# Deep Learning Framework to Detect Face Masks from Video Footage

Aniruddha Srinivas Joshi*, Shreyas Srinivas Joshi†,
Goutham Kanahasabai‡, Rudraksh Kapil§, Savyasachi Gupta¶
B.Tech., Department of Computer Science and Engineering
National Institute of Technology, Warangal, Telangana, India - 506004
aniruddha980@gmail.com*, sj841871@student.nitw.ac.in†,
gauthamkanags@gmail.com‡, rkapil@student.nitw.ac.in§, gsavya10@gmail.com¶

*Abstract*—The use of facial masks in public spaces has become a social obligation since the wake of the COVID-19 global pandemic and the identification of facial masks can be imperative to ensure public safety. Detection of facial masks in video footages is a challenging task primarily due to the fact that the masks themselves behave as occlusions to face detection algorithms due to the absence of facial landmarks in the masked regions. In this work, we propose an approach for detecting facial masks in videos using deep learning. The proposed framework capitalizes on the MTCNN face detection model to identify the faces and their corresponding facial landmarks present in the video frame. These facial images and cues are then processed by a neoteric classifier that utilises the MobileNetV2 architecture as an object detector for identifying masked regions. The proposed framework was tested on a dataset which is a collection of videos capturing the movement of people in public spaces while complying with COVID-19 safety protocols. The proposed methodology demonstrated its effectiveness in detecting facial masks by achieving high precision, recall, and accuracy.

*Index Terms*—Face mask detection, Deep Learning, Computer Vision

## I. INTRODUCTION

With the ever swift development of machine learning algorithms and methodologies in recent times, the task of face detection has been addressed to a large extent. For instance, the face detection model proposed in [1] achieves a precision of 93% even when detecting multiple faces. Due to the advancement of facial detectors, numerous applications such as real-time face recognition systems [2], security surveillance systems [3], etc. have been developed.

Despite the success of such existing techniques, there is an increasing demand for the development of robust and more efficient face detection models. In particular, the detection of masked faces proves to be a challenging and arduous task for existing face detection models due to several reasons. Firstly, traditional face detection algorithms are based on the extraction of handcrafted features. The Viola Jones face detector [4] uses Haar features with the integral images technique to extract facial features. Other feature extraction techniques include the utilisation of the Histogram of Gradients (HOG) [5], Fast Fourier Transform (FFT) and Local Binary Patterns (LBP) [6]. With advancements in the field of deep learning, neural networks can now learn features without utilising prior knowledge for forming feature extractors such as the You Only Look Once (YOLO) algorithm [7].

The pressing concern with the aforementioned approaches when it comes to face mask detection is that the face masks, with their visual diversity and various orientations behave as occlusions and variable noise to the models. This leads to a lack of local facial features, resulting in the failure of even state-of-the-art face detection models. Moreover, there is a lack of large datasets with labeled images of faces with facial masks required in order to analyse the vital characteristics common to masked faces, thus accounting for the low accuracy of existing models. These factors together justify the challenging nature of masked face detection in the field of image processing.

During the COVID-19 pandemic, everyone is advised to wear face masks in public [8]. According to the World Health Organization (WHO), masks can be used for source control (worn by an infected individual to inhibit further transmission) or for the protection of healthy people. At the time of writing, the global pandemic has infected over 11 million people worldwide and has led to over half a million casualties [9]. The wide-scale usage of face masks poses a challenge on public face detection based security systems such as those present in airports, which are unable to detect facial masks. Since the improper removal of masks can lead to contracting the virus, it has become essential to improve facial detectors that rely on facial cues, so that detection can be performed accurately even with inadequately exposed faces.

## II. RELATED WORKS AND LITERATURE

In this section, we review some similar works done in this domain. As elucidated in section I, although research on face detection has been going on for decades and has achieved great success, algorithms and methodologies that are earmarked for face mask detection are limited.

Ge *et al.* [10] developed a deep learning methodology to detect masked faces using LLE-CNNs, which outperforms state-of-the-art detectors by at least 15%. In the given work, the authors introduced a new dataset called MAsked FAces (MAFA), containing 35,806 images of masked faces having different orientations and occlusion degrees. The proposed LLE-CNNs consist of three modules - proposal module, embedding module and verification module. The proposal

module first combines two CNNs to extract candidate facial regions from the input image and represents them with high dimensional descriptors. After that, the embedding module is turns these descriptors into similarity based descriptors using Locally Linear Embedding algorithms and dictionaries trained on a set of faces, comprised of masked and unmasked images. Finally, the verification module is used to identify candidate facial regions and refine their positions with the help of classification and regression tasks.

