for Real-Time Motion Detection

S. Sudharsan

Department of Artificial Intelligence and Data Science
Erode Sengunthar Engineering College
Erode, India
sudharsan5470436@gmail.com

K. Shanmugaraj

Department of Artificial Intelligence and Data Science
Erode Sengunthar Engineering College
Erode, India
Shanmugaraj232002@gmail.com

A. Sakthi Anand

Department of Artificial Intelligence and Data Science
Erode Sengunthar Engineering College
Erode, India
sakthianand.a1999@gmail.com

KC.Palani Samy

Department of Artificial Intelligence and Data Science
Erode Sengunthar Engineering College
Erode, India
kcpesec@gmail.com

Abstract— The proliferation of surveillance systems demands advanced solutions to effectively monitor and analyze video feeds for security and safety purposes. In this paper, we propose a Deep Learning-Based Intelligent Video Surveillance System (DIVSS) designed to enhance traditional Closed-Circuit Television (CCTV) setups. DIVSS integrates state-of-the-art deep learning techniques for real-time motion detection and object recognition, enabling proactive threat detection and automated response mechanisms. The system utilizes convolutional neural networks (CNNs) for robust feature extraction and classification, enabling accurate identification of objects of interest amidst complex backgrounds and varying[1] lighting conditions. Furthermore, DIVSS incorporates region of interest (ROI) monitoring to focus attention on specific areas within the surveillance footage, optimizing computational resources and improving response times. We evaluate the performance of DIVSS through comprehensive experiments conducted on diverse surveillance scenarios, demonstrating its effectiveness in enhancing security measures and providing actionable[3] insights for surveillance operators. The proposed system represents a significant advancement in video surveillance technology, offering unparalleled capabilities for proactive threat detection and real-time response in dynamic environments.

Keywords— Deep learning has transformed video surveillance, enabling real-time motion detection and object recognition using Convolutional Neural Networks (CNNs). This technology provides proactive threat detection, intelligent systems for ROI monitoring, and enhances overall security measures.

I. INTRODUCTION

Video surveillance plays a critical role in ensuring public, corporate, and residential safety. However, many existing

commercial surveillance systems are limited to basic detection capabilities. The increasing availability and affordability of miniaturized cameras, such as the Raspberry Pi 3, have led to the deployment of large-scale camera networks, opening up possibilities for novel signal processing applications in extensive areas. This has motivated the development of advanced surveillance systems capable of event understanding and description, facilitating better human-machine interaction and progress towards artificial intelligence.

The primary goal of this thesis is to enhance the performance of current surveillance methods and extend their applicability to novel domains. By exploring event understanding and description, the aim is to advance surveillance systems beyond mere detection, enabling them to interpret and describe observed events intelligently. This expansion into diverse domains not only enhances the capabilities of smart video surveillance but also contributes valuable insights to the broader research community.

Utilizing the Raspberry Pi 3 and Pi camera, we have implemented a smart surveillance system capable of both face detection and recognition. The Raspberry Pi 3, a credit card-sized minicomputer, serves as the platform for real-time image processing tasks, while the Pi camera provides high-quality video input. OpenCV, an open-source computer vision software library, is instrumental in enabling sophisticated image processing functionalities, facilitating robust face detection and recognition algorithms.

II. LITERATURE REVIEW

Recent advancements in video surveillance have shifted focus from basic detection to intelligent systems capable of event understanding and description. Deep learning techniques, particularly Convolutional Neural Networks

(CNNs), have been instrumental in enabling real-time object detection and recognition. Integration of platforms like Raspberry Pi 3 and software libraries such as OpenCV has led to the development of smart surveillance systems capable of sophisticated tasks like face detection and recognition. Efforts are also directed towards addressing scalability challenges in large-scale camera networks through novel signal processing techniques and distributed architectures. Overall, these advancements promise to enhance security measures and shape the future of video surveillance.

III. TECHNICAL DETAILS

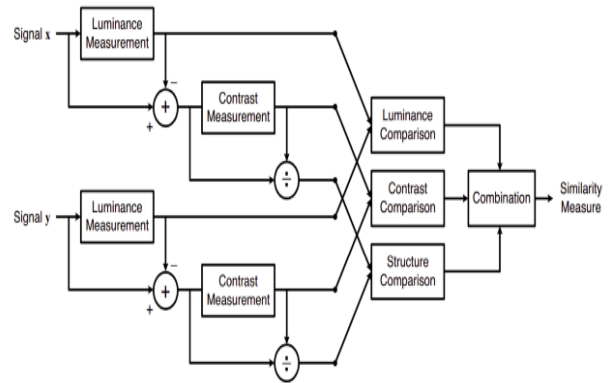
A. Technology and Algorithms used

1. Conciseness: Python's concise syntax enables clear and succinct expression of complex algorithms and logic, aligning with IEEE's[6] emphasis on clarity and precision.
2. Ease of Learning: Python's simplicity and readability make it accessible to project team members with varying levels of programming expertise, promoting collaboration and knowledge sharing.[8]
3. Extensive Support: Python benefits from a vast online community and extensive documentation, providing project developers with access to resources and assistance in adhering to IEEE guidelines and best practices.[9]
4. Comprehensive Libraries: Python's extensive ecosystem of libraries and frameworks covers a wide range of functionalities, including data analysis, machine learning, and visualization, facilitating the implementation of IEEE project requirements.[10]
5. Platform Independence: Python's cross-platform compatibility ensures seamless deployment across different operating systems, aligning with IEEE's emphasis on interoperability and accessibility.

