

Eyes on the City: Leveraging Multi-View Face Detection for Proactive Security in Smart Cities

1st Abdelrhman Yakout

Media engineering and technology

The German University in Cairo

Cairo, Egypt

abdelrhman816atef@gmail.com

2nd Mohammed Salem

Digital media engineering and technology

The German University in Cairo

Cairo, Egypt

mameged2@gmail.com

3rd Yomna Alayary

Digital media engineering and technology

The German University in Cairo

Cairo, Egypt

yomna.alayary@gmail.com

Abstract—With the huge number of video surveillance cameras to this day, it necessitates an efficient system for monitoring and security. This thesis explores the application of real-time face detection and recognition for identifying individuals in facilities and age, gender and race (AGR) detection for data analysis in video surveillance using deep learning techniques and convolutional neural networks(CNN).

Our research aims to create a robust model capable of operating under various conditions such as different lighting conditions, angles and occlusions. When experimenting with 3 different face detection models on surveillance footage, the pre-trained ResNet model. Additionally, trained a YOLOV9 model on WIDER FACE dataset with 30 epochs due to certain limitations, the model had an average precision of 92%. Finally, RetinaFace model was used in the final model.

The model created has 3 stages that each frame pass through: face detection, AGR detection and face recognition. After testing the model on different surveillance footage, the result shown had limited accuracy due to the struggle of the AGR and face recognition model on low resolution faces in the footage as all facial features has to be clearly visible. Otherwise these conditions the AGR and face recognition performed successfully with no issues. Additionally RetinaFace face detection model has an average precision of 96.1%.

I. INTRODUCTION

Due to their numerous applications facial recognition and detection technologies have become more and more popular in recent years. These technologies are very important to a wide range of sectors. Consider the purposes of security and surveillance which include monitoring access control systems and conducting theft investigations. These activities require the identification of individuals present in public spaces. Contributing to law enforcements criminal identification and investigation procedures. Enabling biometric authentication for secure access control systems in a variety of contexts including financial institutions personal devices and border control. Personalizing user experiences in marketing entertainment and social media domains as well.

Face detection and recognition technologies in addition to enhancing convenience and security hold a great promise for enhancing social welfare and public safety. These technologies are particularly useful in helping vulnerable populations such as children and the elderly locate missing individuals. In addition to helping identify accident or disaster victims facial recognition can also speed up emergency response times

and facilitate family reunions. Additionally the possibility of being recognized can discourage criminal activity making the environment safer for everybody. Moreover facial recognition technology can be a useful instrument for expediting emergency response activities. First responders can deliver critical aid more efficiently and possibly save lives by enabling faster and more accurate identification of individuals in need.

Even with these advances there are still a number of issues with face recognition especially when there are multiple views involved. Pose variations make it challenging for conventional methods to reliably detect and recognize individuals because faces are captured from different angles. Illumination variations can also have a substantial impact on how facial features appear making recognition accuracy more difficult. The recognition process can also be made more difficult by partial or total occlusions brought on by hair accessories or other objects. Finally practical applications frequently demand real-time processing capabilities in order to produce effective and instantaneous results. Convolutional neural networks (CNN) a type of deep learning has been suggested as a solution to these problems. The field of face detection and recognition has undergone a revolution. These effective methods have many benefits. High-level feature extraction: in this process facial image data is automatically processed by CNNs to extract robust complex features—even in the face of difficult obstacles like occlusions and pose changes. Increased recognition accuracy because deep learning models outperform conventional techniques in face recognition tasks on a regular basis. The real-time processing potential of CNN-based solutions has been made possible by improvements in hardware and optimized algorithms which qualify them for useful applications. With the use of CNN-powered deep learning and real-time processing this project seeks to create a multi-view face detection and recognition system. This system uses deep learning techniques to improve accuracy and robustness and enable real-time processing capabilities thereby addressing the aforementioned challenges.

II. RELATED WORK

A. Face detection

A collection of articles explores the crucial work of precisely identifying faces in pictures and videos or face de-

tection. An overview of the methods investigated is provided below:

Kumar et.al. [1] and Zhang et.al. [2] provide insightful information about the state of face detection research today. They give thorough summaries of current methods contrasting and comparing their advantages and disadvantages.

