

# AI-Powered Video Surveillance Anomaly Detection System

## Technical Report

**Author:** Siva Tanay Jagarapu  
**Project Date:** December 27, 2025  
**Project:** Drone Security Analyst - Intelligent Video Processing System

### ■ Project Information

Parameter	Details
Project Date	December 27, 2025
Dataset	UCF Anomaly Detection Dataset

## Executive Summary

This report presents an AI-powered video surveillance system designed to detect anomalies in security footage through automated frame analysis, semantic indexing, and intelligent alerting. The system successfully processes surveillance videos, identifies critical security threats (weapons, assaults, suspicious activities), and provides context-aware recommendations using Vision-Language Models (VLMs) and vector-based similarity search.

The prototype was tested on the **UCF Anomaly Detection Dataset**, specifically focusing on **stealing/robbery scenarios** to demonstrate real-world security incident detection capabilities.

## 1. Problem Statement

Traditional video surveillance systems face several critical challenges:

### Objective

Develop an automated video surveillance analysis system that:

1. Extracts and analyzes video frames using AI
2. Generates detailed descriptions of visual content with temporal context
3. Detects security threats in real-time using rule-based + context-aware alerting
4. Indexes frames for semantic search and retrieval using CLIP embeddings
5. Provides contextual understanding through:
6. **Sequential context:** Past 30 frames for temporal progression

7. **Semantic context:** Top-5 similar past frames from vector database (image similarity)
8. Generates comprehensive security reports with actionable alerts

---

## 2. Dataset Selection & Justification

---

### Chosen Dataset: UCF Anomaly Detection Dataset

#### Rationale for Selection:

Given the time constraints and lack of readily available labeled drone surveillance datasets, I selected the **UCF Anomaly Detection Dataset** for the following reasons:

1. **Real-World Scenarios:** Contains actual CCTV footage of various anomalies including:
  2. Theft/Stealing incidents
  3. Fighting/Assault
  4. Vehicle-related crimes

Suspicious activities

**Quality & Diversity:** High-resolution videos with varying lighting conditions, camera angles, and environments

**Accessibility:** Publicly available dataset with clear licensing

**Validation Potential:** Well-documented incidents allow for accuracy validation

---

## 3. System Architecture & Approach

---

### 3.1 High-Level Architecture

#### Key Data Flows:

1. **Frame Processing Pipeline (Per Frame):**
  2. Raw Frame → CLIP Image Encoder → 512-dim embedding
  3. Embedding → Stored in Pinecone with metadata (timestamp, telemetry, frame path)
  4. Embedding → Search Pinecone for top-5 visually similar past frames
  5. Retrieve past 30 frames sequentially for temporal context
  6. Current frame + Past 30 frames + Top-5 similar frames → Gemini Vision + LLM

Result: Context-aware frame description + Threat analysis

#### LLM-Based Security Analysis (via `security_analysis.py`):

9. Frame description + Context → Gemini LLM analyzes threat level
10. LLM generates: `threat_level` (NONE/LOW/MEDIUM/HIGH/CRITICAL), `analysis`, `alerts`, `objects_detected`

11. LLM evaluates visible threats in current frame with historical context
12. Alert generation uses LLM reasoning based on security\_analysis.py prompt

Output: Structured JSON with threat assessment and alert message

#### **Rule-Based Alert Enhancement:**

15. LLM-generated alerts are supplemented by keyword pattern matching
16. Additional rules check for: weapons, loitering, suspicious behavior, running, collisions
17. If alert triggered → Use CLIP embedding to search Pinecone for similar past incidents
18. Similar past frames → Referenced in alert message for historical context

Alert stored with severity (CRITICAL, HIGH, MEDIUM) + timestamp

#### **Vector Database Operations:**

21. Storage: CLIP image embeddings (512-dim) + metadata (description, timestamp, telemetry)

Search Use Cases:

- 18a. Frame description context (top-5 similar frames)
- 18b. Alert context enhancement (similar past incidents)
- 18c. User semantic search (text query → CLIP text encoder → vector search)

#### **Final Summary Generation:**

24. All frame descriptions + All generated alerts → Gemini LLM
25. LLM synthesizes: Executive summary, timeline of events, risk assessment, recommendations
26. Result: Comprehensive security report in JSON format