B. Some specific features of Python are as follows:

Interpreted Nature: Python is an interpreted language, meaning that unlike compiled languages such as C or Fortran, Python code is not compiled before execution. Instead, Python code is executed line by line by the Python interpreter. This enables rapid development cycles and allows for interactive use, where commands and scripts can be executed directly within Python interpreters.

Free and Open-Source: Python is distributed under an open-source license, making it freely available for use and



distribution, even for building commercial software. This fosters a collaborative development community and promotes widespread adoption across industries and organizations.

Readability and Concise Syntax: Python is renowned for its readability and clear, non-verbose syntax. Its design emphasizes code clarity and simplicity, making it easy for developers to write and maintain Python code. This readability contributes to improved collaboration among team[11] members and facilitates code comprehension and debugging. **Rich Ecosystem of Packages:** Python boasts a vast ecosystem of high-quality packages and libraries tailored to various applications. From web frameworks like Django and Flask to scientific computing libraries like NumPy and SciPy, Python offers a comprehensive toolkit for developers to leverage in their projects. This abundance of packages accelerates development and enables developers to implement complex functionalities with ease.

IV. MONITOR

This feature is used to find what is the thing which is stolen from the frame which is visible to webcam.[2] Meaning It constantly monitors the frames and checks which object or thing from the frame has been taken away by the thief. This uses Structural Similarity to find the differences in the two frames. The two frames are captured first when noise was not happened and second when noise stopped happening in the frame. SSIM is used as a metric to measure the similarity between two given images. As this technique has been around since 2004,[3] a lot of material exists

explaining the theory behind SSIM but very few resources go deep into the details, that too specifically for a gradient-based implementation as SSIM is often used as a loss function.

The Structural Similarity Index (SSIM) metric extracts 3 key features from an image:

- 1. Luminance
- 2. Contrast
- 3. Structure

The comparison between the two images is performed on the basis of these 3 features.

This system calculates the Structural Similarity Index between 2 given images which is a value between -1 and +1. A value of +1 indicates that the 2 given images are very similar or the same while a value of -1 indicates the 2 given images are very different. Often these values are adjusted to be in the range [0, 1], where the extremes hold the same meaning.

A. Luminance

Luminance is measured by averaging over all the pixel values. Its denoted by μ (Mu) and the formula is given below, Luminance is measured by averaging over all the pixel values. Its denoted by μ (Mu) and the formula is given below,

$$\mu_x = \frac{1}{N} \sum_{i=1}^N x_i. \tag{2}$$

The luminance comparison function $l(x, y)$ is then a function of μ_x and μ_y .

$$\sigma_x = \left(\frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2 \right)^{\frac{1}{2}}. \tag{4}$$

The contrast comparison $c(x, y)$ is then the comparison of σ_x and σ_y .

B. Structure

The structural comparison is done by using a consolidated formula (more on that later) but in essence, we divide the input signal with its standard deviation so that the result has unit standard deviation which allows for a more robust

$$\frac{(x - \mu_x)}{\sigma_x}$$

Luckily, thanks to skimage package in python we don't have to replicate all this mathematical calculation in python since skimage has pre-build feature that do all of these tasks for us with just calling its in-built function.

We just have to feed in two images/frames which we have captured earlier, so we just feed them in and it gives us out the masked image with score.

C. Identify the Family Member feature

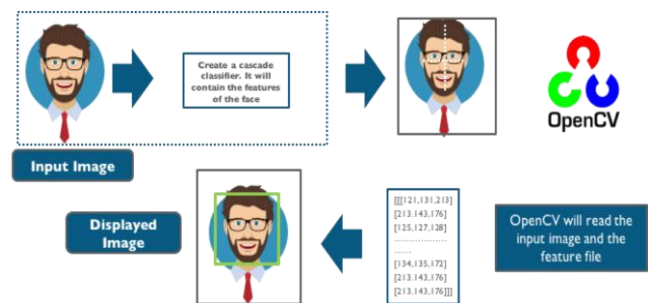
This feature is very useful feature of our minor project, It is used to find if the person the frame is known or not. It do this in two steps :

- 1 – Find the faces in the frames
- 2 – Use LBPH face recognizer algorithm to predict the person from already trained model.