Maximum efficiency and accuracy are given equal priority in the Deng et al. [3] model (RetinaFace). Compared to some two-stage detection techniques it is faster because it uses a single-stage neural network architecture. In order to capture faces at different sizes within an image RetinaFace also uses a unique feature pyramid network. Having an average precision of 96.1% on WIDER FACE dataset. The TinaFace model by Zhu et al. [4] on the other hand places more emphasis on a lightweight model appropriate for devices with limited memory and processing power. Having an average precision of 96.3% on WIDER FACE dataset. Compared to RetinaFace some accuracy may be lost as a result. A Dual Shot Face Detector (DSFD) network is presented by Li et al. [5] that strikes a compromise between excellent accuracy and quickness. With a new Feature Enhancement Module (FEM) to produce richer feature representations and a Progressive Anchor Loss (PAL) to enhance learning across different face scales DSFD employs a two-stage methodology. Furthermore the appropriateness of DSFD for real-time face detection tasks is guaranteed by its lightweight components and optimized network architecture. Having an average precision of 95.0% on WIDER FACE dataset.

RetinaFace utilizes a pyramid feature, which is a multi-level pyramid. First features are extracted from an image and are then fed into the model. Encoding various details at different scales. Larger faces are recognized at the lower level of the pyramid, while smaller faces are recognized at higher levels of the pyramid.

A multi-stage face detection framework called MTCNN is presented by Zhang et al. [6]. To attain great accuracy MTCNN makes use of three convolutional neural networks (CNNs) each of increasing complexity. Candidate bounding boxes are suggested in the first stage refined in the second and facial landmarks such as the mouth nose and eye key points predicted in the third. However FaceBoxes is a single-stage face detection algorithm that is optimized for CPU performance in real-time and is used by Zhang et al. [7] In order to achieve good accuracy and efficiency appropriate for CPU processing FaceBoxes uses a convolutional neural network architecture that has been specially designed for this purpose.

The Chi et al. [8] SRN model uses a two-pronged strategy to address false positives and enhance location accuracy in face detection. Selective Two-step Regression is used to fine-tune possible faces for accurate localization after Selective Two-step Classification removes the majority of negative detections early on. A Receptive Field Enhancement block designed to capture details in difficult poses further improves network performance. With this combination SRN is able to achieve cutting-edge face detection performance. Utilizing

the well-known YOLO object detection framework Qi et al. [9] YOLO5Face modifies it especially for face detection. A face detector architecture that addresses the shortcomings of current techniques is proposed by Liu et al. [10] Three main challenges are addressed: removing false alarms enhancing scale-level data and assigning labels. In order to address these problems MogFace has added three new modules: the Hierarchical Context-Aware Module which lowers false positives the Adaptive Online Incremental Anchor Mining Strategy which improves label assignment and the Selective Scale Enhancement Strategy. Through the resolution of these issues MogFace attains cutting-edge results on multiple face detection benchmarks. By using contextual information Tang et al. [11] address the problem of detecting challenging faces (small blurry occluded). To make use of context it employs Pyramid Anchors in the training process a Context-sensitive Prediction Module ensures precise face location and classification and a Low-level Feature Pyramid Network (LFPN) integrates contextual and facial features. This method enables PyramidBox to recognize faces in difficult real-world situations.

Wang et.al. [12] and Kim et.al. [13] explore the application of deep learning architectures, particularly Fully Convolutional Networks (FCNs), for face detection. FCNs excel at learning spatial features from images, making them well-suited for tasks like identifying and locating faces within an image.

Newer papers in face detection research focuses on the development of models specifically designed for devices with limited resources which was the main focus on Xu et.al. [14] CenterFace model. This is particularly relevant for applications on mobile phones or embedded systems, where computational power and memory usage are critical factors.

B. Face Recognition

This subsection dives into the realm of face recognition, which focuses on identifying individuals from facial images or videos. Here, the papers explore various algorithms and techniques to achieve this task.

By providing thorough surveys of current face recognition systems Kortli et.al. [15] and Parkhi et.al. [16] offer fundamental knowledge. These questionnaires explore various strategies detailing their advantages and disadvantages. You will gain a thorough understanding of the state of the art in face recognition research with this comparative analysis.

