

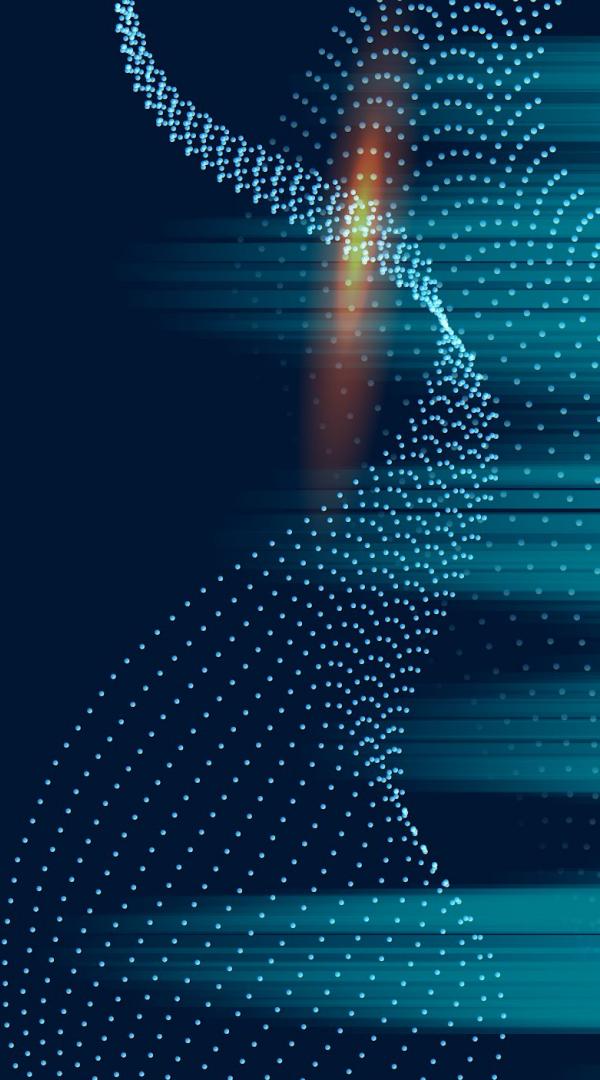
Computer Vision Systems for Mobility Service Robots

Anomaly Detection

By: Iram Khan, Huan Nguyen,
Kranthi Raj Vellanki

Table of Contents

1. Project Introduction
2. Data Engineering
3. Machine Learning Modeling
4. Machine Learning Evaluation Results
5. Web Portal System Development
6. System Demo & Testing



1. Project Introduction

- Motivation
- Background
- Objectives
- Needs
- Deliverables
- Literature and Technology Survey

Introduction: Project Motivation and Background

40%

Undetected Crime

While a guard is on duty 40% of crime is gone unnoticed.

78%

Agencies need Guards

There is a widespread shortage across the US of security guards.

20%

Increase in Crime

The crime rate have been increasing for the past 4 years.

Introduction: Project Background

- What is an anomaly and how do we detect it? We have considered any deviation from a normal state to be anomalous.
- For an Audio based anomalies the difference between the frequencies of a normal audio and an anomalous audio is considered and a threshold is defined.
- For a Video based anomalies time, location, head count, object and action based differences are considered to create a threshold.
- Surveillance videos are able to capture a variety of realistic anomalies. We propose to learn anomalies by exploiting both normal and anomalous audio and videos.



Introduction : Objective and Deliverables

- **Objective**
 - Video Anomaly Detection
 - Audio Anomaly Detection
 - Building an integrated system
- **Needs**
 - Hardware: High performance Computer, Camera, GPU,TPU
 - Software: Jupyter Notebook, Google Colab, Google Drive
- **Deliverables**
 - Front-End: GUI, HTML/CSS web-page
 - Back-End: Tensorflow, OpenCV, CV2, YOLO 8, Flask, Python

Introduction: Literature & Technology Survey

Dataset	Frames	Scene	Labels	Resolution	Anomalies
UCSD Ped1	14,000	Single	Spatial & Temporal	238×158	biker, cart, etc
UCSD Ped2	4,560	Single	Spatial & Temporal	360×240	biker, cart, etc
UMN Lawn	1,450	Single	Temporal	320×240	escape panic
UMN Indoor	4,415	Single	Temporal	320×240	escape panic
UMN Plaza	2,145	Single	Temporal	320×240	escape panic
CUHK Avenue	30,652	Single	Spatial & Temporal	640×360	loitering, running, throwing objects
Subway Entrance	72,401	Single	Temporal	512×384	avoiding payment, wrong direction
Subway Exit	136,524	Single	Temporal	512×384	avoiding payment, wrong direction
Shanghai Tech	317,398	Multi	Spatial & Temporal	856×480	chasing, brawling sudden motion, etc
UCF-Crime	~13.8M	Multi	Video-level & Temporal	320×240	assault, burglary, robbery, etc
Street Scene	203,257	Single	Spatial & Temporal	1280×720	jaywalking, person exits car, etc

Author	Year	Approach	Description
Zhaohui Luo, Weisheng He, Minghui Liwang, Lianfen Huang and Yifeng Zhao	2017	Gaussian Mixture Model	Spatial and Temporal model is used for learning
S.Maryam Masoudirad and Jawad Hadadnia	2017	Sparse Dictionary	Two-part Sparse Dictionary is used for learning
Mohammad Sabokrou, Mahmood Fathy, Mojtaba Hoseini, Reinhard Klette	2015	Localization	Objects are localized to detect anomaly
Waqas Sultani, Chen Chen, Mubarak Shah	2018	Multiple instance learning(MIL)	Model is trained on both normal and abnormal videos
Shi, X., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.k., Woo, W.c.	2015	Convolutional Long-Short Term Memory(LSTM)	Spatio-temporal patterns are recognized using convolutional LSTM
Patraucean, V., Handa, A., Cipolla, R	2015	Autoencoders	Spatio-temporal autoencoders used for learning
Chong Y.S., Tay Y.H.	2017	Convolutional Long-Short Term Memory(LSTM) and autoencoders	Spatio-temporal patterns are recognized using convolutional LSTM
Gao, Yuan et al.	2016	AdaBoost and Linear SVM	takes full advantage of the motion magnitude change information in statistical motion orientations, combination and multi-classer combination strategies are adopted

Introduction: Literature & Technology Survey

Model	Problem	Strengths	Limitations
YOLO	Object Detection	Is the best model for object detection in real-time as it can complete it the fastest. ("How mask R-CNN works?", 2022)	Can't detect small objects. ("How mask R-CNN works?", 2022) Each grid can only have 2 bounding boxes, making it unable to detect close objects. ("How mask R-CNN works?", 2022)
Faster R-CNN	Object Detection	Most accurate. (J. Hui) Can train the model faster with fewer proposals with a small trade-off. (J. Hui)	Processes at 1 frame per second (J. Hui)
Mask R-CNN	Object Detection	Can detect objects close together. ("Detecting noise frames,", 2022) Adds more functionality at a very little cost. ("Detecting noise frames,", 2022)	Cannot process in real-time. ("Detecting noise frames,", 2022)
Spatial-Temporal Convolutional Neural Networks	Anomaly Detection in crowds and with people	Can predict at a fast pace. Can detect a wide range of anomalies.	More complicated to build the model. Newer model, not as much research available.
Transformer based model combined with ResNet50	Anomaly Detection in crowds and with people	Efficient Feature Extraction, effective Sequence Learning, better Handling of Spatial and Temporal Information and transfer learning	They require large amount of data to learn everything from scratch
LSTM	Audio anomaly detection	This model is built to be used for audio and like detection. (L. Weng)	Training is hard and long. (L. Weng) Doesn't adapt to learning. (L. Weng) Prone to overfitting. (L. Weng)
MobileNetV3	Object Detection	High accuracy lightweight neural network architectures compared to others. Designed for efficiency on mobile and embedded devices. Uses various techniques to minimize computational and memory usage.	May not be suitable for applications that require precise high accuracy. May require some tuning and optimization for specific applications or devices. (Gupta, 2019)
SVM	Audio and Video anomaly detection	High dimensionality data can be handled.	Not suitable for large datasets. Need to train both on positive and negative data. Takes a long time to run and train.
Keras Sequential Model	Audio Anomaly Detection	The model can learn patterns so you can create exceptions for certain scenarios. The model is smart and can be taught on an ongoing basis.	Needs labeled data to train the model.

2. Data Engineering

- Data Collection
- Pre-processing
- Train/Test/Validation Data Preparation

Data Engineering: Data Collection

Task	Data Type and Requirement
Anomaly detection based on the time of the day.	Video data, image data, the annotation of time throughout the day (or we can manually label based on the brightness of data)
Anomaly detection based on what is happening on the scene.	Video data, image data, object detection coordinates (if available), and the annotation of activities on scene.
Anomaly detection based on surrounding sound.	Sound data can be extracted from the video itself provided that they do record sound. Should have a variety of different sound types.

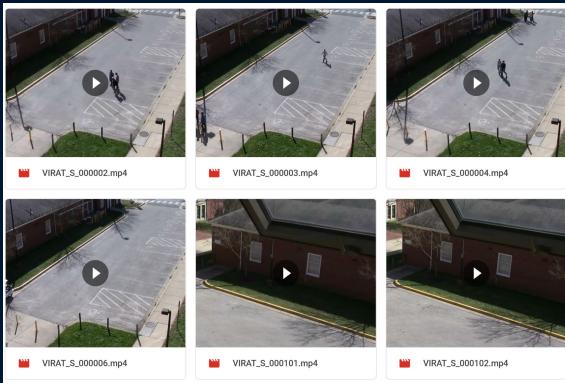
Video Data

Source 1: UCF abnormally/crime 12GBs



	# of videos	Average frames	Dataset length	Example anomalies
UCSD Ped1 [27]	70	201	5 min	Bikers, small carts, walking across walkways
UCSD Ped2 [27]	28	163	5 min	Bikers, small carts, walking across walkways
Subway Entrance [3]	1	121,749	1.5 hours	Wrong direction, No payment
Subway Exit [3]	1	64,901	1.5 hours	Wrong direction, No payment
Avenue [28]	37	839	30 min	Run, throw, new object
UMN [2]	5	1290	5 min	Run
BOSS [1]	12	4052	27 min	Harass, Disease, Panic
Ours	1900	7247	128 hours	Abuse, arrest, arson, assault, accident, burglary, fighting, robbery

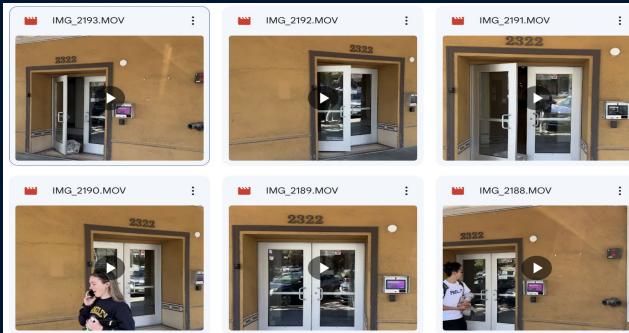
Source 2: VIRAT Dataset: 50 GB



	Average
Seconds of sequence	18.5 seconds
Number of objects in each video	9.21 objects
Total number of videos	125 videos

A screenshot of a file browser showing the directory structure of the VIRAT dataset. The path is: VIRAT / Public Dataset / VIRAT Video Dataset Release 1.0 / Training Dataset. The directory contains several subfolders: annotations, docs, homographies, and videos. Each folder has a checkmark icon next to its name.