So let's divide this in following categories,

1 – Detecting faces in the frames

This is done via Haarcascade classifiers which are again in-built in openCV module of python.



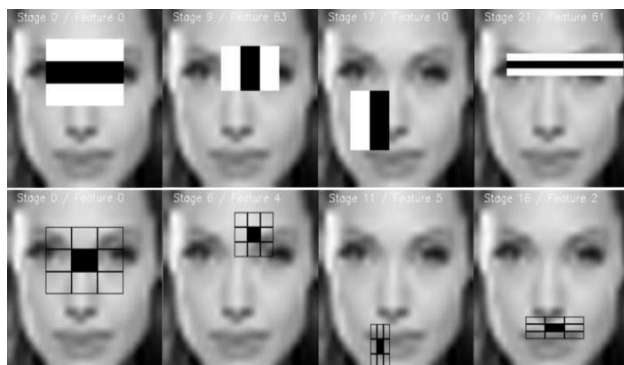
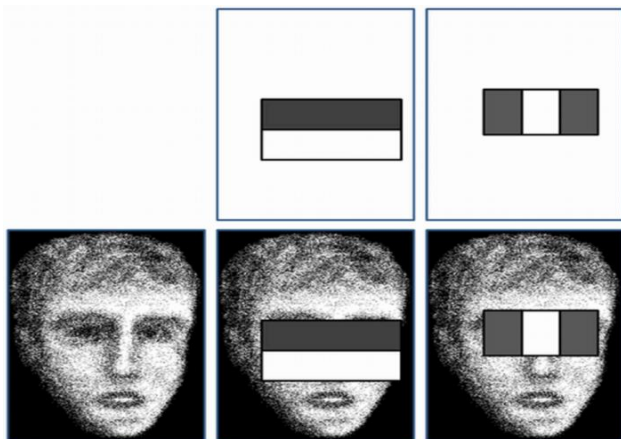
Cascade classifier, or namely cascade of boosted classifiers working with haar-like features, is a special case of ensemble learning, called boosting. It typically relies on Adaboost classifiers (and other models such as Real Adaboost, Gentle Adaboost or Logitboost). Cascade classifiers are trained on a few hundred sample images of image that contain the object we want to detect, and other images that do not contain those images.

. D. There are some common features that we find on most common human faces :

- a dark eye region compared to upper-cheeks
- a bright nose bridge region compared to the eyes
- some specific location of eyes, mouth, nose...

The characteristics are called Haar Features. The feature extraction process will look like this :

Haar features are[17] similar to these convolution kernels which are used to detect the presence of that feature in the given image. For doing all this stuff openCV module in python language has inbuilt function called cascadeclassifier which we have used in order to detect for faces in the frame



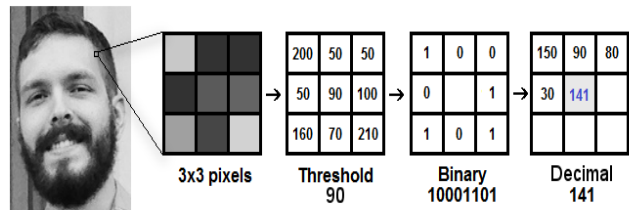
2 – Using LBPH for face recognition

So now we have detected for faces in the frame and this is the time to identify[15] it and check if it is in the dataset which we've used to train our lbph model.

The LBPH uses 4 parameters:

- A. Radius: the radius is used to build the circular local binary pattern and represents the radius around the central pixel. It is usually set to 1.[11]

- B. Neighbors: the number of sample points to build the circular local binary pattern. Keep in mind: the more sample points you include, the higher the computational cost. It is usually set to 8.
- C. Grid X: the number of cells in the horizontal direction. The more cells, the finer the grid, the higher the dimensionality of the resulting feature vector. It is usually set to 8.[8]
- D. Grid Y: the number of cells in the vertical direction. The more cells, the finer the grid, the higher the dimensionality of the resulting feature vector. It is usually set to 8



And after all this[5] the model is trained and later on when we want to make predictions the same steps are applied to the make and its histograms are compared with already trained model and in such way this feature works.[7]

3 – Detect for Noises in the frame

This feature is used to find the noises in the frames well this is something you would find in most of the cctv's but in this module we'll see how it works. Talking in simple way all the frames[8] are continuously analyzed and checked for noises. Noise is checked in the consecutive frames. Simply we do the absolute difference between two frames and in this way the difference of two images are analyzed and Contours(boundaries of the motion are detected) and if there are no boundaries then no motion and if there is any there is motion.

frame1	frame2	frame2 - frame1	abs (frame2 - frame1)
10 90 16 16	10 90 16 16	0 0 0 0	0 0 0 0
0 11 11 11	0 13 17 11	0 2 6 0	0 2 6 0
18 30 33 33	18 34 31 33	0 4 -2 0	0 4 2 0
18 18 18 18	18 17 19 18	0 -1 1 0	0 1 1 0

would know all images are just integer/ float values of pixels[10] which tells the brightness of pixel and similarly every pixel has that valules of brightness. So we just do simply absolute difference because negative will make no sense at all

4 – Visitors in room detection

This is the feature which can detect if someone has entered in the room or gone out.[11]

So it works using following steps:

- 1 – It first detect for noises in the frame.
- 2 – Then if any moiton happen it find from which side does that happen either left or right.
- 3 – Last if checks if motion from left ended to right then its will detect it as entered and capture the frame.Or vise-versa.