GhostFaceNets an effective face recognition model architecture is examined by Alansari et.al. [17]. In order to achieve good recognition accuracy without requiring a lot of computational power this model places a high priority on lightweight design and uses inexpensive operations like depth-wise separable convolutions. using the ResNet-50 backbone model. GhostFaceNets are therefore appropriate for implementation on systems with constrained processing capacity. Having an accuracy of 99.8% on LFW dataset. In addition George et.al. [18] suggest the EdgeFace face recognition network architecture which is effective and lightweight like GhostFaceNets [17]. It does this by fusing the advantages of Transformer models and Convolutional Neural Networks

(CNNs) through a hybrid network architecture that draws inspiration from EdgeNeXt. To further reduce computation in the linear layers without noticeably sacrificing performance a Low Rank Linear (LoRaLin) module is also introduced. Because of this combination EdgeFace can achieve high face recognition accuracy with low computational costs and compact storage which makes it appropriate for deployment on edge devices with limited resources. Having an accuracy of 99.8% on LFW dataset. A Dynamic Class Queue (DCQ) mechanism is also proposed by Li et.al. [19] to address imbalanced class distributions and computational limitations similar to GhostFaceNets [17] and EdgeFace [18]. During each training iteration DCQ reduces computational burden by dynamically selecting a subset of classes for recognition. Furthermore the model is able to handle tail classes with limited training examples because class weights are generated dynamically on the fly. DCQ is able to attain competitive results on extensive facial recognition datasets by using this method despite having restricted computational resources. A model that addresses the performance problems on smaller devices is developed in all three of these papers.

A collection of papers explores loss functions which are an important component in neural network training for face recognition. The goal of these loss functions is to maximize the networks capacity to distinguish between people according to their facial characteristics. Models trained on one data collection method may perform poorly on different evaluation methods this is known as process discrepancy and Kim et.al. [20] DiscFace model addresses this issue. DiscFace improves generalizability by making the model learn features that hold true in various scenarios. The quality-aware mechanism in QMagFace by Terhorst et.al. [21] is integrated into the recognition process. It generates a more robust performance in difficult conditions by estimating the quality of an image and adjusting the recognition confidence score accordingly. For the purpose of training representations of distinct identities existing methods frequently employ a fixed margin penalty. The ElasticFace model by Boutros et.al. [22] suggests a more adaptable strategy. During each training iteration a random margin value is employed which enables the model to acquire a decision boundary that can adjust to fluctuations in the data and could result in enhanced recognition precision. ElasticFace [22] and Deng et.al. [23] ArcFace share the goal of enhancing face recognition models discriminating power. In order to maintain high classification accuracy it uses a particular loss function that promotes a wider margin between the representations of various identities.

Schroff et.al. [24] presented a method for face recognition, it builds a single system rather than learning different models for every task. In a compact space distances signify facial similarity and FaceNet [24] maps face images to this space directly. This makes it possible to carry out common tasks within this new space such as recognition and clustering. This approach is effective because it obtains high accuracy with only a small amount of data (128 bytes) per face.

The complexity of real-world scenarios is addressed by

recent research efforts. The problem of face recognition with masks for example is addressed by Mare et.al. [25] and is becoming more and more problematic as a result of the COVID-19 outbreak. The limitations of occluded facial regions for precise identification are examined in this paper. Furthermore Knoche et.al. [26] concentrates on enhancing face recognition models resistance to image changes. Real-world scenarios may involve capturing faces in varying lighting conditions poses or resolutions. This paper investigates ways to improve the models performance in spite of these changes.

An approach to the problems of face recognition in unrestricted environments is put forth by Liu et.al. [27]. Creating a controllable face synthesis model (CFSM) that can replicate the distribution of target datasets in various settings. The model picks up on a latent space that represents the target datasets stylistic variances. This makes it possible to precisely control the creation of artificial faces with particular traits. The study shows that incorporating this artificial data into a face recognition models training process greatly enhances performance in a range of environmental settings. Differentiating it from every other approach that has been discussed.

C. Discussion

Based on the literature review above in II, there is a gap in the research where they do not experiment on real-time video surveillance systems for both face detection and recognition. By incorporating real-time processing on video surveillance it can be turned into a proactive tool. Enabling immediate detection, could help in identifying the restricted individuals from entering a specified area or even identifying a person from a crowd, who is blacklisted from the environment.

Additionally, implementing an AGR (age, gender and race) detection could be very beneficial for video surveillance systems. Real-time processing of AGR data would allow us to analyze each individual entering the facility, providing a live count of the people inside the facility with how many are males and females and also their age groups and race. Having this data is not only beneficial for keeping track of the number of people in the facility but also for use as a statistical analysis.