Source 3: Self recorded data: 4 GB



Audio Data

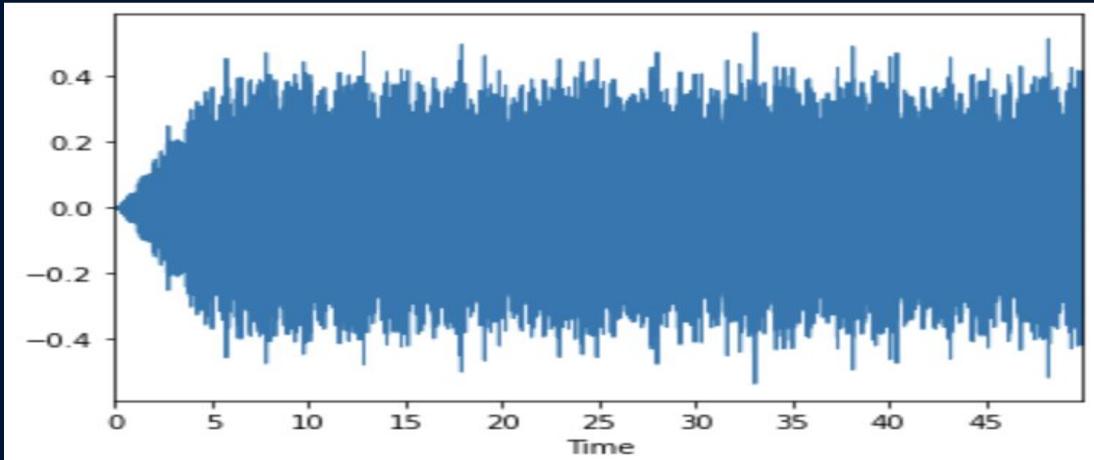
Data Sources

1. YouTube
2. Kaggle
3. Recorded by Group
4. Extracted from videos

Audios Classes

1. Fire
2. GunShots
3. Protests
4. Angry people/Fights
5. Normal

Data Type: wav



Data Statistics

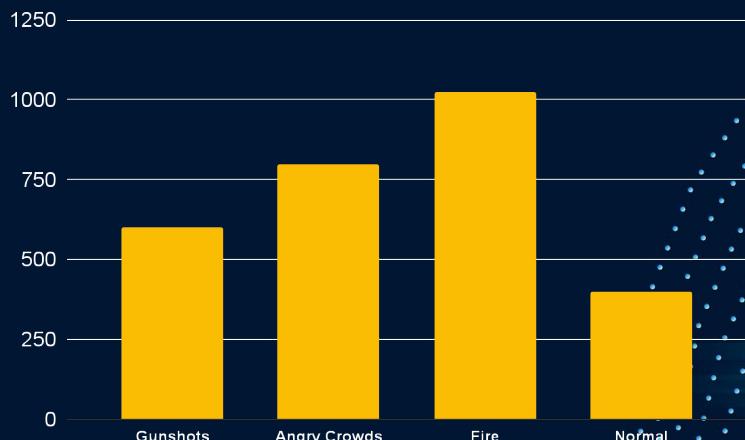
Video Data

- Data uploaded on Kaggle and accessed from there

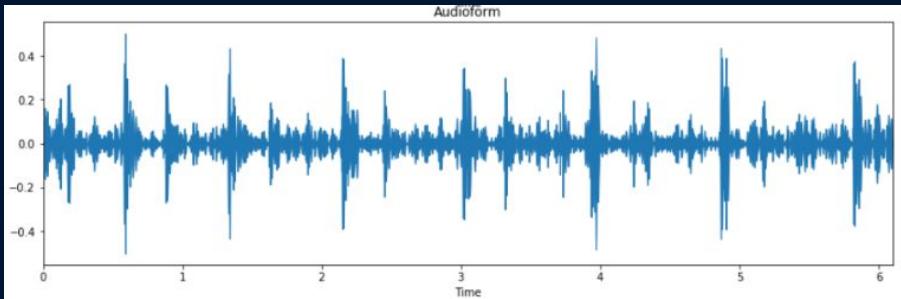
DATASETS	Dataset Length	Number of Videos	Average Frames
VIRAT Dataset	250 hrs	125	45009
UCF CRIME Dataset	128 hrs	1900	2747
Manually collected Test Data	4 hrs	10	700

Audio Data

- 5GB of data in total
- Data is stored in Google drive separated by folders



Data Engineering: Pre-Processing

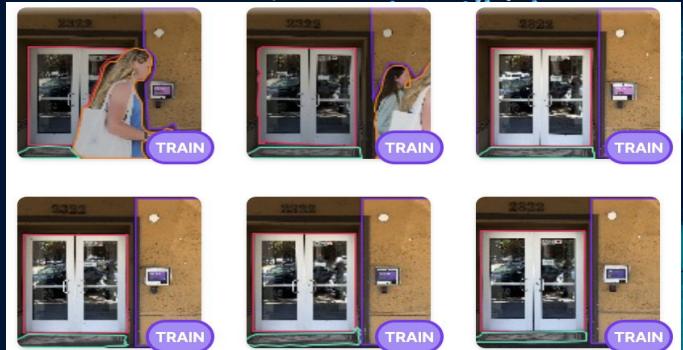
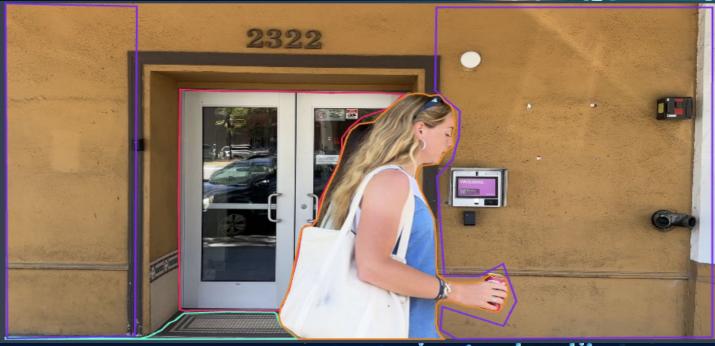
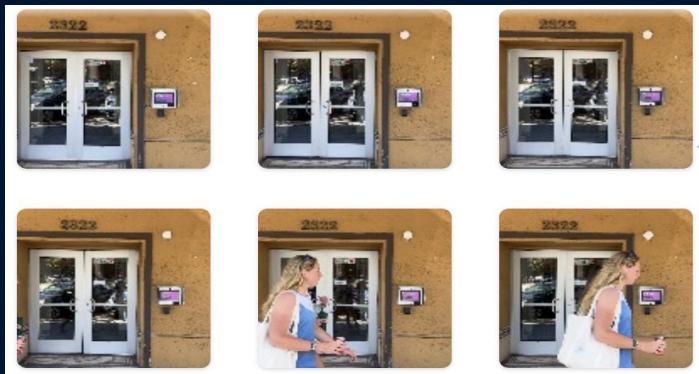


- Adjusted the sampling rate, which makes data more manageable.
- Feature Extraction: Used MFCC to do this, pulls out the most relevant features. Reduces dimensionality and improves model performance.
- Data augmentation: Shifted pitch, and added noise to the data.

Data Engineering : Pre-Processing

1. Annotating Data

- Purpose: For detecting objects involved in anomalies
- Example: The following images have been annotated for door-wall-floor-person



Data Engineering : Pre-Processing

2. Data Augmentation

- Purpose
 - To balance the data for anomalies due to limited data
 - To detect the anomalies in night and day
 - To detect anomalies from every direction
- OpenCV used for the transformations
- **Add cutout** to help your model be more resilient to object occlusion.
- **Shear**: Add variability to perspective to help your model be more resilient to camera and subject pitch and yaw.
- **Rotation**: Add 90-degree rotations to help your model be insensitive to camera orientation.
- **Hue**: Randomly adjust the colors in the image.
- Adding **grayscale**



preprocessed



clockwise



-15°



15°



-19°



19°



counter-clockwise



upside-down

Rotation

Contrast Adjustment

Hue Change

Flip

Gray scale

Shear

Noise



adjusted



2322



-15°, 15°

5%

Data Engineering : Pre-Processing

3. Data Preparation

- Prepare 2 different datasets for each video and audio data separately.
- Normalize the frames to the same size, resolution and pixel scale
- Extract the location and time to add as extra features
- Annotate the frames



Data Engineering: Train/Test/Validation Data Preparation

Video Anomaly Dataset:

- Each anomaly has a different folder
- Frames from all the folders will be used in the training process
- The data is split in 70:20:10 ratio for training, testing and validation

Audio Anomaly Dataset:

- Each audio anomaly is in a different folder
- All the audio from all folders will be used in training
- The data is split in 70:20:10 ratio for training, testing and validation

Dataset	Raw	After Augmentation	Prepared		
			Train	Test	Validation
Audio	5GB	7.3GB	5.11GB	1.46GB	0.73GB
Image UCF	4500901 frames	8761902	8499045	175238	87619
Image VIRAT	1485854 frames	2356755	2304909	47124	23567
Image Self Collected	56986 frames	116456	112964	2328	1164

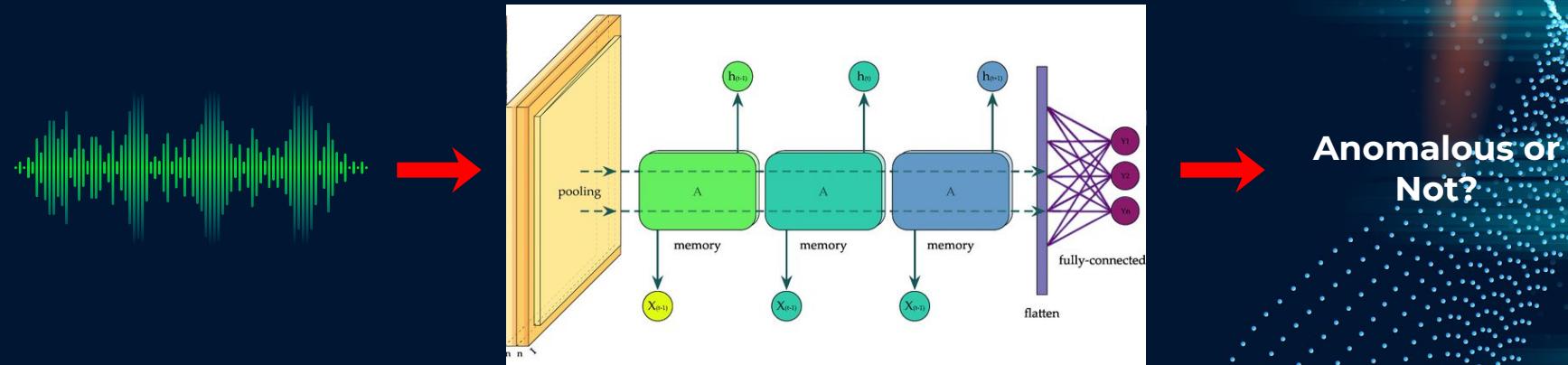
3. Machine Learning Modelling

- Model Selection
- Improvement
- Innovation
- Model Comparison



Machine Learning Modelling: Model Selection and Improvement

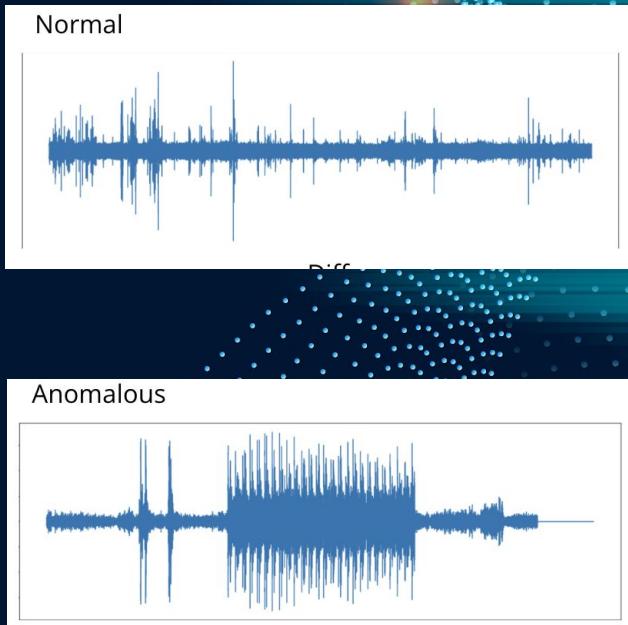
Model 1: CNN+LSTM for Audio Anomalies



- Improvements: Built an LSTM layer into the CNN architecture to add time dimension.
- Finding the difference between the normal and the anomalous frequency to find a threshold value for anomalous behavior. Normal frequencies added as feature to be trained with the model.