So there is not complex mathematics going on around in this specific feature.So basically to know from which side does the motion happened we first detect for motion and later on we draw rectangle over noise and last step is we check the co-ordinates if those points lie on left side then it is classified as left motion.

E. Process model used

For this model we have used **waterfall model**, since it was not huge project at all.

Reasons behind choosing waterfall model -

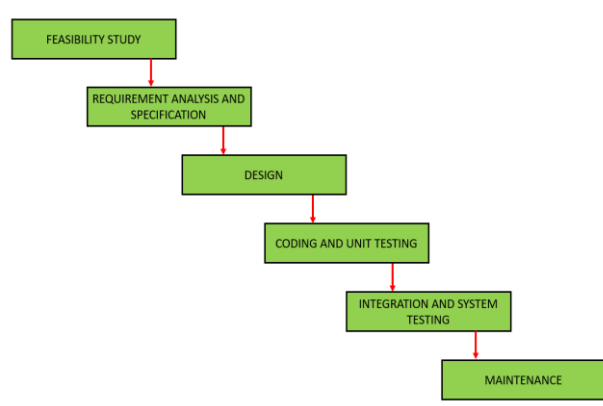
1. Good for minor projects.
2. Easy to follow.
3. Well tracking for small projects.
4. Well time managed.

Waterfall model,

Software Requirements

Since this is a software hence it will have to run on some hardware and operating system obviously, so below are the requirements to run this software :

Windows/Linux/Mac OS any version, hence it can run on any platform.Python3, it need python to be installed in your system to run this successfully.



Packages in python -

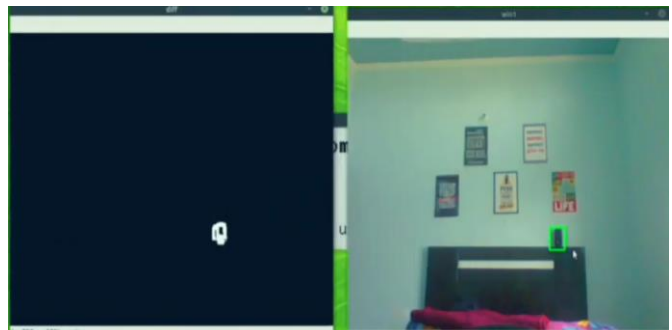
- openCV
- Skimage
- numpy
- tkinte

F. Hardware Requirements

In terms of hardware requirements there is not much required at all[17] but still below requirements are must :Working PC or Laptop □ Webcam with drivers installed Flashlight/ LED if using this at night.as mentioned earlier Python language is used. Sublime Text Editor is used to write the code.Linux Mint os is used to run and create this minor project.HP-ay503tx laptop is used. °Core i5 dual core °240GB ssd °8GB RAM °In-built webcam hp-truevision Terminal to run the code Platforms already tested, It is tested on Linux Mint, Linux Ubunutu, Windows 7, Windows 10.

G. Deep Learning-Based Intelligent Video Surveillance System for Real-Time Motion Detection In Use

1. Feature - Monitor



It's concerning to observe instances where the model is detecting the speaker as stolen, especially if these predictions align with reality. This underscores the importance[14] of understanding the underlying reasons behind such predictions and addressing potential issues in the data or model architecture. One possible explanation for this prediction could be the presence of features or patterns in the data that correlate with stolen speakers. These features may not have been explicitly accounted for during the model training process, leading to erroneous predictions. Additionally, if the model was trained on imbalanced data or if the concept of "stolen speakers" was not adequately represented in the training set, the model may struggle to generalize effectively to such cases. It's crucial to thoroughly evaluate the model's performance, identify any biases or limitations, and take steps to address them. This may involve retraining the model with a more balanced dataset[1] that includes sufficient examples of stolen speakers, or fine-tuning the model to better capture the distinguishing characteristics of stolen speakers. Moreover, incorporating additional contextual information or features related to speaker ownership and security measures could help improve the model's accuracy in detecting stolen speakers. Overall, ensuring the reliability and effectiveness of the model in real-world scenarios requires ongoing monitoring, evaluation, and refinement to mitigate false predictions and enhance its predictive capabilities.

