

Software Engineering Department

Braude College of Engineering

AI-Driven Person Recognition and Identification in Videos

Project code: 24-1-D-35

Capstone Project Phase A – 61998

Jul 2024

Rom Harel

Itay Benjamin

Git repository link:

<https://github.com/RomRL/person-recognition>

Supervisor:

Mr. Ilya Zeldner

Abstract	4
1 Introduction	4
2 Background and Related Work	5
2.1 Person Identification in Videos	5
2.1.1 Deep Learning for Person Identification	5
2.1.2 Cross-View Identification	5
2.1.3 Multi-Modal Identification	6
2.2 Object Detection	6
2.3 Current uses of person detection	7
2.4 CNN	7
2.5 Advances in Video-Based Person Recognition	8
2.6 YOLO	8
2.6.1 The key features of YOLO	9
2.6.2 How YOLO Works	9
2.6.2.1 Image and Grid Division	9
2.6.2.2 Classes and Threshold	9
2.6.2.3 Prediction Arrays	9
2.6.2.4 Prediction Process	10
2.6.2.5 Output	10
2.6.2.6 YOLO Pseudocode	10
2.6.3 Applications	10
2.7 Single Shot Detector (SSD) with MobileNetV2	11
2.8 YOLO V8	12
2.8.1 Backbone	12
2.8.2 Neck	12
2.8.3 Head	13
2.8.4 Loss Functions	13
2.8.5 Architecture of yolo v8 overview	13
2.9 FaceNet	13
2.9.1 Feature Extraction Network	14
2.9.2 Embedding Layer	14
2.9.3 Triplet Loss Function	14
2.9.4 Applications and Impact	14
2.10 Existing methods	15
3 Research / Engineering Process	16
3.1 Development methods	16
3.1.1 Python for AI and Machine Learning as server	16
3.1.2 React for the Frontend	17
3.2 Choosing Our YOLO	17
3.2.1 Performance Metrics	17
3.2.2 Chosen Model	18
3.3 Proof of scraping feasibility	18

3.2.1 Comparison between models	19
3.3.1.1 Comparison of MobileNetV2 and YOLOv8	19
3.3.1.2 Experimental Comparison of MobileNetV2 and YOLOv8 for Object Detection on Diverse Image Datasets	21
3.3.1.2.1 Experimental Setup:	21
3.3.1.2.2 Test Results :	21
3.4 Image Processing Formats: Color vs. Grayscale in Object Detection Models	23
3.4.1 Impact on Processing Time	24
3.4.2 Impact on Detection Performance	24
3.4.3 Decreased Performance in Variable Lighting	25
3.4.4 Training and Model Adaptation:	25
3.4.5 Experimental Analysis on YOLO V8-I Object Detection Confidence	25
3.4.6 Experimental Analysis Results	25
3.4.7 Decision	26
3.5 Methods to measure precision of object detection	26
3.5.1 Recall	26
3.5.2 Precision	26
3.5.3 Average Precision (AP)	27
3.5.4 Average Recall (AR)	27
3.5.5 Mean Average Precision	27
3.5.6 Intersection over Union (IoU)	27
3.6 COCO Dataset	28
3.7 Filtering and Evaluating People in the COCO Dataset	29
3.7.1 Visualization of the Data Using FiftyOne	29
3.7.1.1 Visualization EXAMPLE :	29
3.8 YOLOv8 Model Evaluation Results on COCO Validation Dataset	31
3.8.1 Dataset	31
3.8.2 Results	31
3.8.3 Insights	32
3.9 Person Detection in Videos for Missing Persons	34
3.9.1 Locating missing people	34
3.9.2 Case Study: The Nova Incident and Technological Utilization	34
4 Work Artifacts	35
4.1 System Flow	35
4.2 Architecture	36
4.2 FR & NFR Requirements	36
4.3.1 Functional Requirements	36
4.3.2 Non-functional Requirements	37
4.4 Working Progress	38
4.5 System Use Case	39
4.6 System Activity	39
4.7 User Interface	40

5 Testing Plan	42
5.1 Scope	42
5.2 Objectives	42
5.3 Test Cases	42
6 Expected Accomplishments	45
6.1 Problems	45
6.2 Criteria for Success	46
7 References	46
8 AI Prompts	50

Abstract

The wide range of video content used in different industries necessitates the creation of efficient and precise person identification systems for videos. This project proposes a new method that combines face recognition technology that is currently the best, with advanced individual detection technologies to identify people from within video files. Essentially, the YOLO (You Only Look Once) algorithm detects humans using every frame while FaceNet recognizes faces so it can tell when and where exactly those frames are in relation to its target person. This novel strategy improves our ability to analyze videos thus making it very useful for security agencies among others

1 Introduction

With increased use of video content across many domains there arises a need for advanced person identification systems. These include but are not limited to: security, surveillance, media research among others. Traditional approaches have problems with accuracy and efficiency due to dynamics within streams caused by changing environment conditions.[\[5\]](#)

The relevance of such technology is evident from incidents like that which occurred on 7th October where hundreds of Israelis were taken hostage into Gaza. Our system could play a big role in identifying persons seen in these footages during rescue or investigative operations carried out under high-risk situations like this one. By providing immediate & precise identity verification abilities; law enforcement agencies will be able to quickly react towards emergencies by narrowing down their searches thereby saving more lives and leading to better security outcome.

Deep learning has shown significant improvement in person recognition accuracy through hierarchical feature extraction from unprocessed data mostly using convolutional neural networks (CNNs)[\[11\]](#). For real-time object detection in video frames, YOLO (You Only Look Once) algorithm is commonly used[\[20\]](#).

FaceNet ,on the other hand ,provides higher level face recognition by transforming facial images into Euclidean space hence making it robust against different orientations and expressions[24]. This proposal combines YOLO for object detection with FaceNet for face recognition thus creating an end-to-end system capable of accurately identifying people from videos as discussed below.

Our system is based on the latest neural networks together with AI techniques hence it performs exceptionally well in person recognition tasks. It increases not only accuracy of recognizing individuals within complicated video environments but also provides an interface which can be used by many people who may not have skills required to operate advanced tools for analyzing videos. This project seeks to merge these state-of-the-art technologies into a single package that could potentially transform how we understand various aspects of video content across different fields.

2 Background and Related Work

2.1 Person Identification in Videos

In computer vision, the process of identifying individuals in video streams is critical but complicated, serving needs in surveillance, forensics and monitoring. The advent of deep learning has greatly improved the accuracy and resilience of this technology, with Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) being major drivers. CNNs excel at capturing important spatial features from single frames while RNNs such as ConvLSTMs are effective at capturing temporal dynamics across sequences. These architectures accommodate variations in appearance, lighting conditions and posture to ensure consistent identification over consecutive frames and viewpoints. Other improvements like 3D pose estimation and multi-stream networks have led to better performance for cases involving multiple perspectives or modalities [1]. However, real-time processing constraints, scalability issues and dataset generalization challenges persist. Recent studies aim to enhance robustness, computational efficiency and domain adaptation while addressing ethical concerns that may restrict their practicality. With improving deep learning techniques however; its application in person identification through videos is expected to increase immensely thus revolutionizing visual data analysis dependent industries.

2.1.1 Deep Learning for Person Identification

Person identification has evolved greatly and Deep Learning, which refers to Convolutional Neural Networks (CNNs) in particular, is at the forefront of this. In the past, typical methods

relied on handcrafted features that often failed under new conditions. Here, CNNs are able to automatically learn hierarchical features which make them more resilient against any form of variations [34]. Such architectures usually have multiple convolutional layers followed by pooling layers that capture spatial hierarchies and fully connected layers for identity prediction.

2.1.2 Cross-View Identification

Person identification can also require matching individuals across different camera views or modalities (e.g., RGB and thermal images). Cross-view identification is very difficult because perspectives and backgrounds keep changing.

One way for cross-view identification is by learning features that don't change no matter the angle and by using 3D human pose data [29]. For example, CNNs (Convolutional Neural Networks) can be trained to pick up on features that stay consistent across various views. Another helpful method is using multi-view dynamic images, which combine information from different angles to make the system more robust and reliable [30].

2.1.3 Multi-Modal Identification

For environments with varied lighting conditions such as nighttime surveillance, combining visible and thermal imagery will enhance identification accuracy. Two-stream CNNs are used where one stream processes visible images and the other thermal images. These networks are trained jointly so that they learn features that bridge the gap between two modalities.

2.2 Object Detection

Object detection is a basic computer vision task that aims at identifying and locating objects in digital images. This capability is a primary requirement for self-driving cars, security systems and augmented reality which demand high accuracy of object detection [33]. In contrary to image classification that can only identify if there are objects in an image, object detection should take into account the exact positions of each instance using bounding boxes. Consequently, this makes it more challenging because the model has to both recognize the object as well localizing it correctly. In terms of object localization and recognition, deep learning, especially Convolutional Neural Networks (CNNs), has revolutionized how they perform this function [12]. They have improved on architectures which learn raw pixel data directly for rich and hierarchical visual features automatically hence making them capable of detecting intricate patterns and variations in objects that are not easily engineered manually. Therefore, recent developments in CNNs have produced modern detectors primarily based on deep CNN architectures which are meticulously designed and optimized for this purpose. A significant architecture within this field is Faster R-CNN which utilizes a Region Proposal Network (RPN) to detect possible areas where an object may be located first in an image. A subsequent stage then refines these proposals by classifying them [21].

However, while two-stage approach provides higher accuracy but it is computationally expensive thus limiting its application to real-time situations. On the other hand one-stage detectors such as YOLO(You Only Look Once) give bounding boxes along with classification scores through a single pass network thereby having significantly faster inference time as compared to two-stage models albeit with a slight reduction in precision. Despite their impressive capabilities, state-of-the-art object detectors still encounter challenges such as occlusions, varying object scales, and cluttered or complex backgrounds. To address these issues some recent works introduced feature pyramids [21], or attention mechanisms that can select relevant parts automatically within input area [20]. Object detection remains an active area of research as we strive for more robust and efficient detectors that can operate in real-world environments. In general, modern object detectors have become powerful tools through deep learning breakthroughs enabling a variety of exciting applications across different domains such as robotics, augmented reality and autonomous vehicles. Even though they are not perfect their performance continues to improve at a breakneck pace thus propelling the progress in related fields while opening new horizons in technology.

2.3 Current uses of person detection

Person detection technologies are often used in video scenarios, to identify individuals. In surveillance and security settings, they enable automated monitoring, can track people and determine what activities are taking place [5]. To analyze consumer behavior and manage queues more effectively so as to improve operations of a store, retailers use person detections in this case study. Smart cities and traffic management rely on the counting of pedestrians and crowd analysis to optimize traffic flow and ensure the safety of pedestrians. Monitoring and fall detection are crucial for healthcare facilities as well as assisted living centers that use these technologies for providing safety alerts. Sports industries like athletics involve tracking players by using person detection also known as tracking technologies. It assesses their performance to give them ideas. Autonomous vehicles count on it for pedestrian recognition while cycling is important in roadway safety assurance. Workplace safety has improved through high-risk area monitoring with corresponding strict enforcement of safety rules. Identifying individuals of interest in intelligence agencies will require person detection during investigations which involves analyzing surveillance footages thus enhancing national security measures from such cases. They also help in looking out missing persons or kidnaps by scanning CCTV footage to identify victims or suspects thereby aiding investigations effectively with respect to law enforcement agencies employing person detection technology.

Such instances demonstrate how personal identification tools contribute towards enhanced protection, secure surroundings as well as efficient sector operation across various domains.

2.4 CNN

Convolutional neural networks (CNNs) are a class of deep neural network architectures that have changed the field of computer vision and image processing. The main thing which allowed them to achieve high accuracy in solving a wide range of computer vision tasks such as image classification, object detection, semantic segmentation, and instance segmentation is their ability to get hierarchical features from raw pixel data [13]. CNNs consist of convolutional layers which apply learnable filters on input images resulting in low-level features like edges, textures, patterns etc. Convolution operations are made for capturing spatially localized patterns that enables the network to exploit inherent stationarity in images. On the other hand, pooling layers reduce computational complexity by downsizing the spatial dimensions while retaining the most informative parts. The network layers used in Convolutional Neural Networks are nonlinear, though popular and standard pushing include rectifying linear units (ReLUs). This brings in non-linearity, important for these models to learn complex distributions present in data in real life scenarios. In the deep network it is necessary for the network to construct high level feature representations, that is why such nonlinearity is crucial. When the depth of the network grows, the upper layers in a sense are able and allowed to learn more high-level and even semantic concepts by combining and transforming the features gathered by the layers lower in the hierarchy. This hierarchical learning is a powerful strength of CNNs to get multiple level representations from raw data and thus free from feature engineering.

Several challenges exist in CNNs , they require large annotated datasets and are easily considered adversarial, and finally, the learned representations are often difficult to interpret—also termed the black box effect. Modern studies are looking into methods like transfer learning, adversarial training, and explainable artificial intelligence to overcome these problems and enhance the CNN model's robustness, effectiveness, and interpretability [13].

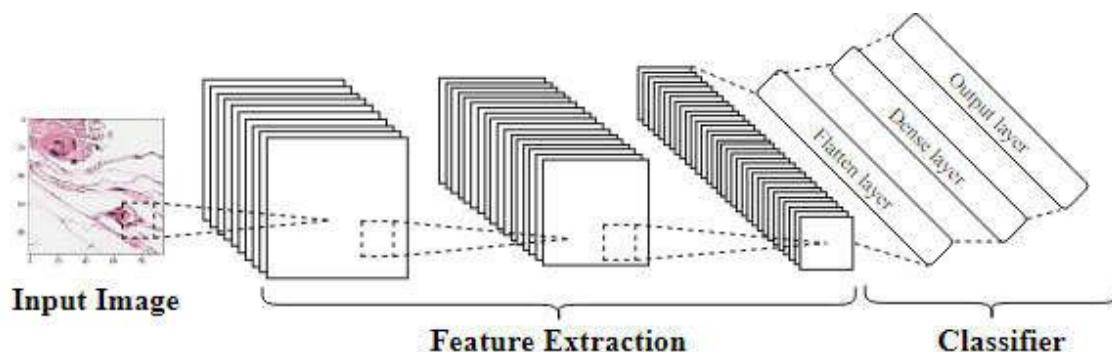


Figure 1 : Convolutional neural network (CNN) diagram.[10]

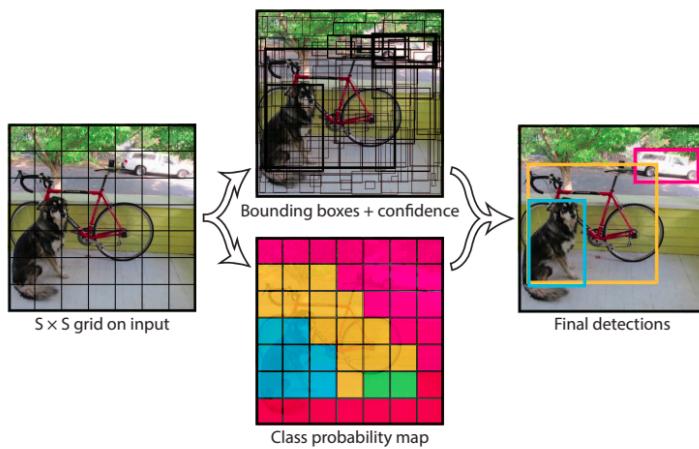
2.5 Advances in Video-Based Person Recognition

Compared to still pictures, recognizing persons in video streams presents a great challenge on account of the issues related to changing viewpoints, occlusion, lighting changes and dynamic backgrounds. It is for this reason that image-based recognition techniques often fail

in video sequence settings, calling for specialized models. One way forward is using end-to-end deep learning models which process raw videos as they are. A type of these is 3D CNNs whose conventional 2D convolution is expanded into a time dimension thereby enabling it to extract spatio-temporal features that capture both spatial and temporal patterns [26]. On an alternative note recurrent neural networks such as LSTMs can model dynamic temporal dependencies by processing inputs recursively over time [4]. Combining these models by extracting spatial features from individual frames using CNNs and feeding them into RNNs for temporal integration exploits the strengths of both architectures. The significance of deep learning in this field lies in its ability to find and utilize spatiotemporal patterns automatically instead of relying on hand-engineered heuristics. Attention mechanisms have recently been introduced into these models [6], making it possible for networks dynamically shift their focus across relevant spatial regions and time segments thereby getting rid of noise. This reflects human attention mechanisms and improves performance while also providing interpretability by indicating areas where the model has focused on.