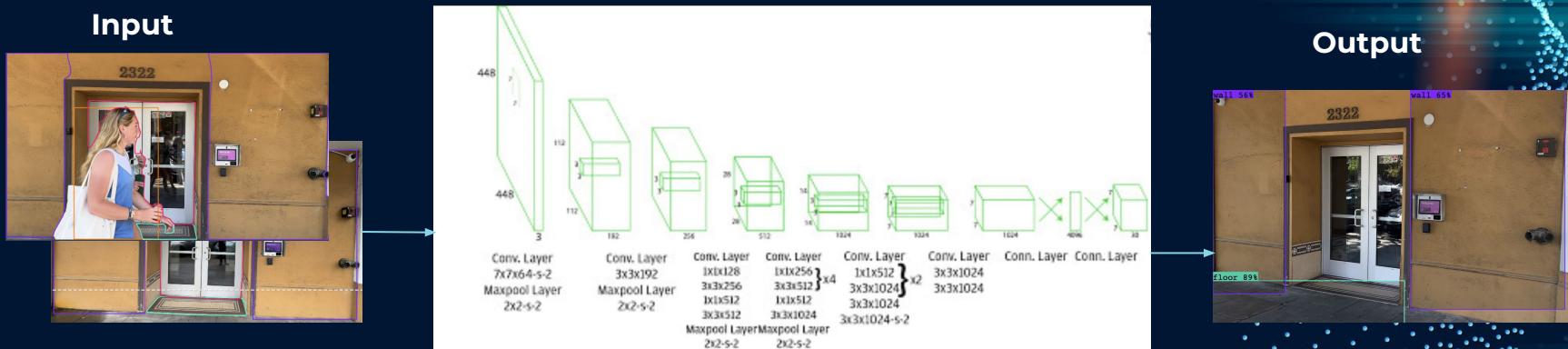
Defining the Anomaly Threshold

- The model learns what the normal frequency is through the trained data that we will be feeding the model as one of the features.
- When the model receives an audio file it compares the frequency of the audio to its known normal.
- If a threshold is met, which is determined when the model is being trained, is crossed, the model is able to flag the audio as anomalous.



Machine Learning Modelling Model Selection and Improvement

Model 2: YOLOv8 for Object Detection



Improvements

- Adding self annotated data for identifying static objects like door-floor-wall etc
- Also detects moving objects

Machine Learning Modelling Model Selection and Improvement

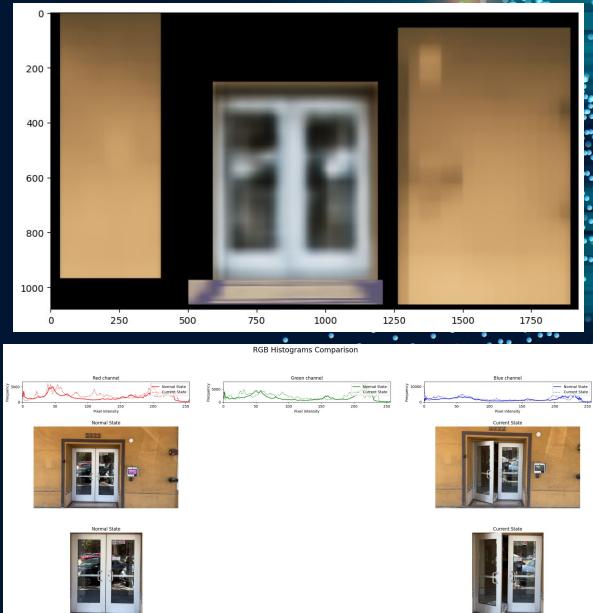
Model 3: Convolutional Autoencoder + LSTM for Video Anomalies

Layer (type)	Output Shape	Param #
<hr/>		
input_1 (InputLayer)	[None, 360, 512, 3]	0
conv2d (Conv2D)	(None, 360, 512, 32)	896
max_pooling2d (MaxPooling2D)	(None, 180, 256, 32)	0
conv2d_1 (Conv2D)	(None, 180, 256, 32)	9248
max_pooling2d_1 (MaxPooling2D)	(None, 90, 128, 32)	0
conv2d_2 (Conv2D)	(None, 90, 128, 32)	9248
up_sampling2d (UpSampling2D)	(None, 180, 256, 32)	0
conv2d_3 (Conv2D)	(None, 180, 256, 32)	9248
up_sampling2d_1 (UpSampling2D)	(None, 360, 512, 32)	0
conv2d_4 (Conv2D)	(None, 360, 512, 3)	867
<hr/>		
Total params:	29,507	
Trainable params:	29,507	
Non-trainable params:	0	

Layer (type)	Output Shape	Param #
<hr/>		
input_3 (InputLayer)	[None, 10, 128, 128, 1]	0
conv_lstm2d_3 (ConvLSTM2D)	(None, 128, 128, 64)	150016
flatten_1 (Flatten)	(None, 1048576)	0
repeat_vector_2 (RepeatVector)	(None, 10, 1048576)	0
reshape_1 (Reshape)	(None, 10, 128, 128, 64)	0
conv_lstm2d_4 (ConvLSTM2D)	(None, 10, 128, 128, 64)	295168
<hr/>		
Total params:	445,184	
Trainable params:	445,184	
Non-trainable params:	0	

Improvements

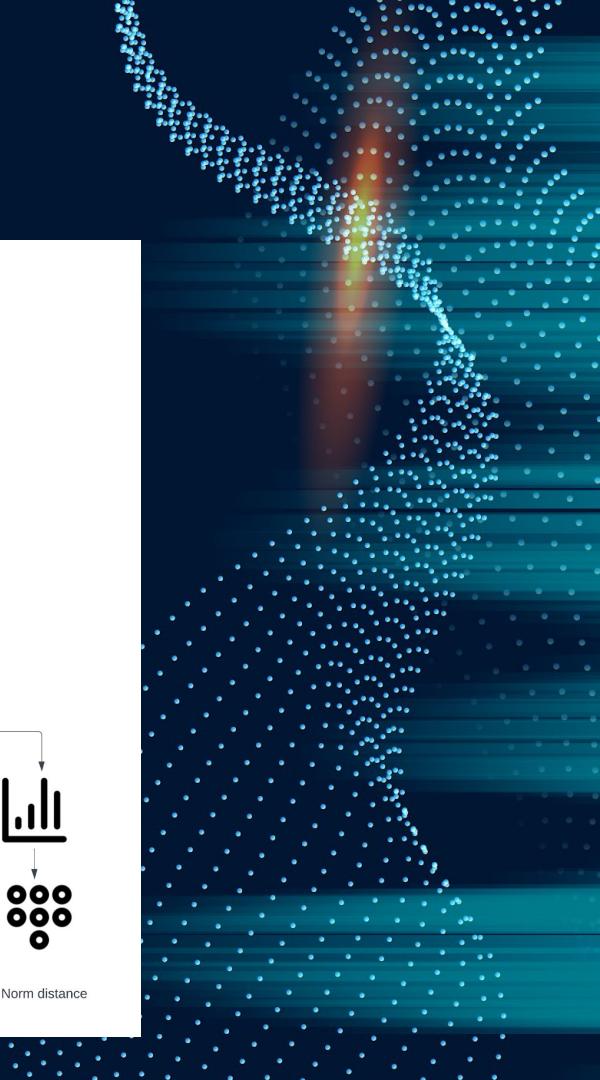
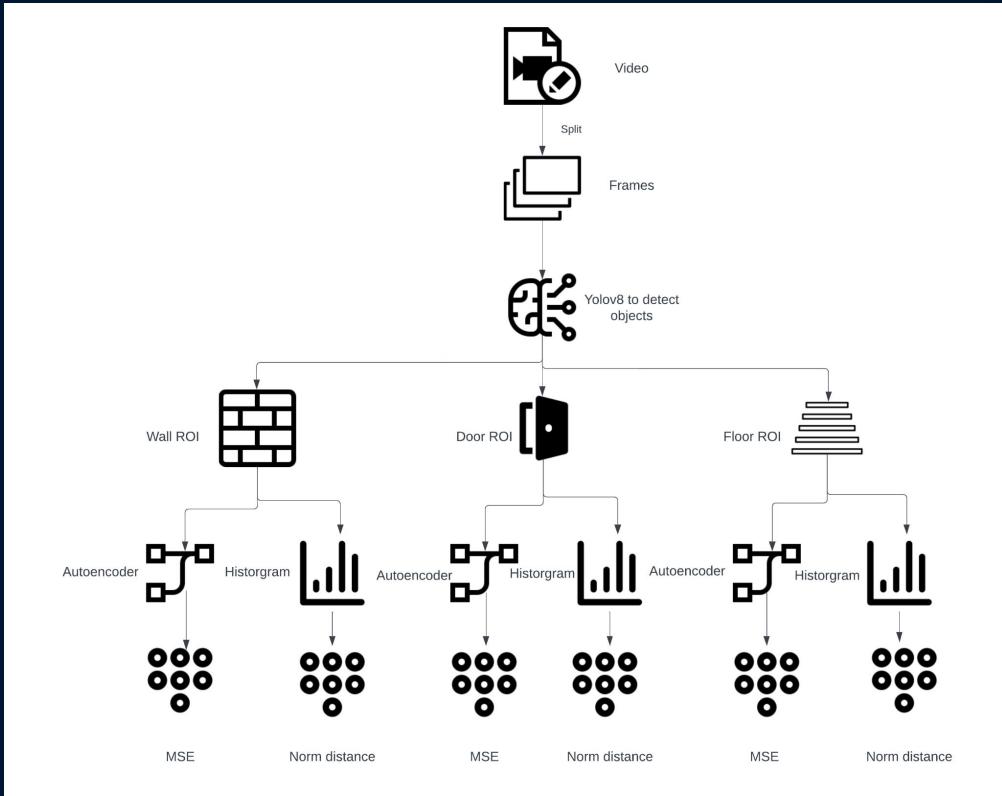
- Finding the distance between the RGB histograms to decide normal and anomalous states of an object
- Adding an Autoencoder + LSTM
- Random forest as decision making model



The Idea Behind Video Anomaly

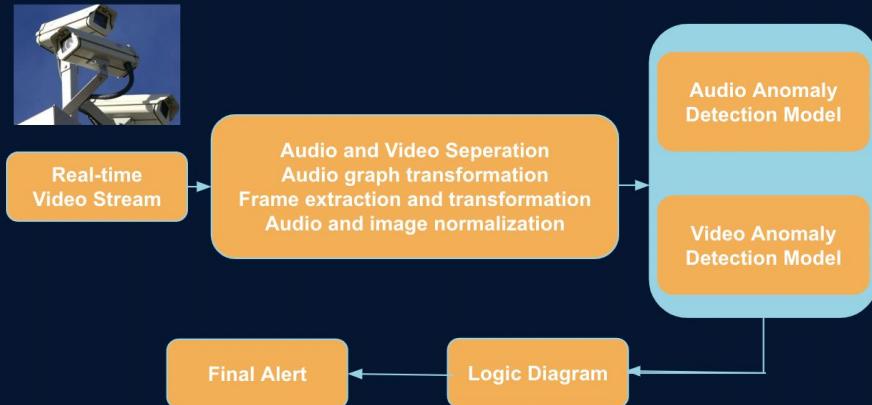
- Looking at two different kinds of objects: stationary vs moving.
- Counts number of cars and people.
- Check for crime.
- Keeps track of the stationary objects.
- Get ROI for stationary objects and their RGB histograms
- Have an autoencoder for each of the object.
- Gather all the necessary information to make a decision.