Nair *et al.* [11] utilised the Viola Jones object detection framework to detect masked faces in surveillance videos. The authors argued that detecting cosmetic components such as face masks takes a significantly longer period than face detection. The framework uses the Viola Jones face detection algorithm to detect the eyes and face of subjects. If eyes are recognised and later the face is recognised as well, it signifies that no face mask was used. However, if eyes are recognised but the face is not, it signifies that a face mask was worn by the person in consideration.

Bu *et al.* [12] built a CNN-based cascaded face detector framework, consisting of three convolutional neural networks. The first CNN, Mask-1 is a very shallow fully convolutional layer network with 5 layers that gives a probability of being a masked face for each detection window, followed by a Non-maximum Supression (NMS) to merge overlapping candidates. Mask-2 is a deeper CNN with 7 layers, which resizes the candidate windows and also sets a detection threshold from the previous CNN. Mask-3 is also a 7 layer CNN which resizes the input windows it receives and gives a likelihood of whether it belongs to a masked face based on a preset threshold. After NMS, the remaining detection windows are the predicted detection results.

Coming to more recent methodologies, Jiang *et. al.* [13] developed RetinaFaceMask, which is a novel framework for accurately and efficiently detecting face masks. The proposed framework is a one-stage detector which consists of a feature pyramid network to combine high-level semantic data with numerous feature maps. The authors propose a novel context attention module for the detection of face masks in addition to a cross-class object removal algorithm that discards predictions with low confidence values. The authors state that their model performs 2.3% and 1.5% more than the baseline result in face and mask detection precision respectively, and 11.0% and 5.9% higher than baseline for recall.

## III. Proposed Approach

In this section, we elucidate our proposed framework, which is illustrated in Figure 1. The proposed framework aims to detect whether people in the video footage of a public area are wearing face masks or not. In order to do so, we first detect the face of the person and then determine if a facial mask is present on the face. It is to be noted that the terms 'face mask' and 'facial mask' are used interchangeably throughout this work.
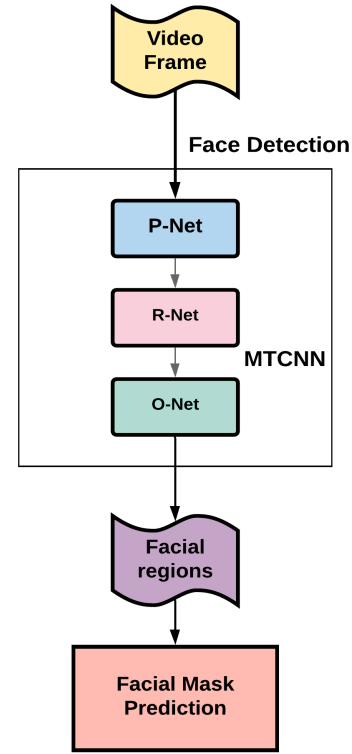


Fig. 1: Workflow of proposed framework

### A. Face Detection

For the task of face detection, we utilized the Multi-Task Cascaded Convolutional Neural Network (MTCNN) [14] as the baseline model. The model is a cascaded structure comprising of three stages of deep convolutional networks that predict the facial landmarks.

The input image is initially resized to different scales in order to build an image pyramid, which behaves as input to the three-staged network elucidated below:

- **Stage 1** consists of a Fully Convolutional Network (FCN) called *Proposal Network (P-Net)* [14], which is used to obtain the potential candidate windows in the input image pyramid and their bounding box regression vectors. In other words, *P-Net* is responsible for proposing candidate facial regions from the input image. These estimated bounding box regression vectors are used to calibrate the candidate windows obtained, after which non-maximum suppression (NMS) is used to combine largely overlapping candidates.

- **Stage 2** consists of a CNN called *Refine Network (R-Net)* [14] to which all the candidate windows obtained from the previous stage are fed. *R-Net* mainly works to filter these candidate windows. This network rejects a large number of false candidates and utilises bounding box regression to calibrate the candidates obtained. For each candidate window, the offset between itself and the nearest ground-truth is predicted, denoted by $L_i^{box}$. The learning task is a regression problem and Euclidean loss

is applied for each sample $x_i$ as:

$$L_i^{box} = \left\| \hat{y}_i^{box} - y_i^{box} \right\|_2^2 \tag{1}$$

where $\hat{y}_i^{box}$ is the target of the network and $y_i^{box}$ is the ground-truth coordinate.