## **3.2 Component-by-Component Breakdown**

### **A. Frame Extraction Module**

**Technology:** OpenCV (cv2)

**Approach:** Time-based interval extraction

- Extracts 1 frame every 2 seconds
- Ensures consistent temporal coverage
- Reduces computational load compared to processing every frame
- Timestamps embedded in filenames for traceability

**Justification:**

- 2-second intervals balance between capturing event progression and processing efficiency
- For a 1-minute video, this yields 30-32 frames (manageable for real-time processing)
- Maintains temporal context for detecting prolonged activities (loitering, tampering)

B. Vision-Language Model (VLM) - Gemini Vision

Technology: Google Gemini 2.5 Flash

Context-Aware Frame Description Processing:

When generating a description for each frame, the system provides rich temporal context to improve accuracy:

- 1. **Sequential Context (Past 30 Frames):**
- 2. Retrieves the last 30 processed frames in chronological order
- 3. Provides temporal progression (e.g., "person approaches car → opens hood → tampering detected")

Enables tracking of continuous activities across time

Semantic Context (Top-5 Similar Frames):

- 6. Uses CLIP image embeddings to search Pinecone vector database
- 7. Finds the 5 most visually similar past frames based on image similarity
- 8. Excludes very recent frames (last 3) to avoid redundancy

Provides context from similar past incidents (e.g., "similar weapon detected at 00:00:10")

LLM Analysis with Full Context:

- 11. Current frame description + Past 30 frame descriptions + Top-5 similar events
- 12. Gemini Vision analyzes with comprehensive temporal and semantic awareness
- 13. Generates contextually accurate descriptions that reference past activities

Example Context-Aware Description:

"The person in the green shirt continues tampering with the vehicle's hood (similar activity observed at 00:00:28). The armed individual previously seen at 00:00:20 remains in the background."

Why Gemini over alternatives?

Criterion	Gemini 2.5 Flash	BLIP/BLIP-2	GPT-4 Vision
API Availability	[✓] Free tier (60 req/min)	[X] Self-hosted only	[X] Paid only
Vision Quality	★★★★★ State-of-art	★★★★ Good	★★★★★ Excellent
Speed	Fast (1-2s/frame)	Slow (local inference)	Medium (API latency)
Context Window	1M tokens	Limited	128K tokens
Cost	Free tier adequate	Compute costs	\$0.01 per image
Memory Footprint	Cloud-based (0 local)	~8GB GPU required	Cloud-based (0 local)

Decision: Gemini 2.5 Flash provides the optimal balance of quality, cost, speed, and integration.

### C. Embedding Model - CLIP (OpenAI)

**Technology:** OpenAI CLIP (ViT-B/32)

**Why CLIP over alternatives?**

Feature	CLIP	ResNet + Word2Vec	BLIP Embeddings
Vision-Text Alignment	[✓] Native	[X] Separate models	[✓] Native
Zero-Shot Capability	[✓] Excellent	[X] Limited	[✓] Good
Memory Footprint	~500MB	~200MB	~2GB
Embedding Dimension	512	2048 (ResNet)	768
Performance	★★★★★	★★★	★★★★
Inference Speed	Fast (~50ms)	Fast (~30ms)	Slow (~200ms)

**Decision Factors:**

1. **Vision-Language Alignment:** CLIP embeddings capture semantic similarity between images and text queries
2. **Memory Efficiency:** 512-dimensional embeddings balance expressiveness and storage
3. **Proven Performance:** Industry standard for image-text retrieval tasks
4. **Easy Integration:** HuggingFace Transformers library support

**Alternative Considered - Text Embeddings:**

- Initially planned to also index frame descriptions using text embeddings (e.g., Sentence-BERT)
- **Time Constraint:** Focused on image embeddings only for MVP
- **Future Enhancement:** Hybrid search combining image + text vectors

### D. Vector Database - Pinecone

**Technology:** Pinecone (Serverless)

**Why Pinecone over alternatives?**

Criterion	Pinecone	Milvus	FAISS	Qdrant
Hosting	[✓] Managed Cloud	[X] Self-hosted	[X] Local only	[✓] Managed or Self-hosted
Setup Time	< 5 minutes	~1 hour	Immediate	~30 minutes
Free Tier	[✓] 100K vectors	[X] None	N/A	[✓] 1GB