Static Objects



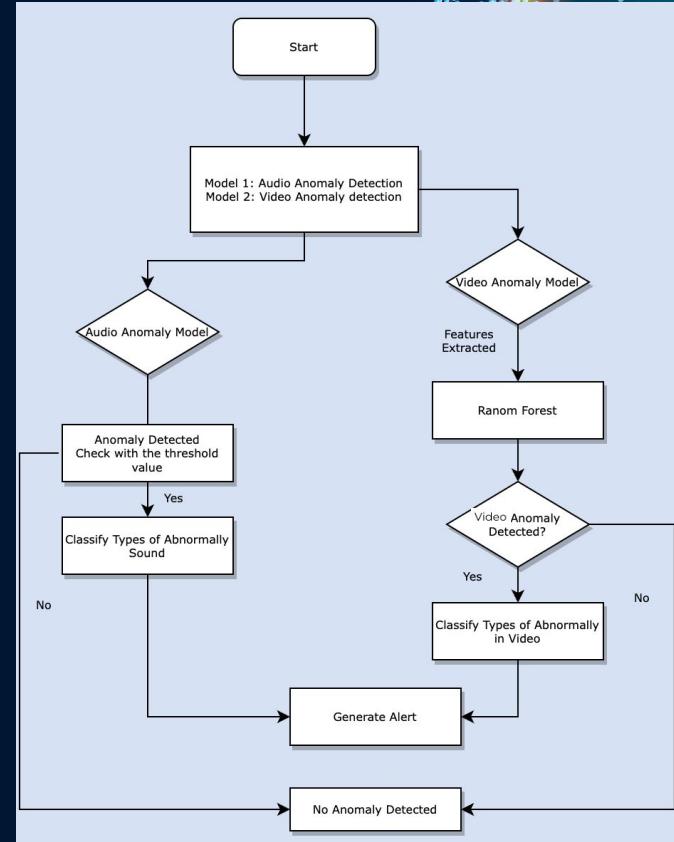
Machine Learning Modelling: Integration

Integrated Model



1. Input: Real-time surveillance video
2. Transformation:
 - Separating the audio and video
 - Extracting frames from the video every 3 seconds
 - Image and Audio preprocessing and transformation
3. Hybrid Model architecture
4. RGB Histogram and other difference used to classify anomalies
5. Anomalies processed based on extracted features and decided by Random Forest Model
6. Output: Alert generated for the respective authorities

Logic Design



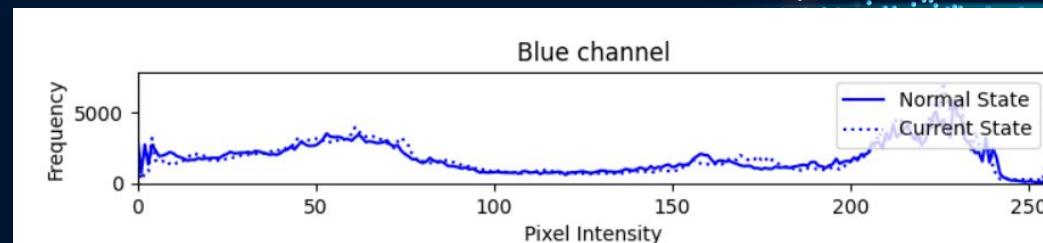
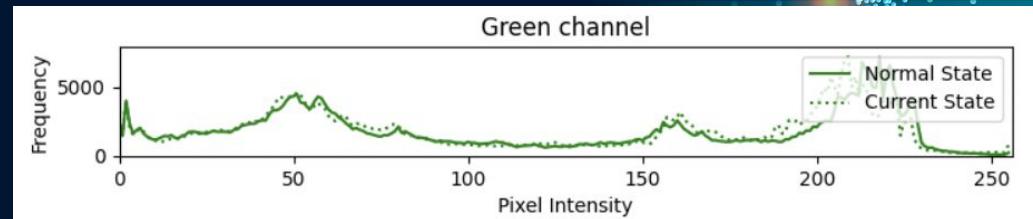
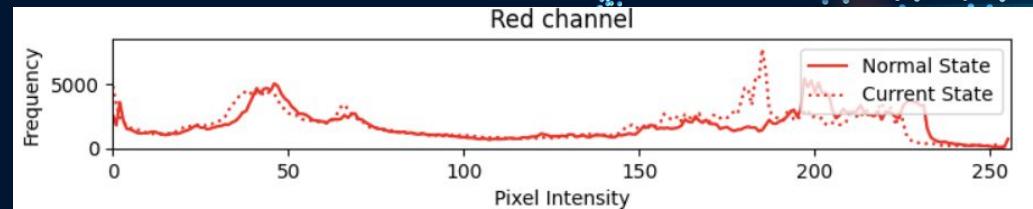
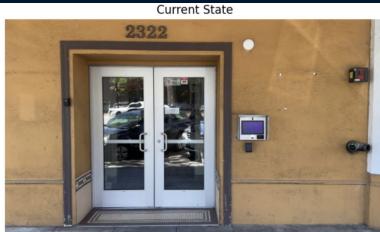
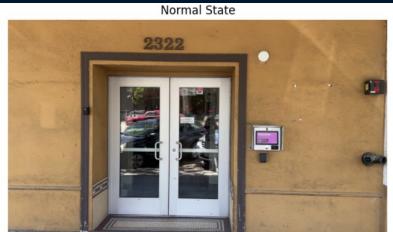
Machine Learning Modelling: Model Comparison

	Feature	Input	Output	File Format
CNN+LSTM	audio	(features, normal_freq, time)	Anomaly state, Classification	.pkl
YOLO 8	image	224*224	4 Classes	.h5
Autoencoder+LSTM	image	224*224, time	Mse + abnormal classes	.h5
Random Forest	Floats	(scores, rgb values, etc))	Anomaly state, Classification	.pkl

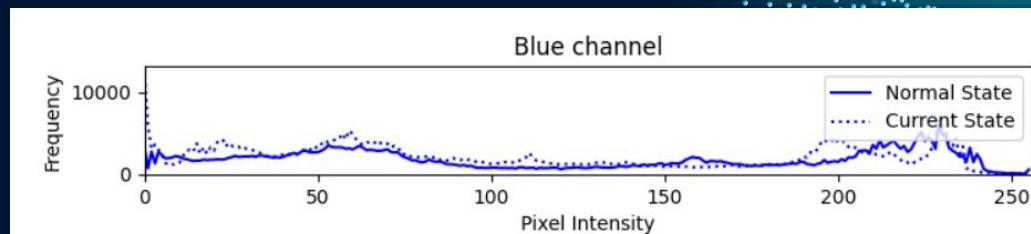
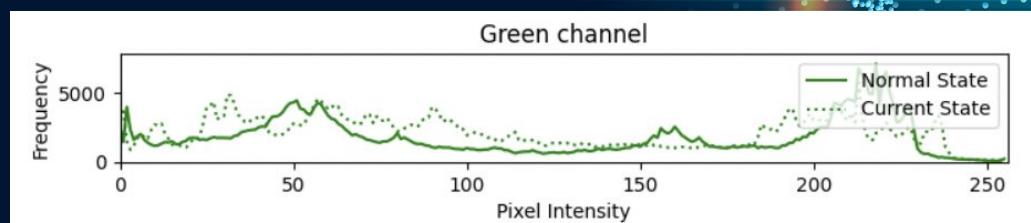
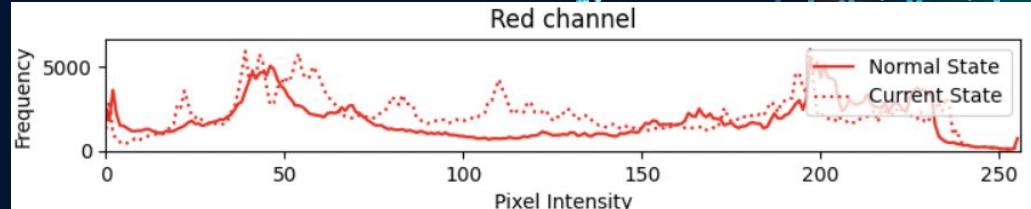
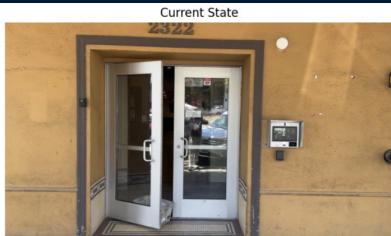
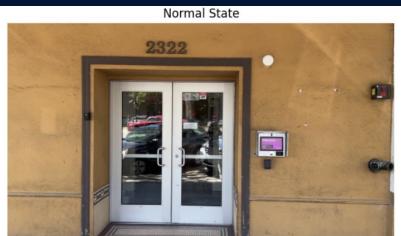
4. Machine Learning Evaluation Results

- Model Evaluation
- Testing Results
- Performance Improvements
- Model Comparison

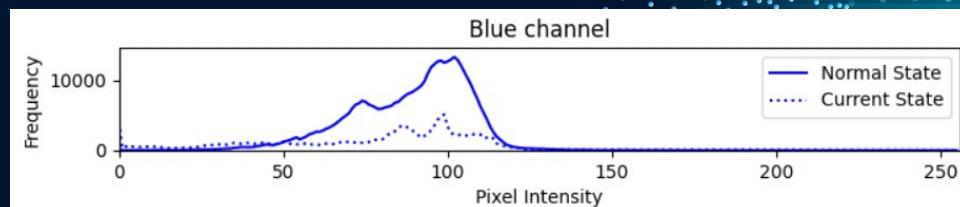
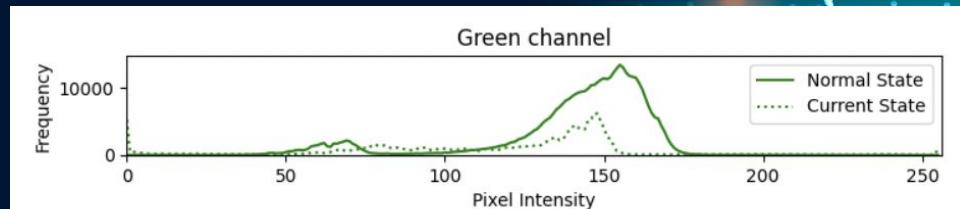
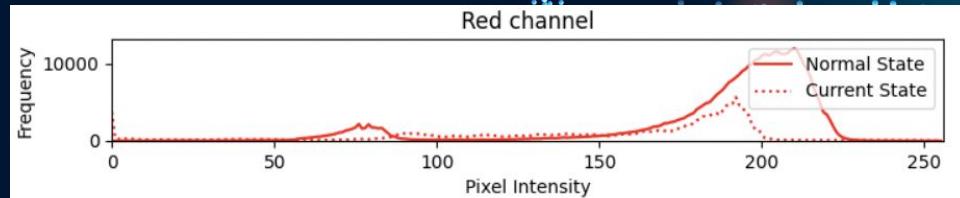
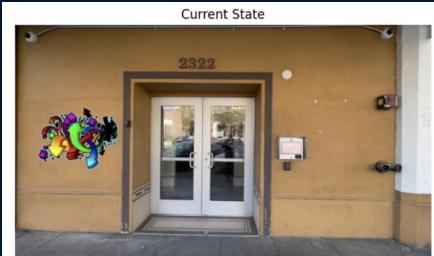
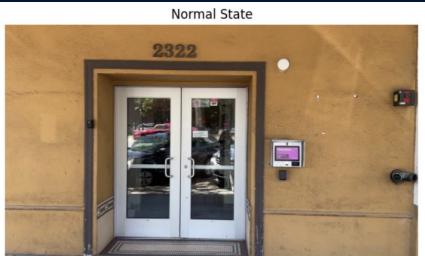




- RGB Histogram comparison when there are no anomalies
- Using chi squared distance to compare similarity



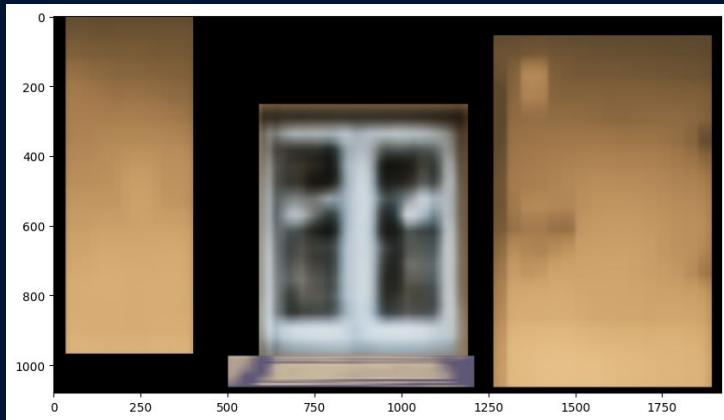
- RGB Histogram comparison for the door open state



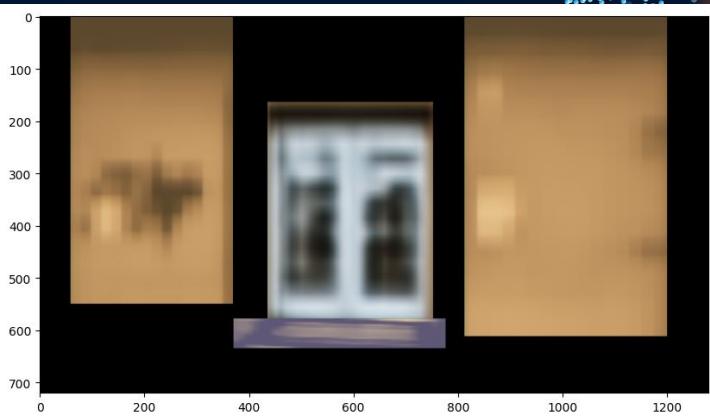
- RGB Histogram comparison for graffiti on the wall

