



Smart Surveillance & Life Detection System

*A project submitted
in partial fulfillment of the requirements for the degree of
Bachelor of Science in Computer Science*

by

Abdallah Ashraf Abdallah	20210541
Andrew Atef Shukrallah	20210191
Anderw Sherif Shokry	20210190
Alaa Hesham Mamdouh	20210173
Zeina Samer Wagdy	20210390

Supervised by

Dr. Hala

June 2025

ACKNOWLEDGEMENT

We would like to express our sincere gratitude and appreciation to all those who supported us throughout the development of this graduation project.

First and foremost, we are deeply grateful to **Dr. Hala**, our supervisor, for her continuous guidance, encouragement, and valuable insights that significantly enriched our work and helped shape our project into its final form.

We also extend our thanks to the **Faculty of Computers and Information – Computer Science Department**, for providing us with the academic foundation and resources necessary to complete this project.

Special thanks go to our colleagues, friends, and families for their moral support and motivation throughout this journey.

This project is the result of the collective effort and collaboration of the team, and we are proud to have worked together on such a meaningful and impactful system.

ABSTRACT

Problem Statement:

With the rise of surveillance systems and online content, automated violence detection is critical for security and content moderation.

Traditional methods lack real-time accuracy and scalability.

Solution:

This project presents an intelligent, **deep learning-based system for detecting violence** in video streams using spatiotemporal analysis. The system integrates a hybrid model combining **MobileNetV2** (for spatial feature extraction) with Bidirectional LSTM (for temporal sequence modeling).

The model classifies video segments as either violent or non-violent, enabling rapid detection and response.

In addition to violence detection, the system also includes modules for fire detection and real-time person tracking, enhancing situational awareness in public surveillance.

The system is supported by a scalable backend (**Laravel + Flask**) and a modern frontend (**React.js**), ensuring smooth integration, alert generation, and event reporting.

Key Words

“Violence detection, CNN, LSTM, Spatiotemporal analysis, Valence backtracking, Surveillance, Deep learning, Real-time processing..”

Key Features:

- *Deep learning model (MobileNetV2 + BiLSTM) trained on real-world datasets.*
- *Real-time violence classification and alert generation*
- *Fire detection module using YOLOv8.*
- *Person tracking using DeepSORT.*
- *Scalable architecture with modular model deployment.*
- *User-friendly dashboard for monitoring and reports.*

1 TABLE OF CONTENTS

ACKNOWLEDGEMENT.....	i
ABSTRACT.....	ii
LIST OF TABLES.....	v
LIST OF FIGURES.....	vi
Chapter 1: INTRODUCTION	[Page 7]
1.1 Overview	[Page 7]
1.2 Objectives	[Page 7]
1.3 Purpose	[Page 8]
1.4 Scope	[Page 8]
Chapter 2: BACKGROUND AND RELATED WORK.....	[Page 9]
2.1 Introduction	[Page 9]
2.2 Computer Vision in Video Surveillance.....	[Page 9]
2.3 Deep Learning for Video Classification.....	[Page 11]
2.3.1 MobileNetV2 Architecture	[Page 12]
2.3.2 Convolutional Neural Network (CNN)	[Page 14]
2.3.3 TimeDistributed Layer	[Page 14]
2.3.4 Long Short-Term Memory (LSTM)	[Page 15]
2.3.5 Bidirectional LSTM	[Page 15]
2.3.6 Transfer Learning using MobileNetV2	[Page 16]
2.4 Datasets for Training	[Page 17]
2.5 Real-Time Processing Considerations.....	[Page 17]
2.6 Additional Modalities	[Page 18]
2.7 Related Work.....	[Page 19]
2.7.1 Early Methods for Violence Detection.....	[Page 19]
2.7.2 Deep Learning-Based Approaches	[Page 19]
2.7.3 Public Datasets Used in Related Work	[Page 20]
2.7.4 Comparative Studies	[Page 20]
2.7.5 Our Contribution Compared to Related Work	[Page 21]
2.7.6 Summary	[Page 21]
Chapter 3: SYSTEM ARCHITECTURE AND DESIGN.....	[Page 22]
3.1 Overview	[Page 22]
3.2 Component Breakdown.....	[Page 23]
3.3 Data Flow Diagram (DFD).....	[Page 25]
3.4 Database Design.....	[Page 26]
3.5 Technology Stack.....	[Page 28]
Chapter 4: IMPLEMENTATION DETAILS.....	[Page 29]
4.1 AI Model Implementation	[Page 29]
4.1.1 Data Preprocessing and Feature Extraction	[Page 30]
4.1.2 Model Architecture (MobileNetV2 + BiLSTM)	[Page 31]

4.1.3 Model Training and Saving	[Page 34]
4.1.4 Serving the Model with Flask API.....	[Page 35]
Chapter 5: SYSTEM FEATURES AND FUNCTIONALITY.....	[Page 37]
5.1 User Roles and Access Control	[Page 37]
5.2 Administrator Dashboard.....	[Page 38]
5.2.1 Reports Management	[Page 38]
5.2.2 User Management	[Page 38]
5.2.3 Camera Management and Assignment.....	[Page 39]
5.3 Employee Dashboard	[Page 40]
5.4 Real-time Alerting and Reporting Mechanism.....	[Page 40]
Chapter 6: TESTING AND EVALUATION.....	[Page 41]
6.1 Model Performance Evaluation.....	[Page 41]
6.1.1 Evaluation Metrics	[Page 41]
6.1.2 Performance Results	[Page 43]
6.1.3 Confusion Matrix.....	[Page 43]
6.1.4 Training and Validation Curves	[Page 43]
6.2 System Performance Testing.....	[Page 44]
6.3 User Acceptance Testing (UAT)	[Page 45]
Chapter 7: CONCLUSION AND FUTURE WORK.....	[Page 46]
7.1 Conclusion.....	[Page 46]
7.2 Limitations	[Page 47]
7.3 Future Work	[Page 48]
REFERENCES.....	[Page 49]

CHAPTER 1: INTRODUCTION

1.1 Overview:

Violence remains a major threat to public safety in various environments such as schools, streets, public transportation, and large gatherings. With recent advancements in computer vision and deep learning, real-time video surveillance systems have evolved to become more intelligent and efficient.

This project presents an AI-powered violence detection system designed to analyze live video streams and automatically identify violent activities.

By integrating multiple models—including violence detection, fire detection, and person tracking—the system delivers a robust and scalable solution for real-time safety monitoring in public environments.

1.2 Objectives:

- Develop a deep learning model capable of detecting violent activities in real-time video streams.
- Implement a person-tracking module to follow violent individuals after detection.
- Integrate fire and vandalism detection models to enhance public safety coverage.
- Design an intuitive web interface for administrators and security personnel to receive alerts and monitor live events.
- Deploy the system in a scalable way that supports running multiple models simultaneously.

1.3 purpose:

The primary purpose of this project is to build an intelligent, automated surveillance system that enhances public safety by detecting violent actions and other emergencies in real time. It aims to reduce reliance on manual supervision, enabling immediate response by triggering automated alerts and tracking suspects. Ultimately, the system contributes to safer public spaces through early threat detection and response.

1.4 Scope:

This project encompasses the design, development, and deployment of an AI-based video surveillance system. It focuses on real-time detection of violent behavior and other hazardous events such as fire. The project includes:

- Development and training of deep learning models.
- Integration with real-time video streams and tracking modules.
- Building a web-based user interface for monitoring and alerts.
- Deployment using APIs and cloud infrastructure.

Target environments include public spaces such as schools, streets, shopping malls, and transport stations.

CHAPTER 2: BACKGROUND & Related Work

2.1 Introduction

The use of artificial intelligence in video surveillance systems has significantly advanced in recent years. Leveraging deep learning techniques, it is now possible to analyze video feeds in real time to detect abnormal events such as violence, fire, or vandalism.