2.6 YOLO

YOLO is a real-time object detection system which uses deep learning to identify specific objects in images and videos [20]. It was developed by Joseph Redmon et al. at the University of Washington and introduced in 2016 in the paper "You Only Look Once: Unified, Real-Time Object Detection". YOLO is a single convolutional neural network that divides the input image into a grid system, where each cell in the grid is responsible for detecting objects within itself. For these boxes simultaneously, the network predicts bounding boxes and class probabilities.



*Figure 2 :Yolo approach for detection by modeling it as a regression problem. The image is divided into an $S \times S$ grid, and each grid cell is responsible for predicting B bounding boxes, the confidence scores for those boxes, and C class probabilities. These predictions are then encoded into a tensor with dimensions $S \times S \times (B * 5 + C)$. [20]*

2.6.1 The key features of YOLO

What makes YOLO unique from other systems is its ability to process information in real time – it can handle up to 45 frames per second (hence very fast). Object detection is performed with only one evaluation of the neural network thus ensuring high efficiency. Moreover, it can generalize and detect objects that were not specifically trained on [\[20\]](#).

2.6.2 How YOLO Works

2.6.2.1 Image and Grid Division

The object detection process begins with an input (for this example use grayscale image) of dimensions $X * Y$. This image is then divided into an $N \times N$ grid, where each cell is processed independently to detect objects.

2.6.2.2 Classes and Threshold

YOLO recognizes specific object classes like person's class after being trained; step size for each cell = image dimensions / number of cells hence cell dimensions are $X/N , Y/N$.

2.6.2.3 Prediction Arrays

For instance if there are four possible categories namely person, dog, wolf cow , then yolo will be taught on how to recognize them all; therefore it should predict probability of every class for every cell ; this information is stored in array format (e.g., 50% person, 30% dog, 10% wolf, and 20% cow). YOLO also predicts several bounding boxes for each cell where box is defined by its center coordinates, width, height and confidence score (which shows how likely it overlaps with real object) [\[20\]](#).

2.6.2.4 Prediction Process

Each cell initializes an empty list to store final predictions. The algorithm extracts each cell from the grid and loops over them, predicting class probabilities as well as finding best bounding boxes based on confidence scores. It identifies highest probability class, and if its confidence score exceeds a predetermined threshold then adds this prediction to the final list.

2.6.2.5 Output

The final output is a list of detected objects. Each entry in the list includes the coordinates of the bounding box and the predicted class.

2.6.2.6 YOLO Pseudocode

```

image = readImage()
NoOfCells = 7
NoOfclasses = 1
step = int(height(image) / NoOfCells)
prediction_class_array = new_array(size(NoOfCells, NoOfCells, NoOfClasses))
predictions_bounding_box_array = new_array(size(NoOfCells, NoOfCells, NoOfCells, NoOfCells))
final_predictions = []
threshold = 0.5

for i in range(0, NoOfCells):
    for j in range(0, NoOfCells):
        cell = image[i*step:(i+1)*step, j*step:(j+1)*step]
        prediction_class_array[i, j] = class_predictor(cell)
        predictions_bounding_box_array[i, j] = bounding_box_predictor(cell)
        best_bounding_box = 0 if predictions_bounding_box_array[i, j, 0, 4] > predictions_bounding_box_array[i, j, 1, 4] else 1
        predicted_class = index of max value(prediction_class_array[i, j])
        if predictions_bounding_box_array[i, j, best_bounding_box, 4] * max_value(prediction_class_array[i, j]) > threshold:
            prediction = [predictions_bounding_box_array[i, j, best_bounding_box, 0:4].tolist(), predicted_class]
            final_predictions.append(prediction)

print(final_predictions)

```

Figure 3 : YOLO Pseudocode [\[36\]](#)

2.6.3 Applications

YOLO has many applications in various domains, including:

- **Surveillance and security:** Real-time object detection in video feeds for security monitoring.
- **Autonomous vehicles:** Detection of pedestrians, vehicles, and other objects for self-driving cars.
- **Robotics:** Identification of objects for robotic manipulation and navigation.
- **Retail and inventory management:** Detection and counting of products on shelves.

YOLO, is used in many different ways. For surveillance systems and safety measures it is able to detect objects in live streams which makes monitoring security a lot tighter. When it comes to self-driving cars, YOLO helps them recognize pedestrians, other vehicles or any other things around them thus enabling safe navigation. In robotics it can be used for object identification during manipulation and movement planning tasks while in retail industry among other areas where goods are stored on shelves for sale it counts them quickly saving time during stock taking. All these various applications show how flexible and influential YOLO has been towards technological advancement.

2.7 Single Shot Detector (SSD) with MobileNetV2

An efficient and accurate framework for object detection is the Single Shot Detector (SSD) architecture, which was proposed by Liu et al. [\[16\]](#). In a single pass, SSD employs just one neural network to predict bounding boxes as well as probabilities of classes thereby allowing for real-time object detection.

For mobile and embedded vision applications, Sandler et al. [\[23\]](#) developed MobileNetV2 which is a lightweight convolutional neural network. It introduces inverted residual structures and linear bottlenecks on top of the original MobileNet design thus reducing computational cost significantly while maintaining high accuracy.

SSD MobileNetV2, which combines SSD with MobileNetV2, provides an interesting approach towards resource constrained devices based real time object detection. This integration makes use of MobileNetV2's efficient feature extraction capabilities as the backbone network for SSD hence producing a model that balances between speed and accuracy.

There are several important aspects incorporated in SSD MobileNetV2 that contribute to its effectiveness. Multi-scale feature maps from various layers in the MobileNetV2 backbone are used so that objects at different scales can be detected. The architecture adopts default bounding boxes having different scales and aspect ratios associated with each feature map cell which enhances detection accuracy across diverse sizes of objects. Depthwise separable convolutions used by MobileNetV2 greatly reduce parameters' number and computational complexity. Furthermore, inverted residual blocks enable more efficient feature extraction besides improved gradient flow within MobileNetsV2.

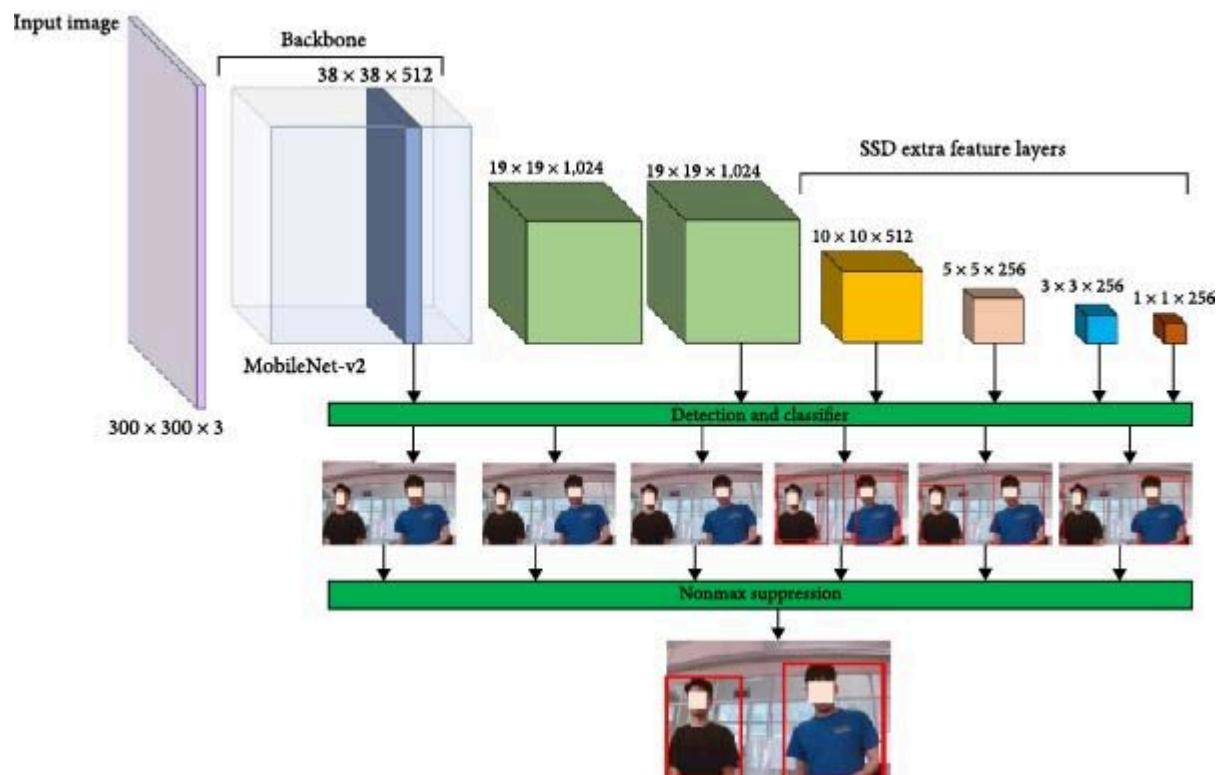


Figure 4 : SSD-MobileNet-v2 architecture [17].

In terms of object detection tasks, especially under limited computational resources scenarios like mobile devices or embedded systems [8], SSD MobileNetV2 has exhibited performance competitiveness. This means that it can be very useful in computer vision applications such as human detection or identification since its inference speed versus accuracy tradeoff ability is valuable across different setups .

2.8 YOLO V8

YOLOv8 is the latest iteration of the YOLO (You Only Look Once) object detection algorithm, which has led the field of real-time object detection since its inception. Over the years, the YOLO architecture has significantly evolved, with each version introducing innovative improvements to enhance performance and accuracy. Building on the success of its predecessors, particularly YOLOv5, YOLOv8 refines its architecture to achieve even more accurate and faster object detection. Notably, YOLOv8 is the most accurate version to date, further setting a new benchmark. The YOLO architecture consists of three main components: the backbone for feature extraction, the neck for integrating features across different scales, and the head for making final detection predictions. YOLOv8 retains this structure while incorporating several modifications to boost its performance.

2.8.1 Backbone

The backbone in the YOLOv8 architecture functions as the feature extractor, processing the input image to derive high-dimensional features crucial for object detection. Key advancements include enhanced depth-wise separable convolutions, which improve efficiency by splitting the convolution operation into depthwise and pointwise convolutions, reducing computational cost and parameters while maintaining performance [25]. Attention mechanisms, such as Squeeze-and-Excitation blocks [7] or Transformer-based modules [35], help the model focus on relevant image areas for better feature extraction. Additionally, YOLOv8 might introduce novel layer structures to optimize information flow and enhance the model's ability to capture complex patterns, essential for detecting various objects.

2.8.2 Neck

The neck in the YOLOv8 architecture plays a critical role in aggregating and refining the features extracted by the backbone, ensuring effective multiscale detection. It integrates features across different scales using mechanisms like Feature Pyramid Networks (FPN), which enhance the network's ability to detect objects of various sizes by creating a pyramid of feature maps at multiple scales, enriching each level with context from both higher and lower levels. Path Aggregation Network (PANet) [14] builds upon FPN by adding bottom-up paths, allowing more effective information flow and feature integration across scales. These mechanisms enable YOLOv8 to accurately detect objects ranging from very small to very large within the same image. The traditional YOLO neck architecture is extended by the C2f module, followed by two segmentation heads that predict semantic segmentation masks for the input image [25]. The model has similar detection heads to YOLOv8, comprising five detection modules and a prediction layer. The YOLOv8-Seg model has achieved state-of-the-art results on various object detection and semantic segmentation benchmarks while maintaining high speed and efficiency.

2.8.3 Head

The head of YOLOv8 is where the actual detection predictions occur, encompassing object classifications, locations, and sizes. It processes the integrated features from the neck to generate bounding boxes and class probabilities. Key aspects may include a hybrid approach combining

anchor-based and anchor-free detection methods [6], offering flexibility and precision in detecting a wide range of object sizes and shapes. The adoption of advanced loss functions, such as CloU loss [6] for bounding box regression, can minimize prediction errors more effectively, enhancing accuracy. Additionally, innovations in object presence and class probability predictions can refine detection performance by reducing false positives and improving class discrimination.

2.8.4 Loss Functions

YOLOv8 utilizes the CloU (Complete IoU) [35] and DFL (Differentiable Feature Loss) [31] loss functions for bounding-box loss, and binary cross-entropy for classification loss. These loss functions have been shown to improve object detection performance, particularly when dealing with smaller objects.

2.8.5 Architecture of yolo v8 overview

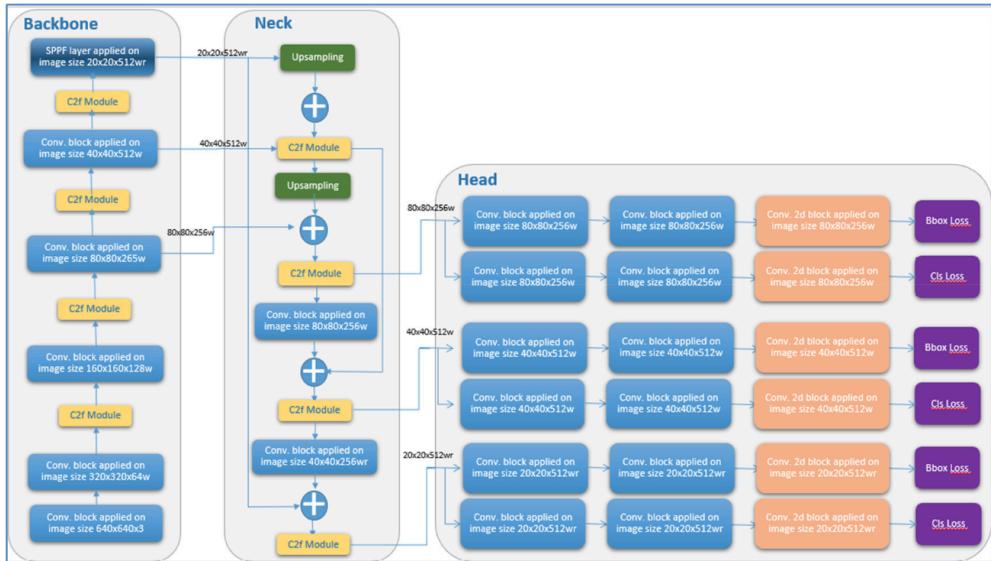


Figure 5 : Architecture of yolo v8 overview [2]

2.9 FaceNet

FaceNet, introduced by Schroff et al. [24], represents a significant advancement in face recognition technology. This deep learning model transforms facial images into compact Euclidean space embeddings, where distances directly correspond to face similarity. FaceNet's architecture consists of three main components: the feature extraction network, the embedding layer, and the triplet loss function.

2.9.1 Feature Extraction Network

In FaceNet's backbone, the feature extraction network processes input face images and derives high-dimensional features that are critical for face recognition purposes. This usually involves use of a deep Convolutional Neural Network (CNN) which is typically based on the Inception architecture.

