# Face recognition in video sequence with pose variation using Neural Networks

Joshi Gajanan Ratnakar rao, Research Student
Computer Science, University -  S.R.T.M.U.Nanded,  joshigajananr@gmail.com

Dr. Prakash Khanale, Professor, DSM College, Parbhani
Computer Science, University -  S.R.T.M.U., Nanded., prakashkhanale@gmail.com

Gautam Kudale, Asst. Proff. VPASC College, Baramati.
Computer Science, University -  S.P.P.U., Pune,  gaukudale@gmail.com

Mahesh Pawar, Asst. Proff. VPASC College, Baramati.
Computer Science, University -  S.P.P.U., Pune, mdpawar77@gmail.com

*Abstract--* Pose invariant face recognition is still an active challenge in face recognition. Face localization and pose invariant feature detection are the crucial steps in robust pose invariant face recognition.  As traditional face detection algorithms like Viola Jones fails to detect face except frontal pose, need of robust algorithms significantly increased. In this paper we are proposing face detection method based on skin color and Earth Movers Distance (EMD) with Particle Swarm Optimization (PSO) for decision making. Pose invariant feature extraction is next challenge and earlier research proved how face alignment degrades quality of face shape and extracted features. This paper proposes Face Feature Prediction Model (FPM) to predict features according to stored face features for min three poses while face registration. Estimated current head pose and current features would be the inputs for FPM. Significant reduction in complexity is possible with use of FPM over traditional face alignment methods. Convolutional Neural Network (CNN) is proposed for accurate face recognition with minor false alarm rate.

*Index Terms*—**Pose Invariant Face Recognition, Face Feature Prediction Model (FPM), Earth Movers Distance (EMD), Particle Swarm Optimization (PSO), Convolutional Neural Network (CNN),**

## I. INTRODUCTION

Face recognition is one of biometric methods identifying individuals by the features of face. Research in this area has been conducted for more than 30 years; as a result, the current status of face recognition technology is well advanced. Many commercial applications of face recognition are also available such as criminal identification, security system, image and film processing. Humans often use faces to recognize individuals and advancements in computing capability over the past few decades now enable similar recognitions automatically.

## II. LITERATURE SURVEY

Jois et al. (2017) have developed an algorithm for face identification using face localization with help of Particle Swarm Optimization (PSO), supported by skin detection algorithm so as to minimize false positive identification. The first requirement of a face recognition system is to select and extract a face from the video frame, which may be static or streaming. The extracted face contains some pixels that represent part of the background, hair, or shoulders, that add noise to data. It is therefore important to be able to identify and expel such pixels from the region of interest in order to capture face-only area from the video frame, so that the recognition process is not misdirected. The authors have used the Least Median Squares (LMS) curve fitting method and derived an ellipsoidal fitness function to classify pixels in "skin" or "no skin" categories. The LMS curve fitting method is superior to the classical Least Squares (LS) method, especially in the case of noisy data. The data is in the form of RGB raster images, and the ellipsoidal fitness function is of the form

$$\frac{(R - r1)^2}{r2^2} + \frac{(G - g1)^2}{g2^2} + \frac{(B - b1)^2}{b2^2} - 1 = 0 \quad (1)$$

Where (B; G; R) represent pixel values in the blue, green, and red bands, respectively. Further, the authors have used the maximal entropy approach to find the face region in a video frame. The principle of maximal entropy maximizes the probability of detecting a face in the given video frame, thus minimizing the proportion of false positives. Detecting a face in the given video frame is called face localization, so that

attention can be focused on the detected face for the purpose of recognition. The quality of a localized face is tested by using the Binary Particle Swarm Optimization (BPSO) metric. A particle in BPSO has two attributes, namely position vector and velocity vector. The BPSO algorithm iteratively updates both of these attributes to obtain the vectors of best fit to achieve optimal performance in face detection.