III. METHODOLOGY

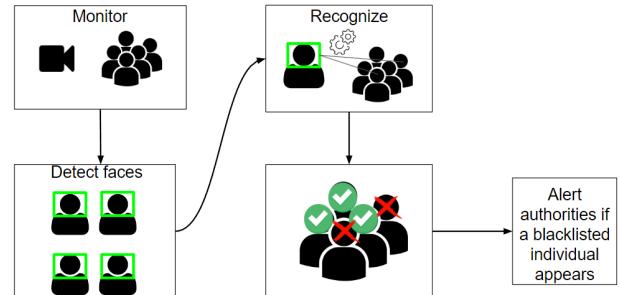


Fig. 1. A visualization of how the model works

Figure 1 shows an overview of how our model would work. Where we would have a video surveillance camera monitoring

a specific facility, always detecting and identifying individuals passing by it. By detecting faces and having a continuous count of people inside the facility and also being able to identify any known blacklisted individuals from the premises, an alert would be sent notifying the authorities of their presence and taking immediate action.

A. Image augmentation and annotations

1) Albumentation tool for image augmentation: A robust Python library called Albumentation was created especially for image augmentation. The process of artificially producing variations of preexisting images is known as image augmentation. The deep learning models benefit greatly from this process because it teaches them to avoid overfitting on the training set and to generalize well. We can create artificial versions of our original face images in this case by using albumentation which covers a greater range of rotations scales lighting and occlusions. Albumentation for instance can be used to modify contrast and brightness to replicate various lighting conditions or to slightly rotate faces to simulate various head positions. These changes expose the model to a wider range of face appearances improving its capacity to recognize faces accurately from different angles in real-world situations.

2) Labelme annotation tool: Labelme is essential to the process of annotating images. It is an open-source free graphical annotation tool that makes labeling images with bounding boxes easier. Labelme in this case enables us to precisely annotate the bounding boxes surrounding every face in the picture regardless of its angle. The ground truth for our model is provided by these annotation data. The augmented image and the bounding box annotations for it are both shown to the model throughout the training process. The model then gains the ability to correlate the features that were extracted from the picture with the existence and positioning of faces. Training our multi-view face detection model requires high-quality labeled datasets which Labelmes user-friendly interface and effective annotation capabilities make possible.

We can create a large and varied training dataset that better prepares our model to handle the challenges of multi-view face detection by combining the strength of Labelme for image annotation and Albumentation for image augmentation. A greater variety of face variations are presented to the model through the augmented images and the carefully labeled annotations supply the ground truth required for efficient training. The development of a reliable and precise multi-view face detection system is made possible by this combined method.

B. Face detection

We examined three face detection techniques and contrasted their respective performances. YOLOv9, a trained ResNet model and the face detection model are well-known models contained in the library of face recognition. Analyzing each models robustness accuracy and efficiency in identifying faces across a range of image datasets was our goal. We sought

to determine which model provided the best face detection abilities by examining the data.

1) YOLO based face detection: We leverage You Only Look Once version 9 (YOLOv9)s innovative features for our real-time multi-view facial detection. Our proposed system is one of the many real-world applications that benefit greatly from this state-of-the-art deep learning object detection models superb balance between accuracy and computational efficiency. That develops. We rely on YOLOv9 and its innovative features for our real-time multi-view face detection. Our suggested system is an excellent option for real-world applications that require real-time performance as this state-of-the-art deep learning object detection model remarkably balances accuracy and computational efficiency. A single integrated neural network is used by YOLOv9 to predict bounding boxes and class probabilities for objects in an image simultaneously in contrast to conventional two-stage detectors that need separate stages for proposal generation and bounding box regression. Achieving real-time processing requires a significant speed advantage which this one-stage detection approach offers.

We go deeper into YOLOv9s advantages and strategically use them to address the problem of multi-view faces. First we use a carefully selected training dataset to tackle the problem of pose variation. This dataset includes a wide range of facial poses including profile views frontal views tilted angles and even partially obscured faces. The YOLOv9 model can learn strong and generalizable features that enable it to correctly identify faces in images regardless of their orientation thanks to the extreme diversity of the data. Since faces are rarely displayed in a perfectly frontal fashion in real-world situations this is especially crucial.

Second we leverage the multi-scale prediction capabilities and anchor boxes of YOLOv9. Anchor boxes are pre-made boxes positioned thoughtfully throughout the image in different sizes and aspect ratios. In order to produce tight-fitting bounding boxes around the identified faces the model refines these anchor boxes based on the features extracted from the extracted image during the prediction stage. Moreover YOLOv9s architecture includes multi-scale prediction. As a result the model can predict at several scales at once and recognize faces of different sizes in the image with accuracy. This is especially useful for multi-view face detection where a faces apparent size can change dramatically based on how it is positioned in relation to the camera. Strengths such as the adaptability of anchor boxes the generalizability from a wide range of training sets and the efficiency of single-stage detection are leveraged by the multi-scale prediction based approach to achieve robust and real-time face detection across a spectrum of facial poses. The subsequent face recognition tasks in our system rely heavily on this reliable detection as their fundamental building block.