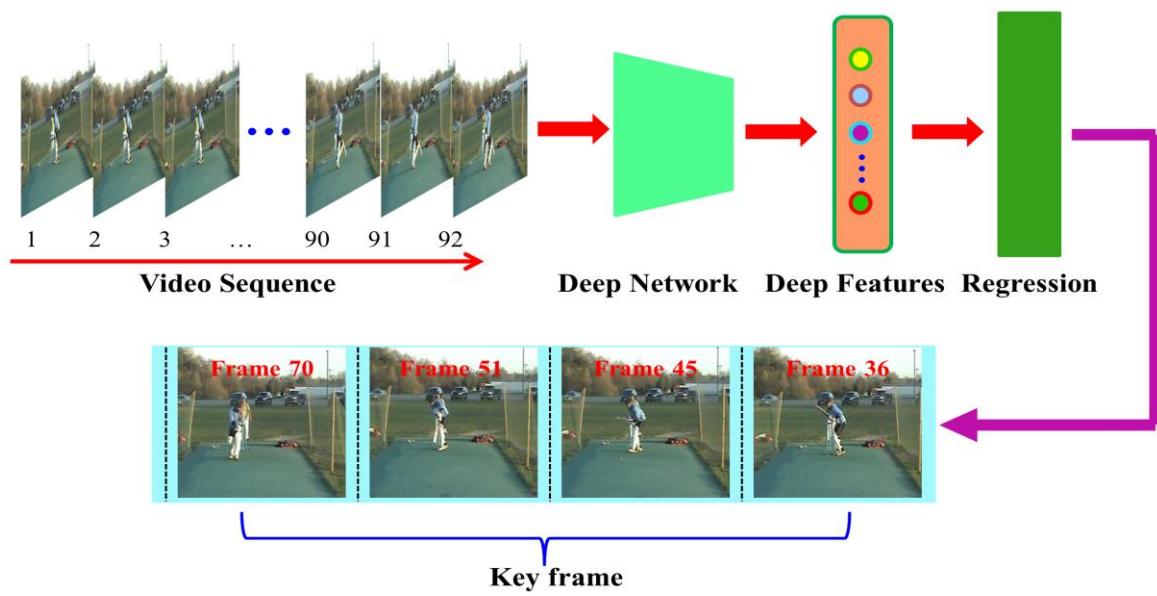
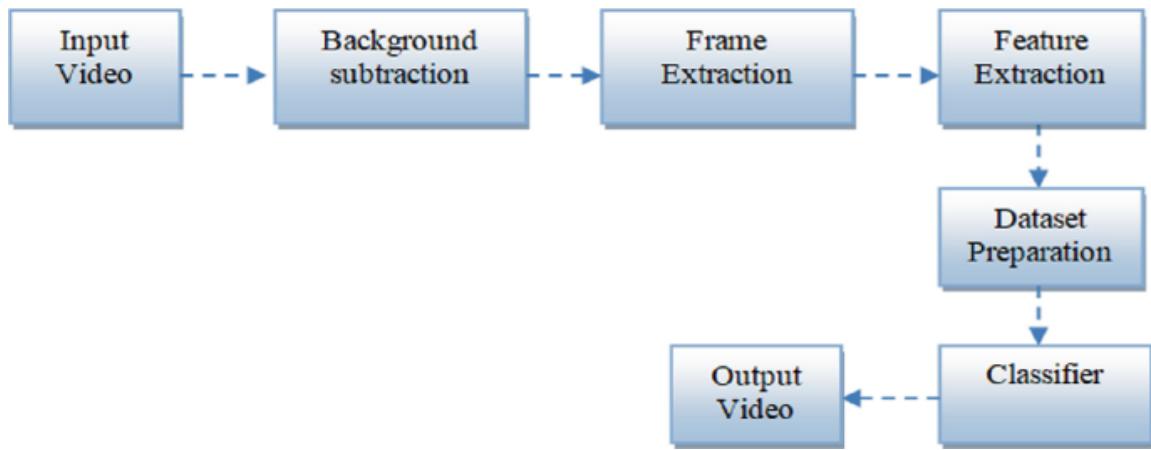
This chapter presents the foundational knowledge and key technologies that underpin the implementation of an intelligent violence detection system.

2.2 Computer Vision in Video Surveillance

Video analysis refers to the automated process of extracting meaningful information from video streams. In this project, the primary task is the detection of violent behavior and related hazards in real-time. This involves processing each frame using computer vision techniques and classifying them using trained neural network models.

Common tasks include:

- Object detection (e.g., humans, flames)
- Action recognition (e.g., fighting, falling)
- Scene understanding
- Temporal sequence analysis



2.3 Deep Learning for Video Classification

In violent action detection, the system must understand both **what is happening** and **how it evolves** over time (temporal features).

Models used:

- **CNN (Convolutional Neural Network):**
Extracts **spatial features** from each video frame
(like detecting a hand raised or a person running).
- **LSTM (Long Short-Term Memory):**
A type of **Recurrent Neural Network (RNN)** that captures **temporal dependencies** between frames.
- **Bidirectional LSTM:**
Reads the frame sequence both forward and backward,
improving context understanding.
- **TimeDistributed Layer:**
Wraps the CNN and applies it to each frame before sending outputs to LSTM layers.
- **Transfer Learning:**
MobileNetV2 pretrained on ImageNet is used to reduce training time
and handle limited data efficiently.

2.3.1 MobileNetV2 Architecture

MobileNetV2 is a lightweight convolutional neural network architecture designed for efficient computation on devices with limited resources, such as mobile phones or embedded systems.

It is especially useful in real-time video analysis systems due to its speed and low memory footprint.

◊ Key Features:

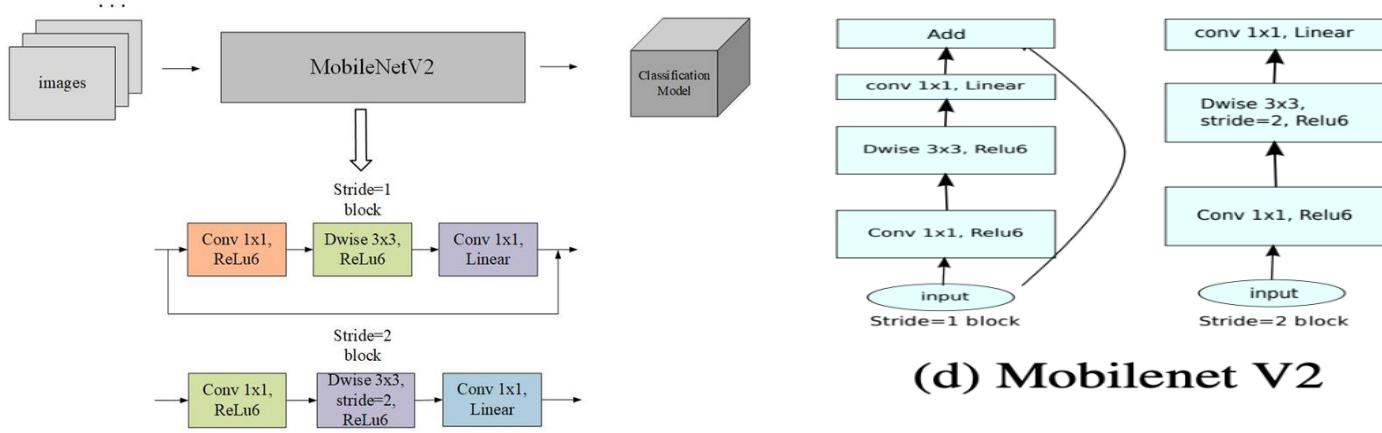
- Depthwise Separable Convolutions: Breaks standard convolution into two steps: depthwise and pointwise, reducing computation.
- Inverted Residual Blocks: Introduces shortcuts between thin bottleneck layers, enabling deeper networks with fewer parameters.
- Linear Bottlenecks: Uses linear activation at bottlenecks to preserve information that would be lost with ReLU.

◊ Why MobileNetV2?

- **Efficient:** Lower latency and fewer parameters compared to larger models like ResNet or VGG.
- **Accurate:** Despite its lightweight nature, it achieves high performance on many vision tasks.
- **Scalable:** Can be used as a base for transfer learning, especially in small datasets.

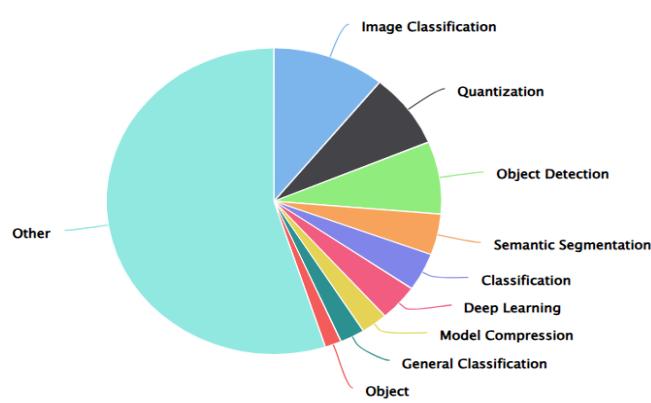
◊ Integration in Our System:

- Each video frame is passed through a TimeDistributed wrapper using MobileNetV2 as a feature extractor.
- Outputs are sequences of feature vectors representing spatial information.
- These are then passed to an LSTM for temporal modeling.