- **Stage 3** comprises of a CNN called *O-Net* [14], which is responsible for proposing facial landmarks from the candidate facial regions obtained from the previous stage. *O-Net* outputs facial landmark locations, namely the eyes, nose, and mouth regions of the face. Similar to the task of bounding box regression, the detection of facial landmarks is a regression problem and the following Euclidean loss is minimised:

$$L_i^{landmark} = \left\| \hat{y}_i^{landmark} - y_i^{landmark} \right\|_2^2 \tag{2}$$

where $\hat{y}_i^{landmark}$ is the facial landmark coordinate predicted by the network and $y_i^{landmark}$ is the ground-truth coordinate.

For the task of face classification, the learning target can be formulated as a binary classification problem.

For each sample $x_i$, cross-entropy loss used was:

$$L_i^{det} = -(y_i^{det} \log p_i + (1 - y_i^{det})(1 - \log p_i)) \tag{3}$$

where $p_i$ is the probability produced by the network that the sample was a face and $y_i^{det} \in \{0, 1\}$ is the ground-truth label.

The output of this stage is the spatial coordinates of the bounding boxes enclosing the facial regions of the subjects in the frame.

### B. Facial Mask Prediction

For the task of identifying faces which are covered by a facial mask, we utilised the MobileNetV2 architecture [15], which is an effective feature extractor for object detection and segmentation. MobileNetV2 was chosen due to its ability to be deployed effortlessly on edge devices.

MobileNetV2 uses depth-wise separable convolutions much like its predecessor, but the main residual block has some key alterations from its predecessor [16]. The new residual block in MobileNetV2, known as the bottleneck residual block is illustrated in Figure 2. There are a total of 3 convolutional layers in a block, where the latter two are: a depth-wise convolution that filters the input and a $1\times1$ point-wise convolution. However, this $1\times1$ convolution is quite different. This projection layer projects input data with a higher number of dimensions (channels) into a tensor with a much lower number of dimensions. As this layer suppresses the amount of data that flows through the network and the output of each block is a bottleneck, it is known as a bottleneck residual block. Hence, the input and output of the block are low-dimensional tensors whereas the filtering that takes place inside the block is on high-dimensional tensors. The other key aspect of MobileNetV2 is the residual connection. This primarily aids with the flow of gradients through the network during backpropagation.

Each layer has batch normalisation and the activation function used is ReLU6. However, an activation function is not applied to the output of the projection layer. Since this layer outputs low-dimensional data, succeeding this layer with non-linearity could destroy valuable information.
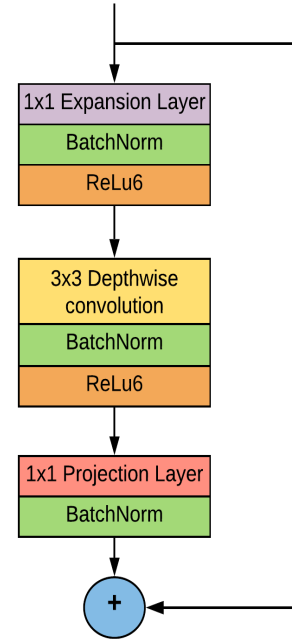


Fig. 2: Bottleneck Residual block

The full MobileNetV2 architecture, as illustrated in Figure 3, comprises of 17 bottleneck residual blocks in a row. This is followed by a regular $1\times1$ convolution. We utilise this base model of the MobileNetV2 architecture as a feature extractor for facial mask detection. We create a *facial mask classifier* using 4 layers, succeeding the earlier mentioned architecture. We downsample each $2\times2$ feature map using the average pooling layer (i.e. they are flattened) to produce a single long feature vector for classification. After passing through a ReLU activation function, we use a softmax function as illustrated in 3 to get the probability distribution over the predicted classifications. This is how the *facial mask classifier* is able to predict whether a subject in a given frame is wearing a facial mask or not.