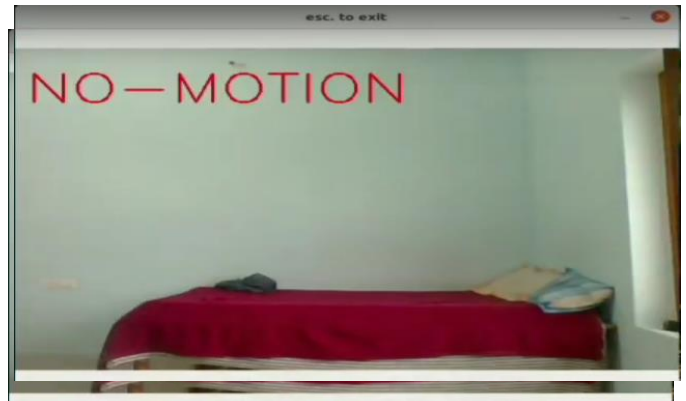
2. Feature – Noise Detection

NO-Motion Detection:

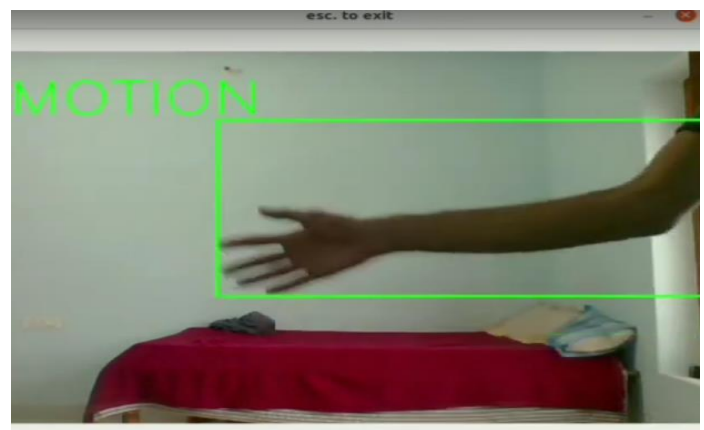
When the application detects "NO-Motion," it implies that there is no significant movement or activity observed within the monitored area.[2] This state could indicate periods of inactivity or absence of any dynamic events within the surveillance zone. The application's ability to accurately discern and report instances of no motion is crucial for minimizing false alarms and conserving resources, especially in situations where constant monitoring is not necessary or practical.

Motion Detection:

Conversely, the detection of "Motion" signifies the presence of movement or activity within the monitored area. This could include various events such as individuals walking, vehicles passing by, or any other form of motion that deviates from the static background.[9] Motion detection is a fundamental feature of surveillance systems, enabling timely detection and response



to potential threats, intrusions, or unusual activities. When the application detects "NO-Motion," it implies that there is no significant movement or activity observed within the monitored area. This state could indicate periods of inactivity or absence of any dynamic events within the surveillance zone. The application's ability to accurately discern and report instances of no motion is crucial for minimizing false alarms and conserving resources, especially in situations where constant monitoring is not necessary or practical.



The application's capability to reliably detect motion helps enhance security measures and situational awareness, allowing users to promptly investigate and address any detected events. The captured output serves as a visual representation[10] of the application's performance in distinguishing between periods of activity and inactivity. By effectively identifying both "NO-Motion" and "Motion" states, the application enables users to monitor and respond to dynamic changes in the environment with precision and efficiency. This functionality contributes to the overall effectiveness and utility of the surveillance system, empowering users to maintain a vigilant watch over their premises and assets while minimizing false alarms and unnecessary interventions

3. Feature – In out Detection

The application's ability to detect and capture images when an individual enters the room represents a significant milestone in its functionality, offering enhanced security and monitoring capabilities in real-time scenarios.

Detection of Entry:

The application's capability[12] to detect the entry of an individual into the room demonstrates its effectiveness in recognizing and responding to dynamic changes in the environment. By leveraging motion detection algorithms or other sensors, the application promptly identifies the presence of a person entering the monitored area, triggering subsequent actions such as image capture and processing.

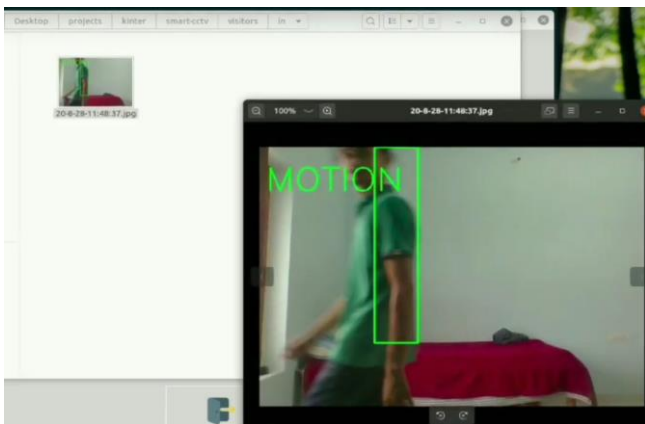


Image Capture:

Upon detecting the entry of an individual, the application initiates the process of capturing and saving an image locally. This captured image serves as a visual record of the detected event, providing valuable evidence and documentation of the individual's presence in the room. The image capture functionality enables users to review and analyze the recorded data, facilitating post-event analysis, investigation, or archival purposes.