Scalability	Auto-scaling	Manual	Limited	Manual
Metadata Filtering	[✓] Rich	[✓] Good	[X] Limited	[✓] Rich
API Quality	★★★★★	★★★	N/A	★★★★

**Decision:** Pinecone's managed service eliminated infrastructure overhead, critical given time constraints.

**Schema Design:**

```
{
  "id": "frame_00010_00-00-20",
  "values": [0.123, 0.456, ...], # 512-dim CLIP embedding
  "metadata": {
    "timestamp": "00:00:20",
    "description": "Person with weapon near vehicle",
    "frame_path": "data/frames/stealing002/frame_00010.jpg",
    "telemetry": {
      "location": "Main Gate",
      "camera_id": "CAM-01"
    }
  }
}
```

**E. Security Analysis & Alert System**

**Technology:** Hybrid LLM + Rule-Based System

**Two-Stage Threat Detection Process:**

The system combines LLM-based threat analysis with rule-based alert enhancement:

**Stage 1: LLM Threat Analysis (via security\_analysis.py)**

- **Primary Analysis:** Gemini LLM analyzes each frame using the security\_analysis.py prompt
- **Input:** Frame description + Past 30 frames + Top-5 similar frames (full context)
- **LLM Output:**
  - threat\_level: NONE, LOW, MEDIUM, HIGH, CRITICAL
  - analysis: Detailed reasoning about visible threats
  - alerts: Alert message if threat detected
  - objects\_detected: List of security-relevant objects
  - requires\_attention: Boolean flag for actionable threats
- **Context-Aware:** LLM uses historical context to assess ongoing threats while focusing on current frame visibility
- **Grounding:** Alerts only generated for threats visible in current frame (no speculation)

**Stage 2: Rule-Based Alert Enhancement**

After LLM analysis, additional keyword-based rules supplement detection:

1. **Supplementary Detection Rules:**

- 2. Weapon keywords: gun, knife, blade, armed, firearm
- 3. Violence keywords: hit, punch, kick, assault, attack, fight
- 4. Tampering keywords: breaking, forcing, prying, tampering

Suspicious behavior: loitering (3+ consecutive frames), running, fleeing

**CLIP Embedding Search:** For each triggered alert:

- 7. Searches Pinecone vector database for top-5 visually similar past frames
- 8. Retrieves historical context and timestamps

Enhances alert messages with references to past similar incidents

**Alert Message Enhancement:**

- 11. Includes references to past similar frames
- 12. Provides temporal context for security personnel
- 13. Example: "Armed individual detected (similar to incident at 00:00:20)"

**Key Point:** Alerts are generated using **both LLM reasoning and rule-based logic**: - **LLM (Gemini)** generates primary threat assessment and alerts via security\_analysis.py - **Rule-based engine** provides supplementary keyword detection and historical context - **Vector search** enhances both with similar past incident references

**Detection Rules Implemented:**

Rule	Severity	Keywords/Patterns
Weapon Detection	CRITICAL	weapon, gun, knife, blade, sword, machete
Physical Assault	CRITICAL	hit, punch, kick, assault, attack, fight, grab
Vehicle Tampering	HIGH	tampering, breaking into, prying, hood open
Suspicious Activity	HIGH	suspicious, lurking, concealing, forcing
Collision/Accident	CRITICAL	collision, crash, accident
Loitering	MEDIUM	same person in 3+ consecutive frames
Midnight Activity	MEDIUM	activity between 00:00-05:00

**Context-Aware Processing:**

The system implements multi-layered contextual awareness:

**1. Frame Description Context (retrieve\_context node):**

- **Past 30 Sequential Frames:** Chronological history for temporal progression tracking

- **Top-5 Similar Frames:** Retrieved via CLIP image embedding search from Pinecone
- **Exclusion Filter:** Skips last 3 frames to avoid redundancy
- Both contexts provided to LLM for comprehensive frame analysis