2) ResNet based face detection: Utilizing a ResNet model that has already been trained is key. One deep learning architecture known for its efficacy in image recognition tasks is called ResNet or Residual Neural Network. Its remaining learning blocks are its main advantage. By adding the input

straight to the convolutional layers output these blocks allow the network to learn complex features avoiding the vanishing gradient problem a common training challenge for deep networks. This method makes it easier to learn feature hierarchies from images in an efficient manner.

We use a pre-trained ResNet for multi-view face detection. The purpose of this pre-training is to teach the ResNet generic image features such as edges textures and basic shapes using a large dataset of labeled images. These qualities provide a solid basis for the tasks that follow. This is where pre-training becomes useful. Because of the large dataset it was trained on the pre-trained ResNet already possesses strong feature extraction capabilities. Main stage goes to the pre-trained ResNet. The input picture is fed into the network possibly with multiple faces in different orientations. Features are gradually extracted from the image as it passes through the convolutional layers capturing progressively more intricate details. These characteristics serve as the fundamental components of face detection.

After being extracted the features are fed into more layers that are made expressly for face detection. In these layers classifiers that have been trained to detect faces in an image regardless of their orientation may be used. Bounding boxes enclosing the faces found in the image would be the stages output. Our advantage is substantial because we use a pre-trained ResNet. Our system is able to manage the complexity of multi-view face detection thanks to the networks pre-learned feature extraction capabilities. With the help of the pre-trained network faces in different orientations can be recognized and key features that are useful for face recognition and other later tasks can be extracted.

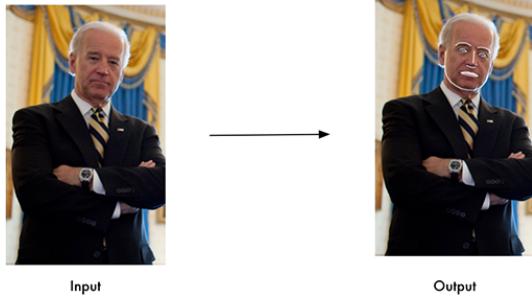


Fig. 2. Feature extraction detection example

3) RetinaFace based face detection: RetinaFace [3] is a deep learning face detection model that is able to find faces in a frame with accuracy and speed, due to its focus on being lightweight model. Unlike most face detection models discussed in II that utilizes a multi-layered face detection, RetinaFace does this in one shot.

ResNet forms the backbone of the model. At every level of the pyramid it effectively creates these feature maps based on the image. Deformable convolutional networks (DCN) are a technique that is used to further improve the models comprehension of the context within the image. With the help of DCN the model can concentrate on particular regions within

the features giving regions that probably contain faces more attention.

RetinaFace applies a cascade regression strategy after feature extraction and enrichment. This entails a number of adjustments to determine the precise position and dimensions (bounding box) of the face in the picture. Furthermore the model predicts the positions of important facial landmarks like the mouth nose and eyes in addition to bounding boxes. Applications for image editing and face recognition benefit greatly from this capability. RetinaFace uses a multi-task loss function to guarantee that these detections are accurate. This function allows the model to simultaneously improve its performance on both landmark localization and bounding box prediction by combining their errors.

After reviewing the different models for face detection, I came to a conclusion to use the RetinaFace model. That conclusion is due to the results of the different models above. The ResNet model could not detect faces unless image had high resolution, whole face is visible and face is close to the camera not faraway. The YOLO model was far better than the ResNet. Thus coming to the conclusion of using the RetinaFace model as the face detection model, as the model did not face the same issues as YOLO and ResNet.

C. Face Recognition

Various data must be processed effectively for the model to work as there is a growing need for real-time facial recognition applications. This covers jobs like recognizing faces in a live video stream and creating encodings for recognized faces.

1) Pre-processing: Creating face encodings is a crucial step in the pre-processing stage of facial recognition. The fundamental features of a persons face can be represented numerically by a face encoding. This compact representation acts as a faces unique identifier within the system it is usually a vector of numbers. A key component of our systems face encoding process is the face-recognition library which uses deep learning models that have already been trained for face recognition.

Large-scale labeled facial image datasets are used to train these deep learning models. The model gains the ability to recognize and extract a set of important features that set one face apart from another during the training phase. These characteristics can be the separation between the eyes the contour of the jawline the prominence of the cheekbones or any other special traits. Averaging 99.38% accuracy on Labeled faces in the wild dataset (LFW) was also achieved. A deep learning model can be used to create face encodings for previously undiscovered faces once it has been trained. This is accomplished by using the built in face encoding function.