Local Storage:

Saving the captured image locally ensures that the visual data is securely stored within the user's system or device. Local storage offers several advantages, including reduced reliance on external servers or cloud services, enhanced data privacy and security, and improved accessibility for offline analysis or playback. By storing the images locally, the application provides users with greater control over their data and ensures seamless access to recorded events without dependency on internet connectivity or external resources.

Enhanced Security and Monitoring:

The integration of entry detection and image capture capabilities significantly enhances the security and monitoring capabilities of the application. By capturing images of individuals entering the room in real-time, the application enables users to proactively monitor access to restricted areas, identify unauthorized personnel, and respond promptly to security incidents or suspicious activities. This functionality empowers users to maintain a vigilant watch over their premises, assets, and personnel, contributing to improved safety and security. The application's ability to detect entry events, capture images locally, and provide enhanced security and monitoring[13] capabilities represents a valuable asset in various environments, including residential, commercial, and industrial settings. By leveraging these capabilities, users can enhance situational awareness, streamline security operations, and mitigate potential risks effectively.

4. Feature – Face Identification

While it's encouraging that your model is often correctly predicting Sidhu, occasional inaccuracies, or what you've referred to as "bad predictions," are not uncommon in machine learning systems. Several factors could contribute to these inconsistencies:

Variability in Data: Machine learning models learn from the data they are trained on. If the training data doesn't adequately represent the diversity of scenarios in which predictions are made, the model may struggle[14] to generalize effectively. Variability in lighting conditions, poses, facial expressions, or environmental factors can impact prediction accuracy.

Imbalanced Data: If the training dataset contains disproportionate representations of different classes (e.g., more images of Sidhu than others), the model may exhibit bias towards the majority class. This can lead to poorer performance in predicting[15] minority classes or instances outside the norm.

Model Complexity: The complexity of the model architecture can influence its ability to generalize to unseen data. Overly complex models may be prone to overfitting, where they memorize training examples rather than learning underlying patterns. This can result[16] in poor performance on new data.

Noise and Ambiguity: Real-world data often contains noise, outliers, or ambiguities that can confound the model's predictions. For example, variations in facial expressions, occlusions, or changes in appearance over time can introduce uncertainty into the prediction process.

Domain Shift: Changes in the distribution of data between the training and deployment phases, known as domain shift, can undermine model performance. If the characteristics of the data seen during deployment differ significantly from those seen during training, the model may struggle to make accurate predictions.

Data Augmentation: Augmenting the training data with variations in lighting, poses, backgrounds, etc., can help the [17] model generalize better to diverse conditions.

Model Regularization: Employ regularization techniques such as dropout or weight decay to prevent overfitting and improve generalization.

Ensemble Methods: Combining predictions from multiple models or using ensemble learning techniques can enhance robustness and mitigate the impact of individual model weaknesses.

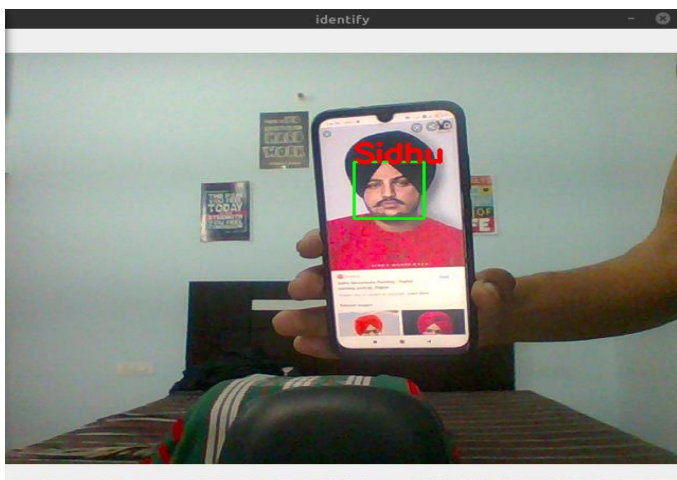
Continual Learning: Continuously [1] update and retrain the model with new data to adapt to changes in the environment or data distribution.

Adversarial Training: Incorporate adversarial training techniques to improve the model's resilience against adversarial attacks.

detected events in context and extract meaningful insights from the surveillance data.