Ding and Tao (2017) have addressed some of the difficulties with human face recognition in surveillance videos that occur due to blur, pose variations, and occlusion. The authors propose to use convolutional neural networks (CNN) for developing a comprehensive framework for facing the challenges in video-based face recognition (VFR). The authors experimented by artificially blurring images in training data and training a convolutional neural network to learn blur-insensitive features. The authors further propose a Trunk-Branch Ensemble CNN (TBE-CNN) model to enhance robustness of the CNN against pose variations and occlusion. TBE-CNN extracts features efficiently with help of convolutional layers at the low and middle levels that are shared between trunk and branch networks. Finally, the authors propose an improved triple loss function to promote the discriminative power of the TBE-CNN. The authors have established the efficiency of their proposed trunk-branch ensemble convolutional neural networks for video-based face recognition by applying it to three popular video face databases, namely PaSC, COX Face, and YouTube Faces. The triplet loss was initially proposed by Schroff et al. (2015) to train CNNs for face recognition. Since the input face representation image is l2-normalized for the triplet loss function, the input face representations for the triplet loss lie on a unit hypersphere. The performance of pairwise and triplet loss functions depends on sampling from image pairs or triplets from possible combinations. Still images are treated as single-frame videos. The similarity between two images is measured by the cosine metric. The authors have reported results of their experiments, where sample images have been selected from three popular databases, namely PaSC, COX Faces, and YouTube Faces.

An et al. (2017) have developed a new face alignment method to achieve pose-invariant face recognition and have named the new method adaptive pose alignment (APA). The authors make an observation that face recognition is significantly affected by three factors, namely wide variations in pose, variation in illumination, and changing facial expressions. Out of these three causes of variation in face representation in video images, pose problem is still considered to be the most challenging task. The process of face recognition consists of four stages, namely face detection, detection of facial landmarks and face alignment, feature extraction or feature learning, and feature comparison, The authors state that, while several researchers prefer to align all faces to the frontal pose, which is pre-defined and uniform, the new adaptive method learns alignment templates according to facial poses, so

that every new face can be aligned to a related template. The authors also propose a feature normalization method that can enhance discriminative feature representation of faces, when combined with the APA method. The authors introduce a sequential process, where pose-invariant face recognition is achieved through sequential application of APA-based image alignment, deep representation by a loss function, and feature normalization. Finally, the authors have shown that the proposed method achieves state-of-the-art performance on popular and challenging face datasets like IJB-A, IJB-C, and CPLFW datasets.

Khan et al. (2017) have used convolutional neural network and edge computing to develop deep unified model for face recognition. The common tools used by the research community for face recognition include Scale-Invariant Feature Transform (SIFT) and Speeded Up Robust Feature (SURF). The authors claim that their proposed algorithm for face detection and recognition based on Convolutional Neural Networks (CNN) outperform the traditional techniques. The proposed system achieved 97.9% accuracy on the test data. The authors finally remark that the proposed method is secure, reliable, and easy to use. The proposed system does not require any additional hardware or software.

Sang et al. (2016) have addressed the problem of face recognition through three-dimensional (3D) face models, rather than traditional 2D based models. The authors propose a pose-invariant face recognition method with help of RGB-D images. The proposed method can handle occlusion and deformation, both challenging problems in two-dimensional face recognition, by employing depth. The authors have also used depth for developing a similarity measure by way of frontalization and symmetric filling. The authors demonstrate that the additional information on depth has led to an improved performance of face recognition with large pose variations by applying the proposed method to Bosphorus, CurtinFaces, Eurecon, and Kiwi databases.

Cootes et al. (2001) have based their study on their own earlier research where they showed that a statistical model of appearance is used for matching two images in two stages. First, an Active Shape Model is matched to boundary features in the image. Second, an eigenface model is used for reconstructing the texture in a shape-normalized frame. However, this approach may not produce an optimal fit of the model to the image. The authors have proposed to use a direct optimization approach by matching shape and texture simultaneously to develop a rapid, accurate and robust algorithm. This has been achieved by developing a parametric model for appearance. The model is as follows.

$$x = \_x + Qsc$$
$$g = \_g + Qgc;$$

where _x is the mean shape, g is the mean texture in a mean shaped patch, Qs; Qg are matrices describing the modes of variation, and c is the parameter, controlling the shape and

texture. The authors have discussed gradient-based optimization method to achieve the best possible performance. Ramadan and Abdel-Kader (2009) have developed a feature selection algorithm based on Particle Swarm Optimization (PSO) and applied it to the problem of face recognition. The computational paradigm Particle Swarm Optimization (PSO) is based on the idea of collaborative behavior inspired by the social behavior of flocks of birds or schools of _sh. Feature selection is important because it reduces the number of features by removing irrelevant, noisy, and redundant data and improves recognition accuracy. The two feature extraction techniques by the authors are discrete cosine transforms (DCT) and the discrete wavelet transform (DWT).The authors have shown through experiments that the PSO-based feature selection algorithm generates excellent recognition results with minimal set of selected features.