The facial regions obtained from the face detection model discussed in Face Detection (Section III-A) are passed as input to the aforementioned *facial mask classifier* and the output is a bounding box over each face region, with the label 'Mask' indicating the presence of a face mask or 'No Mask' when no face mask is worn by the subject in consideration. This output is illustrated in Figure 6.

## IV. EXPERIMENTAL EVALUATION

In this section, we discuss the dataset used for conducting this study and the results obtained by the proposed approach. The experiments were conducted on Google Colab [17] with
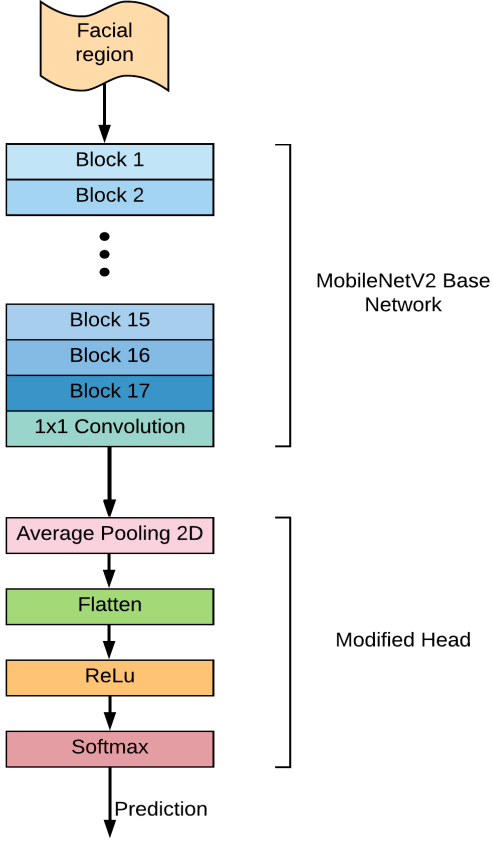
Fig. 3: Facial mask classifier constructed using MobileNetV2 architecture

Intel(R) Xeon(R) 2.00 GHz CPU, NVIDIA Tesla T4 GPU, 16 GB GDDR6 VRAM and 13 GB RAM. All programs were written in *Python* - 3.6 and utilised *OpenCV* - 4.2.0, *Keras* - 2.3.0 and *TensorFlow* - 2.2.0.

### A. Dataset Used

The dataset used in this work is a collection of footage videos of public places from multiple geographical locations, compiled from YouTube. There are a total of 15 video samples in the dataset, each with an average duration of 1 minute. The videos capture the movement of people in public areas after the imposition of various safety rules and regulations in wake of the COVID-19 pandemic. The videos showcase people from multiple ethnicities and also capture different types of face masks worn by the public. Our dataset contains videos captured using different specifications of cameras and has a multitude of camera angles, varying illumination conditions, noise, and an average frames per second (FPS) of 30. Figure 5 illustrates a few sample videos present in this dataset.

### B. Experimental Results and Statistics

The proposed approach has been evaluated by measuring the precision, recall, and accuracy metrics of the face detection model and *facial mask classifier* respectively.



Fig. 4: Visualisation of the results obtained by the proposed approach



Fig. 5: Some samples from the video dataset used in this work

$$Precision = \frac{TP}{TP + FP} \times 100\% \qquad (4)$$

$$Recall = \frac{TP}{TP + FN} \times 100\% \qquad (5)$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \qquad (6)$$