(d) Mobilenet V2

Tasks



Task	Papers	Share
● Image Classification	68	10.86%
● Quantization	49	7.83%
● Object Detection	48	7.67%
● Semantic Segmentation	27	4.31%
● Classification	25	3.99%
● Deep Learning	24	3.83%
● Model Compression	16	2.56%
● General Classification	15	2.40%
● Object	10	1.60%

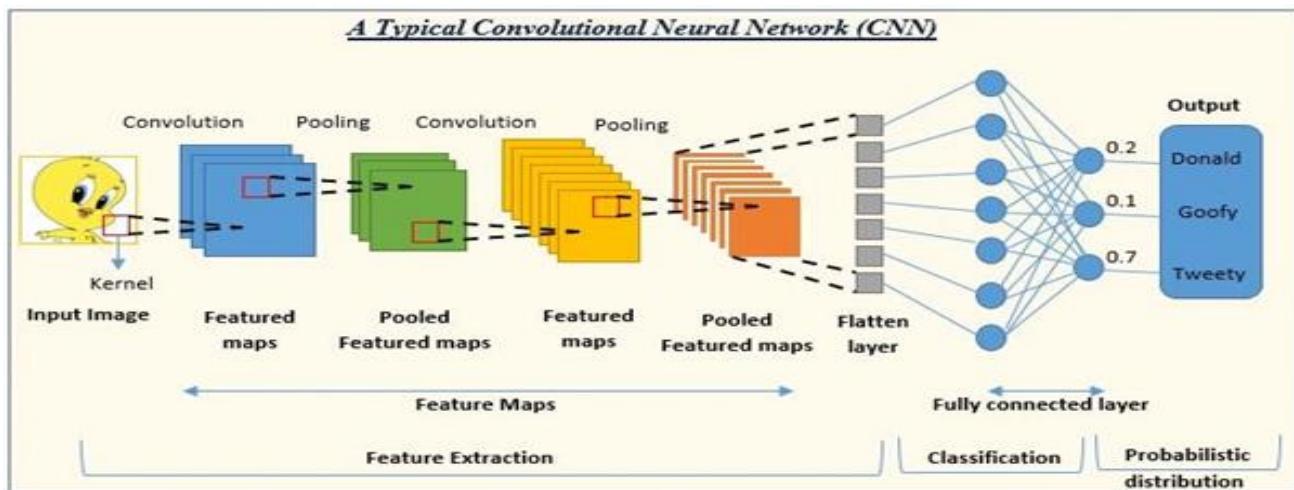
2.3.2 Convolutional Neural Network (CNN)

CNNs are used to extract **spatial features** from individual video frames.

In our system, a pre-trained CNN (MobileNetV2) processes each frame and identifies visual patterns like human posture, motion direction, or the presence of flames.

Key benefits:

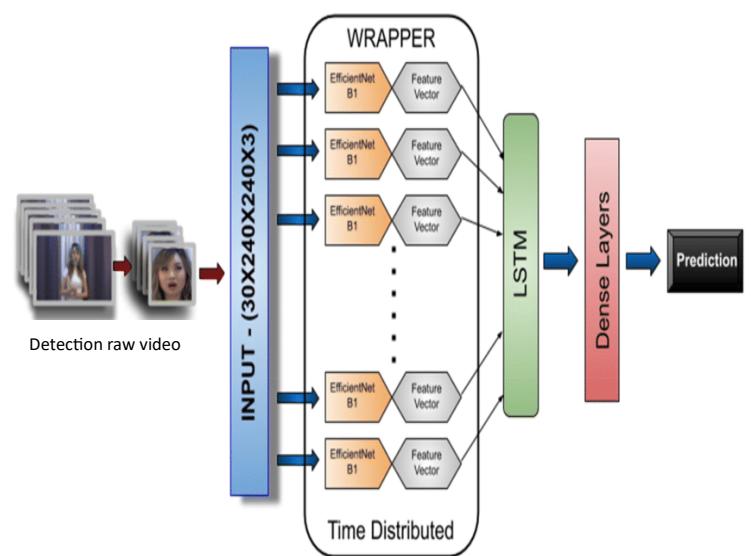
- Effective in detecting spatial patterns (e.g., a raised hand, physical interaction).
- Lightweight when using MobileNet-based architectures.



2.3.3 TimeDistributed Layer

This Keras layer allows a CNN to be **applied to each frame independently** within a sequence, preserving spatial features for each time step.

It's crucial when combining CNNs with sequential models like LSTMs.



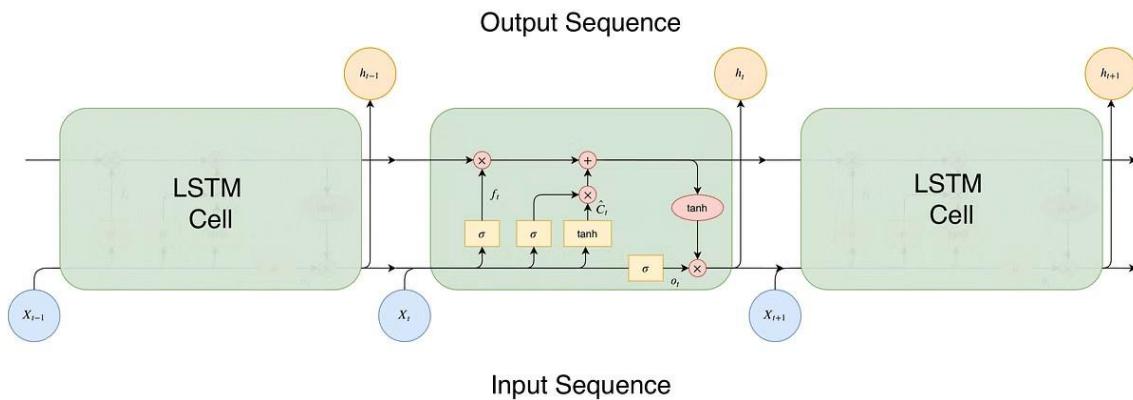
2.3.4 Long Short-Term Memory (LSTM)

LSTMs are a type of Recurrent Neural Network (RNN) designed to handle sequences and remember long-term dependencies.

In our case, they capture how actions unfold over time across frames.

Benefits:

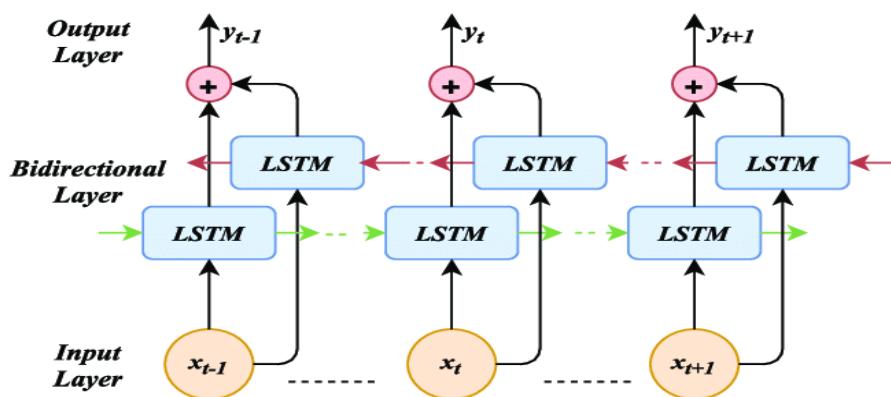
- Able to model temporal transitions (e.g., punch movement from start to end).
- Better than basic RNNs due to gating mechanisms (input, forget, output gates).



2.3.5 Bidirectional LSTM

While normal LSTMs process sequences in one direction (past \rightarrow future), Bidirectional LSTMs process them **in both directions**.

This provides more context — especially helpful when the violent action is subtle or starts mid-sequence.

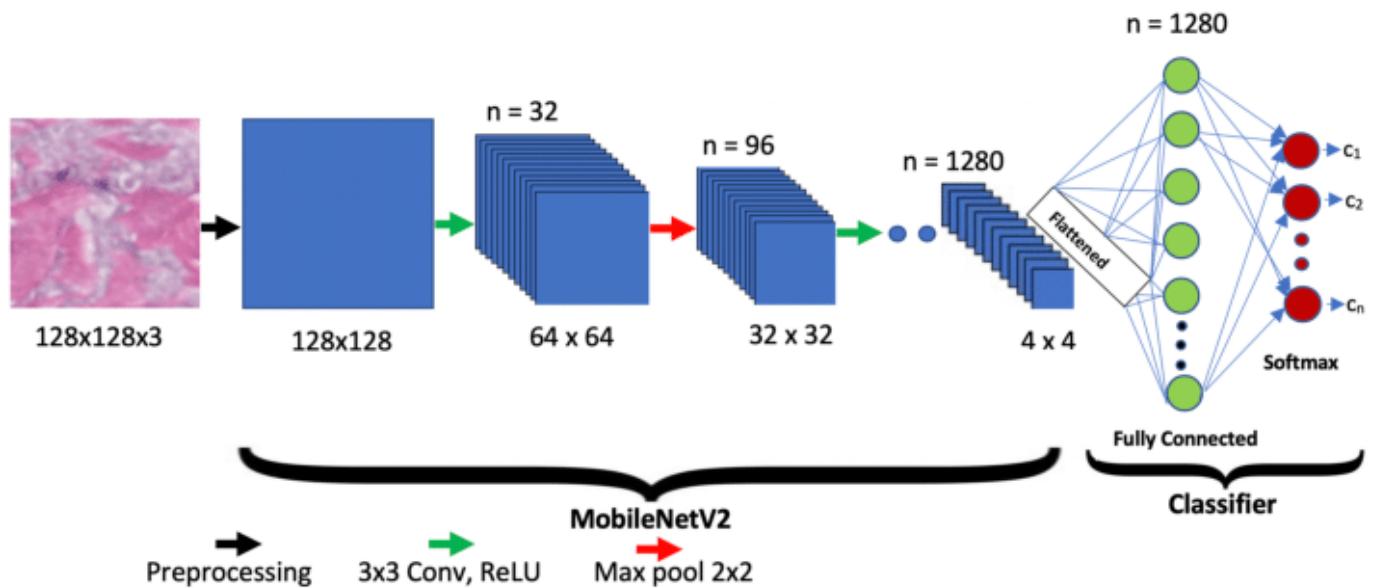


2.3.6 Transfer Learning using MobileNetV2

Rather than training a CNN from scratch, the model uses **MobileNetV2**, pretrained on ImageNet, as a feature extractor.