Bichwe and Shende (2015) have considered the problem of multi-view face recognition. The authors claim that multi-view recognition system is insensitive to pose variations, and hence is robust. The authors have considered the cases of still-image recognition and video-based recognition separately and then added multi-view-based recognition as an extension. The authors use spherical harmonics to define the robust feature. Spherical harmonics involve a set of orthogonal basis functions defined on the unit hypersphere and are used for expanding expand square integrable functions linearly. Variation in pose causes self-occlusion of facial feature, creating challenges to designing robust face recognition algorithms. Use of multi-view data provides a promising approach to handle pose variation and its inherent challenges.

Mudunuri and Biswas (2016) have proposed an automatic face recognition algorithm for low resolution face images captured in an uncontrolled environment. A common transformation matrix is learnt through multidimensional scaling for the entire face. Facial features of low resolution and high resolution images are simultaneously transformed in such a way that the distance between the two is approximates the distance between them if both the sets of images had been captured under identical conditions. The authors propose an algorithm that matches facial images across pose, illumination, and resolution. The transformation matrix is learned from high resolution frontal and low resolution non frontal images.

Kasar et al. (2016) have reviewed the literature on the use of a neural network for face recognition. The major challenges in face recognition are due to variation in illumination_, changes in pose, orientation, and expression. Researchers have tried a variety of neural network models, including feed-forward and recurrent neural networks. This is a review paper and provides a collection of results reported by several researchers on the use of artificial neural networks for face recognition.

### III. PROPOSED METHOD

Proposed method for pose variant face recognition starts with video sequence acquisition as we are aiming face recognition in video sequence. For testing and evaluation purpose we are using NCR-IIT Facial Video Database which contains RGB video with different users and continuous face movements. One of many objectives is tracking a face in video sequence which is continuous in nature. Such kind of tracking that is achieved is against pose variations. Skin color detection along with EMD will make this possible. After tracking the face, the current head pose is estimated and features are extracted.

Novel idea in proposed pose variant face recognition is predicting face features for target pose stored in training database while registering user faces. Features will be predicted according to current extracted features and current estimated pose in order to match features with best possible 3 pose features stored in database. Finally CNN will play an important role in face recognition with good accuracy.

### A. Face Tracking

Skin color detection, EMD calculation and PSO are the three main steps involved in face tracking.

### 1. Skin Color Detection

Approach is based assumption of maximum skin like pixels in an image belongs to human face. Once image acquisition and preprocessing is done next step is to classify each pixel from image between two classes as Skin pixels and Non Skin pixels. We have formulated an Ellipsoidal equation based on earlier study [1] to find out probability of each pixel belongs to face skin. The Ellipsoidal equation for skin pixel detection is[1],

$$\frac{(R-166)^2}{13^2} + \frac{(G-100)^2}{37^2} + \frac{(B-71)^2}{182^2} - 1 = 0 \quad (1)$$
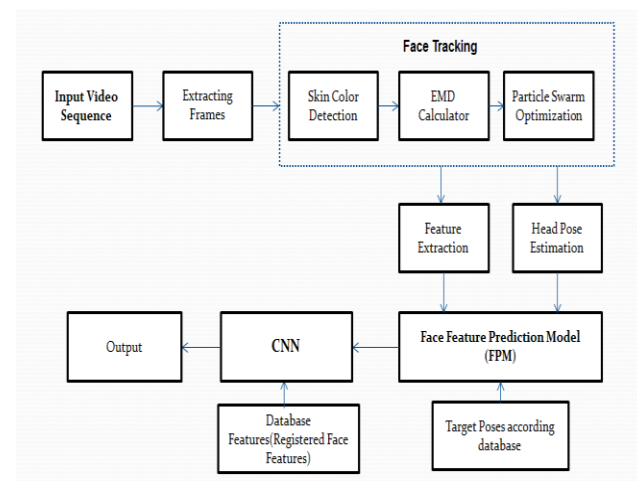
Eq.1 Ellipsoidal Equation



Fig.1. Block Diagram of Proposed Method

*2. EMD*

EMD is calculated as a fitness function to classify suspected kernels into face or non-face kernels. EMD is robust and reliable measure of spectral similarity between images. Here once pixels are classified in to skin and non-skin pixels, image is divided into small non overlapping boxes i.e. kernels and then EMD can be calculated between database face kernel and all image kernels. Kernel with best EMD value will be selected as face kernel. Size of kernel is decided by trial and error method while evaluating the performance of tracking algorithm.