It recognizes faces in images and extracts facial features using the pre-trained deep learning model in the library. After that a condensed numerical representation of the extracted features is created by encoding them in the face. The generated face encoding for a given person remains largely consistent across images of that same person as long as their lighting and facial expressions are reasonably similar. The feature

mentioned above allows the system to compare the face encodings of people who are known to be in the database with the face encodings of faces that are detected in the video stream. By using methods to compute the similarity between encodings the system determines whether the detected face corresponds to a known individual.

The recognition process and the raw picture data are essentially connected by face encoding. By reducing the complex information present in an image to a more manageable and comparable format it enables the system to perform real-time face recognition effectively.

2) Real-time processing: Getting a steady stream of video frames from a camera is the foundation of real-time processing. Our code makes use of cv2 in OpenCV in order to access the desired camera. Serving as a bridge this function connects to the camera to allow frames to be retrieved at a predetermined rate. A snapshot of the scene taken by the camera at a specific moment is represented by each retrieved frame. The live video stream that the system is analyzing for faces is made up of this series of continuously recorded and processed frames.

Our system extracts a face encoding using the face recognition function for each detected face location obtained from the face detection model. This function makes use of the pre-trained deep learning model in the face recognition library much like the method used during pre-processing for known faces. The bounding box encloses a portion of the frame that is analyzed by the model this region most likely contains a face. Then it extracts a condensed numerical representation that captures the key facial features of the person that was detected—the face encoding.

The comparison of detected faces face encodings with the encodings of known people kept in the systems database forms the basis of real-time face recognition. By repeatedly going through each extracted face encoding our code achieves this. Every encoding is compared against all known face encodings kept in the face-encodings list using the compare faces function. Whether a match is found between the current encoding and any of the known encodings is indicated by the list of boolean values that this function returns.

At last the user is shown the processed frame which may contain faces the user has identified with names above them. Our code makes use of OpenCVs feature to show the frame on the screen. Through the use of bounding boxes drawn around detected faces and names displayed for those who are recognized the user is able to view the live video stream. The system produces a smooth real-time facial recognition experience by continuously taking pictures processing them and displaying them. Through a series of coordinated actions real-time processing essentially creates a dynamic system that can recognize faces in a live video stream from a static database of recognized faces. Real-time face recognition is made possible by an intricate interplay of functions from capturing video frames to extracting face encodings and comparing them against a database.

D. DeepFace AGR based detection

DeepFace is a Python library designed for face analysis tasks, developed by the OpenAI research team. It leverages deep learning techniques and pre-trained neural networks to perform facial analysis, including age estimation, gender prediction, emotion recognition, and facial recognition. The first step involves loading an image into the program for analysis. The input image is read using the OpenCV library, which is a popular computer vision library in Python. The analyze function is the key component responsible for performing facial analysis on the provided image. Such as age, gender and race (AGR).

There are certain limitations associated with deepface that must be acknowledged. The accuracy of age, gender and race prediction heavily depends on the quality of the input image and the performance of the pre-trained models used by DeepFace. The success of the analysis is dependant upon the presence of detectable faces in the image. The pre-trained models may be biased towards certain age or gender distributions in the training data, potentially leading to inaccuracies in some cases.

Finally, for the final draft of our model, there are 3 stages that each frame pass through. First the frame is passed through the RetinaFace for identifying where each face is. After identifying and saving the coordinates of all the known face locations we loop over each face, where each loop passes a cropped image of the face through the deepface model served as a preliminary step for feature extraction for the final stage, then we save the individuals characteristics (age, gender and race). In the end the face goes to the final stage for recognizing any face from our dataset. Only if a face is recognized do we show is the frame his/her name, otherwise we do not write unknown to avoid over crowding of information in a face.

IV. RESULTS

For my face detection model I used WIDER FACE dataset for both training and evaluating the model. In addition, LFW dataset was utilized for face recognition model used for only evaluating the model. Furthermore, I incorporated a custom dataset consisting of images of me and other figures with masks on to further increase people wearing masks detection. This dataset was augmented using albumentation tool (III-A1) and annotated using labelme tool (III-A2). This custom dataset enhances models performance on recognizing me in different lighting conditions and occlusions thanks to albumentation tool for image augmentation to simulate my face in various conditions.