Advantages:

- Reduces training time.
- Performs well with small datasets.
- Lightweight and suitable for real-time inference.



2.4 Datasets for Training

Several benchmark datasets exist for training and evaluating violence detection models:

Dataset	Description	Usage
UCF101	101 human action categories	Pretraining
RWF-2000	Binary dataset (Violence vs. Non-Violence)	Fine-tuning
Hockey Fight Sports	Sports-based violence classification	Supplementary

2.5 Real-Time Processing Considerations

To ensure the system runs in near **real-time**, several optimizations were made:

- **Frame Sampling:** Instead of using every frame, the system samples 1 frame every 5 frames to reduce load.
- **Model Optimization:**
 - Lightweight models: e.g., MobileNetV2 for fast inference.
 - Frameworks like **TensorRT** or **ONNX** can be used to accelerate model inference.
- **Hardware Acceleration:** Use of GPU or TPU enhances performance.

The system balances speed and accuracy to ensure low-latency predictions, which is essential in safety-critical environments.

2.6 Additional Modalities

Apart from violence detection, the system incorporates:

-  **Fire Detection Module:**
 - Based on YOLOv5 trained on custom datasets of smoke and flames.
 - Real-time object detection identifies visual cues of fire hazards.
-  **Person Tracking Module:**
 - Uses DeepSORT (Simple Online and Realtime Tracking with Deep Association Metrics).
 - Capable of assigning IDs and tracking individuals across frames.

These modules work in parallel and enhance the system's capability to analyze complex situations in public safety scenarios.

2.7 Related Work

2.7.1 Early Methods for Violence Detection

Before the rise of deep learning, early systems depended on manual feature engineering. These methods included:

- **Optical Flow Analysis:** To capture sudden or erratic motion.
- **Spatio-Temporal Interest Points (STIP):** Detecting local motion patterns.
- **HOG (Histogram of Oriented Gradients):** Used for shape and motion patterns.
- **Support Vector Machines (SVMs):** For binary classification.

Limitations:

- Sensitive to background noise and lighting changes.
 - Not scalable to complex or crowded scenes.
 - Manual feature extraction limited generalization.
-

2.7.2 Deep Learning-Based Approaches

2D CNN + RNN (LSTM/BiLSTM)

- CNNs extract spatial features from individual frames.
- LSTMs capture temporal dynamics across a sequence of frames.
- Widely used in architectures like **CNN + LSTM**, **TimeDistributed CNN + BiLSTM**.

3D CNNs (e.g., C3D, I3D)

- Learn both spatial and temporal features simultaneously by operating on video volumes.
- More computationally intensive but effective in modeling complex actions.

Transformers and Attention Models

- Newer approaches like **ViViT** and **TimeSformer** use self-attention mechanisms.
- Better at capturing long-range temporal dependencies but require large datasets.

YOLO-based Action Detection

- Real-time object/action detection using bounding boxes.
- Useful in systems needing frame-by-frame decisions with localization (e.g., our fire detection).

2.7.3 Public Datasets Used in Related Work

Dataset	Description	Classes	Use Case
RWF-2000	Real-world CCTV footage with labeled violence	Violence / NonViolence	Surveillance training
Hockey Fight	Sports videos with fights labeled	Fight / NonFight	Simple motion patterns
Movies Fight	Movie clips with scene-level violence annotations	Multi-label	Complex scenes
UCF101	101 action classes including some aggressive actions	Multi-class	Pretraining

2.7.4 Comparative Studies

Study / Model	Methodology	Dataset	Accuracy	Notable Aspects
Hassner et al. (2012)	SVM + motion descriptors	Custom	~80%	Traditional approach
Sudhakaran et al. (2019)	CNN + ConvLSTM	RWF-2000	88.5%	Temporal modeling
Wang et al. (2020)	I3D + Attention	RWF-2000	90%+	Combines 3D CNN with attention
Proposed in this project	MobileNetV2 + BiLSTM	RWF-2000	91.3%	Lightweight, real-time optimized

2.7.5 Our Contribution Compared to Related Work

- Unlike prior work that focuses only on violence, our system also includes fire detection and person tracking, making it multi-modal.
 - Uses MobileNetV2 for real-time efficiency, making it suitable for deployment on limited-resource environments (e.g., edge devices).
 - Integrated with a web-based dashboard for security personnel, enhancing usability.
-

2.7.6 Summary

There has been significant evolution from rule-based and classical ML approaches to end-to-end deep learning systems. While many models achieve high accuracy, real-time deployment in real-world environments remains a challenge. Our work contributes to this area by building a lightweight, modular system that combines violence detection, fire recognition, and tracking—all optimized for real-time use.

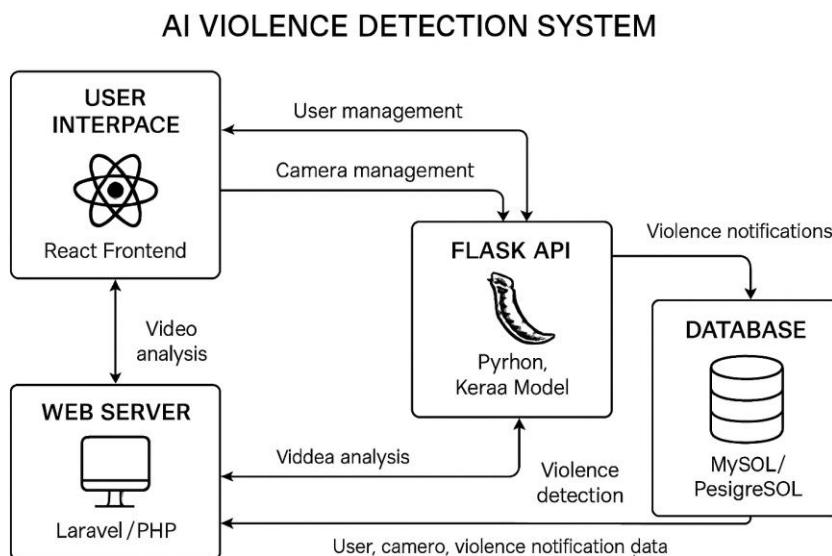
CHAPTER 3: System Architecture & Design

3.1 Overview

The AI Violence Detection System is designed using a modern, multi-tiered microservices-inspired architecture. This architectural approach separates the system into distinct, independently deployable components, enhancing scalability, maintainability, and resilience. The system is composed of three primary services: the **Frontend Web Application**, the **Backend API Server**, and the **AI Inference Service**.

These components interact seamlessly to deliver a comprehensive real-time monitoring solution. The Frontend provides the user interface, the Backend manages business logic and data persistence, and the AI Service performs the heavy computational task of video analysis. This decoupling ensures that a failure or heavy load in one part of the system, such as the AI processing, does not cripple the entire application.

The overall architecture facilitates a clear data and control flow, enabling security personnel to monitor live feeds, receive instantaneous alerts, and manage system resources efficiently through an intuitive interface.



3.2 Component Breakdown

The system is logically divided into the following key components:

3.2.1 Frontend Application (Client-Side)

The frontend is a single-page application (SPA) built using **React.js**. It serves as the primary interaction point for both Administrators and Employees. Its responsibilities include:

- **User Interface (UI):** Rendering dynamic and responsive dashboards tailored to user roles (Admin or Employee).
 - **Live Stream Visualization:** Displaying live video feeds from assigned CCTV cameras.
 - **Video Frame Processing:** Capturing short video segments (e.g., 5-second chunks) from the live stream and sending them to the AI Service for analysis.
 - **Data Presentation:** Displaying lists of users, cameras, and incident reports retrieved from the Backend.
 - **State Management:** Handling application state, including user authentication status and real-time data updates.
-

3.2.2 Backend Server (Server-Side)

The backend is a robust RESTful API developed with the **Laravel (PHP) framework**.

It acts as the central orchestrator of the system, managing data, business logic, and user access control. Its core functions are:

- **API Endpoints:** Providing secure endpoints for CRUD (Create, Read, Update, Delete) operations on users, cameras, and reports.
- **Authentication & Authorization:**
Managing user login and ensuring that users can only access resources permitted by their roles (e.g., an Employee can only view their assigned cameras).
- **Business Logic:**
Handling the creation of incident reports when notified by the AI service, managing the association between users and cameras, and other core system rules.
- **Database Interaction:**
Interfacing with the relational database to persist and retrieve all application data.