EMD can be formulized as linear programing problem as follows,

$$EMD(P,Q) = \frac{\sum_{i=1}^{m}\sum_{j=1}^{n} d_{ij}F_{ij}}{\sum_{i=1}^{m}\sum_{j=1}^{n} F_{ij}} \qquad (2)$$

Where $P$ and $Q$ are the two suppliers, $F$ is flow matrix and $F_{ij}$ is the flow between $Pi$ to $Qj$ and vice versa. $D$ is the ground distance matrix $d_{ij}$ is the ground distance between $Pi$ and $Qj$.



Fig2.EMD Value for Different Kernel Positions.

*3. Particle Swarm Optimization-*

Is a technique that optimizes an issue by iteratively attempting to improve an applicant arrangement with respect to a given measure of value. It takes care of a problem by having a population of candidate solution, the kernel particles, and moving these particles around in the search-space as indicated by the fitness function over the position and speed of the particles.[1]

*B. Pose Estimation*

Pose estimation corresponds to finding out alignment and orientation of detected face. This is achieved by several morphological operations and extracting few geometric features on detected face image. Geometric features include centroid, orientation, eccentricity, convex hull etc. The experimental work is made simpler by passing the traditional methods of pose estimation like detection of landmarks such as like eyes, nose, lips etc. Because this was requiring appropriate face pose to identify those landmarks. Such kind of manual land marking is purposefully avoided.  Proposed method of estimating pose from geometrical features is found very simpler and promising.
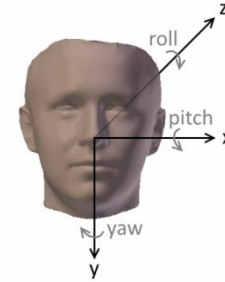


Fig3.Face Pose in yaw, pitch and roll angle.

Face pose estimation is the estimation of 3 popular Euler angles yaw ($\alpha$), pitch ($\beta$) and roll ($\gamma$). Those three angles will give clear idea about orientation of face in all three dimensions, however images used for processing will be 2D images only.

*C. Feature Extraction*

Feature extraction is vital step in robust face recognition. This stage requires meaningful feature subset extraction from input face image for reliable face recognition. We have studied various face feature extraction methods but we have found geometrical face features extraction method more promising for face recognition. Geometrical features are easier to predict and scale for expected face pose, so we found them easily compatible with our proposed face feature prediction model.

*Geometrical Features:*

Geometrical feature extraction involves detection of key landmarks on human face like eye, nose, mouth etc. The feature-based or analytic approach computes a set of geometrical face features of eyes, a mouth, and a nose. In this representation, outline of the face and positions of the different facial features form a feature vector. Usually, for good extraction process, the feature points are chosen in terms of their reliability. To compute the geometrical relationships, the location of these points are used. Such system is insensitive to position variations in the image. Nevertheless, Geometric features present the shape and locations of facial components, which are extracted to form a feature vector that represents the face[]. Following are the different geometrical face features,
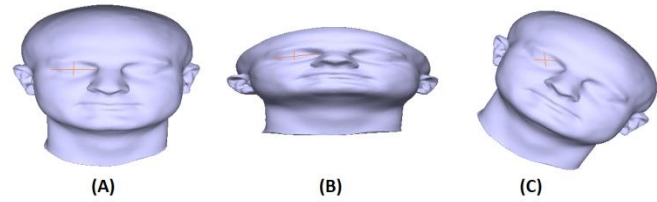
Fig4. Variation in eye height and width for different angles.

*1. Eye Height:* distance between eyelids in terms of pixels.

*2. Eye width:* horizontal width of eye in terms of pixels and extracted mainly when frontal face is detected otherwise this features will be predicted.

*3. Eye to Nose Tip Distance:* can be measured from center of two eyes to nose tip in pixels.

*4. Face Width from Nose:* can be measured from center of nose to side edge of face in pixels.

*5. Mouth to Nose Tip Distance:* can be measured from center of moth to nose tip in pixels.

*6. Mouth Width:* width between two corner points of lips.

*7. Face Height:* it is complete face height irrespective of face pose.

### D. Face Feature Prediction Model (FPM)

Face Feature Prediction Model is core of our proposed pose variant face recognition. The proposed model will make a proper prediction of features even though the input face is differently posed. Inputs expected by this model are estimation of current pose, features extracted in current pose and expected pose angle. Common formulation for each feature is difficult because of their geometrical independence. Each feature is separately formulated.

3D face model rendering in Matlab helped us to formulate for each feature.