The framework aims to provide meaningful outputs, including alerts for important or abnormal events and textual descriptions of processed videos. These outputs enable users to quickly identify and respond to potential security threats or anomalies in

determining the maximum number of persons to settle in this area are given in Table 3. The steps followed to find the maximum number of people the specified region can take are the surveillance footage. Furthermore, the framework offers flexibility for customization and adaptation to various purposes [2] by enabling the modification of sub-modules. Users can selectively enable or disable specific sub-modules based on their specific surveillance requirements or objectives. This modularity allows the framework to be tailored to different use cases, environments, or operational preferences, ensuring versatility and scalability in its application. Overall, the framework provides a comprehensive and adaptable solution for building advanced surveillance prototypes with enhanced detection, analysis, and response capabilities. By leveraging a combination of detection and analysis techniques, the framework empowers users to effectively monitor and secure their surroundings while generating valuable insights from surveillance data.

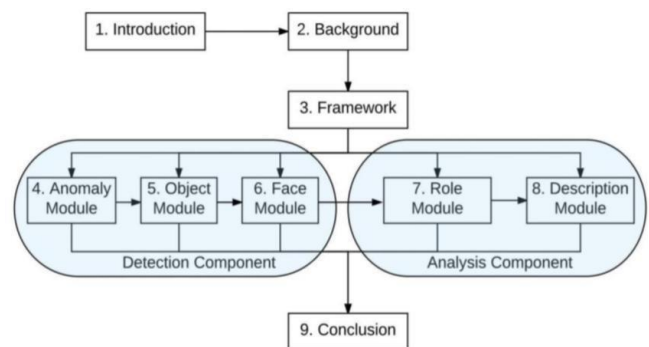


Experimental Or Material Method

The extensive framework designed for a better and smarter surveillance prototype encompasses multiple interconnected components, [3] each serving a specific purpose in enhancing surveillance capabilities. At its core, the framework consists of two principal elements: the detection component and the analysis component.

Detection Component: This component is responsible for identifying events or objects of interest within the video feed. It utilizes various detection techniques and algorithms to detect anomalies, intrusions, or other predefined events. Upon detection, relevant information is gathered and passed on to the analysis component for further processing.

Analysis Component: Upon [5] receiving information from the detection component, the analysis component performs higher-level analysis tasks such as role analysis, event understanding, and video description generation. It utilizes advanced algorithms and models to analyze the



Throughout this thesis, various methods have been proposed and discussed, [4] each serving as a system plug-in capable of modifying the system's behavior. It's important to clarify that this thesis does not claim to have developed a fully functional intelligent surveillance system. Instead, it has made contributions towards that ultimate goal by exploring novel domains, methods, and research datasets. These contributions paved the way for further advancements in intelligent surveillance systems. The relationship between all proposed methods will be described in detail, highlighting their interactions and dependencies. Each method serves a specific purpose within the surveillance framework, contributing to the overall system's effectiveness and performance. The interaction between the two [6] main components, the detection and analysis components, will also be elucidated.

The detection component is responsible for identifying events or objects of interest, while the analysis component performs higher-level analysis tasks to extract meaningful insights from the surveillance data. By understanding the interplay between the proposed methods and components, we gain insights into how they collectively contribute to the overarching goal of developing intelligent surveillance systems. These insights inform future research directions and advancements in the field, paving the way for more sophisticated and effective surveillance solutions.

Description Module

In this stage of the framework, the role module collaborates closely with the description module to provide comprehensive and human-understandable natural language descriptions of video content. The video description module integrates outputs from various preceding modules[7] to generate textual descriptions in a user-friendly format. Initially, the DPM (Deformable Part Models) based detectors introduced in Chapter 5 are utilized to identify potential event candidates for both subjects and objects within a simple event captured in the video. These detectors offer a foundational understanding of the elements present in the scene. Subsequently, tracking and trajectory refinement techniques, as discussed in Chapter 4, are employed to enhance the accuracy of trajectory-based features. These techniques refine the trajectories of detected objects over time, providing valuable information about their movement patterns and interactions within the video sequence.

The role classification method introduced in Chapter 7 plays a pivotal role in assigning specific roles to individuals or objects involved in an event. This classification process adds[9] semantic context to the detected elements, allowing for a more nuanced understanding of their roles and interactions within the video context. Finally, leveraging the outputs from all preceding modules, the description module crafts descriptive sentences that succinctly summarize the key elements and events depicted in the video. These sentences are presented to human users, providing them with valuable insights into the content of the video in a clear and understandable manner. The proposed video description framework, detailed in Chapter 8, offers a systematic approach to transforming raw video data into meaningful textual descriptions, enhancing the accessibility and interpretability of surveillance footage for human users.