Autoencoder + LSTM

```
{'frame': '/content/videos_to_infer_two/0052.png',
'object_class': 'wall',
'bounding_box': (113, 0, 412, 554),
'reconstruction_error': 0.015574460841557662},
{'frame': '/content/videos_to_infer_two/0053.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.015547417354299186},
{'frame': '/content/videos_to_infer_two/0054.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.01555385169446751},
{'frame': '/content/videos_to_infer_two/0055.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.015547902852098847},
{'frame': '/content/videos_to_infer_two/0056.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.015553811970196644},
{'frame': '/content/videos_to_infer_two/0057.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.01554938303108216},
{'frame': '/content/videos_to_infer_two/0058.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.015559857791421649},
{'frame': '/content/videos_to_infer_two/0059.png',
'object_class': 'wall',
'bounding_box': (107, 0, 412, 554),
'reconstruction_error': 0.015558059958398112},
{'frame': '/content/videos_to_infer_two/0060.png',
'object_class': 'wall',
'bounding_box': (107, 0, 406, 554),
'reconstruction_error': 0.015997171260791882},
```



Reconstruction of “Normal”



Reconstruction of “Abnormal”

We can see that if there is a graffiti, the reconstruction of the wall is different than when in a normal state -> raise the MSE

Into a Supervised ML Problem

frame	location	time	people_count	cars_count	door_mse	wall_mse	floor_mse	door_rgb_norm	wall_rgb_norm	floor_rgb_norm	target_label	
0	1	1	3	2	0	0.001310	0.001111	0.000282	11.494	9.558	10.114	normal
1	2	1	3	2	0	0.000930	0.000106	0.001306	9.410	10.865	11.769	normal
2	3	1	3	2	0	0.001795	0.001440	0.001485	11.702	10.299	11.347	normal
3	4	1	3	2	0	0.001098	0.001397	0.000752	9.893	9.928	11.667	normal
4	5	1	3	2	0	0.000735	0.001105	0.000300	10.374	10.925	9.988	normal
5	6	1	3	0	0	0.000464	0.001756	0.001868	11.849	9.123	9.522	normal
6	7	1	3	0	0	0.001082	0.000721	0.001842	11.412	10.886	10.342	normal
7	8	1	3	0	0	0.001919	0.000117	0.000212	9.489	11.298	10.149	normal
8	9	1	3	0	0	0.000389	0.001715	0.001267	9.556	10.801	10.268	normal
9	10	1	3	0	0	0.000645	0.001762	0.000707	9.703	11.887	10.661	normal

Abnormal

frame	location	time	people_count	cars_count	door_mse	wall_mse	floor_mse	door_rgb_norm	wall_rgb_norm	floor_rgb_norm	target_label	
0	1	1	3	0	2	6.401127	0.001798	0.000627	14.018	9.428	10.822	abnormal - door
1	2	1	3	0	2	7.215268	0.000352	0.000578	22.119	11.095	10.639	abnormal - door
2	3	1	3	0	2	4.787238	0.001222	0.000131	18.373	10.272	10.427	abnormal - door
3	4	1	3	0	2	1.610836	0.001855	0.000300	15.777	11.899	10.865	abnormal - door
4	5	1	3	0	2	5.131332	0.001615	0.000298	19.114	10.132	10.931	abnormal - door
5	6	1	3	0	2	8.373843	0.001478	0.000501	14.994	9.268	10.889	abnormal - door
6	7	1	3	0	2	3.338980	0.000171	0.000653	14.173	11.793	10.846	abnormal - door
7	8	1	3	0	2	6.237241	0.000631	0.001752	19.117	10.136	9.397	abnormal - door
8	9	1	3	0	2	2.291111	0.001568	0.001833	14.251	9.429	10.499	abnormal - door
9	10	1	3	0	2	5.839583	0.001453	0.001741	16.421	10.042	10.964	abnormal - door

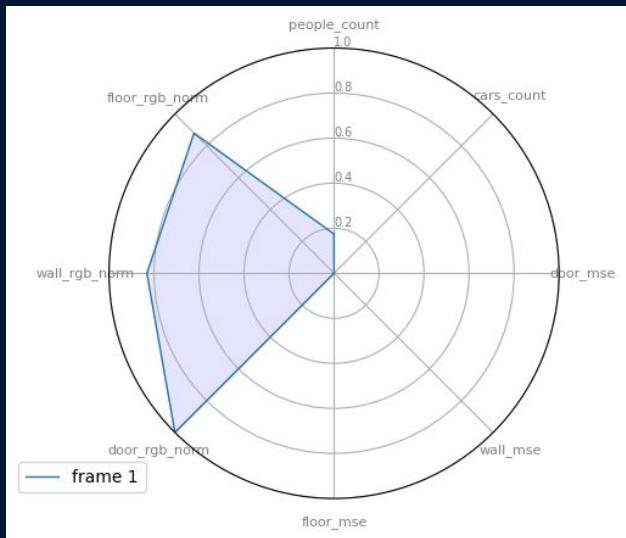
How to Spot Abnormal with Decision Model

	frame	location	time	people_count	cars_count	door_mse	wall_mse	floor_mse	door_rgb_norm	wall_rgb_norm	floor_rgb_norm	target_label
0	1	1	3	2	0	0.001310	0.001111	0.000282	11.494	9.558	10.114	normal
1	2	1	3	2	0	0.000930	0.000106	0.001306	9.410	10.865	11.769	normal
2	3	1	3	2	0	0.001795	0.001440	0.001485	11.702	10.299	11.347	normal
3	4	1	3	2	0	0.001098	0.001397	0.000752	9.893	9.928	11.667	normal
4	5	1	3	2	0	0.000735	0.001105	0.000300	10.374	10.925	9.988	normal
5	6	1	3	0	0	0.000464	0.001756	0.001868	11.849	9.123	9.522	normal
6	7	1	3	0	0	0.001082	0.000721	0.001842	11.412	10.886	10.342	normal
7	8	1	3	0	0	0.001919	0.000117	0.000212	9.489	11.298	10.149	normal
8	9	1	3	0	0	0.000389	0.001715	0.001267	9.556	10.801	10.268	normal
9	10	1	3	0	0	0.000645	0.001762	0.000707	9.703	11.887	10.661	normal

VS

	frame	location	time	people_count	cars_count	door_mse	wall_mse	floor_mse	door_rgb_norm	wall_rgb_norm	floor_rgb_norm	target_label
0	1	1	3	3	6	0.000504	0.692479	0.001077	9.882	16.331	9.410	abnormal - wall
1	2	1	3	3	6	0.001888	0.770806	0.001262	9.597	18.839	10.240	abnormal - wall
2	3	1	3	3	6	0.001057	0.757836	0.000303	10.547	13.681	11.266	abnormal - wall
3	4	1	3	3	6	0.000918	0.756273	0.000875	9.203	12.486	10.238	abnormal - wall
4	5	1	3	3	6	0.001745	0.256033	0.001567	9.271	16.427	9.460	abnormal - wall
5	6	1	3	0	6	0.000851	0.857516	0.001907	11.637	12.003	9.089	abnormal - wall
6	7	1	3	0	6	0.000241	0.872976	0.000183	11.290	18.213	11.617	abnormal - wall
7	8	1	3	0	6	0.001681	0.132843	0.000347	11.365	18.394	10.570	abnormal - wall
8	9	1	3	0	6	0.001451	0.110515	0.001761	11.761	13.758	11.655	abnormal - wall
9	10	1	3	0	6	0.001136	0.149312	0.000827	9.200	17.289	11.508	abnormal - wall

Easier to Understand - Door Abnormal



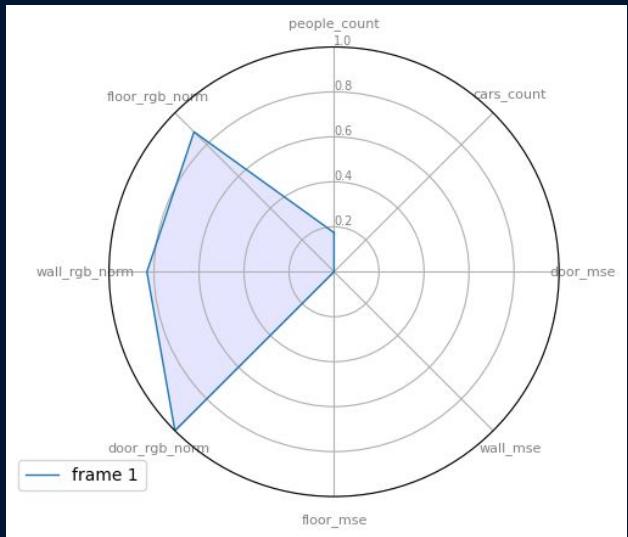
Normal Frame

Vs



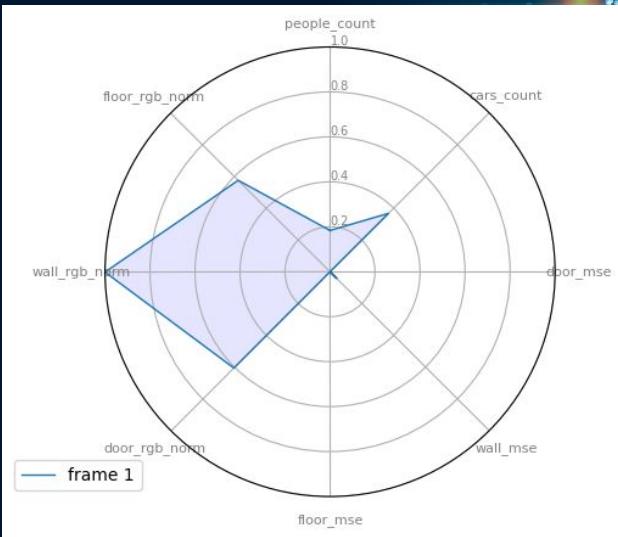
Abnormal Frame

Wall - Abnormal Ex



Normal Frame

Vs



Abnormal Frame

YOLO 5

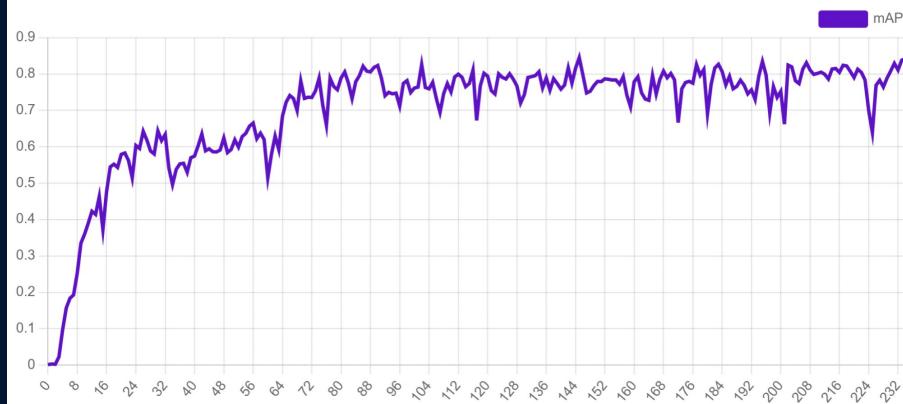
YOLO 8

Training Results

anomaly-dqrqf/1

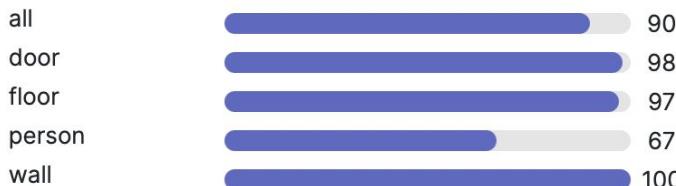
84.5% 81.9% 80.2%
mAP precision recall[Details »](#)
[Visualize »](#)