1. Eye Height

Let's assume eye height and eye width extracted in current pose is $E_h, E_w$ respectively and predicted eye height and eye width for input pose angles $E'_h, E'_w$ respectively. For simplicity we formulated each predicted feature as a function of single Euler angles yaw $(\alpha)$ or pitch $(\beta)$ or roll $(\gamma)$. Assuming frontal face features as reference features and assuming yaw $(\alpha)$, pitch $(\beta)$, roll $(\gamma)$ equals to Zero degree.

$$E_h = Cos(\beta) * E_{hr}$$

$$E_{hr} = \frac{E_h}{Cos(\beta)}$$

$$E'_h = Cos(\beta') * E_{hr} \qquad (3)$$

Eye height is predicted based on detected height, current head pitch $(\beta)$, and targeted head pitch $(\beta')$. Only pitch angle is considered while calculating height features as they are sensitive to pitch mainly. $E_{hr}$ in equation(3) is the reference eye height when face is completely in frontal direction and all three Euler angles are equals to Zero.

2. Eye Width.

Width based features are mainly sensitive to change in yaw $(\alpha)$. Similar to Eye height, Eye width can be predicted based on current head yaw $(\alpha)$, extracted width and target head yaw $(\alpha')$.

$$E_w = \left[\frac{Cos(\alpha)}{2} + k\right] * E_{wr}$$

$$E_{wr} = \frac{E_w}{\left[\frac{Cos(\alpha)}{2} + k\right]}$$

$$E'_w = \left[\frac{Cos(\alpha')}{2} + k\right] * E_{wr} \qquad (4)$$

Similar to equation (3) and (4) all height and width based features can be predicted based on input Euler angles and expected Euler angles and their current value.

**Algorithm:** Face Feature Prediction Model (FPM)

**Input:** Detected Face Pose Euler Angles yaw $(\alpha)$, pitch $(\beta)$ and roll $(\gamma)$, Extracted Features$(X)$, target Face Pose angles yaw $(\alpha')$, pitch $(\beta')$ and roll $(\gamma')$
**Output:** Predicted Features$(X')$,
*Initialization*
1. Input all Extracted Features$(x)$ and three Euler angles of current head pose, yaw $(\alpha)$, pitch $(\beta)$ and roll $(\gamma)$,

2. Determine target head pose angles yaw ($\alpha'$), pitch ($\beta'$) and roll ($\gamma'$).

### LOOP Prediction

3. For **i=1** to **total number of features**.
4. classify geometric feature in height and width based feature.

5. **If** (Extracted Features($x_i$) = Height Based)
6.     **Then** predict features using pitch ($\beta$) angle.
7.     **Find** reference feature value with pitch ($\beta$) angle =0°,

$$X_{hr} = \frac{X_h}{Cos(\beta)}$$

8.     Predict feature with target pitch ($\beta'$) angle,

$$X_h' = Cos(\beta') * X_{hr} \qquad (3)$$

9. **Else if** (Extracted Features($x_i$) = Width Based)
10.     **Then** predict features using yaw (α) angle.
11.     **Find** reference feature value with yaw (α) angle =0°,

$$X_{wr} = \frac{X_w}{\left[\frac{Cos(\alpha)}{2} + k\right]}$$

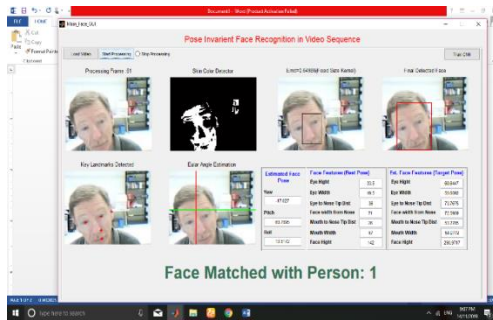12.     Predict feature with target yaw ($\alpha'$) angle,

$$X_w' = \left[\frac{Cos(\alpha')}{2} + k\right] * X_{wr} \qquad (4)$$



13. **End if.**
14. **End for**.
15. **Return** Predicted Features($X'$),

### E. Face Recognition by CNN.

Convolutional Neural Network is a type of deep neural network mainly used for analyzing visual images. CNN is very popular in various computer vision applications like Face recognition etc. The Convolutional Layer makes use of a set of learnable filters. A filter is used to detect the presence of specific features or patterns present in the original image (input).

Elements of a Neural Network includes Input, Hidden, Output layers. Input Layer accepts input features. It provides information from the outside world to the network, no computation is performed at this layer, node here just pass on the information (features) to the hidden layer. Hidden Layer nodes are not exposed to the outer world; they are the part of the abstraction provided by any neural network. Hidden layer performs all sort of computation on the features entered through the input layer and transfer the result to the output layer. Last layer is Output Layer which brings up the information learned by the network to the outer world.