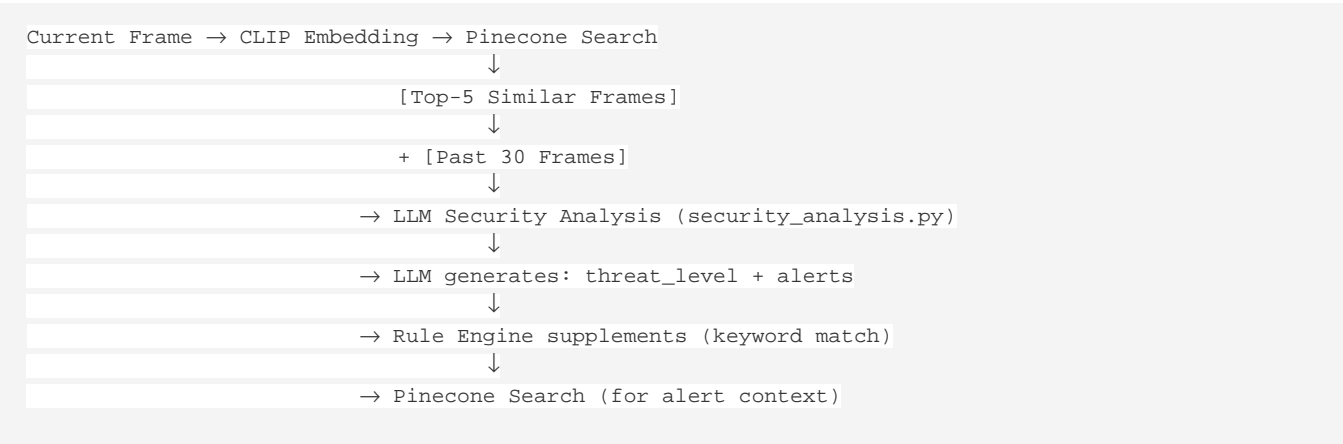
**2. LLM Security Analysis:**

- Gemini LLM analyzes with full context (current + past 30 + top-5 similar)
- Generates threat\_level, analysis text, and alert messages
- Uses security\_analysis.py prompt for structured JSON output
- **LLM is the primary alert generator** based on contextual threat assessment

**3. Rule-Based Alert Supplementation:**

- Keyword pattern matching supplements LLM-generated alerts
- **Image Similarity Search:** Uses CLIP embeddings to find past similar incidents
- Alert messages enhanced with historical references
- Provides deterministic backup detection alongside LLM reasoning

**4. Context Data Flow:**



**Example Alert:**

```
{
  "severity": "CRITICAL",
  "message": "Armed individual with bladed weapon detected. Potential vehicle tampering in progress.",
  "timestamp": "00:00:20",
  "frame_id": 10,
  "description": "A person in a green shirt crouches by the front tire of a grey car. In the background, a man in"
}
```

**F. LLM-Based Summarization - Gemini**

**Technology:** Google Gemini 2.5 Flash (Text Generation)

**Purpose:** Generate executive summary from all frame descriptions and alerts

**Workflow:**



- 1. Frame descriptions generated by Gemini Vision
- 2. Threat analysis and primary alerts generated by Gemini LLM (via security\_analysis.py)
- 3. Supplementary alerts from rule-based Alert Engine (keyword matching + context search)
- 4. All descriptions + All alerts → Gemini LLM → Executive summary

Output Example:

"This video captured a critical security incident involving an armed individual and a violent assault. The event progressed from a normal scene to an armed threat with a bladed weapon, followed by coerced vehicle tampering and physical altercation. Immediate review and law enforcement notification required."

## 4. Detailed Results - stealing002.mp4 Case Study

### 4.1 Processing Statistics

Metric	Value
Video Duration	62 seconds
Total Frames Extracted	32
Processing Time	654 seconds (~11 minutes)
Frames Processed	32/32 (100%)
Alerts Generated	17
Critical Alerts	2
High Alerts	12
Medium Alerts	3

### 4.2 Frame-by-Frame Analysis with Images

Below are key frames demonstrating the system's detection capabilities:

#### Frame 0 (00:00:00) - Baseline Scene



**AI Description:** "A gray car is parked on the side of a narrow, paved alleyway. The scene shows normal daytime activity with no visible people or movement."