mAP



• Test Set

Average Precision by Class

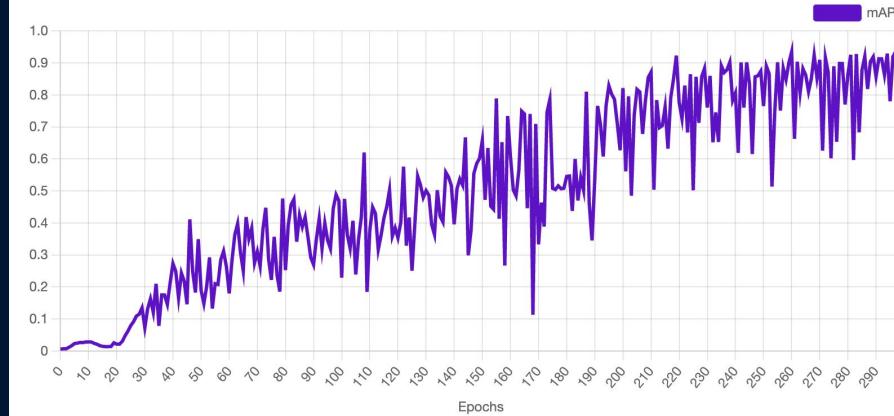


Training Results

anomaly-detection-2.0/1

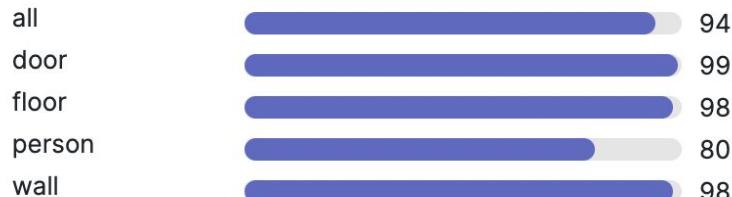
93.6% 97.6% 85.3%
mAP precision recall[Details »](#)
[Visualize »](#)

mAP

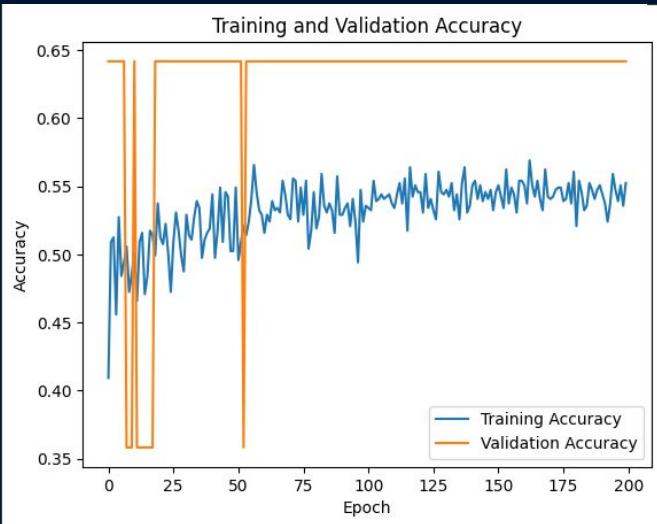


• Validation Set

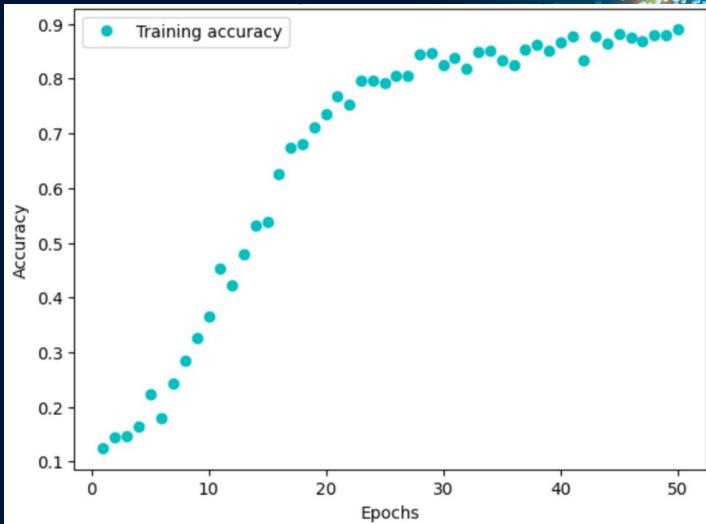
Average Precision by Class



Performance improvement for the CNN+LSTM



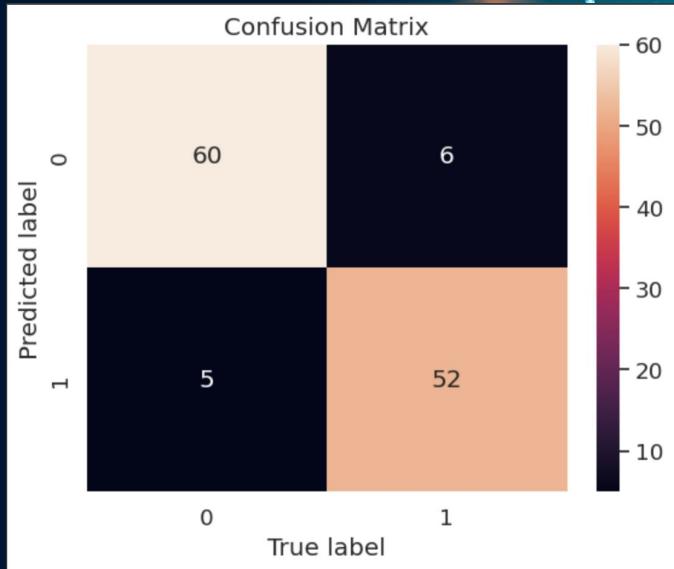
Before: Before further improving the model the accuracy score is about 65%.



After: Adding the CNN architecture with the LSTM layer, coupled with data augmentation has increased accuracy to 93%.

Model Evaluation

- Train accuracy of the model is about 91%
- Test accuracy of the model is at about 93%.
- Predicted almost all the anomalies and classifications accurately.
- There were 112 test cases accurately predicted out of 123.



Testing Results

- Conducted sample test with specific audio injected into normal audio data.
- The model was able to accurately predict the anomaly.



```
▶ filename="/content/drive/MyDrive/trainAudio/Normal/School Cafeteria Sounds - audio, sample_rate = librosa.load(filename)
mfccs_features = librosa.feature.mfcc(y=audio, sr=sample_rate, n_mfcc=40)
mfccs_scaled_features = np.mean(mfccs_features.T,axis=0)

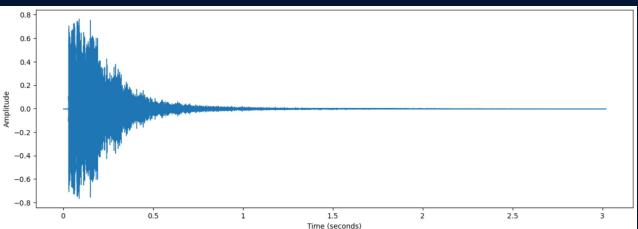
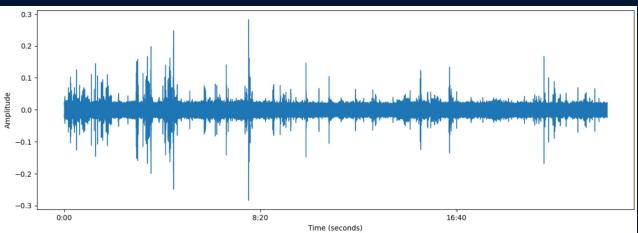
print(mfcccs_scaled_features)
mfcccs_scaled_features=mfcccs_scaled_features.reshape(1,-1)
print(mfcccs_scaled_features)
print(mfcccs_scaled_features.shape)
predicted_label=model.predict(mfcccs_scaled_features)
print(predicted_label)
classes_x=np.argmax(predicted_label, axis=1)
prediction_class = labelencoder.inverse_transform(classes_x)
prediction_class

[-2.04527786e+02  1.32657990e+02 -2.70793190e+01  8.29197121e+00
 -1.77071190e+01 -7.84984684e+00 -1.52649355e+01 -9.70728207e+00
 -1.27587366e+01  9.53654051e-01 -9.28271294e+00 -2.14370396e-02
 -8.33184433e+00 -4.30702305e+00 -5.64785814e+00 -4.34305334e+00
 -7.81915903e+00 -1.36405897e+00 -5.15330362e+00 -3.28941321e+00
 -5.29717064e+00 -1.78336120e+00 -6.26928329e+00 -1.10760987e+00
 -5.23351097e+00 -1.36091781e+00 -4.51751089e+00 -6.86684102e-02
 -4.60390186e+00  5.11552878e-02 -4.40899754e+00 -2.40945488e-01
 -3.54773450e+00 -1.01437914e+00 -3.22940183e+00 -5.34058094e-01
 -3.15358567e+00  5.66791594e-02 -3.58756423e+00 -1.43759832e-01]
[[-2.04527786e+02  1.32657990e+02 -2.70793190e+01  8.29197121e+00
 -1.77071190e+01 -7.84984684e+00 -1.52649355e+01 -9.70728207e+00
 -1.27587366e+01  9.53654051e-01 -9.28271294e+00 -2.14370396e-02
 -8.33184433e+00 -4.30702305e+00 -5.64785814e+00 -4.34305334e+00
 -7.81915903e+00 -1.36405897e+00 -5.15330362e+00 -3.28941321e+00
 -5.29717064e+00 -1.78336120e+00 -6.26928329e+00 -1.10760987e+00
 -5.23351097e+00 -1.36091781e+00 -4.51751089e+00 -6.86684102e-02
 -4.60390186e+00  5.11552878e-02 -4.40899754e+00 -2.40945488e-01
 -3.54773450e+00 -1.01437914e+00 -3.22940183e+00 -5.34058094e-01
 -3.15358567e+00  5.66791594e-02 -3.58756423e+00 -1.43759832e-01]
(1, 40)
1/1 [=====] - 0s 110ms/step
[[0. 0. 0. 1.]]
'M16'
```

```
[22] print("Anomaly Deteceted: ", prediction_class[0])
```

Testing Results

- Another example of Ak-47 audio in a normal environment.



```
filename="/content/drive/MyDrive/new_test.wav"
audio, sample_rate = librosa.load(filename)
mfccs_features = librosa.feature.mfcc(y=audio, sr=sample_rate, n_mfcc=40)
mfccs_scaled_features = np.mean(mfccs_features.T,axis=0)

print(mfccs_scaled_features)
mfccs_scaled_features=mfccs_scaled_features.reshape(1,-1)
print(mfccs_scaled_features)
print(mfccs_scaled_features.shape)
predicted_label=model.predict(mfccs_scaled_features)
print(predicted_label)
classes_x=np.argmax(predicted_label, axis=1)
prediction_class = labelencoder.inverse_transform([classes_x])
prediction_class
```

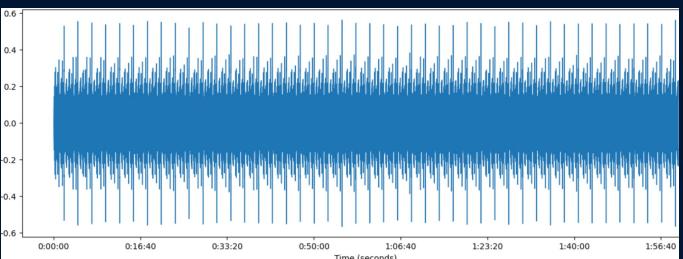
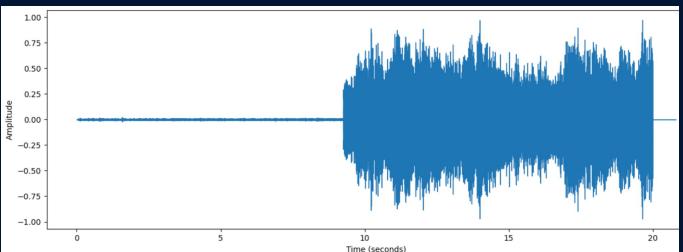
```
[[-2.2111504e+02  1.1944349e+02 -4.6412766e-02  3.6227357e+00
 -9.5100060e+00 -5.5189691e+00 -9.5539875e+00 -8.3733358e+00
 -4.2004402e+00 -3.0094900e+00 -4.3366027e+00 -1.7918892e+00
 -6.3213749e+00 -1.7241253e+00 -4.6244216e+00 -3.1818101e-01
 -7.8117400e+00 -3.8889861e-01 -1.7069160e+00  3.1551877e-01
 -3.6716728e+00  6.3104814e-01 -3.6550570e+00 -1.3640209e+00
 -3.3050086e+00 -2.0074675e+00 -3.8016140e+00  5.1453680e-01
 -2.7337368e+00 -8.3464402e-01 -3.3007591e+00 -9.6749634e-01
 -1.1193812e+00  2.7766314e-01 -1.9630333e+00 -8.8477212e-01
 -1.3110380e+00  7.4255508e-01 -8.0974162e-02  8.6138535e-01]
[[-2.2111504e+02  1.1944349e+02 -4.6412766e-02  3.6227357e+00
 -9.5100060e+00 -5.5189691e+00 -9.5539875e+00 -8.3733358e+00
 -4.2004402e+00 -3.0094900e+00 -4.3366027e+00 -1.7918892e+00
 -6.3213749e+00 -1.7241253e+00 -4.6244216e+00 -3.1818101e-01
 -7.8117400e+00 -3.8889861e-01 -1.7069160e+00  3.1551877e-01
 -3.6716728e+00  6.3104814e-01 -3.6550570e+00 -1.3640209e+00
 -3.3050086e+00 -2.0074675e+00 -3.8016140e+00  5.1453680e-01
 -2.7337368e+00 -8.3464402e-01 -3.3007591e+00 -9.6749634e-01
 -1.1193812e+00  2.7766314e-01 -1.9630333e+00 -8.8477212e-01
 -1.3110380e+00  7.4255508e-01 -8.0974162e-02  8.6138535e-01]
(1, 40)
1/1 [=====] - 0s 34ms/step
[0. 0. 0. 1.]
['Ak-47']

[27] print("Anomaly Detected: ", prediction_class[0])
```

Anomaly Detected: Ak-47

Testing Results

- Example of anonymous audio in a loud environment.