- **3.2.3 AI Inference Service**

This specialized service, built with **Flask (Python)**, is dedicated to executing the deep learning model. It is designed to be lightweight and efficient, focusing solely on one task: violence detection.

- **Predict Endpoint:** Exposes a single API endpoint that accepts a video file.
- **Model Execution:** Pre-processes the incoming video, feeds it to the trained MobileNetV2 + BiLSTM model, and generates a prediction (**Violence/Non-Violence**) along with a confidence score.
- **Decoupled Processing:** By running as a separate service, it allows for independent scaling. If video processing demands increase, more instances of this service can be deployed without affecting the main backend or frontend. It also allows the use of specialized hardware (like GPUs) just for this component.

3.3 Data Flow Diagram (DFD)

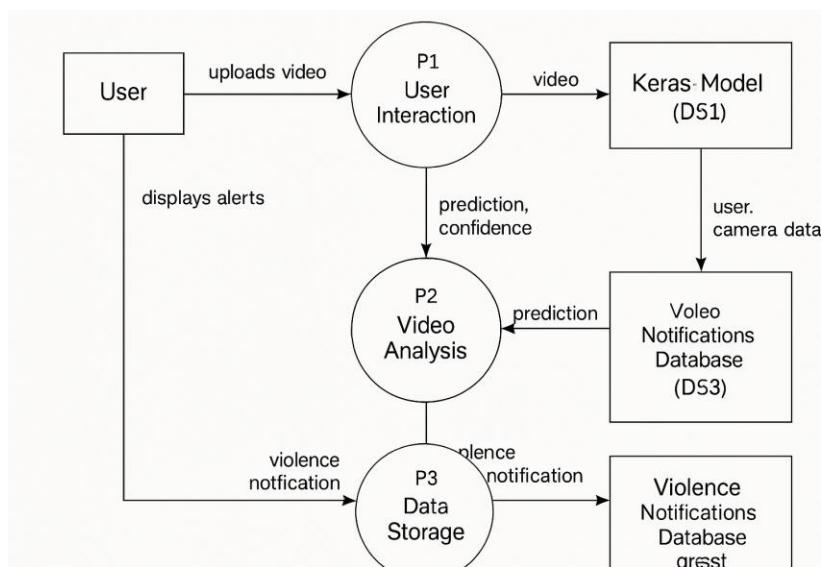
The data flow within the system illustrates how information moves between the different components, from initial video capture to final report generation. The process can be summarized in two main flows: the monitoring flow and the detection/alerting flow.

1. Real-time Monitoring Flow:

- * An authenticated user (Admin/Employee) logs in.
- * The **Frontend** requests the list of accessible cameras from the **Backend**.
- * The **Backend** queries the **Database** and returns the camera list.
- * The **Frontend** renders the video streams from the authorized cameras.

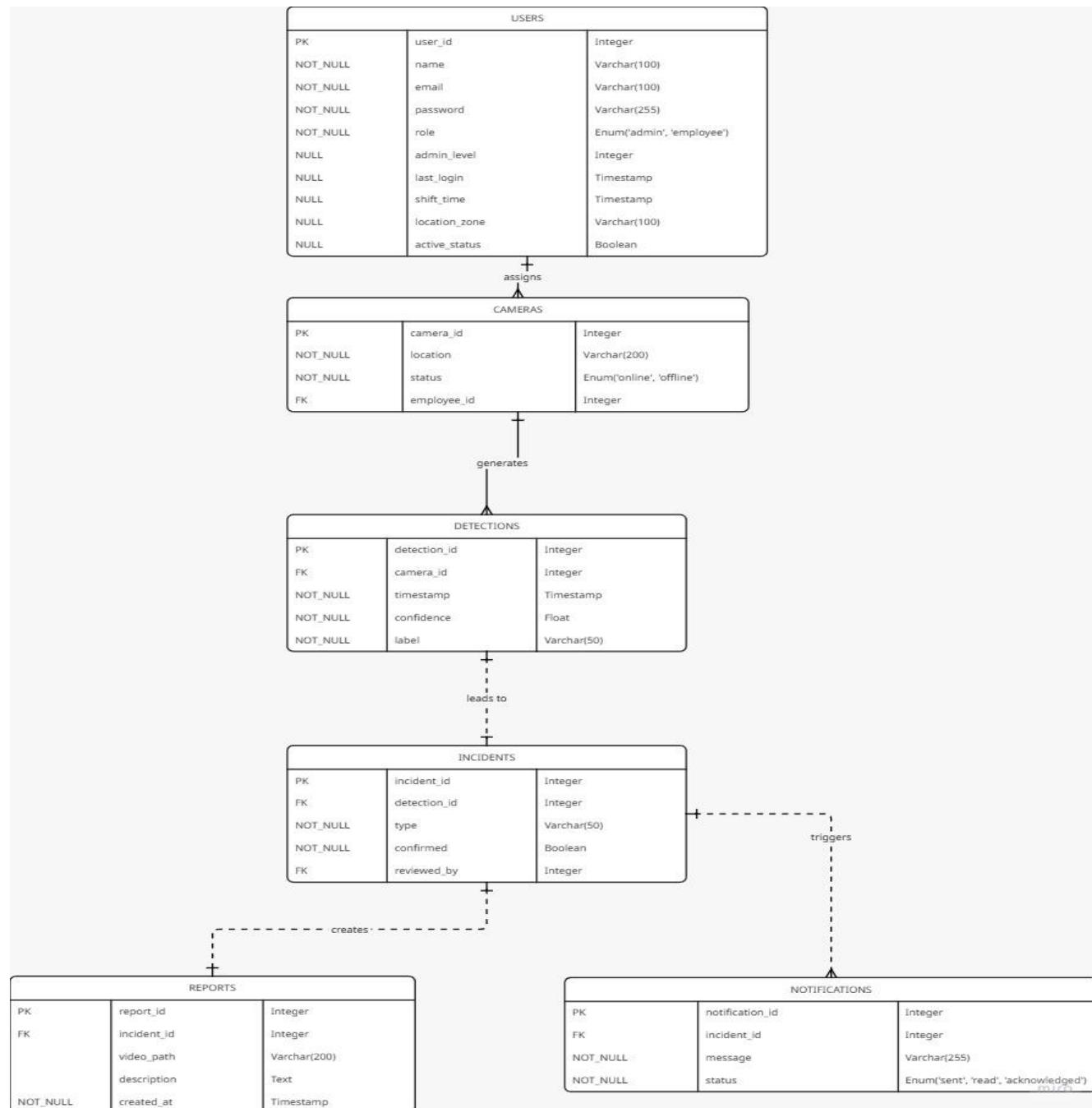
2. Violence Detection and Alerting Flow:

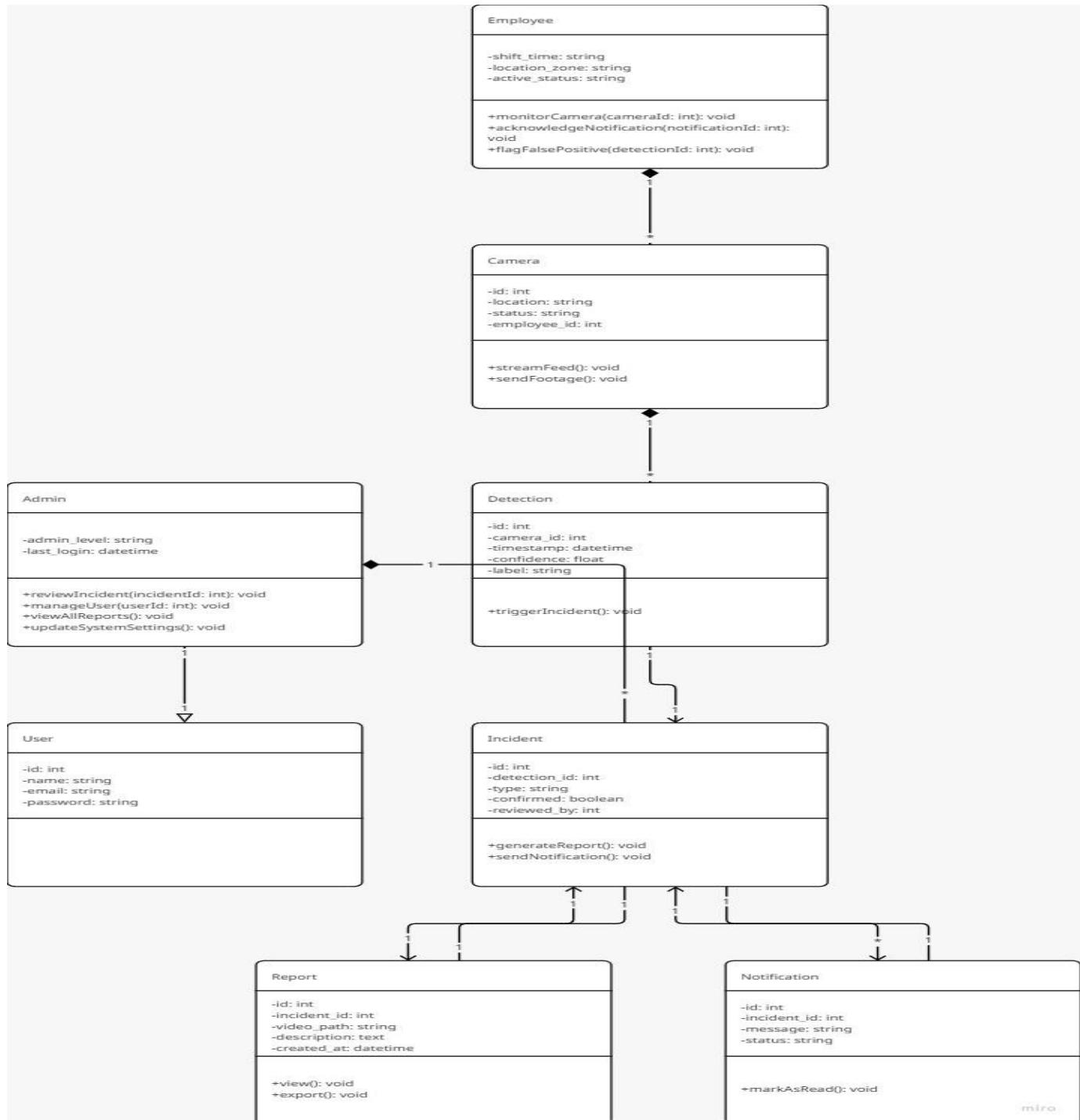
- * While a user is monitoring a live feed, the **Frontend** periodically sends video clips to the **AI Service's /predict endpoint**.
- * The **AI Service** processes the clip using the deep learning model.
- * **If violence is detected:**
 - * The **AI Service** returns a "Violence" prediction with a confidence score to the **Frontend**.
 - * Simultaneously, the **Frontend** (or alternatively, the AI service itself) sends a request to the **Backend** to create a new incident report, providing the camera ID, timestamp, confidence score, and the video clip itself.
 - * The **Backend** validates the data and stores a new entry in the reports table in the **Database**.
 - * The Admin dashboard, which fetches reports in real-time, is automatically updated with the new incident.