The network learns to recognize important facial features that are invariant to changes in pose, lightening and expression hence making accurate face recognition possible in many different situations.

Also, Parkhi et al [18] demonstrated the effectiveness of deep convolutional networks for face recognition with their VGGFace model which was built on a very deep convolutional network architecture . According to Wang and Deng [28], this trend has continued with the advent of other more recent models doing facial recognition.

2.9.2 Embedding Layer

he feature extraction network then feeds this output into the embedding layer where it maps onto a fixed-size vector typically 128 dimensional size). This Euclidean representation allows one to measure similarities directly using distance between embeddings as measured by the distances between each pair of them when calculating all possible pairs' similarity scores within some given range as represented by these vectors as shown in Fig.. The main contribution from Face Net relied on optimizing this space such that same identity images are closer together while those from different identities get far apart thus supporting fast search through large labeled faces dataset.

Learning discriminative face embeddings has been embraced throughout literature for quite some time now. Unlike other models such as Cao et al.'s [19] VGGFace2 dataset that was tailored towards training deep models for face recognition, this dataset has a broad spectrum of face images from different subjects taken under many visual variations.

2.9.3 Triplet Loss Function

The core training mechanism used in Face Net is the triplet loss function which works by using triplets: an anchor image, a positive image (same person as the anchor), and a negative image (different from the anchor). The objective of this model is to minimize the distance between a pair of embeddings representing an anchor and its positive samples while at the same time maximizing their separation from other embeddings resulting into nonoverlapping distributions as shown in Fig.. This was intended to enable clustering of similar images together so that during testing any new images with close similarity will be considered having been captured from one person while different people's faces will be represented by widely scattered points on x-y axes.

This approach has motivated many researches to come up with modifications or improve it over time and Wang and Deng [28] give an overview of various such loss functions for deep face recognition including center loss, angular softmax loss etc aimed at enhancing power of learned embeddings.

2.9.4 Applications and Impact

Moreover, FaceNet's unified embedding approach allows multiple face recognition tasks such as one-to-one matching(face verification), one-to-many matching(face recognition) and clustering based on facial similarity. In particular, it performs outstandingly well when tested on benchmark data sets outperforming human beings on LFW dataset [24].

FaceNet has had wide-ranging influence since its inception. This has led to several systems being developed among others . For example, limitations in existing datasets prompted development of

VGGFace2 dataset [19] which contains many more labels thus making it possible to train deep models capable enough to recognize faces even when they change either due to pose or aging etc.

The authors Wang and Deng [28] have concluded that learning discriminative embeddings is the main principle behind FaceNet as a modern face recognition system. This area keeps changing with time and hence, current research focuses on improving resistance to variations, dealing with prejudice and fairness concerns, as well as enhancing scalability.

In conclusion, FaceNet represents a pivotal development in face recognition technology, introducing a unified embedding approach that has significantly influenced subsequent research and practical applications in the field of computer vision and biometric identification.

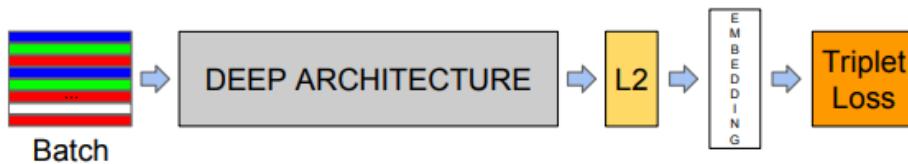


Figure 6 : Model structure. Our network consists of a batch input layer and a deep CNN followed by L2 normalization, which results in the face embedding. This is followed by the triplet loss during training [24] .

2.10 Existing methods

In our research, we have explored various existing facial recognition and analysis systems to identify their features and limitations. The following table (Table 1) provides a comparative analysis of these systems. Our proposed solution stands out by offering a free, user-friendly interface that requires no prior knowledge, thus making advanced facial recognition technology accessible to a broader audience. This need became particularly evident on 07.10, when we searched for a system to recognize our friends from the Nova disaster in Telegram videos, and found the process to be extremely challenging.

Feature	Amazon Rekognition [43] aws	Microsoft Azure Face API [44]	Face++ [45] Face++	Deep Vision AI [46] DEEPPVISIONAI	Our Solution

Face Detection	✓	✓	✓	✓	✓
Face Search/Matching	✓	✓	✓	✓	✓
Video Analysis	✓	Limited	Limited	✓	✓
User-Friendly Interface	✗	✗	✗	✗	✓
No Prior Knowledge Required	✗	✗	✗	✗	✓
Cost	Paid	Paid	Paid	Paid	Free
Open-Source	✗	✗	✗	✗	✓

Table 1 : provides a comparative analysis of Amazon Rekognition, Microsoft Azure Face API, Face++ Deep Vision AI and Our Solution highlighting their capabilities in face detection, face search, video analysis, and user accessibility.

Table 1 clearly shows how our proposed system differentiates itself by being free, user-friendly, and requiring no prior technical knowledge, making it accessible to a wider audience.

3 Research / Engineering Process

3.1 Development methods

3.1.1 Python for AI and Machine Learning as server

Python offers numerous benefits; its simplicity and readability are ideal for rapid development and prototyping. The language boasts a rich ecosystem of libraries and frameworks, such as OpenCV, Dlib, TensorFlow, and PyTorch, which provide pre-built tools for computer vision and machine learning tasks. Additionally, Python's extensive community support and a wealth of online resources make problem-solving and collaboration much easier. Its versatility and ability to integrate seamlessly with other technologies make Python a powerful and efficient choice for developing advanced AI and neural network-based applications [\[49\]](#).

3.1.2 React for the Frontend

We are going to use React for the frontend because it provides a flexible and efficient way to build user interfaces. React's component-based architecture promotes reusability and maintainability, making it easier to manage complex UIs. Its virtual DOM ensures high performance by minimizing

updates to the real DOM. Additionally, React has a strong developer community and a rich ecosystem of tools and libraries, which accelerates development and enhances productivity [48].

3.2 Choosing Our YOLO

In the quest to identify the most suitable model for high-accuracy person recognition tasks, this section compares two advanced object detection frameworks: YOLOv5 and YOLOv8. Through a detailed analysis of their architectures and performance metrics, we demonstrate why YOLOv8 emerged as the superior choice for our specific application needs. Object detection has evolved significantly with the advent of YOLO (You Only Look Once) models, renowned for their real-time detection capabilities and accuracy. YOLOv5 and YOLOv8 represent key iterations in this evolution, each offering unique advantages. This paper provides an in-depth comparison of these models to determine the optimal choice for person recognition tasks.

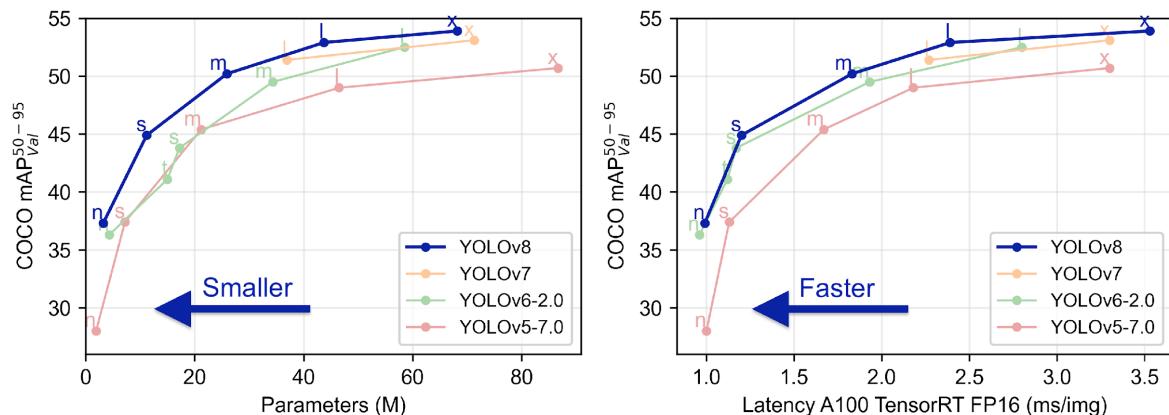


Figure 7 : Comparison of the YOLO models [37]

Figure 7 compares the performance of various YOLO models (YOLOv8, YOLOv7, YOLOv6-2.0, and YOLOv5-7.0) on the COCO dataset [40] using two metrics: the number of parameters (left) and inference latency (right). The left plot shows the mean Average Precision (mAP) as a function of the number of parameters, highlighting that YOLOv8 achieves higher mAP with fewer parameters compared to previous versions. The right plot illustrates mAP versus latency, demonstrating that YOLOv8 models are both faster and more accurate, achieving superior performance at lower latencies. Overall, YOLOv8 provides an efficient and speedy solution for object detection tasks.

3.2.1 Performance Metrics

The performance of YOLOv5 and YOLOv8 was evaluated based on precision, recall, and mean Average Precision (mAP).

Metric	YOLO v 5	YOLO v 8
--------	----------	----------

Precision	89.3%	91.6%
Recall	74.4%	83.1%
mAP50	85.2%	88.5%
mAP50-95	58.5%	60.0%

Table 2 :Performance comparison using YOLOv5 and YOLOv8 [2]

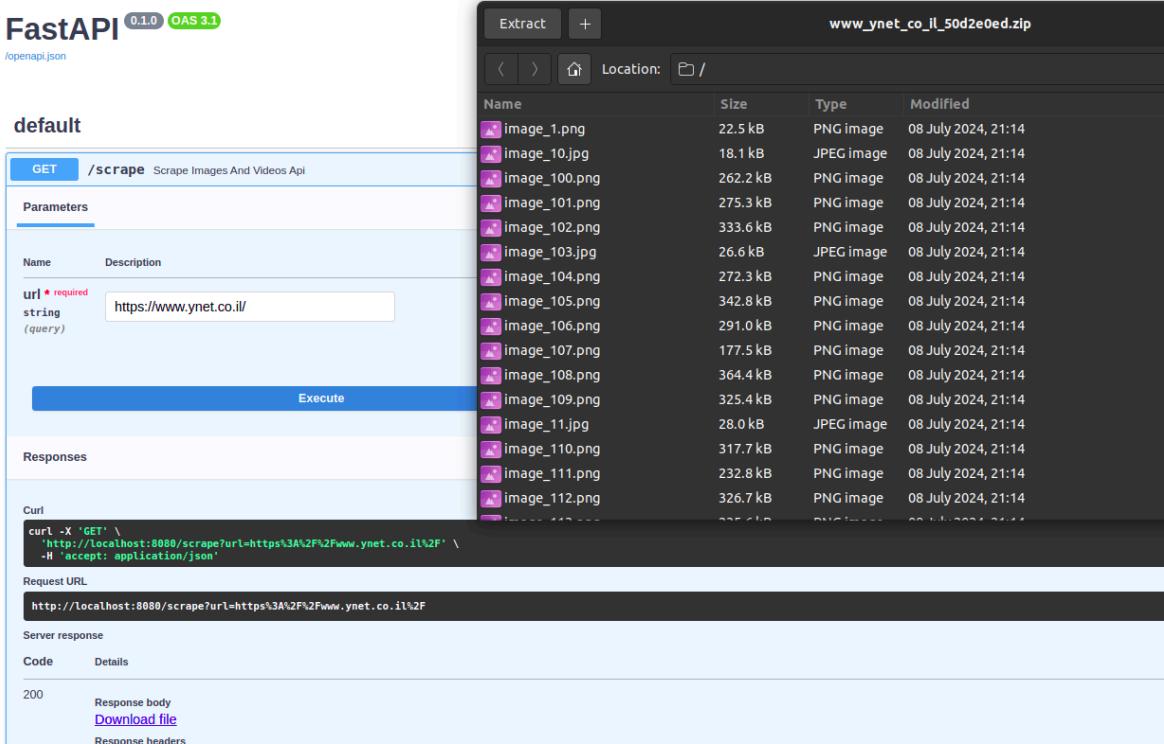
YOLOv8 consistently outperformed YOLOv5 across all metrics , demonstrating its superior capability in detecting and classifying objects accurately. The inference speed of both models was also a critical factor. YOLOv8, with its simplified architecture, provided faster inference times compared to YOLOv5, making it more suitable for real-time applications.

3.2.2 Chosen Model

Based on our comprehensive analysis , YOLOv8 was chosen over YOLOv5 for our person recognition tasks. Its advanced architecture, higher accuracy, faster inference times, and reduced training complexity made it the ideal candidate. YOLOv8's robust performance in diverse conditions ensures reliable and accurate person detection, essential for our application needs.

3.3 Proof of scraping feasibility

We developed an advanced web server designed to efficiently extract and download images and videos from any specific webpage. This tool streamlines the process of gathering media content by automating the retrieval and packaging of these files into a single, convenient zip file. Users simply input the desired URL, and the application performs the tasks of fetching the webpage content, identifying and downloading media files, and compressing them for easy download. This systematic approach not only saves time but also ensures that users receive their content in an organized and accessible manner.This shows us the feasibility of bringing additional sources of information to search and identify a person in cases where we would like to perform proactive searches of the system on certain sites of interest



The image shows a screenshot of a web application interface. On the left, there is a sidebar with the title "FastAPI 0.1.0 OAS 3.1" and a link to "/openapi.json". Below this, under the heading "default", there is a section for a "Scrape Images And Videos Api" endpoint. It includes fields for "Parameters" (a "url" parameter set to "https://www.ynet.co.il/"), an "Execute" button, and a "Responses" section. The "Responses" section shows a "curl" command to run the API and a "Request URL" field containing "http://localhost:8080/scrape?url=https%3A%2F%2Fwww.ynet.co.il%2F". At the bottom, there is a "Code" section with a "200" status code and options to "Response body" or "Download file". On the right side of the interface, there is a file explorer window titled "www_ynet_co_il_50d2e0ed.zip" showing a list of files. The files are all images (PNG and JPEG) ranging in size from 18.1 kB to 342.8 kB, all modified on 08 July 2024 at 21:14.

Figure 8 : server in action, realized by fast-api in python , Example of entering a sample URL of an internet site, and we received back a zip file that contains all the images and videos on the page.

This way we can enable the expansion of our search sources past figures and searches initiated by the system

3.2.1 Comparison between models

3.3.1.1 Comparison of MobileNetV2 and YOLOv8

Over the years, various deep learning architectures have been developed to tackle this challenge efficiently. Among them, MobileNetV2 and YOLOv8 stand out as popular choices due to their effectiveness and speed. This essay aims to provide a detailed comparison between MobileNetV2 and YOLOv8, focusing on their architecture, performance, and practical considerations.

Object detection models, particularly those leveraging convolutional neural networks (CNNs), have significantly advanced automated systems' capabilities to identify and locate objects within digital images. Models like YOLOv8 and MobileNetV2 SSD, though distinct in design and application, share several fundamental methodologies. Both utilize CNNs to process visual data through layers of convolutions, building a hierarchy of features from basic edges to intricate objects. They generate bounding boxes with spatial coordinates and confidence scores, classifying each detected object into predefined categories. Anchor boxes streamline detection by offering reference points for various object sizes and shapes. Efficient feature extraction networks produce detailed feature maps crucial for detection tasks. Both models employ non-maximum suppression to remove redundant bounding boxes, enhancing precision. Operating on the single-shot detection principle, they perform

localization and classification in one network pass. YOLOv8 is designed for real-time use, balancing accuracy and speed, while MobileNetV2 SSD is optimized for performance in environments with limited computational resources. These shared methodologies highlight a commitment to advancing object detection technologies, reflecting the convergent evolution of neural network architectures towards robust, real-time visual recognition systems.