To test our model performance of real video surveillance, three videos varying in different conditions. First video surveillance in eye level footage consisting of 460 frames showing the facial features clearly of all the faces. Second video surveillance is positioned as a birds eye view in the morning and consists of 6832 frames. Last surveillance video is same as second video, however at night and consisting of 10871 frames.

A. Face detection

1) YOLO:

X-axis here represents the number of epochs used to train the model (30 epochs).

Y-axis here represents the magnitude of the loss or metric being measured.

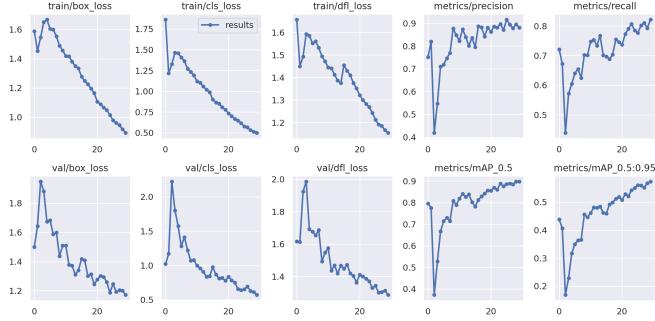


Fig. 3. YOLOV9 results after training.



Fig. 4. YOLOV9 results on video surveillances using confidence of 0.4.

2) ResNet: ResNet model was able to detect faces in the eye level video surveillance due to the face features being visible unlike the other 2 videos. However, also missing multiple faces, due to small face sizes and occlusion blocking the face. Additionally not able to detect any face at all in the other 2 videos due to its angles and more severe conditions.



Fig. 5. ResNet result on eye level video surveillance.

3) RetinaFace: RetinaFace had the best result out of the other 2 models being able to detect faces in all three videos even the faces with occlusions and small sizes.

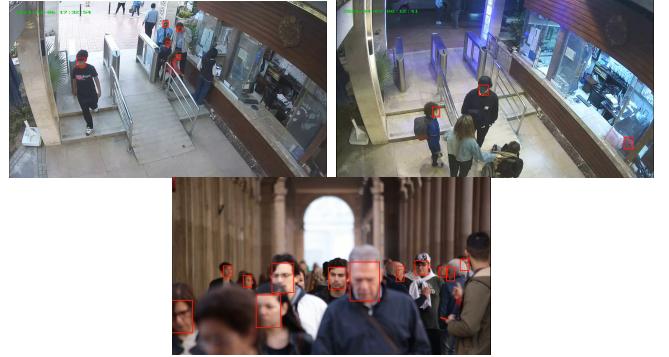


Fig. 6. RetinaFace results on video surveillances.

Video Surveillance	Total Frames	RetinaFace	YOLOV9	ResNet
Video 1	460	99.3%	75%	34.7%
Video 2	6832	98.1%	62.8%	0%
Video 3	10871	93%	12.5%	0%

TABLE I
AVERAGE PRECISION OF THE 3 MODELS ON VIDEO SURVEILLANCE.

The table above I supports our decision, because of RetinaFace average precision across our video surveillance dataset.

B. Face recognition



Fig. 7. Face recognition on eye level video surveillance

Figure 7 show the result of a eye level video surveillance. The reason for choosing eye level rather the traditional birds eye view is because in the eye level we can see the whole facial features, while the other view can not completely view all facial landmarks when encoding the face detected.

C. AGR detection

Figures 8 showcase the results of our final model on video surveillance footage. By clearly being able to detect faces and



Fig. 8. Results from our models final version

computing their AGR as well in every Frame and waiting to find a match if any known individual in our database appears.

ACKNOWLEDGMENT

My deepest gratitude to Prof. Mohamed Salem for their most appreciated support and guidance. Thanks also to my mentor Eng. Yomna Islam, for their time in aiding me with the bachelor.

I would also like to thank my family for their love and belief in me. To my friends, I thank you for your constant support and motivation.