Activation function in CNN decides, whether a neuron should be activated or not by calculating weighted sum and further adding bias with it. The purpose of the activation function is to introduce non-linearity into the output of a neuron.
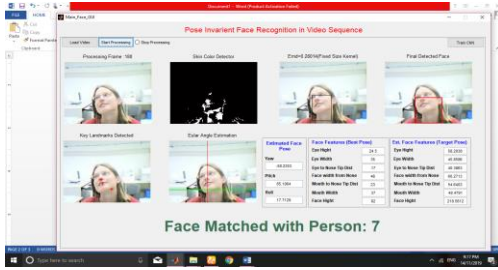
## IV. EXPERIMENTAL RESULTS

We conducted experimental study on NRC-IIT Facial Video Database. This database contains video sequence of faces of different persons with different poses. Our proposed skin color detector and EMD worked well on almost every video sequence to track face irrespective head pose. We evaluated accuracy of pose invariant face recognition with proposed Feature Prediction Model(FPM). Following results proves proposed FPM along with CNN classifier robust against face pose variation.

Face Matched with Person: 7

Transaction on Pattern Analysis and Machine Intelligence, Vol 38. N. 5, May 2016.

[10] Manisha Kasar, Debnath Bhattacharyya and Tai-hoon Kim, "Face Recognition Using Neural Network: A Review", International Journal of Security and Its applications Vol. 10, No. 3, 2016, pp. 81-100.

[11]

## V. CONCLUSION

Skin color detector with EMD calculation found very robust for face tracking with varying poses. It is found that proposed Face Feature Prediction model along with CNN has given very good results. After comparing face features with (best pose) the target features, it is found that almost all of the parameters like eye height, eye width, eye to nose distance etc. are nearly equal. Our experimental work shows that FPM has reduced the software complexity and improved classifier accuracy. FPM can accurately predict features for expected head pose. FPM output is dependent on head pose estimation and feature extraction. So more focus is still required for accurate head pose estimation and feature extraction.

## VI. REFERENCES

[1] Subramanya Jois, Rakshit Ramesh, Anoop K, "Face Localization using Skin colour and Maximal Entropy based Particle Swarm Optimization for Facial Recognition," 2017 4th IEEE Uttar Pradesh Section International Conference on Electrical, Computer and Electronics (UPCON) GLA University, Mathura, Oct 26-28, 2017

[2] Changxing Ding, Dacheng Tao, "Trunk-Branch Ensemble Convolutional Neural Network for Video-based Face Recongnition", 2016, IEEE Transactions on Pattern Analysis and Machine Intelligence.

[3] Zhanfu AN, Weigong Deng, Jiani Hu, Yaoyao Zhong, Yuying Zhao, "APA: Adaptive Pose Alignment for Pose-Invariant Face Recognition", 2019, DOI 10.1109/Access.2019.2894162, IEEE Access

[4] Muhammad Zeeshan Khan, Sadd Harous, Saleet-Ul-Hussan, Muhanmmand Usman Ghani Khan, razi Iqbal, Shahid Mumtaz, "Deep Unified Model for Face Recognition based on Convolution Neural Network and Edge Computing", 2019, IEEE Access, DOI 10.1109/ACCESS.2019.2918275, IEEE Access.

[5] Gaoli Sang, Jing Li, and Qujun Zhao, "Pose-Invariant Face Recogniton via RGB-D images", Hindawi Publishing Corporation Computational Intelligence and Neuroscience Volume 2016, Article ID 3563758.

[6] Timothy F. Cootes, Gareth J. Edwards, and Christopher J. Taylor, "Active Appearance Models", IEEE Transaction on Pattern Analysis and Machine Intelligence, Vol. 23, No. 6 June 2001.

[7] Rabab M. Ramdan and Rehab F. Adbel – Keder, " Face Recognition Using Particle Swarm Optimization- Based Selected Features", International Journal of Signal Processing, Image Processing and Pattern Recognition, Vol. 2, No. 2, June 2009.

[8] Ms. Madhavi R. Bichwe, "Face Recognition in Video by Pose variations", IEEE International Conference on Computer, Communication and Control (IC4-2015).

[9] Sivaram Prasad Mudunuri and Soma Biswas,"Low Resolution Face Recognition Across variations in Pose and Illumination", IEEE