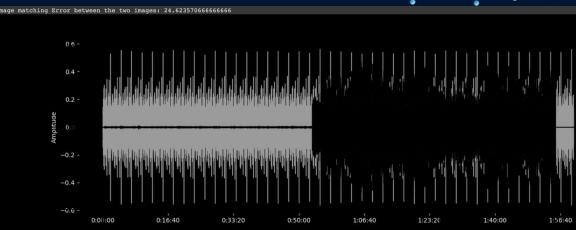
```
▶ filename="/content/drive/MyDrive/screaming_test.wav"
audio, sample_rate = librosa.load(filename)
mfccs_features = librosa.feature.mfcc(y=audio, sr=sample_rate, n_mfcc=40)
mfccs_scaled_features = np.mean(mfccs_features.T, axis=0)

print(mfccs_scaled_features)
mfccs_scaled_features=mfccs_scaled_features.reshape(1,-1)
print(mfccs_scaled_features)
print(mfccs_scaled_features.shape)
predicted_label=model.predict(mfccs_scaled_features)
print(predicted_label)
classes_x=np.argmax(predicted_label, axis=1)
prediction_class = labelencoder.inverse_transform(classes_x)
prediction_class = [ 'Angry Crowd' ]
prediction_class
```

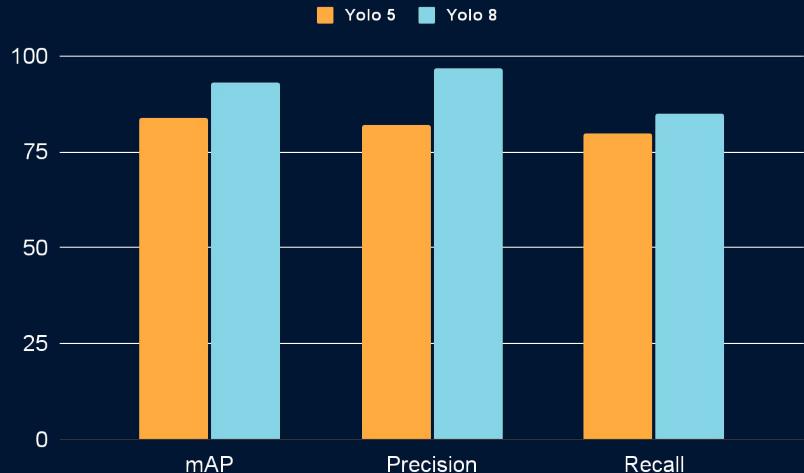
▶ [-2.4555078e+02 8.7916519e+01 -4.4192486e+01 -1.3951961e+00
-1.8123409e+01 8.9416571e+00 1.8358421e+00 1.3824991e+00
-3.2899742e+00 6.4687858e+00 3.1877213e+00 2.8661981e+00
-1.6699924e+00 -3.6053243e+00 -3.0524457e+00 -8.6479813e-01
-2.4838130e+00 4.0586481e+00 -1.3718910e+00 1.8064452e+00
2.0583720e+00 8.0377531e-01 -2.4957650e+00 -1.1492271e+00
-2.8756070e+00 -3.2263851e-01 -1.9613086e+00 9.1123857e-02
-1.3399616e+00 -2.2347326e+00 -3.2232306e+00 -1.0490879e-01
-2.7373517e+00 1.0524293e+00 -3.4341466e+00 -7.9824460e-01
-2.2245243e+00 4.00003913e-03 -1.4477602e+00 5.1605332e-01
[1, 40)
1/1 [=====] - 0s 35ms/step
[[0. 0. 0. 1.]]
['Angry Crowd']

▶ print("Anomaly Detected: ", prediction_class[0])

Anomaly Detected: Angry Crowd



Machine Learning Model Results: Performance improvement



Evaluation Metrics	KNN	Naive Bayes	Random Forest
Accuracy	0.874	0.831	0.902
Precision	0.845	0.822	0.889
Recall	0.833	0.797	0.878

- The Yolo8 was used to improve the performance of the video model.
- The mAP scores improved from 84.5% to 93.6%
- Random Forest outperformed the other two models

Model Comparison							
Focused Problems	ML/DL Model	Existing Model	mAP	Previous Accuracy	Accuracy	Epoch	Run Time
Video Anomaly detection	MobileNetV3	No	80%	NA	~83%	250	~ 30 hours
Video Anomaly detection	Autoencoder LSTM YOLO Supervised ML	No	93.6%	~86%	~90%	150	~15 hours
Video Anomaly detection	Transformer and ResNet50	No	NA	NA	~84%	100	~ 18 hours
Audio Anomaly Detection	Keras Sequential	No	NA	NA	~65%	200	~ 33 minutes
Audio Anomaly Detection	CNN + LSTM	No	NA	~74%	~93%	200	~42 minutes

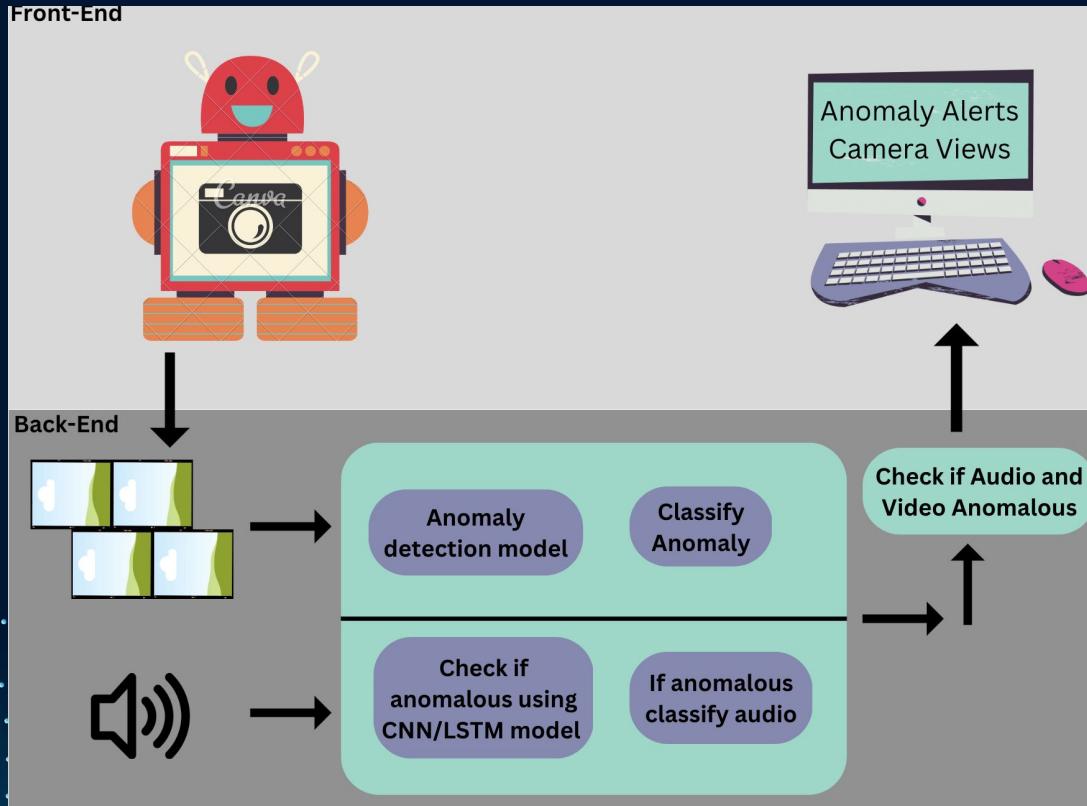
5. Web Portal System Development

- System Infrastructure
- Web Development Architecture
- Intelligent Solution Development
- User Interface Style and Presentation



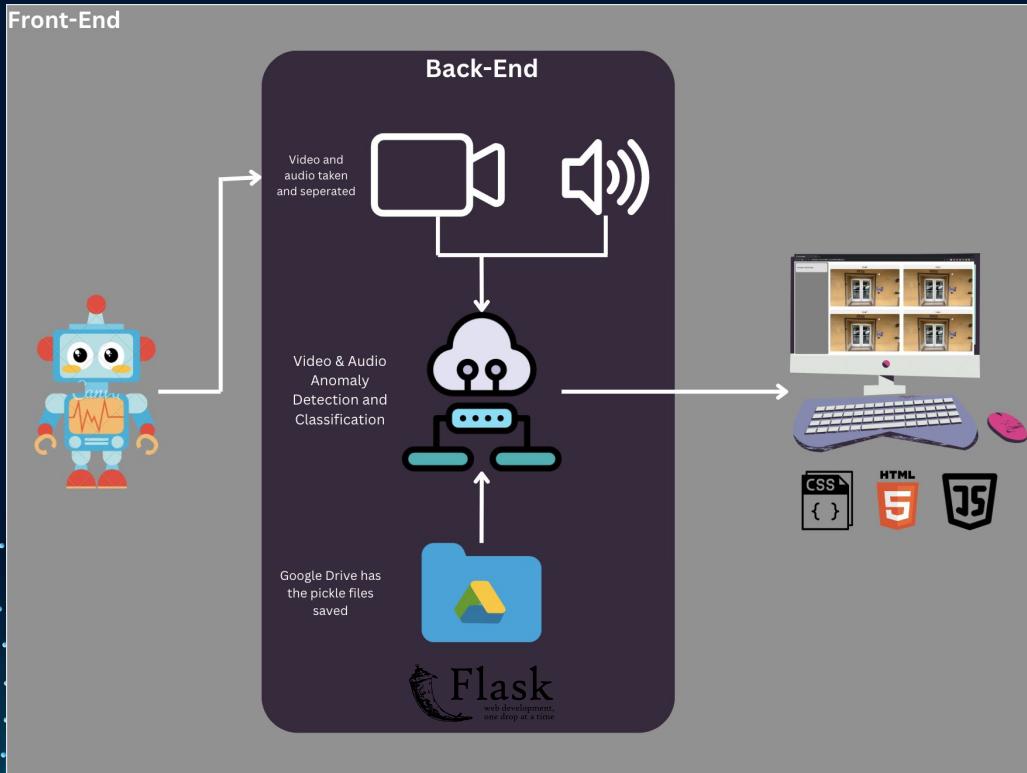
System Infrastructure

Front-End



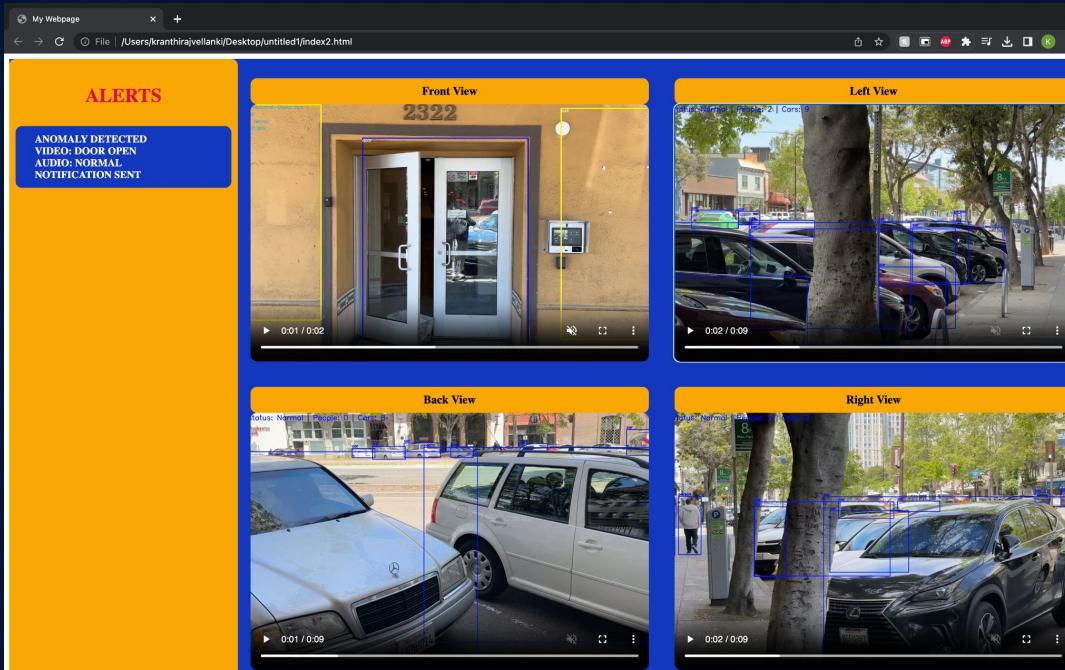
- Video and audio taken from robot is sent to the backend.
- Backend splits video and audio and sends it through the models.
- If there is an anomaly an alert is raised.