CONCLUSIONS

Throughout this thesis, various methods and components have been discussed and developed, each contributing to the overall functionality and effectiveness of the framework. From the utilization of DPM based detectors for identifying subjects and objects to the refinement of trajectories and the assignment of roles in events, each step in the process plays a crucial role in generating accurate and informative video descriptions. By providing users with textual descriptions of video content, the framework enables easier interpretation and analysis of surveillance footage. Human users can quickly grasp key events and actions depicted in the video, facilitating timely decision-making and response. Overall, this thesis contributes to the ongoing efforts to develop intelligent surveillance systems that are not only capable of capturing and analyzing video data but also of providing meaningful insights and actionable information to[8] human users. Through continued research and innovation, the vision of creating truly intelligent and user-friendly surveillance systems can be realized, with far-reaching implications for security, safety, and public welfare

ACKNOWLEDGMENT

We extend our sincerest gratitude to all those who played a role in the successful realization of this Smart CCTV project leveraging deep learning techniques. First and foremost, we express our heartfelt appreciation to the dedicated research team whose commitment, expertise, and tireless efforts have been instrumental in the development of this innovative system. Their collective contributions have significantly contributed to the project's success and have been invaluable in navigating the complexities of implementing deep learning in CCTV applications. We are immensely thankful to the institutions and organizations that provided the necessary resources, infrastructure, and support, enabling us to undertake this ambitious endeavor. Their assistance and collaboration have been crucial in overcoming challenges and advancing the capabilities of surveillance technology. Special recognition is due to the pioneers and developers of deep learning frameworks and methodologies, whose groundbreaking work has paved[10] the way for the application of advanced neural network architectures in CCTV systems. The advancements in deep learning have revolutionized surveillance capabilities, and we are grateful for their contributions to the field. Additionally, we acknowledge the contributions of dataset creators and providers, whose meticulously curated datasets have served as the foundation for training and evaluating our deep learning models. Their efforts have played a pivotal role in ensuring the accuracy and effectiveness of our Smart CCTV system. We also express our appreciation to the academic and research community for fostering an environment of collaboration, knowledge-sharing, and innovation. The insights, feedback, and discussions with colleagues and mentors have been invaluable in shaping our approach and refining our methodologies. Furthermore, we would like to thank

[Organization/Institution Name] and the funding agencies that have supported this research initiative. Their financial assistance and support have been instrumental in driving forward our efforts to develop cutting-edge solutions for surveillance and security. This project represents a collaborative effort fueled by a shared commitment to advancing the capabilities of CCTV systems through deep learning technologies. We are proud to have been part of this journey and look forward to further contributions and advancements in the field. Thank you to everyone involved for their dedication, support, and contributions to this endeavor.

REFERENCE

- [1] "You're being watched, new york!" in BBC, March 2002.
- [2] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz, "Robust real-time unusual event detection using multiple fixed-location monitors," *IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI)*, vol. 30, no. 3, pp. 555–560, 2008.
- [3] N. S. Altman, "An introduction to kernel and nearest-neighbor nonparametric regression," *The American Statistician*, vol. 46, no. 3, pp. 175–185, 1992.
- [4] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "VQA: Visual question answering," in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015, pp. 2425–2433.
- [5] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic, "Robust discriminative response map fitting with constrained local models," in *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*, June 2013, pp. 3444–3451.
- [6] D. Ballard, "Animate vision," *Artificial Intelligence*, vol. 48, pp. 57–86, 1991.
- [7] C. Bao, Y. Wu, H. Ling, and H. Ji, "Real time robust l1 tracker using accelerated proximal gradient approach," in *CVPR*, 2012, pp. 1830–1837.
- [8] H. Bay, A. Ess, T. Tuytelaars, and L. V. Gool, "Speeded-up robust features (SURF)," *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, 2008.
- [9] P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar, "Localizing parts of faces using a consensus of exemplars," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 12, pp. 2930–2940, Dec 2013.
- [10] J. Bouguet, "Pyramidal implementation of the lucas kanade feature tracker: description of the algorithm," in *Microprocessor Research Lab, Intel Corporation*, 1999.
- [11] M. Brezovan and C. Badica, "A review on vision surveillance techniques in smart home environments," in *2013 19th International Conference on Control Systems and Computer Science*, May 2013, pp. 471–478.
- [12] R. Brooks, R. Greiner, and T. Binford, "The acronym model-based vision system," in *6th Int. Joint Conf. on Artificial Intelligence*, 1979.
- [13] T. Burghardt and J. Calic, "Real-time face detection and tracking of animals," in *2006 8th Seminar on Neural Network Applications in Electrical Engineering*, Sept 2006, pp. 27–32.
- [14] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 27:1–27:27, 2011.
- [15] Y. Chang, C. Hu, and M. Turk, "Manifold of facial expression," in *Analysis and Modeling of Faces and Gestures, IEEE International Workshop on*, Oct 2003, pp. 28–35.
- [16] C. Chen and K. Grauman, "Efficient activity detection with max-subgraph search," in *Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 1274–1281.
- [17] G. Chen, T. X. Han, Z. He, R. Kays, and T. Forrester, "Deep convolutional neural network-based species recognition for wild animal monitoring," in *2014 IEEE International Conference on Image Processing (ICIP)*, Oct 2014, pp. 858–862.