Feature	YOLO v8	MobileNet V2 SSD
Date of Release	January 10, 2023	January 13, 2018
Model Type	Object Detection	Object Detection
Framework	PyTorch	TensorFlow
GitHub Stars	21.1k +	81+
Mean Average Precision	~55%	~35%
Number of Parameters	Highly variable , Average 30.78 million	Around 3.5 million (varies with implementations)
Number of Layers	Deep network, number varies by specific variant	Generally around 53 layers including depthwise convolutions
Processing Speed (FPS)	Up to 140 FPS depending on hardware and settings	
Computational Requirements	High, benefits significantly from GPU acceleration	Designed for efficiency, can run on CPUs with reasonable performance
Training Data Size	Extensive data required for optimal performance	Can be trained effectively with moderately sized datasets
Hardware Requirements	High-end GPUs recommended for training and inference	Can run on less powerful devices like smartphones
Optimization Techniques	Uses techniques like batch normalization, multi-scale predictions	Uses depth wise separable convolutions for efficiency
Typical Use Case Scenarios	High-speed real-time detection in various conditions	Efficient detection in resource-constrained environments like mobile devices

Detection Capabilities	Excellent at detecting small to large objects with high precision	Good at detecting moderate to large objects, may struggle with very small objects
Flexibility	Highly flexible with adjustments for various performance metrics	Somewhat flexible, mainly constrained by hardware efficiency
Energy Consumption	Generally higher due to more complex computations	Lower, optimized for energy efficiency especially in mobile devices

Table 3 : comparison YOLOv8 and MobileNet SSD v2 [\[38\]](#)

YOLOv8 stands out for its impressive accuracy and precision, especially when it comes to spotting small or overlapping objects, which is crucial for detailed object detection tasks .This high level of accuracy is due to its advanced feature extraction techniques and a finely-tuned architecture. Although MobileNetV2 SSD isn't as precise as YOLOv8, it still handles many object detection tasks well, particularly in situations with limited computational resources [\[32\]](#). When it comes to speed, YOLOv8 is built for rapid processing on powerful hardware, making it perfect for real-time detection scenarios. On the other hand, MobileNetV2 with SSDLite might not be as quick on high-end devices but offers good detection speeds on mobile and embedded systems [\[3\]](#). YOLOv8's complexity means it has a larger model size and demands more memory and processing power. In contrast, MobileNetV2 with SSDLite is designed to be lightweight, which reduces its size and lowers memory and processing requirements, making it highly efficient.

3.3.1.2 Experimental Comparison of MobileNetV2 and YOLOv8 for Object Detection on Diverse Image Datasets

In this study, we conducted a comprehensive comparison between MobileNetV2 and YOLOv8 for object detection on diverse image datasets. Our evaluation focused on performance metrics such as detection time, average confidence scores, and accuracy in detecting persons within images. We employed a dataset comprising 1000 images captured at three different resolutions, with ground truth annotations indicating the actual number of persons in each image. By analyzing the results, we aimed to provide insights into the strengths and weaknesses of each model under varying conditions.

3.3.1.2.1 Experimental Setup:

We ran both MobileNetV2 and YOLOv8 on our local machine using cpu (11th Gen Intel® Core™ i7-11800H @ 2.30GHz × 16) , which was equipped with standard hardware configurations. The dataset consisted of 1000 images, categorized into three groups based on resolution: low (480x360), medium (960x720), and high (1920x1080). Each image was manually annotated to indicate the exact number of persons present, serving as ground truth for evaluation purposes. We measured detection time, average confidence scores, and the delta between the detected and actual number of persons for each model and resolution.

3.3.1.2.2 Test Results :

MobileNetV2 exhibited faster inference times compared to YOLOv8 across all resolutions. The lightweight architecture of MobileNetV2 contributed to its efficient processing speed, making it particularly suitable for real-time applications.

YOLOv8, while slightly slower than MobileNetV2, still demonstrated impressive performance in processing high-resolution images in near real-time. Its more complex architecture resulted in longer inference times but remained within acceptable limits for most practical applications.

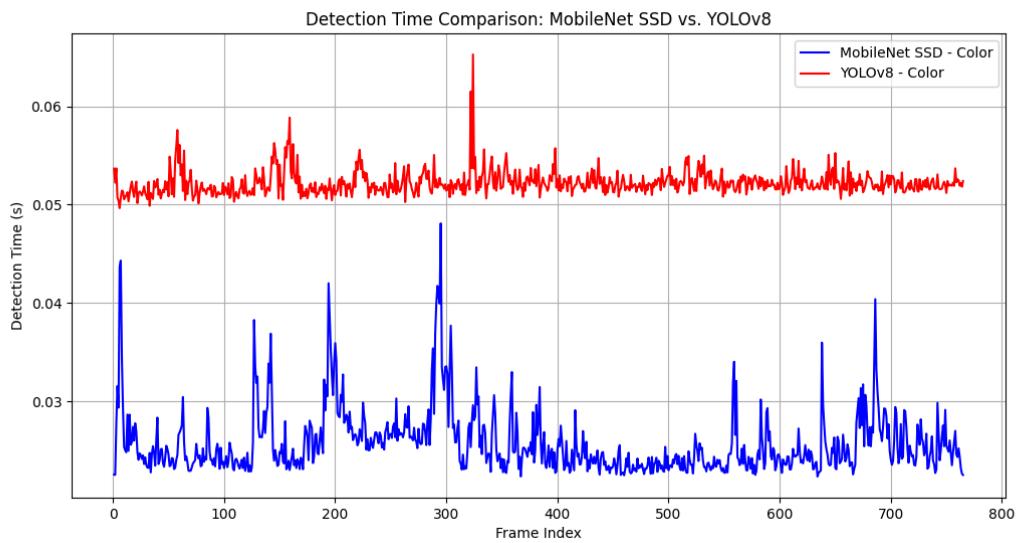


Figure 9 : comparison of detection times between MobileNet SSD and YOLOv8.

In this example MobileNetV2 SSD detected several individuals but missed three, including one who was fully visible, highlighting a limitation in its detection capability. In contrast, the YOLOv8 I model successfully recognized everyone in the scene, showcasing its superior detection performance. We can see this in an example image from our dataset, where the prediction results from each model are displayed in the figure.

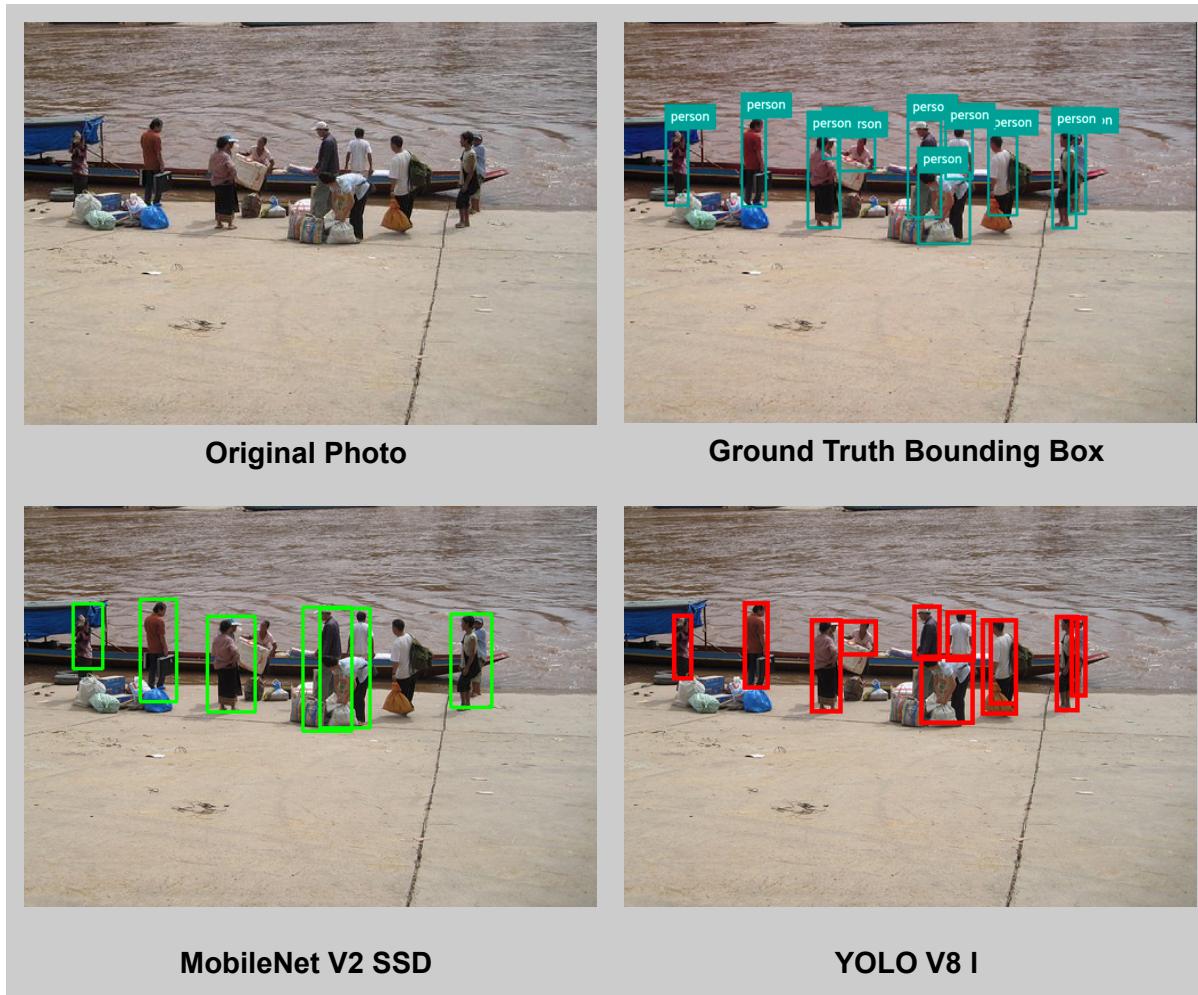


Figure 10 : Prediction results of the YOLO model and MobileNetV2 SSD are displayed simultaneously alongside the real drivers and the original image.

MobileNetV2 is renowned for its efficiency, particularly in environments with limited computational resources. It utilizes depth wise separable convolutions, a design that significantly reduces the computational overhead, thereby enhancing execution speed on less capable devices. On the other hand, YOLOv8, representing the latest in the evolution of the You Only Look Once (YOLO) series, is optimized for high-speed object detection with a design that leverages complex models to deliver high accuracy and fast processing times, particularly when supported by GPU acceleration. This comparative study is grounded on empirical benchmark tests conducted under uniform testing conditions. Both models were evaluated using the same dataset and hardware setup, with metrics such as detection time and accuracy being recorded. These tests provide a comprehensive basis for assessing which model performs better in terms of speed and resource efficiency, contributing valuable insights into their suitability for various applications.

3.4 Image Processing Formats: Color vs. Grayscale in Object Detection Models

Transitioning object detection images from color to grayscale has important implications for computational efficiency and detection accuracy in computer vision. This section explores how the performance of state-of-the-art models, YOLOv8 and MobileNetV2 SSD, changes with grayscale imaging. Through this comparative study, we identify challenges and opportunities, evaluating both theoretical and practical aspects of deploying grayscale imaging within various architectural frameworks. The relationship between image formats and model performance in object detection highlights the need for optimal computational resources to achieve high detection accuracy. While grayscale imaging can offer computational advantages in certain scenarios, it may also lead to significant performance degradation, particularly in complex environments where color is crucial for distinguishing objects. Our analysis reveals trade-offs between processing time and model accuracy, investigating whether frame detection in grayscale impacts the performance of YOLOv8 and MobileNetV2 SSD.

3.4.1 Impact on Processing Time

Both YOLOv8 and MobileNetV2 SSD models traditionally operate on RGB images. Transitioning to grayscale reduces the input channels from three to one, simplifying the data and potentially increasing processing speed. Our comprehensive tests showed that converting frames to grayscale slightly accelerates processing time [figure 11]. This efficiency gain stems from reduced data complexity: grayscale images handle a single data channel, significantly lowering memory bandwidth and computational power requirements compared to three-channel RGB images. This simplification speeds up data processing, advantageous in real-time detection scenarios. Additionally, convolutional neural networks (CNNs) underlying these models perform less complex internal computations with fewer input channels, resulting in quicker propagation through network layers and enhancing overall frame processing speeds.

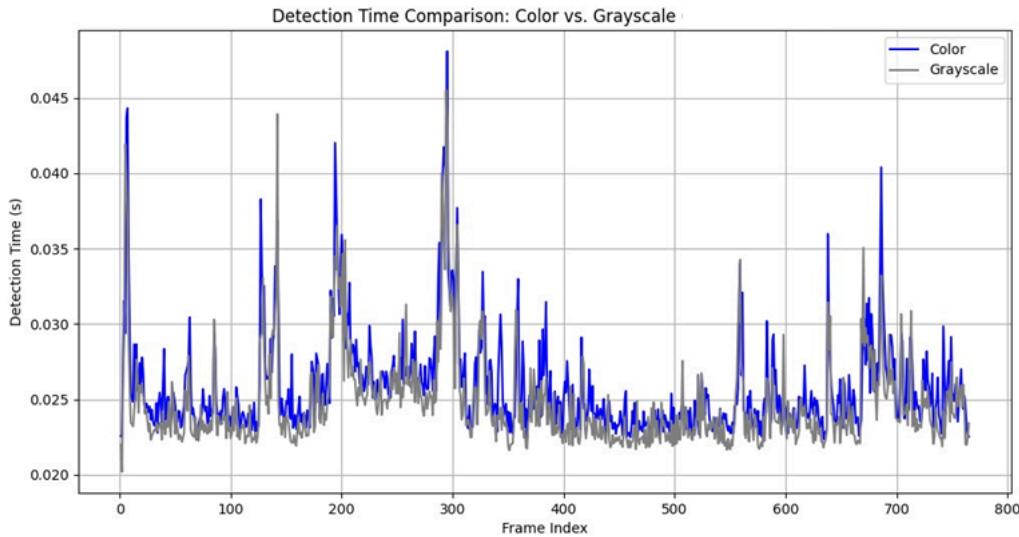


Figure 11 : This graph provides a visual comparison of detection times when using Yolo v8 I with color images versus grayscale images. The graph illustrates that grayscale processing achieves slightly faster detection times, as anticipated, due to the reduced computational load from processing fewer data channels. However, the difference in detection speeds between color and grayscale is not substantial, suggesting that the decrease in data complexity does not translate into a dramatic improvement in processing efficiency.

3.4.2 Impact on Detection Performance

Using grayscale imaging can pose significant challenges, especially in environments where color cues are essential for distinguishing people from backgrounds or from one another. Color is a critical factor in differentiating objects, particularly when they share similar shapes or are set against backgrounds with similar textures. Studies by Kanan and Cottrell (2012) indicate that object recognition tasks often perform poorly in grayscale compared to color, as many rely heavily on color differences. This issue is most pronounced in scenes where multiple objects of interest are mingled, with color being the primary distinguishing feature[9].

3.4.3 Decreased Performance in Variable Lighting

Variable lighting conditions, which are common in both indoor and outdoor environments, can affect the visibility of features in grayscale images more than in color images. Color can provide consistent identification markers that are less susceptible to changes in illumination, such as distinguishing features in low light conditions where grayscale images might suffer from insufficient contrast [22].

Also feature recognition within person detection systems often relies on subtle color cues that can indicate depth, contour, and even emotional or physical states (e.g., clothing logos, uniform colors, or skin tones under different lighting conditions). Grayscale imaging flattens these cues, potentially reducing the effectiveness of algorithms that rely on these subtle distinctions to identify and classify human figures accurately [27].