Web Development Architecture



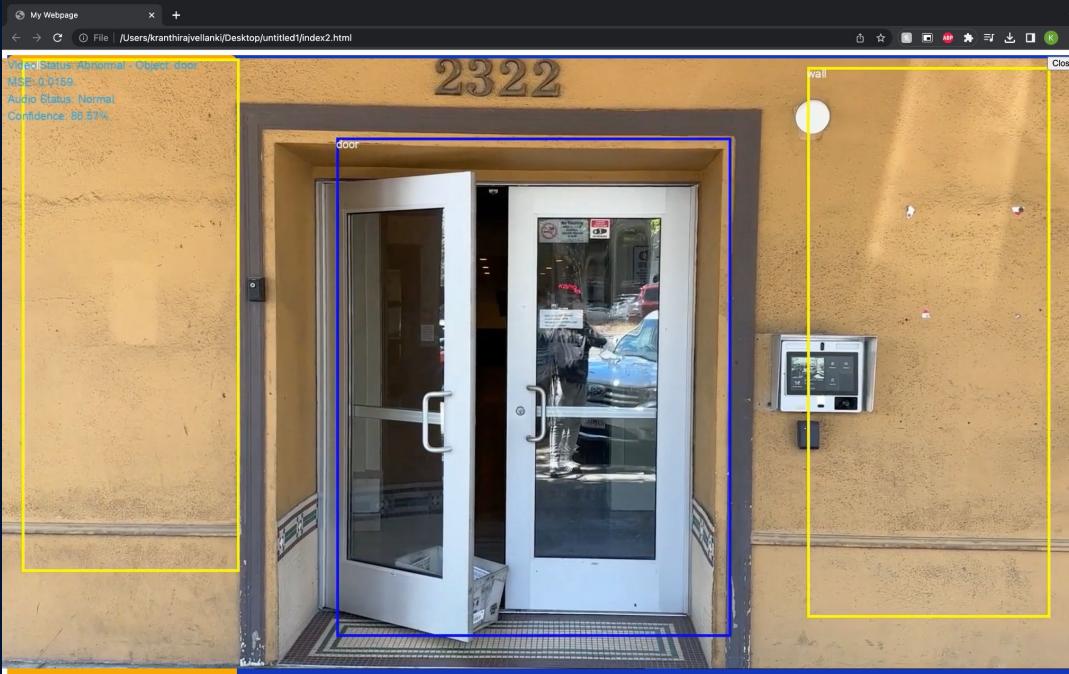
- Back end built with flask, this is where all the processing is taking place.
- Connected to Google Drive using an API to access stored files.
- If an alert is needed it is passed to the frontend to JavaScript where the alert is raised on the front end.

GUI



- Application shows video of all 4 views.
- There is a box on the left side for notifications of anomalies detected.

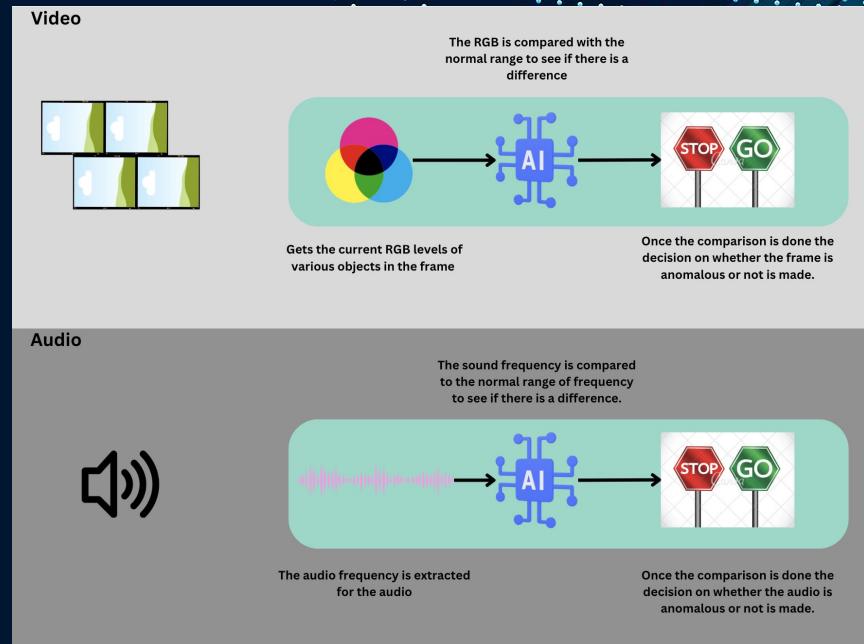
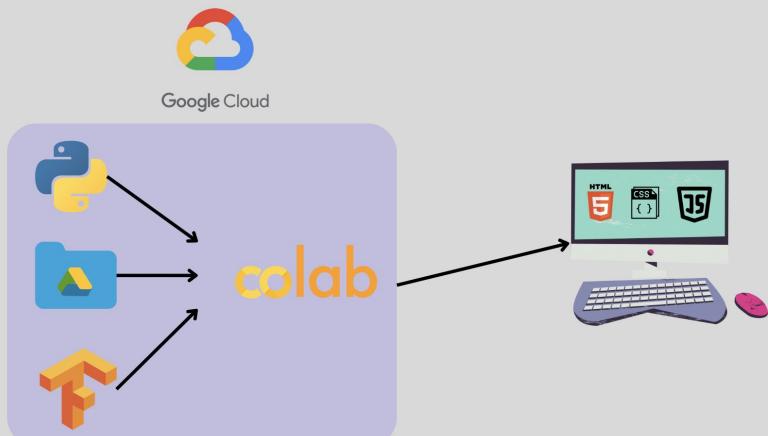
GUI



- Once the anomaly notification is clicked on, the video of the camera where an anomaly is detected is opened up and overlaid on the application.
- The video is rewinded to where the anomaly is first detected and played from there.

Intelligent Solution Development

Train & Deploy



- Trained and deployed using Google Cloud, which is connected to the web app.
- Developed an intelligent system to utilize RGB and audio frequencies in real time.

User Interface Development: Final Web GUI Deployment

- **Camera Module:** 4 View Camera Mounted to Robot
 - Example: Uses pre-recorded video from iPhone
- **Algorithm Module:** Uses Pickle files to load models and label encoder.
- **Video Acquisition Model:** Draw bounding box, along with creating alert of specific anomaly.
- **Audio Model:** Scans audio to determine if an anomaly is found.
- **Web Framework:** Utilized Flask for real time predictions
- **GUI:** HTML/CSS and Javascript for front-end



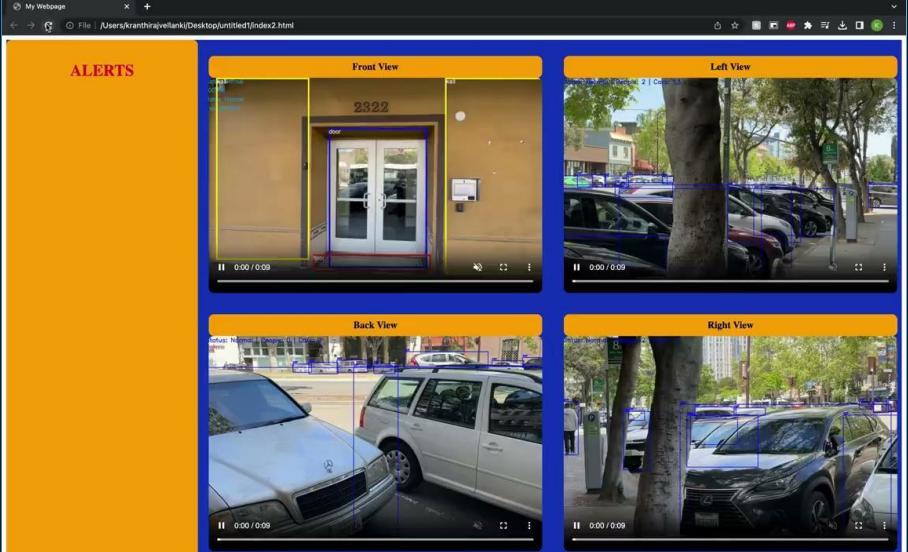
GUI Runtime Comparison

Task	Model	Runtime
Video Anomaly Detection	Layered Model	2 Minutes
Audio Anomaly Detection	LSTM/CNN Model	16 Seconds
Video & Audio Anomaly Detection	All of Above	2.5 Minutes

6. System Demo

- System Demo
- Future Work

System Demo and Testing

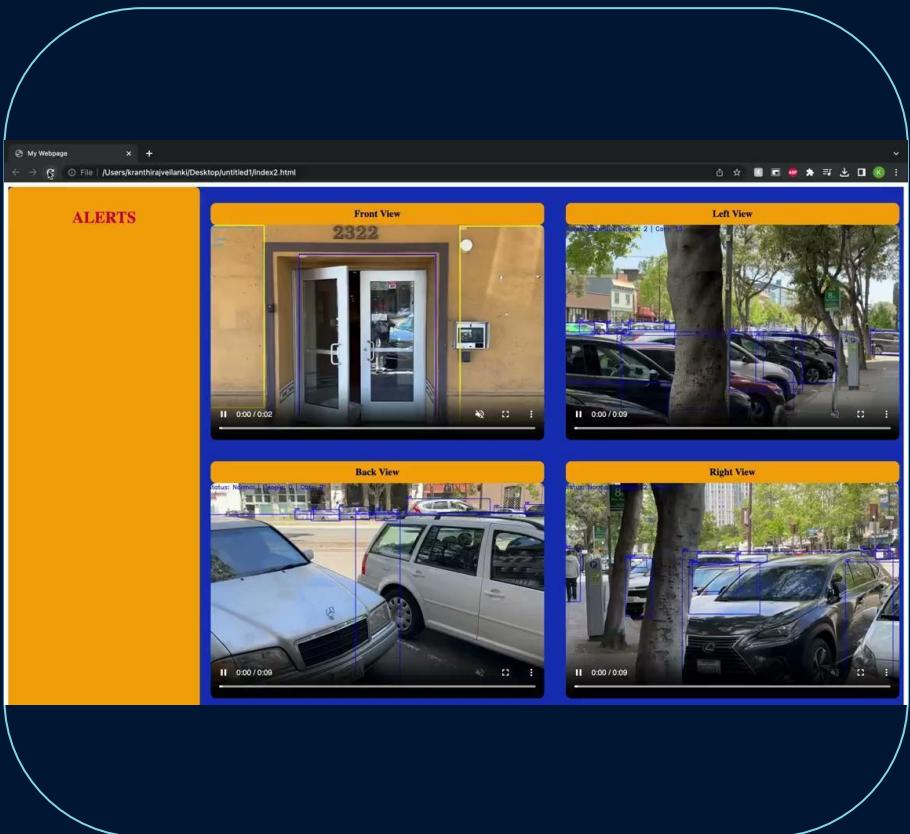


Scenario 1:

Appearance of Graffiti on
the wall

Audio Gunshot Sounds

System Demo and Testing



Scenario 2:

Front door left open

Future Work

- Add support for classification of additional anomalies.
- Make GUI runtimes faster.
- Update the model built with YOLO 8 to YOLO NAS
- Add clickable alert for sound alerts as well.

Questions?