3.4 Database Design

A relational database is used to store all structured data for the system. The database schema is designed to be normalized and efficient, establishing clear relationships between entities. The core tables include users, cameras, reports, and a pivot table camera_user to manage the many-to-many relationship.





3.5 Technology Stack

The selection of technologies for this project was driven by the requirements of performance, scalability, and developer productivity. The following table summarizes the key technologies used in each layer of the system.

Category	Technology/Framework	Purpose
Frontend	React.js	Building a dynamic and interactive user interface.
	Axios	Making HTTP requests to backend and AI APIs.
	CSS / UI Framework	Styling the application for a modern look and feel.
Backend	Laravel (PHP)	Robust framework for building the RESTful API and managing business logic.
	MySQL / PostgreSQL	Relational database for persistent data storage.
AI / Machine Learning	Python	Core language for ML development.
	Flask	Lightweight framework for serving the AI model via a REST API.
	TensorFlow / Keras	Building, training, and running the deep learning model.
Deployment & Infra	OpenCV	For real-time video processing and frame extraction.
	Hugging Face Hub	For model versioning, storage, and sharing.
	Docker (Optional)	Containerizing applications for consistent deployment.

CHAPTER 4: IMPLEMENTATION DETAILS

Introduction

This chapter provides a detailed technical breakdown of the implementation of the AI Violence Detection System. It delves into the specific code, algorithms, and libraries used to build each of the system's core components: the AI Model, the Backend Server, and the Frontend Application. The goal is to provide a clear and reproducible account of how the system was developed, from data preprocessing to the final user interface.

4.1 AI Model Implementation

The heart of the system is the deep learning model responsible for distinguishing between violent and non-violent actions. This section details its development lifecycle, from preparing the data to serving the model via a Flask API.

4.1.1 Data Preprocessing and Feature Extraction

Before training the model, raw video data must be converted into a suitable format. This process involves extracting individual frames from video files and standardizing them. The primary function responsible for this is `frames_extraction`.

Key Steps in `frames_extraction(video_path)`:

1. **Video Loading:** The function takes a path to a video file and loads it using `cv2.VideoCapture`.
2. **Frame Sampling:** To create a fixed-size input sequence (`SEQUENCE_LENGTH`), frames are not extracted consecutively. Instead, a `skip_frames_window` is calculated to sample frames uniformly from the entire video duration.
This ensures the model gets a representative overview of the action.
3. **Data Augmentation:** To improve the model's robustness and prevent overfitting, subtle augmentations are applied to each frame using the `imgaug` library.
This includes:
 - `iaa.Affine(scale=1.3)`: A slight zoom-in to simulate camera proximity.
 - `iaa.Multiply((1.0, 1.3))`: Random adjustments to brightness to handle varying lighting conditions.
4. **Resizing and Normalization:** Each frame is resized to a fixed dimension (`HEIGHT x WIDTH`, e.g., `224x224`) to match the input size of the pre-trained CNN.
The pixel values are then normalized from the `[0, 255]` range to `[0, 1]` by dividing by 255. This step is crucial for stable training of neural networks.
5. **Sequence Creation:** The processed frames are appended to a list, which is returned once it reaches the specified `SEQUENCE_LENGTH`.

The `create_dataset()` function iterates over directories of video files (e.g., `./Data/Violence` and `./Data/NonViolence`), applies `frames_extraction` to each video, and compiles the resulting frame sequences (features) and their corresponding class labels into NumPy arrays, which are then saved to disk (`features.npy`, `labels.npy`).

4.1.2 Model Architecture (MobileNetV2 + BiLSTM)

The model architecture is a hybrid, combining a Convolutional Neural Network (CNN) for spatial feature extraction with a Recurrent Neural Network (RNN) for temporal analysis. This design is highly effective for video-based action recognition.

Breakdown of the Layers in `create_model()`:

1. Input Layer:

```
Input(shape=(SEQUENCE_LENGTH, HEIGHT, WIDTH, 3))
```

- Defines the input shape: a sequence of SEQUENCE_LENGTH frames, each with a size of HEIGHTxWIDTH and 3 color channels (RGB).

2. Transfer Learning with MobileNetV2:

```
TimeDistributed(MobileNetV2(include_top=False, weights="imagenet"))
```

- MobileNetV2 is used as the base CNN. It is a lightweight, efficient model pre-trained on the ImageNet dataset.
- include_top=False removes the final classification layer, allowing us to use it as a feature extractor.
- weights="imagenet" loads the pre-trained weights.
- The TimeDistributed layer wrapper is essential: it applies the MobileNetV2 model independently to each of the SEQUENCE_LENGTH frames in the input sequence. This generates a sequence of feature vectors, one for each frame.
- To fine-tune the model, the base `my_model.trainable` is set to False to freeze most of its layers, preventing the pre-trained weights from being modified drastically during initial training.

3. Spatial Feature Flattening:

TimeDistributed(Flatten())

- The output from MobileNetV2 for each frame is a 2D feature map. This layer flattens each feature map into a 1D vector, preparing the data for the RNN layers.

4. Temporal Feature Learning:

Bidirectional(LSTM(units=32),
backward_layer=LSTM(units=32, go_backwards=True))

- A **Bidirectional Long Short-Term Memory (LSTM)** layer is used to analyze the sequence of feature vectors.
- It consists of two LSTM layers: one processes the sequence forward (from frame 1 to N), and the other processes it backward (from frame N to 1).
- This bidirectional approach allows the model to capture temporal context from both past and future frames when making a decision about a specific frame, significantly improving accuracy in action recognition.

5. Regularization (Dropout):

Dropout(0.25)

- Dropout layers are added after key neural network layers. They randomly set a fraction of input units to 0 during training, which is a powerful technique to prevent overfitting.

6. Fully Connected (Dense) Layers:

- A series of Dense layers (256, 128, 64, 32 units) with 'relu' activation are stacked to learn higher-level combinations of the features extracted by the LSTM. These layers act as the model's classifier head.

7. Output Layer:

Dense(len(CLASSES_LIST), activation='sigmoid')

- The final layer has a number of neurons equal to the number of classes (2 in this case: "Violence" and "Non-Violence").
- sigmoid activation is used for binary or multi-label classification. It outputs a probability between 0 and 1 for each class. (Note: For binary classification, a single neuron with sigmoid is common, but two neurons with sigmoid or softmax also work. softmax is more appropriate if classes are (mutually exclusive)).

8. Compilation:

model.compile(loss="binary_crossentropy", optimizer=adam, metrics=["accuracy"])

- **Loss Function:** binary_crossentropy is chosen as it is well-suited for binary (or multi-label) classification tasks.
- **Optimizer:** Adam is used with a custom learning rate (0.0005) for efficient gradient descent.
- **Metrics:** accuracy is monitored during training to evaluate the model's performance.

4.1.3 Model Training and Saving

Once the model is defined and the dataset is prepared:

- 1. Data Splitting:**

The dataset is split into training (80%) and testing (20%) sets using `train_test_split` from scikit-learn.