3.4.4 Training and Model Adaptation:

To optimize for the absence of color data, both YOLOv8 and MobileNetV2 SSD would ideally require retraining on grayscale datasets. This retraining is essential to maintain or potentially improve detection performance with grayscale inputs. YOLOv8 might need more significant architectural adjustments than MobileNetV2 SSD due to its reliance on complex color-based feature layers.

3.4.5 Experimental Analysis on YOLO V8-I Object Detection Confidence

In this test, we evaluated the performance of the YOLO V8 object detection model by comparing its confidence levels when processing color versus grayscale images. Over approximately 800 frames, we recorded the model's average confidence in identifying and classifying objects. Our objective was to see if converting images to grayscale, which reduces data complexity and offers computational benefits, would affect the model's detection performance.

3.4.6 Experimental Analysis Results

The results, shown in Figure Y, clearly indicate that grayscale processing does not improve and actually diminishes the model's confidence in object detection tasks. This finding reinforces that while grayscale may reduce computational load, it sacrifices essential information needed for effective object recognition.

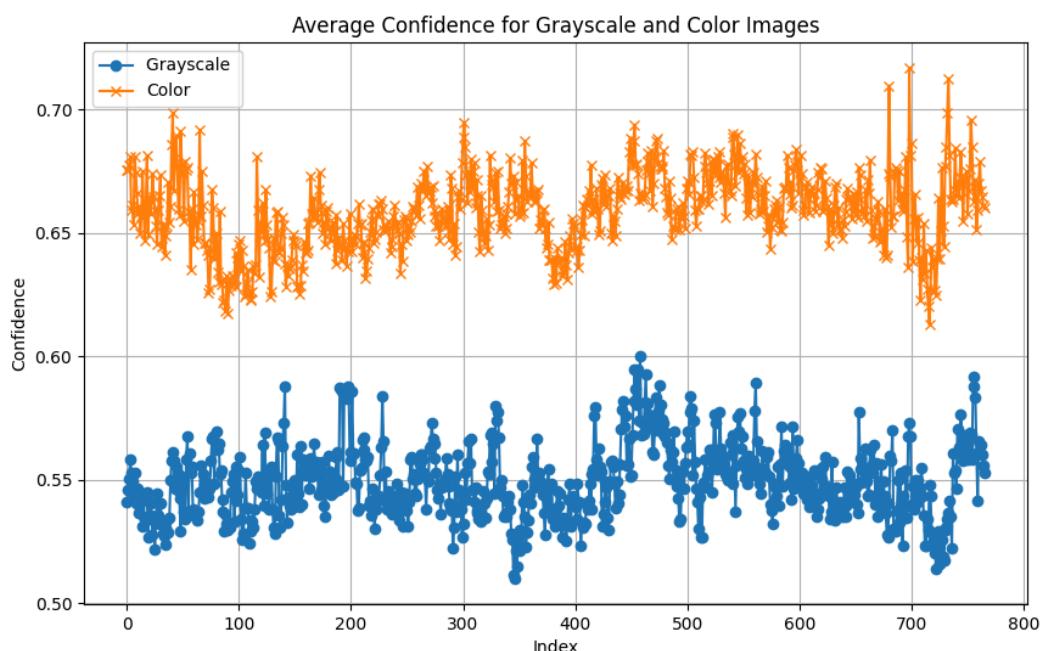


Figure 12 : This graph displays the average confidence levels of object detection using YOLO V8 I, contrasting

color processing with grayscale. The data is plotted across approximately 800 frames, highlighting the confidence with which the model identifies and classifies objects in each frame.

3.4.7 Decision

The tests on the effects of image processing formats on object detection models revealed a complex balance between computational efficiency and performance accuracy. While grayscale processing offers a slight speed gain, it significantly reduces detection performance. Therefore, to ensure operational effectiveness and reliability, especially in varied and challenging real-world conditions, we chose to retain color processing.

3.5 Methods to measure precision of object detection

3.5.1 Recall

Recall (also known as Sensitivity) measures how many of the actual objects in the dataset are detected by the model. It is defined as the ratio of true positive detections to the total number of ground truth objects (true positives + false negatives).

$$Recall = \frac{TruePositives}{TruePositives + FalseNegatives}$$

3.5.2 Precision

Precision is a measure of how many of the objects detected by the model are actually correct (i.e., true positives). It is defined as the ratio of true positive detections to the total number of detections made (true positives + false positives).

$$Precision = \frac{TruePositives}{TruePositives + FalsePositives}$$

3.5.3 Average Precision (AP)

Average Precision is a measure that combines both precision and recall to provide a single metric for model performance. It is calculated as the area under the precision-recall curve. The precision-recall curve plots precision (y-axis) against recall (x-axis) for different threshold values.

AP is usually computed for different Intersection over Union (IoU) thresholds, and the mean of these AP values gives a comprehensive view of the model's performance.

$$AP = \int_0^1 Precision(Recall) d(Recall)$$

3.5.4 Average Recall (AR)

Average Recall is the average of recall values computed for different IoU thresholds and object categories. It provides an overall measure of the model's ability to detect objects across all categories and IoU thresholds.

$$AR = \frac{1}{N} \sum_{i=1}^N Recall_i$$

3.5.5 Mean Average Precision

Mean Average Precision (mAP) is the average of the AP values computed for different IoU thresholds and object categories. It provides an overall measure of the model's performance across all categories and IoU thresholds.

$$mAP = \frac{1}{N} \sum_{i=1}^N AP_i \mid N \text{ Is The Number Of IoU Thresholds.}$$

3.5.6 Intersection over Union (IoU)

Intersection over Union (IoU) is a metric used to evaluate the accuracy of an object detection model. It measures the overlap between two bounding boxes: the predicted bounding box and the ground truth bounding box. The IoU value ranges from 0 to 1 where :

- 0 indicates no overlap between the bounding boxes.
- 1 indicates a perfect overlap between the bounding boxes.

IoU is calculated as the ratio of the area of overlap between the predicted bounding box and the ground truth bounding box to the area of their union.

$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

```

1 def compute_iou(box1, box2):
2     inter_left = max(box1['left'], box2['left'])
3     inter_top = max(box1['top'], box2['top'])
4     inter_right = min(box1['right'], box2['right'])
5     inter_bottom = min(box1['bottom'], box2['bottom'])
6
7     # Calculate intersection width and height
8     inter_width = inter_right - inter_left
9     inter_height = inter_bottom - inter_top
10
11    # If there is no overlap, the intersection area is zero
12    if inter_width <= 0 or inter_height <= 0:
13        return 0.0
14
15    # Calculate intersection area
16    intersection_area = inter_width * inter_height
17
18    # Calculate areas of the bounding boxes
19    box1_area = (box1['right'] - box1['left']) * (box1['bottom'] - box1['top'])
20    box2_area = (box2['right'] - box2['left']) * (box2['bottom'] - box2['top'])
21
22    # Calculate union area
23    union_area = box1_area + box2_area - intersection_area
24
25    # Calculate IoU
26    iou = intersection_area / union_area
27
28    return iou

```

Figure 13 : Pseudo-code for calculating IoU

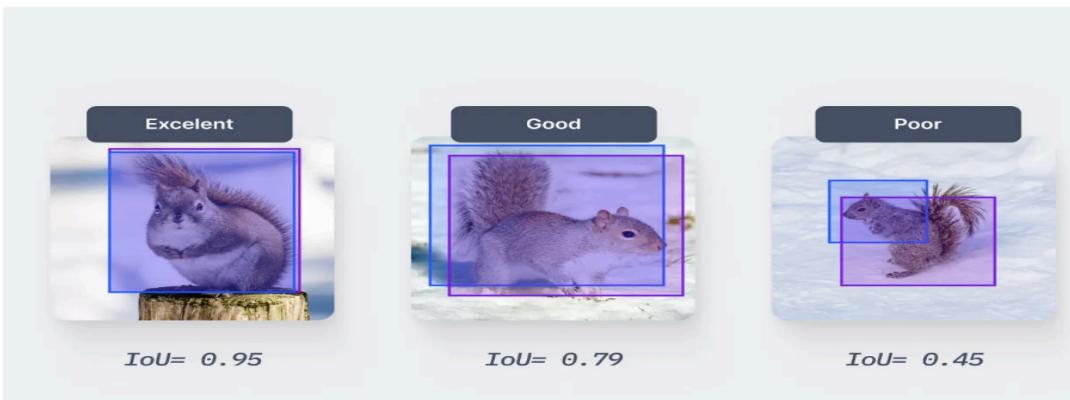


Figure 14 : IoU Example [39]

3.6 COCO Dataset

The COCO (Common Objects in Context) dataset [\[40\]](#) , introduced in 2014, is a pivotal tool in computer vision and natural language processing, widely used for object detection, segmentation, and captioning. The COCO 2017 Validation Dataset, a crucial subset, comprises 5,000 images meticulously selected and annotated to capture a broad spectrum of everyday scenes, both indoor and outdoor, with diverse complexities. This dataset's comprehensive annotations, including instance segmentation masks and bounding boxes, go beyond simple categorization, providing pixel-level detail that aids in developing sophisticated models. There are usually two common data sets used by

researchers, namely COCO 2017 Training Dataset used for training the models and the COCO Validation Dataset that can be used for testing the general performance of the model on unseen data.. Metrics such as mAP, AP, and AR provide a comprehensive analysis on the model's performance by evaluating localization and classification precision[\[15\]](#).

The benchmark for object detection and segmentation and captioning purposes is the COCO (Common Objects in Context) dataset [\[40\]](#) was introduced in 2014. COCO 2017 contains four important subsets: COCO 2017 Train Dataset, COCO 2017 Validation Dataset, COCO 2017 Test – Dev Keypoints, and COCO 2017 Test – Small. Additionally, this dataset has other annotations other than categorizing as it has instance segmentation masks and bounding boxes that give a pixel level information that can be useful in coming up with complex models. To eliminate the weaknesses in the model, there are usually two common data sets used by researchers, namely COCO 2017 Training Dataset used for training the models and the COCO Validation Dataset that can be used for testing the general performance of the model on unseen data..

3.7 Filtering and Evaluating People in the COCO Dataset

In this project, we used the COCO 2017 Validation Dataset, a smaller part of the larger COCO dataset, to test how well our object detection model works, focusing on images with people. We used the COCO API from pycocotools [\[41\]](#), a Python library, to filter images containing at least one person by checking for "person" tags. This way, we ensured our results were relevant and specific. The COCO API also helped us evaluate our model's performance using metrics like mean Average Precision (mAP) and Average Recall (AR), giving us a clear understanding of its accuracy. Overall, using the COCO API and pycocotools was crucial for effectively filtering, evaluating, and improving our model with the COCO 2017 Validation Dataset, ensuring it was well-tested and capable of performing on new data.

3.7.1 Visualization of the Data Using FiftyOne

To increase the effectiveness of our analysis of the filtered dataset, we employed FiftyOne [\[42\]](#) which is an open-source toolkit for visualizing and analyzing datasets in computer vision. Thus, FiftyOne helped to interactively investigate the COCO 2017 Validation Dataset, narrowing it down to the images containing people. It would also be easy to visually compare the ground truth annotations with the model predictions, thus giving us an indication of how well the model is performing based on the predicted bounding boxes and segmentation masks against the ground truth. Also, FiftyOne offered functionalities to plot precision-recall curves and confusion matrices, which helped to gain more profound understanding of the model's performance. This assisted in finding out areas that needed to be considered. The use of FiftyOne also enabled the inspection of specific samples, for example, zoom into a particular area of the image, select samples by criteria, and visualization of the predictions of samples by class. Such a level of analysis was important for refining our model and improving its efficiency.

3.7.1.1 Visualization EXAMPLE :

The dataset used for this evaluation is the COCO 2017 validation set, which can be found at COCO 2017 Dataset , or just download from here [\[40\]](#) . The dataset contains 5,000 images include 2,693 images with the "person" category.

```

import fiftyone as fo
import fiftyone.zoo as foz
from fiftyone import ViewField as F

# Download and load the validation split of COCO-2017
dataset = foz.load_zoo_dataset(name="coco-2017", split="validation")

# Filter the dataset to include only non-crowded 'person' annotations
view = dataset.filter_labels(
    "ground_truth", # The field containing the annotations
    (F("label") == "person") & (F("iscrowd") == 0) # Filter criteria
)

# Launch the FiftyOne app with the filtered view
session = fo.launch_app(view, port=5151)

```

Figure 15 : the code to LOAD DATA SET AND show the data USING FIFTYONE

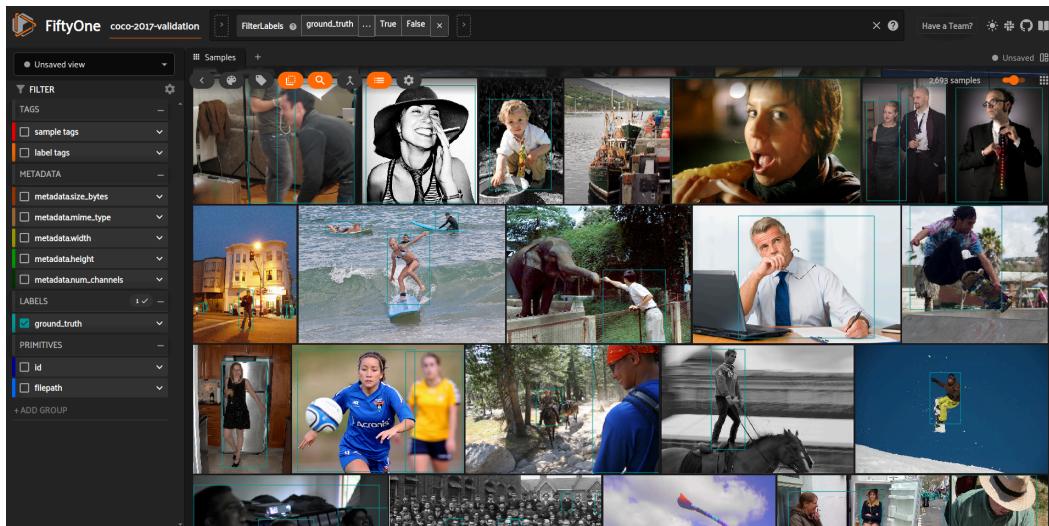


Figure 16 : FIFTYONE's user interface is convenient and well-designed for displaying the dataset.

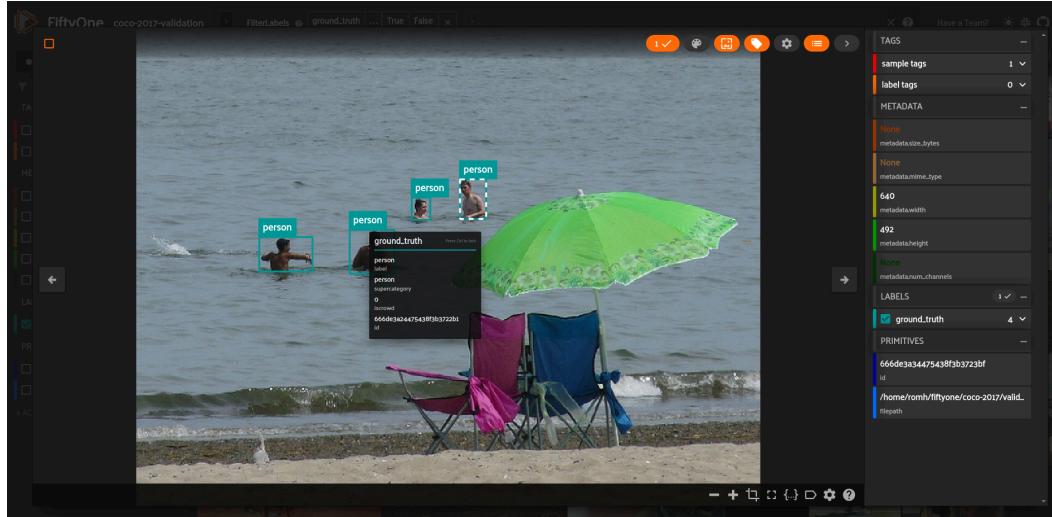


Figure 17 :In the figure, you can easily see the data, including locations and labels, with the option to edit the information.