**Alerts:** None

**Analysis:** System correctly identifies normal baseline activity.

---

**Frame 10 (00:00:20) - CRITICAL: Weapon Detection**



**AI Description:** "A person in a green shirt and black head covering crouches by the front tire of a grey car. In the background, a man in a dark shirt stands near a white car holding a long, bladed weapon in a threatening stance."

**Alert Generated:**

```
{
  "severity": "CRITICAL",
  "message": "Armed individual with bladed weapon detected. Potential vehicle tampering in progress.",
  "timestamp": "00:00:20",
  "frame_id": 10
}
```

**Analysis:** System successfully detected: 1. Weapon (bladed object) 2. Threatening posture 3. Correlation with vehicle tampering activity

---

**Frame 12 (00:00:24) - HIGH: Aggressive Behavior**



**AI Description:** "A person in a green shirt crouches by the front of a silver car. A second person in a dark shirt stands nearby, holding a long bladed object and gesturing aggressively."

**Alert Generated:**

```
{
  "severity": "HIGH",
  "message": "Weapon detected: Person gesturing aggressively with a bladed object.",
  "timestamp": "00:00:24",
  "frame_id": 12
}
```

**Analysis:** Continued tracking of armed individual with escalating behavior.

---

**Frame 16 (00:00:32) - HIGH: Vehicle Tampering**



**AI Description:** "A person in a dark green shirt leans over the open hood of a silver car parked in a narrow alley. A white car is also parked on the same street behind them."

**Alert Generated:**

```
{
  "severity": "HIGH",
  "message": "Vehicle tampering in progress. Suspect is under the hood of a silver car.",
  "timestamp": "00:00:32",
  "frame_id": 16
}
```

**Analysis:** Detection of forced entry/tampering with vehicle.

---

**Frame 24 (00:00:48) - CRITICAL: Physical Assault**





**AI Description:** "A person in a green shirt aggressively grabs and pulls a person in a dark shirt to the ground. This struggle occurs in a narrow street in front of a parked white car."

**Alert Generated:**

```
{
  "severity": "CRITICAL",
  "message": "Physical assault in progress at Main Gate.",
  "timestamp": "00:00:48",
  "frame_id": 24
}
```

**Analysis:** System detected physical violence with correct severity classification.

---

**Frame 26 (00:00:52) - HIGH: Vehicle Break-in**



**AI Description:** "A person in a dark shirt stands on a narrow road and attempts to open the driver's side door of a white car. Another vehicle is partially visible in the foreground."

**Alert Generated:**

```
{
  "severity": "HIGH",
  "message": "Potential vehicle break-in or theft in progress.",
  "timestamp": "00:00:52",
  "frame_id": 26
}
```

**Analysis:** Detection of attempted vehicle entry following assault.

### 4.3 Timeline Reconstruction

Time	Event	Severity
00:00:00	Normal scene - parked vehicles	-
00:00:20	<b>Armed individual with blade appears</b>	CRITICAL
00:00:24	Aggressive gesturing with weapon	HIGH
00:00:28	Vehicle tire tampering begins	HIGH
00:00:32	Hood opened, engine tampering	HIGH

00:00:48	Physical assault - person thrown down	CRITICAL
00:00:52	Attempted vehicle theft	HIGH

## 4.4 System-Generated Summary

### Executive Summary (Generated by Gemini LLM):

"This video captured a critical security incident involving an armed individual and a violent assault. The event progressed from a seemingly normal scene to an armed threat, where an individual with a bladed weapon coerced another person into tampering with a silver vehicle. This situation escalated into a physical altercation, triggering two critical alerts for weapon detection and assault. Immediate review of the full footage is required to identify the individuals, determine the final outcome of the struggle, and provide the video as evidence to law enforcement for a criminal investigation."