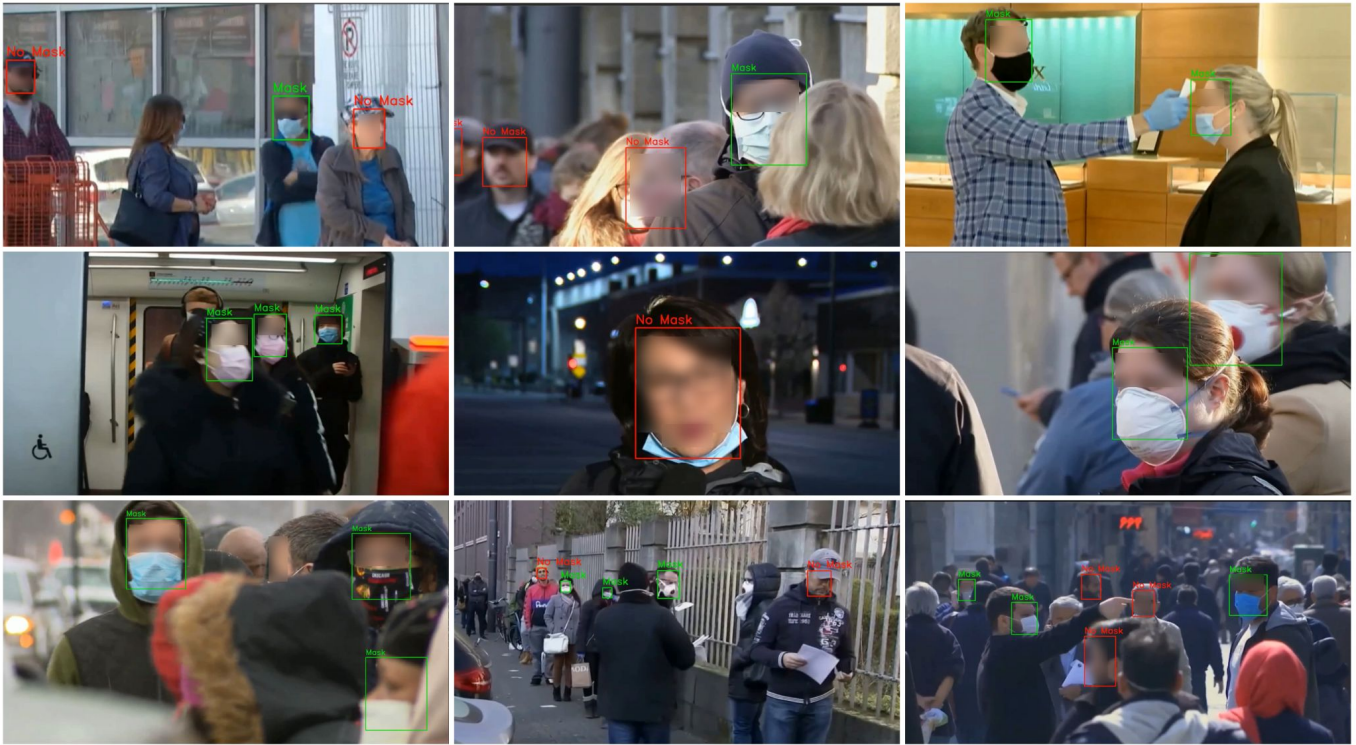where *TP, TN, FP and FN* denote the true positive, true

Fig. 6: Some instances of the results obtained by the proposed approach

negative, false positive, and false negative observations respectively.

*1) Face Detection:* The face detection model mentioned in Section III-A achieved a *precision* of 94.50%, *recall* of 86.38%, and *accuracy* of 81.84% on the chosen dataset.

*2) Facial Mask Prediction:* The *facial mask classifier* mentioned in Section III-B achieved a *precision* of 84.39%, *recall* of 80.92%, and *accuracy* of 81.74% on the chosen dataset.

TABLE I: Comparison of proposed framework with Cascaded framework for mask detection [12]

| Approach | Accuracy | Recall |
|---|---|---|
| Proposed Framework | 81.74% | 80.92% |
| Cascaded framework for mask detection | 86.6% | 87.8% |

TABLE II: Comparison of proposed framework with RetinaFaceMask [13]

| Approach | Face | | Mask | |
|---|---|---|---|---|
| | Precision | Recall | Precision | Recall |
| Proposed Framework | 94.50% | 86.38% | 84.39% | 80.92% |
| RetinaFaceMask with MobileNet | 83.0% | 95.6% | 82.3% | 89.1% |

Table I compares our proposed framework to the cascaded framework used in [12]. The higher accuracy of the cascaded framework is due to the fact that it was designed to work on images rather than videos. Also, the "MASKED FACE" dataset [12] which was used to test the cascaded framework comprises of people wearing head gear. On the other hand, the dataset used to evaluate our proposed framework captures the various types of face masks worn by the public as a precautionary measure for disease control.

Table II compares our proposed framework to RetinaMask [13]. It can be observed that our proposed framework achieves a higher precision value in detecting masks and faces as compared to RetinaMask. However, RetinaMask achieves a higher recall as the dataset it was evaluated on comprises of images of a close-up of people's faces which accounts for their better recall figures in detecting masks and faces. Also, the authors of RetinaMask do not mention the effectiveness of their model in detecting multiple faces at once, while our model works well in detecting multiple faces, as illustrated in Figure 6.