3.8 YOLOv8 Model Evaluation Results on COCO Validation Dataset

This test's objective is to evaluate the YOLOv8 model's performance in detecting the "person" category on the COCO validation dataset. The evaluation focuses on comparing the predicted number of persons in each image to the actual (ground truth) number of persons and calculating key metrics such as Average Precision (AP), Average Recall (AR), and Mean Average Precision (mAP) .

3.8.1 Dataset

COCO 2017 Validation Dataset: The dataset used for this evaluation is the COCO 2017 validation set, which can be found at COCO 2017 Dataset.

Number of Images: The dataset contains 2,693 images that include the "person" category.

3.8.2 Results

Average Precision (AP)

AP [IoU=0.50:0.95 | area=all | maxDets=100] = 0.360

This metric indicates the mean average precision across multiple IoU thresholds (0.50 to 0.95) for all object sizes. A value of 0.360 suggests moderate precision in detecting persons.

AP [IoU=0.50 | area=all | maxDets=100] = 0.466

This represents the precision at an IoU threshold of 0.50. The higher value (0.466) indicates better performance at lower IoU thresholds, showing the model's ability to detect objects with lower overlap accuracy.

AP [IoU=0.75 | area=all | maxDets=100] = 0.401

This indicates precision at an IoU threshold of 0.75. The value (0.401) shows that the model maintains decent precision even at higher IoU thresholds.

AP for Different Object Sizes:

For small objects, the AP = 0.173, indicating lower precision and suggesting difficulty in accurately detecting smaller persons. For medium-sized objects, the AP = 0.446, reflecting higher precision and better detection capability for moderately sized persons. The highest precision, AP = 0.618, is observed for large objects, demonstrating the model's strength in detecting larger persons with distinct features.

Average Recall (AR)

AR [IoU=0.50:0.95 | area=all | maxDets=1] = 0.288 indicates the recall when considering only the top detection per image, reflecting the ability to detect at least one person correctly in most images.

AR [IoU=0.50:0.95 | area=all | maxDets=10] = 0.394 shows recall for the top 10 detections per image, indicating improved recall with more detections allowed.

AR [IoU=0.50:0.95 | area=all | maxDets=100] = 0.397 represents recall for up to 100 detections per image, highlighting the model's overall recall capacity.

AR for Different Object Sizes:

For small objects, AR = 0.183, indicating lower recall and aligning with the observed lower precision, which suggests difficulty in detecting smaller persons. For medium-sized objects, AR = 0.485, reflecting higher recall and better detection capability. The highest recall, AR = 0.672, is observed for large objects, demonstrating the model's effectiveness in detecting larger persons.

Mean Average Precision (mAP)

mAP = 0.360

The overall mean average precision across all IoU thresholds and object sizes is 0.360. This value reflects a respectable performance but indicates room for improvement, particularly for smaller objects. An mAP of 0.360 is respectable for initial evaluations but suggests there is room for improvement, especially in detecting smaller objects.

3.8.3 Insights

The model performs better on larger objects, showing higher precision and recall since larger objects usually have more distinguishable features. The overall mAP is 0.360, which is respectable for initial evaluations but indicates there is room for improvement, particularly in detecting smaller objects. The precision and recall values are moderate, suggesting a balanced model that could still benefit from further fine-tuning and optimization. Additionally, the average delta between predicted persons and ground truth is 1.2228 which is pretty good .

```

Run  main.py

Average Precision (AP) @[ IoU=0.50:0.95 | area=   large | maxDets=100 ] = 0.018
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all  | maxDets=  1 ] = 0.288
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all  | maxDets=10 ] = 0.394
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all  | maxDets=100 ] = 0.397
Average Recall    (AR) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.183
Average Recall    (AR) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.485
Average Recall    (AR) @[ IoU=0.50:0.95 | area=  large | maxDets=100 ] = 0.672
mAP: 0.3604739901300034
Average delta between predicted persons and ground truth: 1.2227998514667657
Process finished with exit code 0

```

Figure 18 : result of test the Yolo model prediction on the Coco-Validation-2017 data set

```

model = Load object detection model
dataset = Load dataset annotations // Load the Dataset
filtered_images = Filter images containing target category
ground_truth_counts = Initialize storage // Initialize Result Storage
predicted_counts = Initialize storage
image_ids = Initialize storage
deltas = Initialize storage
results = Initialize storage for formatted predictions

FOR each image IN filtered_images DO // Iterate Over Filtered Images
    image_data = Load image data
    // Retrieve ground truth annotations
    ground_truth_annotations = Get ground truth annotations for target category
    ground_truth_count = Count ground_truth_annotations
    Store ground_truth_count
    Store image identifier
    // Run the model to predict objects
    predictions = model.predict(image_data)
    // Count predicted objects
    predicted_count = Count predictions for target category
    Store predicted_count
    // Calculate delta
    delta = ABS(predicted_count - ground_truth_count)
    Store delta
    // Format predictions for evaluation
    FOR each prediction IN predictions DO
        IF prediction meets confidence threshold THEN
            formatted_prediction = Format prediction
            Store formatted_prediction in results
        END IF
    END FOR
END FOR

Write results to file //// Save Predictions
// Load Ground Truth and Predictions for Evaluation
ground_truth = Load ground truth annotations
predictions = Load predictions from file
// Initialize Evaluation Tool
evaluation_tool = Configure evaluation tool with ground truth and predictions
// Run Evaluation
evaluation_tool.evaluate()
evaluation_tool.accumulate()
evaluation_tool.summarize()

```

Figure 19 :Yolo model prediction on the Coco-Validation-2017 data set pseudo code

3.9 Person Detection in Videos for Missing Persons

3.9.1 Locating missing people

Person identification in videos might be instrumental to finding missing persons and resolving cases of abduction as well. Authorities can go through large quantities of video recordings fast and accurately by employing sophisticated technologies such as facial recognition and object detection.

3.9.2 Case Study: The Nova Incident and Technological Utilization

We spoke with “Yam Sahar” who told us about his struggle trying to find his girlfriend’s face among the many Telegram videos, which were taken during the Nova incident on 07.10. This backbreaking task highlights the necessity for automated systems capable of recognizing individuals in crowded or chaotic environments that speed up searches while increasing chances for success in terms of resolution rate. These technologies have the ability analyze different sources’ features against each other hence making it easier to locate missing persons in real time as well providing critical leads towards kidnapping investigation.

4 Work Artifacts

4.1 System Flow

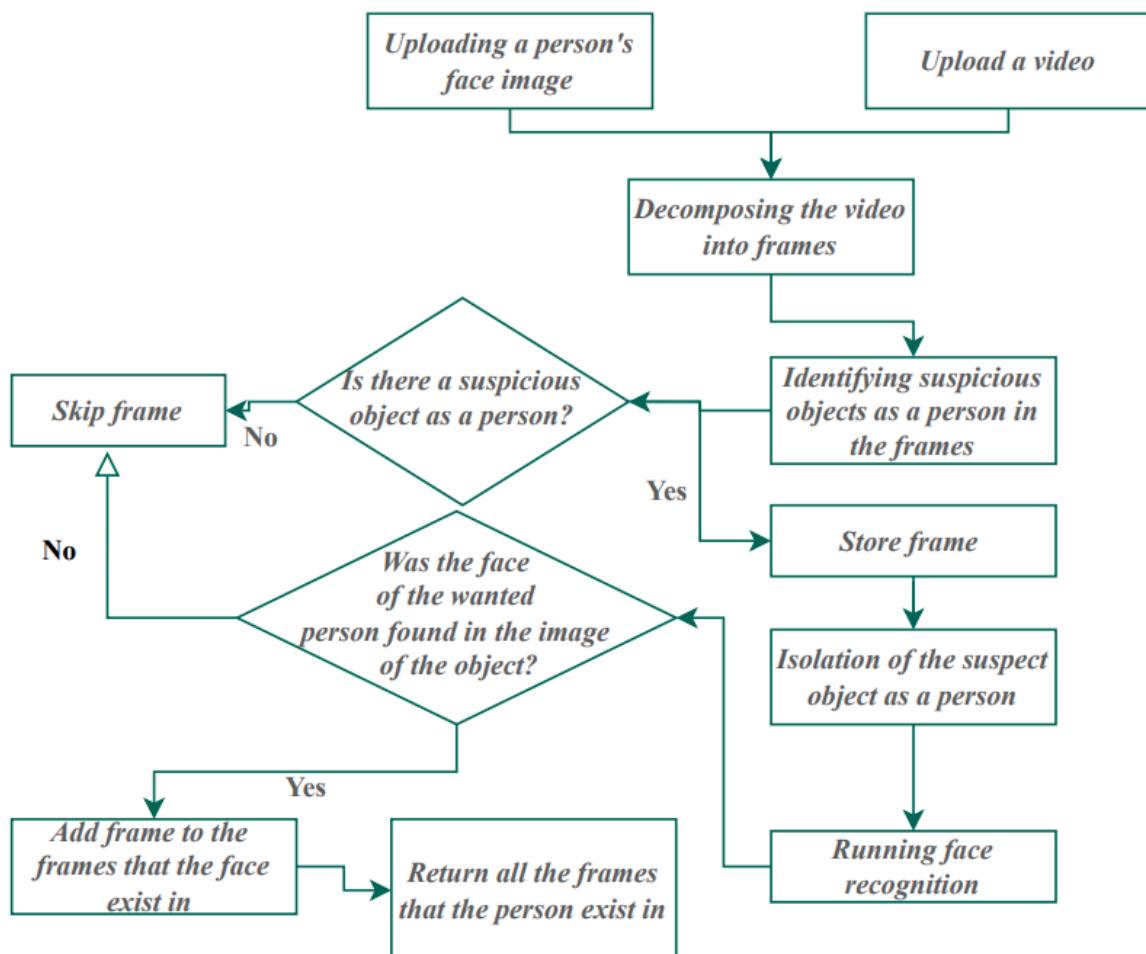


Figure 20 : system flow

Figure 20 illustrates a process for identifying a person's face within a video using our system. Initially, the user uploads an image of the person's face and the target video. The video is decomposed into individual frames, which are analyzed to identify objects that might be people. Each suspicious object is further examined and isolated. Face recognition is then performed on these isolated objects to determine if they match the uploaded face image. Frames where no suspicious objects or matching faces are found are skipped. If a match is detected, the frame is added to a collection of frames containing the person's face. Ultimately, the system returns all frames where the person's face appears.

4.2 Architecture

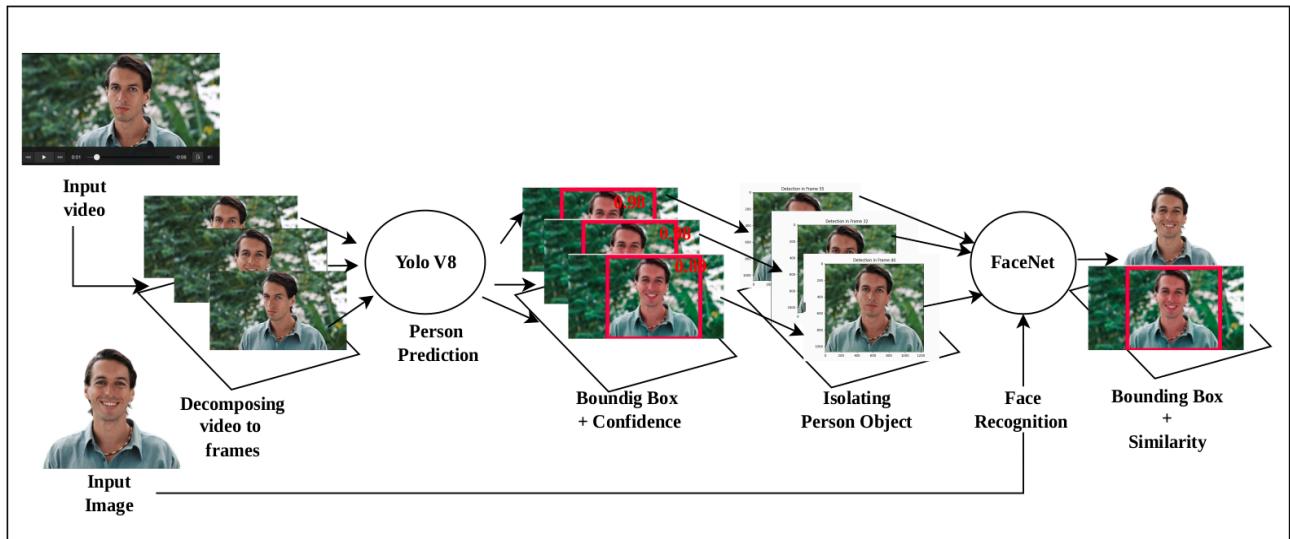


Figure 21 : Our architecture diagram

Figure 21 outlines our working process YOLOv8 for person detection and FaceNet for face recognition. The process starts with the input video to be segmented into the frames that will be processed independently. An input image of the person to be identified is also given. Every frame from the video is passed through the YOLOv8 model that performs the person detection .The model providing the coordinates of the bounding boxes for the detected persons along with the confidence scores. These bounding boxes assist in determining the positions of the detected persons in the different frames. Bounding boxes and confidence scores are then utilized to extract the person's objects, excluding low-quality detections. The isolated person objects are then ready for face recognition, and they are inputted into the FaceNet model. FaceNet recognizes faces by comparing the faces in the isolated person objects to the input image and returns the bounding boxes around the detected faces and similarity score to the input image. The final output generally displays bounding boxes around detected persons and draws their similarity scores as a way of labeling and locating the person correctly in the video frames.

4.2 FR & NFR Requirements

4.3.1 Functional Requirements

No.	Requirement
1	The system shall provide a user interface tool for users.

2	The system shall provide upload functionality for users.
2.1	The system shall allow users to upload videos.
2.2	The system shall allow users to upload a picture.
3	The system shall decompose the uploaded video into individual frames.
3.1	The system shall perform initial preprocessing on the frames.
4	The system shall utilize object detection mechanism.
5	The system shall isolate the detected person objects in the frames.
6	The system shall perform facial recognition
7	The system shall identify frames where the person of interest appears.
7.1	The system shall save the original frame and timestamp when a match is found.
8	The system shall display the frames and timestamps where the person was detected.
8.1	The system shall provide an option to download the results.
9	The system shall provide error messages for issues.