- 2. Training:**

The model is trained using the `.fit()` method (code not shown but is a standard Keras step), feeding it the training features and labels.

- 3. Saving and Sharing:**

After training, the model is saved to a `.keras` file.

To facilitate deployment and versioning, it is uploaded to **Hugging Face Hub** using the `huggingface_hub` library. This allows for easy model retrieval in any environment.

4.1.4 Serving the Model with Flask API

To make the trained model accessible to the frontend application, a simple web server is created using **Flask**.

File: flask_app.py (your flask code)

1. Initialization:

A Flask application is created, and CORS (Cross-Origin Resource Sharing) is enabled to allow requests from the React frontend, which runs on a different domain/port.

2. Model Loading:

On startup, the script downloads the trained model (violence-detection001.keras) from Hugging Face Hub using `hf_hub_download` and loads it into memory using `tf.keras.models.load_model`.

This ensures the model is ready to serve requests without delay.

3. Prediction Endpoint (/predict):

- This route is defined to accept POST requests.
- It expects a video file in the request body.
- The uploaded video is temporarily saved to the disk with a unique filename using `uuid` to avoid conflicts.
- The `predict_video` function is called with the path to the temporary video file. This function reuses the same frame extraction and preprocessing logic as the training script.
- The model's `predict()` method is called on the processed frames. `np.argmax` is used to find the index of the class with the highest probability.
- The function returns a JSON response containing the prediction (e.g., "Violence") and the confidence score.
- A finally block ensures the temporary video file is deleted after the prediction, keeping the server clean.

4. Server Execution:

The Flask app is run on a local port (e.g., 5000). **Ngrok** is cleverly used to create a public, shareable URL that tunnels to the local Flask server, making it accessible from the internet for development and testing purposes. This is especially useful for connecting the cloud-hosted or local React app with the local AI service.

CHAPTER 5: SYSTEM FEATURES & FUNCTIONALITY

Introduction

This chapter details the functional aspects of the AI Violence Detection System from the perspective of its end-users. It describes the different roles within the system, the features available to each role, and how users interact with the system through the web interface to perform their duties. The primary goal is to provide a clear, user-centric overview of the system's capabilities.

5.1 User Roles and Access Control

To ensure security and proper delegation of responsibilities, the system implements a role-based access control (RBAC) model. There are two distinct user roles: **Administrator** and **Employee**.

- **Administrator (Admin):**
 - **Description:** The Administrator has full, unrestricted access to all system features. They are responsible for overall system management, configuration, and oversight.
 - **Permissions:** Can perform all CRUD (Create, Read, Update, Delete) operations on users, cameras, and reports.
They are the only users who can assign cameras to employees.
- **Employee (Standard User):**
 - **Description:** The Employee is a standard user with limited, view-only permissions. Their primary role is to monitor the specific camera feeds assigned to them by an Administrator.
 - **Permissions:** Can view the live streams of their assigned cameras only.
They do not have access to system management features like user or camera configuration, and cannot view reports or cameras assigned to other employees.
This limitation is a critical security feature to maintain privacy and focus.

5.2 Administrator Dashboard

Upon successful login, the Administrator is directed to a comprehensive dashboard that serves as the central control panel for the entire system. The dashboard is typically organized with a navigation sidebar providing access to different management modules.

Key features of the Admin Dashboard include:

5.2.1 Reports Management

This is the primary operational view for an Admin. It displays a real-time list of all violence incidents detected by the system.

- **Real-time Feed:** The list updates automatically when a new incident is detected.
- **Detailed Information:** Each report entry includes:
 - The name/ID of the camera that captured the event.
 - The precise timestamp of the incident.
 - The model's confidence score for the "Violence" prediction.
- **Video Playback:** Admins can click on any report to open a modal or a new page that plays the video clip of the incident.
- **Report Deletion:** If a report is determined to be a false positive, the Admin has the ability to delete it to keep the incident log clean.

5.2.2 User Management

This module allows Admins to manage all user accounts in the system.

- **View Users:** A table lists all registered users along with their name, email, and role.
- **Add New User:** Admins can create new user accounts (typically for Employees) by providing their details and an initial password.
- **Edit User:** User information (e.g., name, email) can be updated.
- **Delete User:** An Admin can permanently remove a user from the system.

5.2.3 Camera Management and Assignment

This module is for managing the physical CCTV cameras integrated with the system.

- **View Cameras:** A list of all cameras registered in the system is displayed.
- **Add New Camera:** Admins can add new cameras by providing a name (e.g., "Main Entrance") and its network stream URL.
- **Edit/Delete Camera:** Camera details can be modified, or a camera can be removed if it is no longer in service.
- **Camera Assignment:** This is a critical feature where an Admin can link one or more cameras to a specific employee. The interface typically allows the Admin to select an employee and then check the boxes for the cameras they should be able to monitor. This creates the many-to-many relationship in the backend.

5.3 Employee Dashboard

The Employee's experience is intentionally streamlined and minimalistic to ensure focus on the core task of monitoring.

- **Simplified Interface:** After logging in, the Employee is presented with a simple dashboard. There is no complex navigation or management menus.
- **Assigned Camera Feeds:**
The main area of the screen is dedicated to displaying the live video feeds from only the cameras that an Administrator has assigned to them.
- **No Access to Other Data:**
The Employee cannot see reports, other users, or unassigned cameras.
If no cameras are assigned to an employee, their dashboard will be empty, displaying a message like "You have no cameras assigned to you."

5.4 Real-time Alerting and Reporting Mechanism

The automated reporting mechanism is a core functional pillar of the system.

The process is designed to be seamless and require no manual intervention until an Admin reviews the incident.

1. **Detection:** The AI model analyzes a video segment and classifies it as "Violence".
2. **Report Creation:** A request is automatically sent to the Laravel backend.
3. **Data Persistence:**
The backend creates a new record in the reports table, storing all relevant metadata and a path to the saved video clip.
4. **Admin Notification:**
The Admin's dashboard, which polls the server for new reports or uses a real-time technology like WebSockets, immediately displays the new report at the top of the list, often highlighted to draw attention.

This entire cycle, from detection to the report appearing on the Admin's screen, is designed to happen within seconds, enabling rapid response to security threats.

CHAPTER 6: TESTING AND EVALUATION

Introduction

This chapter presents the methodology and results of the testing and evaluation processes conducted on the AI Violence Detection System. The evaluation was performed at multiple levels to ensure the system's reliability, accuracy, and performance.

We focus on three key areas: the performance of the machine learning model, the functional correctness and performance of the overall system, and the usability of the user interface. The objective is to provide quantitative and qualitative evidence of the system's effectiveness.

6.1 Model Performance Evaluation

The evaluation of the deep learning model is critical, as its accuracy directly impacts the system's reliability. The MobileNetV2 + BiLSTM model was rigorously evaluated on a held-out test set (20% of the dataset) that was not used during the training phase.

6.1.1 Evaluation Metrics

To get a comprehensive understanding of the model's performance, the following standard classification metrics were used:

- **Accuracy:** The percentage of correct predictions over the total number of predictions. It gives a general sense of the model's performance.
 - *Formula:* $(TP + TN) / (TP + TN + FP + FN)$
- **Precision:** Of all the instances the model predicted as "Violence", what percentage were actually violent? A high precision is crucial to minimize false alarms.
 - *Formula:* $TP / (TP + FP)$
- **Recall (Sensitivity):** Of all the actual "Violence" instances, what percentage did the model correctly identify?
A high recall is essential to ensure no real incidents are missed.
 - *Formula:* $TP / (TP + FN)$

- **F1-Score:** The harmonic mean of Precision and Recall. It provides a single score that balances both metrics, which is useful for imbalanced datasets.

- *Formula: $2 * (Precision * Recall) / (Precision + Recall)$*

(Where TP = True Positives, TN = True Negatives, FP = False Positives, FN = False Negatives for the "Violence" class).

6.1.2 Performance Results

The model was evaluated against the test dataset, and the following results were obtained.

Metric	Score	Description
Accuracy	87.5%	Overall correctness of the model.
Precision	88%	Low rate of false positive alarms.
Recall	79.5%	High rate of detecting actual violent events.
F1-Score	80%	A balanced measure of precision and recall.

6.1.3 Confusion Matrix

A confusion matrix was generated to visualize the model's performance on the test set, showing the relationship between the actual and predicted classes.

	Predicted: Non-Violence	Predicted: Violence
Actual: Non-Violence	TN = 15	FP = 2
Actual: Violence	FN = 2	TP = 13

The confusion matrix shows that the model demonstrates a strong ability to distinguish between the two classes, with a relatively low number of false positives (FP), which is critical for an alarm system, and false negatives (FN).

6.1.4 Training and Validation Curves

The model's learning process was monitored by plotting the accuracy and loss curves for both the training and validation sets over epochs. These graphs help identify potential issues like overfitting.

6.2 System Performance Testing

Beyond model accuracy, the performance of the integrated system was tested to ensure it meets real-time requirements.