Finally, our proposed framework has also been tested on a video dataset unlike the aforesaid approaches which deal with image datasets. The video dataset used to evaluate the proposed framework contains videos taken using different specifications of cameras and has a multitude of camera angles, varying illumination conditions and noise. Thus, the proposed approach will perform well on real world camera captures.

*C. Analysis of the proposed approach*

From the earlier discussion, it can be observed that the effectiveness of the *facial mask classifier* depends on the effectiveness of the face detection model. If the face detection

model fails to detect a face or incorrectly identifies an object as a face, the performance of the *facial mask classifier* is affected.

The following key observations were made about the effectiveness of the proposed approach:

1) It is able to detect facial masks on subjects present at a considerable distance from the camera.
2) It performed well even in scenarios where the public areas captured were crowded.
3) It satisfactorily detected the presence of facial masks on subjects not directly facing the camera (i.e. only a side profile of the face was visible) in most cases.
4) It was able to identify subjects who were incorrectly wearing a facial mask (i.e. the mask was not covering their mouth and nose) and labeled them as 'No Mask'.

These observations are illustrated in Figure 6.

## V. CONCLUSIONS AND FUTURE WORK

In this work, a new approach for detecting face masks from videos is proposed. A highly effective face detection model is used for obtaining facial images and cues. A distinct facial classifier is built using deep learning for the task of determining the presence of a face mask in the facial images detected. The resulting approach is robust and is evaluated on a custom dataset obtained for this work. The proposed approach was found to be effective as it portrayed high *precision*, *recall*, and *accuracy* values on the chosen dataset which contained videos with varying occlusions and facial angles. The effectiveness of the facial mask classifier largely confides on the ability of the face detection algorithm to accurately identify faces in the video frames. This could be the subject of future research in this direction.

## VI. ACKNOWLEDGEMENT

## REFERENCES

[1] Lijing Zhang and Yingli Liang, "A fast method of face detection in video images," in *Proc. of International Conference on Advanced Computer Control*, vol. 4, 2010, pp. 490–494.

[2] N. R. Borkar and S. Kuwelkar, "Real-time implementation of face recognition system," in *Proc. of International Conference on Computing Methodologies and Communication (ICCMC)*, 2017, pp. 249–255.

[3] Z. Jian and S. Wan-juan, "Face detection for security surveillance system," in *Proc. of International Conference on Computer Science Education*, 2010, pp. 1735–1738.

[4] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, 2001, pp. I–511 – I–518.

[5] M. Murugappan and S. Murugappan, "Human emotion recognition through short time electroencephalogram (eeg) signals using fast fourier transform (fft)," in *Proc. of IEEE International Colloquium on Signal Processing and its Applications*, 2013, pp. 289–294.

[6] F. A. Alomar, G. Muhammad, H. Aboalsamh, M. Hussain, A. M. Mirza, and G. Bebis, "Gender recognition from faces using bandlet and local binary patterns," in *Proc. of International Conference on Systems, Signals and Image Processing (IWSSIP)*, 2013, pp. 59–62.

[7] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 779–788.

[8] "Coronavirus disease (covid-19) advice for the public: When and how to use masks," Apr 2020. [Online]. Available: https://www.who.int/emergencies/diseases/novel-coronavirus-2019/advice-for-public/when-and-how-to-use-masks

[9] "Coronavirus cases count." [Online]. Available: https://www.worldometers.info/coronavirus/

[10] S. Ge, J. Li, Q. Ye, and Z. Luo, "Detecting masked faces in the wild with lle-cnns," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 426–434.

[11] A. Nair and A. Potgantwar, "Masked face detection using the viola jones algorithm: A progressive approach for less time consumption," *International Journal of Recent Contributions from Engineering, Science & IT (iJES)*, vol. 6, pp. 4–14, 12 2018.

[12] W. Bu, J. Xiao, C. Zhou, M. Yang, and C. Peng, "A cascade framework for masked face detection," in *Proc. of IEEE International Conference on Cybernetics and Intelligent Systems (CIS) and IEEE Conference on Robotics, Automation and Mechatronics (RAM)*, 2017, pp. 458–462.

[13] M. Jiang, X. Fan, and H. Yan, "Retinamask: A face mask detector," 2020.

[14] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.

[15] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[16] A. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," 04 2017.

[17] "Google colab." [Online]. Available: https://colab.research.google.com/