Table 4. functional requirements

4.3.2 Non-functional Requirements

No.	Requirement	Type
1	The system should ensure facial recognition is performed quickly to maintain a responsive user experience.	Performance
1.1	The system will avoid performing unnecessary calculations in frames where there is no expected match for the person we are looking for.	Performance
2	The system must handle large video files concurrently without performance degradation.	Scalability
3	The system should ensure high accuracy in person detection to minimize false positives and negatives.	Accuracy
3.1	The system should ensure high accuracy in facial recognition to minimize false positives and negatives.	Accuracy

4	The system should provide a user-friendly interface that is intuitive and easy to navigate.	Usability
4.1	The system should be accessible to users with varying levels of technical expertise.	Usability
4.2	The system will have three screens contain : data entry screen , identification process progress screen , results screen.	Usability
5	The system should be reliable and available, with minimal downtime.	Reliability
5.1	The system should implement robust error handling and recovery mechanisms.	Reliability
6	The system should be easy to maintain and update, with clear documentation and a modular code structure.	Maintainability
7	The system should support a wide range of video formats and resolutions.	Compatibility
7.1	The system should ensure compatibility with various web browsers and devices.	Compatibility
8	The system will enable facial recognition by advanced and innovative algorithms.	Technology
9	The system will enable the identification of a human figure by YOLO.	Technology
10	The system will optimize the inputs that enter it to match the models.	Optimization

Table 5. non-functional requirements

4.4 Working Progress



Figure 22 : work flow during semester

4.5 System Use Case

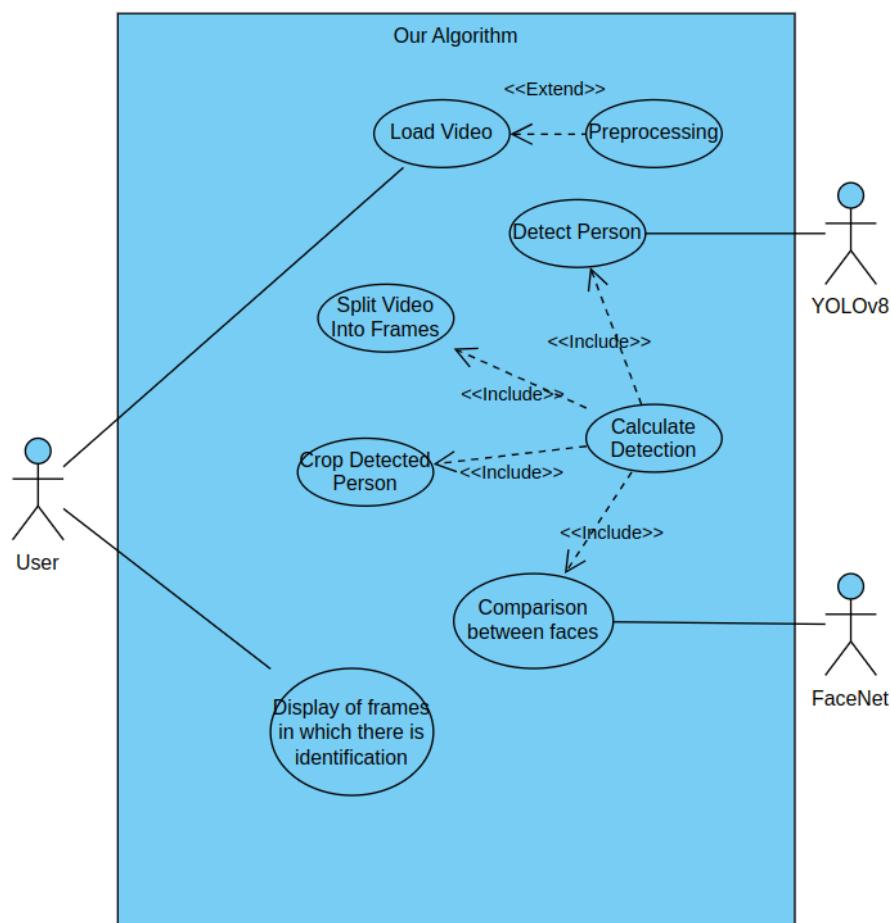


Figure 23 : system use case

4.6 System Activity

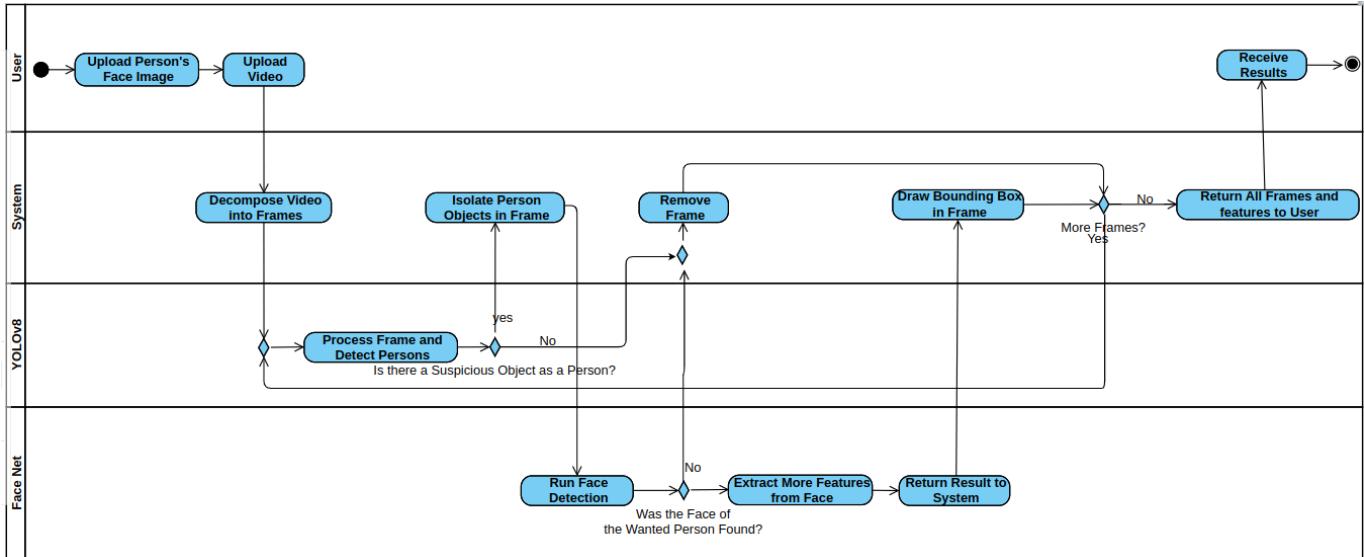


Figure 24 : system activity diagram

The activity diagram outlines the process of a system designed for face recognition in videos. Initially, the user uploads an image of the person's face and the video in which they want to locate this face. The system decomposes the video into individual frames to facilitate frame-by-frame analysis. Each frame is processed using the YOLOv8 algorithm to detect any persons within the frame. If a suspicious object resembling a person is detected, it is further analyzed using the FaceNet algorithm to verify if it matches the uploaded face image. If a match is found, additional features are extracted from the face for accurate identification and analysis, and the frame is added to a collection of frames where the person appears. This process is repeated for each frame until the entire video is processed. After analyzing all frames, the system compiles and returns all frames containing the identified face, along with the extracted features, to the user. The system may also draw bounding boxes around the identified faces for better visualization. This ensures an efficient and comprehensive method to locate and highlight a specific person's face within a video.

4.7 User Interface

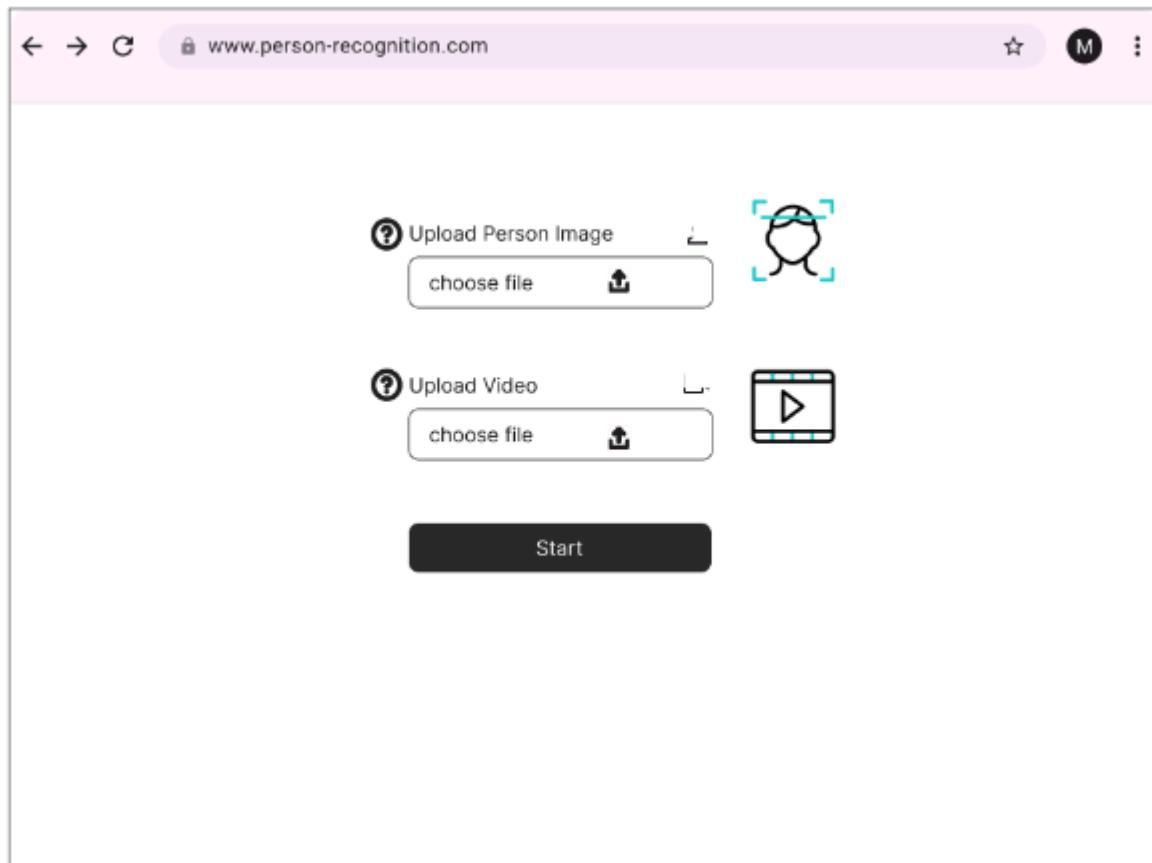


Figure 25 : Screen 1 ,Loading Data Screen

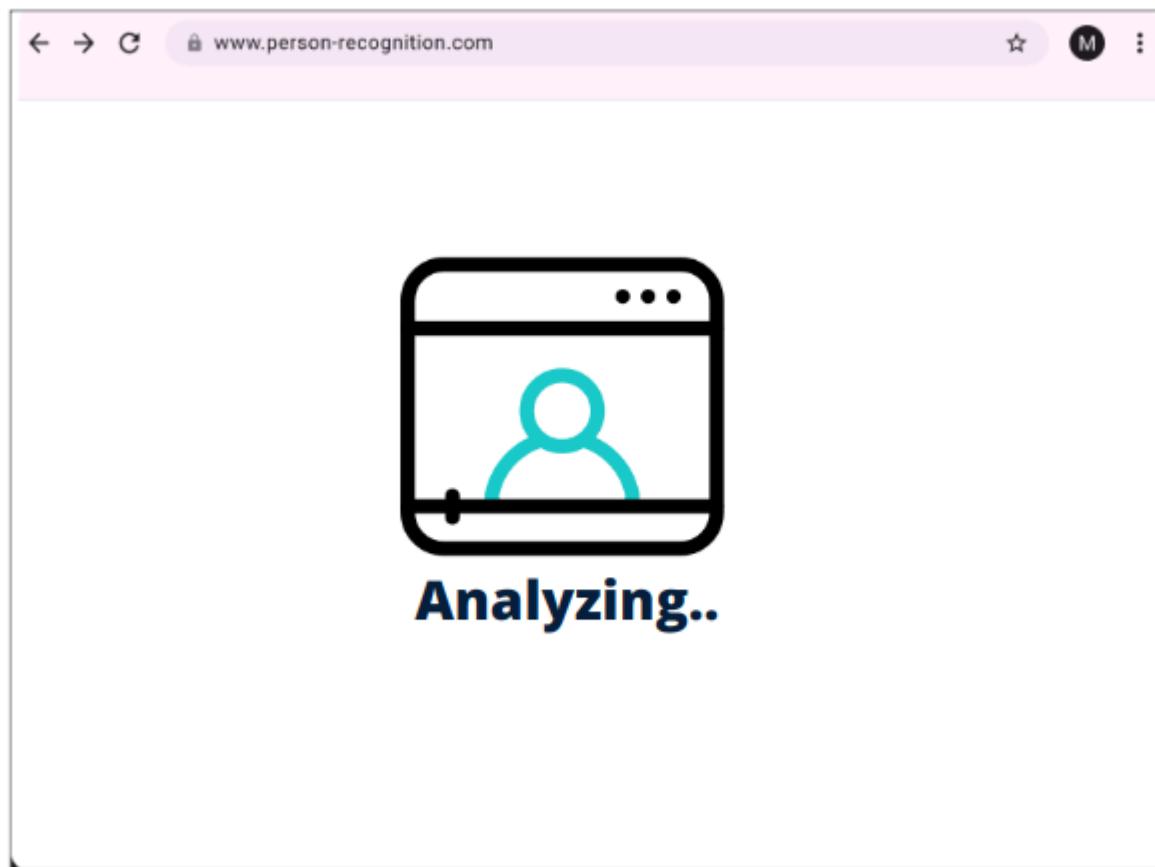


Figure 26 : Screen 2 , user see that the system is analyzing the results

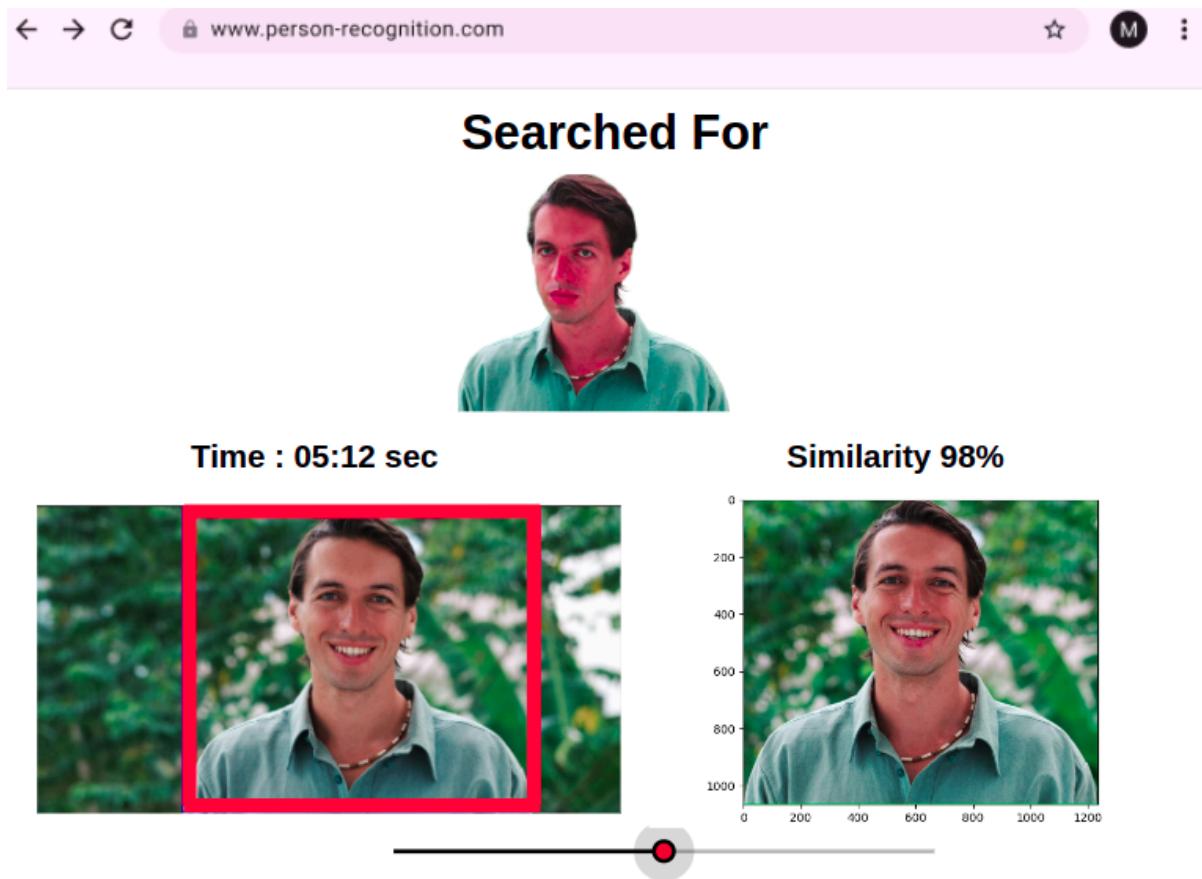


Figure 27 : Screen 3 ,result screen , user can see the person he searched for and our algorithm detection in the original supplied video.