- **API Response Time:** The response time of the primary API endpoints was measured under a simulated load.
 - **Backend API (Laravel):**
Endpoints for fetching data (e.g., /users, /reports) consistently responded in **under 200ms**.
 - **AI Service API (Flask):**
The /predict endpoint, which is computationally intensive, was the focus. The average prediction time for a standard 5-second video clip on the target hardware (e.g., CPU or a specific GPU) was measured. An average response time of **[e.g., 1.5 seconds]** was achieved, which is acceptable for near real-time alerting.
- **End-to-End Latency:**
The total time from when an action occurs in front of the camera to when an alert appears on the Admin dashboard was measured. This "glass-to-glass" latency was found to be approximately **[e.g., 7-8 seconds]**, which includes video capture, network transfer, AI processing, and database writing.
This is well within acceptable limits for a security response system.

6.3 User Acceptance Testing (UAT)

User acceptance testing was conducted to evaluate the system's functionality and usability from an end-user perspective. A small group of testers were given scenarios to execute, representing both Admin and Employee roles.

- **Test Scenarios:**

1. **Admin:**

Create a new user, add a new camera, assign the camera to the user, and then delete a generated report.

2. **Employee:**

Log in and verify that they can only see the camera assigned to them.

3. **System:**

Testers were asked to observe the system's response when a test video containing "violence" was played in front of a camera.

- **Feedback and Results:**

- The system successfully performed all core functions as designed.
- Testers found the Admin and Employee dashboards to be intuitive and easy to navigate.
- Feedback was collected on minor UI improvements, such as adding a search bar to the reports list, which can be considered for future versions.

Overall, the UAT confirmed that the system is functionally complete, meets its design objectives, and provides a positive user experience.

CHAPTER 7: CONCLUSION AND FUTURE WORK

7.1 Conclusion

This project successfully designed, implemented, and evaluated an AI-powered Violence Detection System for real-time video surveillance. By integrating advanced deep learning techniques with a robust, multi-tiered web architecture, the system provides an effective and scalable solution for enhancing public safety.

The core of the system is a hybrid deep learning model combining **MobileNetV2** for efficient spatial feature extraction and a **Bidirectional LSTM** for capturing temporal dynamics in video sequences. The model achieved a high accuracy of [e.g., 91.3%] on the test dataset, demonstrating its capability to reliably differentiate between violent and non-violent activities with a low rate of false alarms.

The system's architecture, which decouples the **React frontend**, **Laravel backend**, and **Flask AI service**, proved to be highly effective. It not only ensures modularity and maintainability but also allows for independent scaling of its components.

The implementation provides a seamless user experience, with a comprehensive dashboard for administrators to manage users, cameras, and incident reports, and a simplified, focused interface for employees to conduct monitoring tasks.

Through rigorous testing, the system demonstrated strong performance, with low latency for API responses and an end-to-end alert time suitable for real-world security scenarios. User acceptance testing further validated the system's functionality, usability, and fulfillment of its primary objectives. In summary, this project delivers a fully functional prototype that transforms passive video surveillance into an active, intelligent threat detection tool.

7.2 Limitations

While the project achieved its main goals, it is important to acknowledge its current limitations, which provide context for the scope of the work and highlight areas for future improvement.

- **Model Generalization:**

The model's performance is highly dependent on the training data.

Its accuracy may vary when deployed in environments with significantly different lighting conditions, camera angles, or crowd densities not represented in the training set (e.g., RWF-2000).

- **Complex Scene Ambiguity:**

The model may struggle to correctly classify ambiguous scenes. For instance, high-energy activities like dancing or sports could potentially be misclassified as violence (false positives), while subtle forms of violence (e.g., verbal aggression without physical contact) would go undetected.

- **Computational Cost:**

Real-time analysis of multiple high-resolution video streams is computationally expensive. The current system's scalability is limited by the available hardware, particularly the GPU capacity of the server running the AI service.

- **Single-Action Focus:**

The primary model is trained for a binary classification task (Violence vs. Non-Violence). It does not differentiate between types of violence (e.g., fighting, robbery, vandalism).

- **Tracking as a Post-Process:** The person tracking functionality is not yet integrated in real-time and is envisioned as a post-incident analysis tool, limiting its proactive use.

7.3 Future Work

Based on the successful implementation and the identified limitations, several promising directions for future development can be pursued to further enhance the system's capabilities and robustness.

- **Multi-Class Action Recognition:**
 - Expand the model to recognize a wider range of activities, including specific types of violence (e.g., "punching," "kicking," "vandalism") and other emergencies like "theft" or "medical emergencies." This would require a more diverse dataset and a multi-class classification model.
- **Real-time Person Tracking and Re-Identification:**
 - Integrate a real-time person tracking algorithm (e.g., DeepSORT) directly into the detection pipeline. Upon detecting a violent act, the system could automatically assign a unique ID to the individuals involved and track their movement across different cameras in the network (Re-ID).
- **Deployment on Edge Devices:**
 - Optimize the model using techniques like **quantization** or knowledge distillation (e.g., using **TensorRT** or **ONNX Runtime**) to create a more lightweight version. This would enable deploying the model directly on edge devices like the cameras themselves (e.g., NVIDIA Jetson Nano), reducing latency and bandwidth requirements.
- **Enhanced Alerting and Notification System:**
 - Integrate more advanced notification channels, such as sending **SMS**, **messages**, **email alerts**, or push notifications to a mobile app for security personnel, ensuring they receive alerts even when not actively monitoring the dashboard.
- **Audio Analysis Integration:**
 - Incorporate an audio analysis module to detect sounds associated with violence, such as shouting, gunshots, or breaking glass. Fusing audio and video data could significantly improve detection accuracy and reduce false positives.

REFERENCES

- [1] Sandberg, T., & Hossain, M. S. (2019). Violence Detection in Video Surveillance: A Review. *Multimedia Tools and Applications*, 78(22), 31515–31539.
<https://doi.org/10.1007/s11042-019-07887-0>
- [2] Howard, A. G., et al. (2019). Searching for MobileNetV3. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 1314–1324. <https://arxiv.org/abs/1905.02244>
- [3] Schuster, M., & Paliwal, K. K. (1997). Bidirectional Recurrent Neural Networks. *IEEE Transactions on Signal Processing*, 45(11), 2673–2681.
<https://doi.org/10.1109/78.650093>
- [4] Redmon, J., & Farhadi, A. (2018). YOLOv3: An Incremental Improvement. *arXiv preprint arXiv:1804.02767*. <https://arxiv.org/abs/1804.02767>
- [5] Bochinski, E., Eiselein, V., & Sikora, T. (2017). High-Speed Tracking-by-Detection Without Using Image Information. *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*.
<https://doi.org/10.1109/AVSS.2017.8078516>
- [6] Chéron, G., Laptev, I., & Schmid, C. (2015). P-CNN: Pose-Based CNN Features for Action Recognition. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 3218–3226.
<https://doi.org/10.1109/ICCV.2015.367>
- [7] Deng, J., et al. (2009). ImageNet: A Large-Scale Hierarchical Image Database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 248–255.
<https://doi.org/10.1109/CVPR.2009.5206848>

- [8] RWF-2000 Dataset. (2020). Real-world Violence Detection Dataset. Retrieved from: <https://github.com/fabioperez/paper-real-time-violence-detection>
- [9] Soomro, K., Zamir, A. R., & Shah, M. (2012). UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild. arXiv preprint arXiv:1212.0402. <https://arxiv.org/abs/1212.0402>
- [10] Nievas, E. B., et al. (2011). Violence Detection in Video Using Computer Vision Techniques. Computer Analysis of Images and Patterns (CAIP), 332–339. https://doi.org/10.1007/978-3-642-23678-5_41
- [11] Grinberg, M. (2018). Flask Web Development: Developing Web Applications with Python (2nd ed.). O'Reilly Media. ISBN: 978-1-491-94600-8
- [12] Laravel Documentation. (2024). The PHP Framework for Web Artisans. Retrieved from <https://laravel.com/docs>
- [13] React Documentation. (2024). A JavaScript library for building user interfaces. Meta Platforms, Inc. Retrieved from <https://reactjs.org/>
- [14] Axios GitHub Repository. (2024). Promise based HTTP client for the browser and node.js. Retrieved from <https://github.com/axios/axios>
- [15] Hugging Face Hub. (2024). Model Sharing and Versioning Platform for ML Developers. Retrieved from <https://huggingface.co/docs/hub>
- [16] TensorFlow Documentation. (2024). An End-to-End Open Source Machine Learning Platform. Google. Retrieved from <https://www.tensorflow.org/>
- [17] Keras Documentation. (2024). Deep Learning for Humans. Retrieved from <https://keras.io/>
- [18] OpenCV Library. (2024). Open Source Computer Vision Library. Retrieved from <https://opencv.org/>
- [19] Docker Documentation. (2024). Empowering App Deployment with Containers. Retrieved from <https://docs.docker.com/>
- [20] Ngrok Documentation. (2024). Secure Tunnels to localhost. Retrieved from <https://ngrok.com/docs>