REFERENCES

- [1] A. Kumar, A. Kaur, and M. Kumar, “Face detection techniques: a review,” *Artificial Intelligence Review*, vol. 52, pp. 927–948, 2019.
- [2] F. Zhang, X. Fan, G. Ai, J. Song, Y. Qin, and J. Wu, “Accurate face detection for high performance,” *arXiv preprint arXiv:1905.01585*, 2019.
- [3] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia, and S. Zafeiriou, “Retinaface: Single-stage dense face localisation in the wild,” *arXiv preprint arXiv:1905.00641*, 2019.
- [4] Y. Zhu, H. Cai, S. Zhang, C. Wang, and Y. Xiong, “Tinaface: Strong but simple baseline for face detection,” *arXiv preprint arXiv:2011.13183*, 2020.
- [5] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang, “Dsfid: dual shot face detector,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5060–5069, 2019.
- [6] B. Zhang, J. Li, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, Y. Xia, W. Pei, and R. Ji, “Asfd: Automatic and scalable face detector,” *arXiv preprint arXiv:2003.11228*, 2020.
- [7] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, “Faceboxes: A cpu real-time face detector with high accuracy,” in *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pp. 1–9, IEEE, 2017.
- [8] C. Chi, S. Zhang, J. Xing, Z. Lei, S. Z. Li, and X. Zou, “Selective refinement network for high performance face detection,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, pp. 8231–8238, 2019.
- [9] D. Qi, W. Tan, Q. Yao, and J. Liu, “Yolo5face: why reinventing a face detector,” in *European Conference on Computer Vision*, pp. 228–244, Springer, 2022.
- [10] Y. Liu, F. Wang, J. Deng, Z. Zhou, B. Sun, and H. Li, “Mogface: Towards a deeper appreciation on face detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4093–4102, 2022.
- [11] X. Tang, D. K. Du, Z. He, and J. Liu, “Pyramidbox: A context-assisted single shot face detector,” in *Proceedings of the European conference on computer vision (ECCV)*, pp. 797–813, 2018.
- [12] Y. Wang, X. Ji, Z. Zhou, H. Wang, and Z. Li, “Detecting faces using region-based fully convolutional networks,” *arXiv preprint arXiv:1709.05256*, 2017.
- [13] M. Kim, A. K. Jain, and X. Liu, “Adaface: Quality adaptive margin for face recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18750–18759, 2022.
- [14] Y. Xu, W. Yan, G. Yang, J. Luo, T. Li, and J. He, “Centerface: joint face detection and alignment using face as point,” *Scientific Programming*, vol. 2020, pp. 1–8, 2020.
- [15] Y. Kortli, M. Jridi, A. Al Falou, and M. Atri, “Face recognition systems: A survey,” *Sensors*, vol. 20, no. 2, p. 342, 2020.
- [16] O. Parkhi, A. Vedaldi, and A. Zisserman, “Deep face recognition,” in *BMVC 2015-Proceedings of the British Machine Vision Conference 2015*, British Machine Vision Association, 2015.
- [17] M. Alansari, O. A. Hay, S. Javed, A. Shoufan, Y. Zweiri, and N. Werghi, “Ghostfacenets: Lightweight face recognition model from cheap operations,” *IEEE Access*, 2023.
- [18] A. George, C. Ecabert, H. O. Shahreza, K. Kotwal, and S. Marcel, “Edgeface: Efficient face recognition model for edge devices,” *arXiv preprint arXiv:2307.01838*, 2023.
- [19] B. Li, T. Xi, G. Zhang, H. Feng, J. Han, J. Liu, E. Ding, and W. Liu, “Dynamic class queue for large scale face recognition in the wild,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3763–3772, June 2021.
- [20] I. Kim, S. Han, S.-J. Park, J.-W. Baek, J. Shin, J.-J. Han, and C. Choi, “Discface: Minimum discrepancy learning for deep face recognition,” in *Proceedings of the Asian conference on computer vision*, 2020.
- [21] P. Terhörst, M. Ihlefeld, M. Huber, N. Damer, F. Kirchbuchner, K. Raja, and A. Kuijper, “Qmagface: Simple and accurate quality-aware face recognition,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3484–3494, 2023.
- [22] F. Boutros, N. Damer, F. Kirchbuchner, and A. Kuijper, “Elasticface: Elastic margin loss for deep face recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1578–1587, 2022.
- [23] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4690–4699, 2019.
- [24] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- [25] T. Mare, G. Duta, M.-I. Georgescu, A. Sandru, B. Alexe, M. Popescu, and R. T. Ionescu, “A realistic approach to generate masked faces applied on two novel masked face recognition data sets,” *arXiv preprint arXiv:2109.01745*, 2021.
- [26] M. Knoche, M. Elkadeem, S. Hörmann, and G. Rigoll, “Octuplet loss: Make face recognition robust to image resolution,” in *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*, pp. 1–8, IEEE, 2023.
- [27] F. Liu, M. Kim, A. Jain, and X. Liu, “Controllable and guided face synthesis for unconstrained face recognition,” in *European Conference on Computer Vision*, pp. 701–719, Springer, 2022.