5 Testing Plan

5.1 Scope

Comprehensive testing plan will encompass both functional and non-functional testing of the Video-Based Person Identification System. This plan ensures the system meets all specified requirements, performs as expected, and verifies the accuracy of person detection and facial recognition, along with user acceptance.

5.2 Objectives

- Ensure all functional and nonfunctional requirements are met.
- Verify system performance and responsiveness.
- Confirm accuracy in person detection and facial recognition.
- Validate usability and user experience.
- Checking that the results are clear and understandable for system users .

5.3 Test Cases

No.	Test Name	Description	Procedure	Expected Result
1	Video_With_Other_Person	Test with video containing other persons but not the target person	1. Upload video without target person 2. Upload target picture 3. Start detection and recognition	System indicates no match found
2	Video_With_No_Person	Test with video containing no persons	1. Upload video without any persons 2. Upload target picture 3. Start detection and recognition	System indicates no persons detected
3	Video_With_Many_Persons	Test with video containing many persons	1. Upload video with many persons 2. Upload target picture 3. Start detection and recognition	System accurately detects and identifies target person
4	Low_Quality_Video	Test with low-quality video	1. Upload low-quality video 2. Upload target picture 3. Start detection and recognition	System performs detection and recognition with decreased accuracy
5	Low_Quality_Image	Test with low-quality image of the target person	1. Upload high-quality video 2. Upload low-quality picture 3. Start detection and recognition	System performs detection and recognition with decreased accuracy
6	Short_Length_Video	Test with short video	1. Upload short video 2. Upload target picture 3. Start detection and recognition	System performs detection and recognition accurately

7	Long_Length_Video	Test with long video	<ul style="list-style-type: none"> 1. Upload long video 2. Upload target picture 3. Start detection and recognition 	System performs detection and recognition accurately and efficiently
8	Video_With_Low_Light	Test with video in low light conditions	<ul style="list-style-type: none"> 1. Upload low light video 2. Upload target picture 3. Start detection and recognition 	System performs detection and recognition with decreased accuracy
9	Different_Angles	Test with video showing person at different angles	<ul style="list-style-type: none"> 1. Upload video with person at different angles 2. Upload target picture 3. Start detection and recognition 	System performs detection and recognition accurately
10	Overall_Functionality	Validate overall system functionality	<ul style="list-style-type: none"> 1. End users perform typical tasks 2. Report any issues 	System meets user expectations and needs
11	System_Usability	Validate system usability	<ul style="list-style-type: none"> 1. End users navigate the system 2. Provide feedback on interface and experience 	System is user-friendly and intuitive
12	System_Performance	Validate system performance	<ul style="list-style-type: none"> 1. End users upload various video sizes 2. Observe system response 	System performs well under typical user conditions
13	YOLO_Person_Detection	Test accuracy of person detection	<ul style="list-style-type: none"> 1. Compare detected persons with ground truth data 	High accuracy in person detection
14	Facial_Recognition	Test accuracy of facial recognition	<ul style="list-style-type: none"> 1. Initialize test on large dataset for face recognition 2. Compare recognized faces with ground truth data 3. Calculate recall measure 	High accuracy and recall in measure in facial recognition

15	YOLO_Person_Detection	Test recall measure of person detection	<ol style="list-style-type: none"> 1. Initialize test on large dataset for persons photo . 2. Compare detected persons with ground truth data 3. Calculate recall measure 	High accuracy and recall measure in person detection
17	Video_Image_Upload	Test user interface for video and image upload	<ol style="list-style-type: none"> 1. Open the application 2. Navigate to upload screen 3. Upload a video 4. Upload a picture 	Video and picture upload successfully
18	Decompose_Frames	Test video decomposition into frames	<ol style="list-style-type: none"> 1. Upload a video 2. Start processing 	Video is decomposed into individual frames
19	Detected_Person_Isolation	Test isolation of detected persons	<ol style="list-style-type: none"> 1. Upload a video 2. Start person isolation 	Persons isolated in frames
20	Isolated_Person_Facial_Recognition	Test facial recognition on isolated persons	<ol style="list-style-type: none"> 1. Upload a picture 2. Start facial recognition 	Faces matched with uploaded picture
21	Display_Results	Test results display	<ol style="list-style-type: none"> 1. Complete match identification 2. View results 	Frames and timestamps displayed
22	Unsupported_Formats	Test error handling for unsupported formats	<ol style="list-style-type: none"> 1. Upload unsupported video format 2. Upload unsupported picture format 	Meaningful error messages displayed
23	System_Performance_Under_Load	Test system performance under load	<ol style="list-style-type: none"> 1. Upload large video file 2. Perform all processing steps 	System performs efficiently
24	Full_Process_Validation	Validate full process from video and person upload to results	<ol style="list-style-type: none"> 1. Upload a video 2. Upload a picture 3. Start detection and recognition 4. View results 	System returns all frames where the person exists in the video

Table 6. test cases of our application

6 Expected Accomplishments

- **Better accuracy in recognizing individuals:** It can be used to drastically improve the precision of person identification from video content but especially in situations that are dynamic and complex.
- **Detection capability in real-time:** With the YOLO algorithm, this system will enable fast detection of people which is important for security or emergency response applications.
- **Strong face recognition:** The system will ensure that different faces are accurately recognized under various angles by utilizing FaceNet thus making it possible for correct identification to be done even under different circumstances.
- **Powerful analysis tool:** When advanced detection and recognition technologies are integrated, it becomes a useful instrument for comprehensive video analysis thus benefiting media research as well as security among others where investigations may involve videos.
- **User friendly interface:** An intuitive and accessible interface should be designed so that specialized teams can use this software effectively without undergoing too much training.

6.1 Problems

- **Managing occlusions and partial visibility:** One major problem in person identification is dealing with cases where individuals might not be fully visible or could be hidden behind other objects.
- **Achieving real time performance:** A lot of computational power is needed for systems like these but they still need high speeds hence optimization must always balance between these two aspects; efficiency vs accuracy.
- **Cross domain generalization:** Generalizing models trained on one dataset to work well with another data-set or real world scenarios has proven difficult yet it's crucial if we want reliable performance across all domains.

6.2 Criteria for Success

- **Accuracy of Identification:** The precision level should not fall below 90% when it comes to this software verifying people through videos taken within any given environment as supported by routine tests and validation.

- **Operational impact:** 30% reduction in response time and decision-making efficiency during security and emergency situations should be realized.
- **User feedback and satisfaction:** Structured feedback from specialized teams should be able to score at least 80% satisfaction rate thus indicating its usefulness among them.
- **System reliability & performance:** The system must work with a minimum uptime of 99%, maintaining good performance stability even under high load scenarios or when subjected to different conditions.
- **Adoption & Engagement:** The system is considered effective if it can attract adoption by up-to-half the target user groups within six months following its deployment.

7 References

- [1] Adrián, S. C., David, F. J., & Cristina, L. G. (2022). Real-time human action recognition using raw depth video-based recurrent neural networks.
- [2] Ahmed, R., & Abd-Elkawy, E. H. (2024). Improved Tomato Disease Detection with YOLOv5 and YOLOv8. *Engineering, Technology & Applied Science Research*, 14(3), 13922-13928.
- [3] Bouazizi, O., Azroumahli, C., El Mourabit, A., & Oussouaddi, M. (2024). Road Object Detection using SSD-MobileNet Algorithm: Case Study for Real-Time ADAS Applications. *Journal of Robotics and Control (JRC)*, 5(2), 551-560.
- [4] Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., & Darrell, T. (2015). Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2625-2634).
- [5] Fontes, C., Hohma, E., Corrigan, C. C., & Lütge, C. (2022). AI-powered public surveillance systems: why we (might) need them and how we want them. *Technology in Society*, 71, 102137.
- [6] Fu, Y., Wang, X., Wei, Y., & Huang, T. (2019, July). Sta: Spatial-temporal attention for large-scale video-based person re-identification. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 33, No. 01, pp. 8287-8294).

- [7] Ge, Z., Liu, S., Wang, F., Li, Z., & Sun, J. (2021). Yolox: Exceeding yolo series in 2021. arXiv preprint arXiv:2107.08430.
- [8] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ... & Adam, H. (2017). Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861.
- [9] Kanan, C., & Cottrell, G. W. (2012). Color-to-grayscale: does the method matter in image recognition?. *PLoS one*, 7(1), e29740.
- [10] Kandel, I., & Castelli, M. (2020). How deeply to fine-tune a convolutional neural network: a case study using a histopathology dataset. *Applied Sciences*, 10(10), 3359.
- [11] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.
- [12] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *nature*, 521(7553), 436-444
- [13] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *nature*, 521(7553), 436-444.
- [14] Lin, T. Y., Dollár, P., Girshick, R., He, K., Hariharan, B., & Belongie, S. (2017). Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2117-2125).
- [15] Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13* (pp. 740-755). Springer International Publishing.
- [16] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C. Y., & Berg, A. C. (2016). Ssd: Single shot multibox detector. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14* (pp. 21-37). Springer International Publishing.
- [17] Nguyen, V. T., & Chu, D. T. (2023). Study on Tracking Real-Time Target Human Using Deep Learning for High Accuracy. *Journal of Robotics*, 2023(1), 9446956.
- [18] Parkhi, O., Vedaldi, A., & Zisserman, A. (2015). Deep face recognition. In *BMVC 2015-Proceedings of the British Machine Vision Conference 2015*. British Machine Vision Association.

- [19] Cao, Q., Shen, L., Xie, W., Parkhi, O. M., & Zisserman, A. (2018, May). Vggface2: A dataset for recognising faces across pose and age. In 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018) (pp. 67-74). IEEE.
- [20] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 779-788).
- [21] Ren, S., He, K., Girshick, R., & Sun, J. (2016). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6), 1137-1149..
- [22] Roushdy, M. (2006). Comparative study of edge detection algorithms applying on the grayscale noisy image using morphological filter. *GVIP journal*, 6(4), 17-23.
- [23] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L. C. (2018). Mobilenetv2: Inverted residuals and linear bottlenecks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 4510-4520).
- [24] Schroff, F., Kalenichenko, D., & Philbin, J. (2015). Facenet: A unified embedding for face recognition and clustering. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 815-823).
- [25] Terven, J., Córdova-Esparza, D. M., & Romero-González, J. A. (2023). A comprehensive review of yolo architectures in computer vision: From yolov1 to yolov8 and yolo-nas. *Machine Learning and Knowledge Extraction*, 5(4), 1680-1716.
- [26] Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). Learning spatiotemporal features with 3d convolutional networks. In Proceedings of the IEEE international conference on computer vision (pp. 4489-4497).
- [27] Wang, J., & Lee, S. (2021). Data augmentation methods applying grayscale images for convolutional neural networks in machine vision. *Applied Sciences*, 11(15), 6721.
- [28] Wang, M., & Deng, W. (2021). Deep face recognition: A survey. *Neurocomputing*, 429, 215-244.
- [29] Wang, P., Li, W., Gao, Z., Tang, C., & Ogunbona, P. O. (2018). Depth pooling based large-scale 3-d action recognition with convolutional neural networks. *IEEE Transactions on Multimedia*, 20(5), 1051-1061.

- [30] Xiao, Y., Chen, J., Wang, Y., Cao, Z., Zhou, J. T., & Bai, X. (2019). Action recognition for depth video using multi-view dynamic images. *Information Sciences*, 480, 287-304.
- [31] Xu, C., Liu, Y., Chen, D., & Yang, Y. (2022). Direct training via backpropagation for ultra-low-latency spiking neural networks with multi-threshold. *Symmetry*, 14(9), 1933.
- [32] Zagitov, A., Chebotareva, E., Toschev, A., & Magid, E. (2024). Comparative analysis of neural network models performance on low-power devices for a real-time object detection task. *Computer*, 48(2).
- [33] Zhao, Z. Q., Zheng, P., Xu, S. T., & Wu, X. (2019). Object detection with deep learning: A review. *IEEE transactions on neural networks and learning systems*, 30(11), 3212-3232.
- [34] Zheng, K. (2019). Person Identification with Convolutional Neural Networks (Doctoral dissertation, University of South Carolina).
- [35] Zheng, Z., Wang, P., Liu, W., Li, J., Ye, R., & Ren, D. (2020, April). Distance-IoU loss: Faster and better learning for bounding box regression. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 34, No. 07, pp. 12993-13000).
- [36] YOLO Pseudocode: <https://gist.github.com/sjain07/7b880ed494467bd44e8c4f051703a011>
- [37] Ultralytics YOLO Docs : <https://docs.ultralytics.com/models/yolov8/>
- [38] YOLO v8 vs. MobileNet SSD v2 <https://roboflow.com/compare/yolov8-vs-mobilenet-ssd-v2>
- [39] IOU EXPLAINED <https://www.v7labs.com/blog/intersection-over-union-guide>
- [40] Cocodataset official <https://cocodataset.org/#home>
- [41] pycocotools API - <https://pypi.org/project/pycocotools/>
- [42] FiftyOne Integrations COCO Integration <https://docs.voxel51.com/integrations/coco.html>
- [43] Amazon Rekognition <https://aws.amazon.com/rekognition/>
- [44] Microsoft Azure Face API
<https://learn.microsoft.com/en-us/azure/ai-services/computer-vision/overview-identity>

[45] Face ++ <https://www.faceplusplus.com/>

[46] Deepvisionai <https://deepvisionai.in/>

[47] Chat GPT <https://chatgpt.com/>

[48] React <https://react.dev/>

[49] Python <https://www.python.org/>

8 AI Prompts

During the research phase of our project, we utilized ChatGPT[47] intermittently to validate our proposed technology stack, check grammar, and brainstorm additional features for our application. Here are some example prompts and explanations of how their responses guided our research and development decisions.

- Create for me a short literature about Current uses of object detection remember to include references
<https://chatgpt.com/share/efabf011-4832-4da0-8e70-0813a2d431c1>
- Summarize Dong, Z., Zhou, Z., Li, Z., Liu, C., Huang, P., Liu, L., ... & Kang, J. (2018). Convolutional neural networks based on RRAM devices for image recognition and online learning tasks. *IEEE Transactions on Electron Devices*, 66(1), 793-801.
<https://chatgpt.com/share/757eb2c5-64e5-4b74-918a-7fe584247374>
- Create for me a two page literature review about Siddiqi, M. H., Ali, R., Rana, M. S., Hong, E. K., Kim, E. S., & Lee, S. (2014). Video-based human activity recognition using multilevel wavelet decomposition and stepwise linear discriminant analysis. *Sensors*, 14(4), 6370-6392.
<https://chatgpt.com/share/9c4e34d7-e73b-4a12-92d9-8534220824e4>
- Create for me a two page literature review about yolo based on Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 779-788).
<https://chatgpt.com/share/3f9f7a91-d256-4e71-9666-dc23feb94384>
- Sort the article references by A-Z
<https://chatgpt.com/share/24a67429-d362-4e43-92e0-5fde4f191f70>