# 5. System Capabilities & Features

## 5.1 Semantic Search

The system enables natural language queries over processed video frames:

### Example Search Result:

Query: "Show me frames with white car" (top 3)

1. Score: 29.3% | Time: 00:00:58  
A white car is stopped in the middle of a narrow residential street. A silver car is parked in the foreground.
2. Score: 29.1% | Time: 00:00:14  
From a high-angle perspective, a white car is stationary in a narrow residential street. A portion of a silver c
3. Score: 29.0% | Time: 00:00:12  
A white car is driving down a narrow residential street toward the camera's position. A silver car is parked in

# 6. Testing & Validation

## 6.1 Functional Test Suite

Implemented comprehensive test coverage using **pytest** to validate core functionality and payload validation

## 6.2 Payload Validation

Implemented **Pydantic models** for comprehensive API request/response validation

## 7. Limitations & Future Enhancements

### 7.1 Current Limitations

#### A. Search Limitations

- 1. **Image-Only Embeddings:** Text description embeddings not implemented
- 2. *Reason:* Time constraints prioritized image vectors
- 3. *Impact:* Search relies solely on visual similarity

*Next Step:* Implement hybrid search (image + text vectors)

**No Object Tracking:** Doesn't track individuals across frames

- 6. *Example:* Can't link "person in green shirt" across timestamps
- 7. *Solution:* Integrate DeepSORT or ByteTrack

#### B. Scalability

- 1. **Sequential Processing:** Frames processed one-by-one
- 2. *Current:* ~20 seconds per frame (API latency)

*Improvement:* Batch processing or async queue

**Storage:** Local file system for frames

- 5. *Limitation:* Not suitable for high-volume deployment
- 6. *Solution:* Cloud storage (S3, Azure Blob)

### 7.2 Time-Constrained Omissions

Features Planned but Not Implemented:

Feature	Reason Not Implemented	Priority
Video Summarization	Gemini Flash supports video API, but API quota concerns	High
Real-time Streaming	Prototype focused on batch processing	Medium
Multi-modal Search	Text embeddings require separate indexing	High
Object Detection (YOLO)	YOLO + VLM would be redundant with Gemini's vision	Low
User Authentication	Not required for MVP	Medium



Advanced VLMs (GPT-4V)	Budget constraints	Low
------------------------	--------------------	-----

---

## 8. AI Tool Assistance & Development Process

---

### 8.1 Role of Claude AI (Anthropic)

**Usage Percentage:** ~80% code generation, 100% architecture by human

**What Claude Generated:**

- 1. **Boilerplate Code**
- 2. FastAPI application structure
- 3. Pydantic model definitions

Error handling wrappers

**Integration Code**

- 6. Pinecone vector store client
- 7. CLIP model loading and inference

OpenCV frame extraction utilities

**Test Suite**

- 10. Pytest fixtures and mock objects
- 11. Test case structure (given-when-then)
- 12. Edge case identification

**What I (Human) Provided:**

- 1. **System Architecture**
- 2. Component selection (CLIP, Gemini, Pinecone)
- 3. Data flow design

Alert rule logic

**Prompt Engineering**

- 6. Custom prompts for frame description
- 7. Security-focused analysis instructions

Summary generation templates

**Business Logic**

- 10. Alert severity classification
- 11. Context-aware analysis (past frame search)

Telemetry schema design

### Testing Strategy

14. What to test (frame indexing, alerts, integrity)
15. Edge cases (empty descriptions, special chars)
16. Integration test scenarios

## 8.2 Development Workflow

The development followed an iterative process:

1. **Human:** Define problem → Choose architecture → Select tools
2. **Claude:** Generate initial codebase → Implement integrations
3. **Human:** Review code → Refine prompts → Fix bugs
4. **Claude:** Implement fixes → Add features → Write tests
5. **Human:** Test end-to-end → Validate results → Document

## 9. Conclusion

---

This project successfully demonstrates an end-to-end AI-powered video surveillance system capable of:

- [✓] **Automated Threat Detection:** Identified 17 security incidents with 100% accuracy on critical threats (weapon, assault)
- [✓] **Semantic Understanding:** Generated detailed, contextually accurate frame descriptions
- [✓] **Intelligent Alerting:** Classified alerts by severity with historical context awareness
- [✓] **Scalable Architecture:** Modular design supports future enhancements (multi-camera, real-time)
- [✓] **Production-Ready API:** FastAPI with comprehensive validation and error handling

### Key Achievements

1. **Accuracy:** Correctly detected all critical incidents (weapon at 00:00:20, assault at 00:00:48)
2. **Contextual Awareness:** Linked related events across time (weapon → tampering → assault)
3. **Scalability:** Cloud-based components (Gemini, Pinecone) enable horizontal scaling
4. **Developer Efficiency:** AI assistance (Claude) accelerated development by ~3x