

# Face Pose Classification Method using Image Structural Similarity Index

Hari C.V. and Praveen Sankaran

Department of Electronics & Communication Engineering  
National Institute of Technology Calicut, Calicut, Kerala, India - 673 601  
Email:haricv@gmail.com, psankaran@nitc.ac.in

**Abstract**—Face pose estimation methods try to identify/classify the position and orientation of human faces present in an image. This paper proposes a new method of face pose classification based on the structural similarity index. It is based on the measure of the similarity between the facial image with a facial pose with a set of images in the database with different poses.

**Index Terms**—Face Pose Classification, Structural Similarity Index (SSIM)

## I. INTRODUCTION

Face pose classification is one of the major research areas in the field of computer vision and robotics. Human computer interface, driver assistance. . . etc are some of the major applications of pose estimation. Pose estimation is the process of identifying the position and orientation of a human face with respect to a reference coordinate system.

Chutorian and Trivedi [1] gives a survey on different head pose identification methods. According to them, head pose estimation is the process of inferring the orientation of a human head from digital imagery as in a computer vision context. Almost all of the head pose estimation methods are based on some rigid model with inherent limitations. Head pose estimation problem is complicated with varying conditions like lighting, background and camera geometry.

Niyogi and Freeman [2] proposed a head pose identification method based on a non-linear mapping from the input image to an output parametric description. This mapping through the examples from a training set, which gives the pose as the nearest neighbour of the input. It is a modified vector quantization which stores an output parameter code with each quantized input code with a tree structured vector quantizer for efficient indexing.

Beymer [3] proposed pose estimation by finding the eyes and nose lobe features. This method is also template-based, with tens of facial feature templates covering different poses and different people. The recognizer also applies an affine transform to the input to bring the three feature points into correspondence with the same points on the model.

Sherrah et al. [4] proposed pose estimation with a similarity-to-prototypes philosophy. The similarities are calculated robustly with the help of principal component analysis (PCA) which provide an identity invariant representation. Here, the differences in poses are enhanced by the orientation selective Gabor filters. Different filter

orientations are claimed to be optimal at different poses.

More recently, Goudeliset. al. [5] proposed a method for automatic pose extraction in head-and-shoulder videos. Mutual information between the frames is the key idea behind this technique. Mutual information evaluates the information content of each facial image (contained in a video frame) of facial poses in comparison to a given ground truth image. This method is able to find any pose required with a good accuracy.

In this paper, we extend the idea of image quality comparison with a ground truth image and propose a new face pose classification technique from a video sequence containing head-and-shoulder videos. This method is based on the structural similarity of facial image from the input video (which forms the test input image) to the facial images in the database (which are the ground truth images for the different classes).

## II. IMAGE SIMILARITY MEASUREMENT

Common image similarity measuring techniques include Mean Square Error (MSE), Peak Signal to Noise Ratio (PSNR), Root Mean Square Error (RMSE), Mean Absolute Error (MAE) and Signal to Noise Ratio (SNR). Another type of image quality measurement is based on the models of Human Visual System (HVS) [6]. Structural similarity index (SSIM) [7] is a full reference metric, in other words, the measuring of image quality based on an initial uncompressed or distortion-free image as reference. SSIM was designed to improve on traditional methods like peak signal-to-noise ratio (PSNR) and mean squared error (MSE), which have proved to be inconsistent with human eye perception. Since SSIM provides a comparative similarity between two images as shown in Fig.(1), its use can be expanded beyond the intended image quality analysis. We propose a novel SSIM based face pose classification method that takes into consideration the fact that closer poses should yield more similar images and hence should provide higher SSIM values.

### A. Universal Image Quality Index

Universal Image Quality Index [8] measures the similarity as a factor of distortion using a combination of three different factors: loss of correlation, luminance distortion and contrast distortion. If  $X = \{X_i | i = 1, 2, \dots, N\}$  and  $Y = \{Y_i | i = 1, 2, \dots, N\}$  be the original and the test signals respectively, the Universal Image Quality Index  $Q$  between  $X$  and  $Y$  is defined as

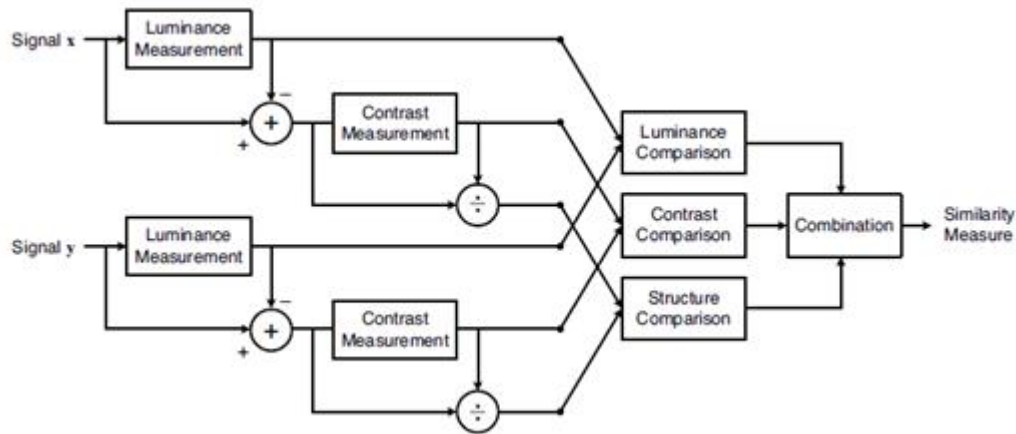


Fig.1. SSIM Measurement System [7]

$$Q = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \frac{2\mu_x \mu_y}{(\mu_x^2 + \mu_y^2)} \frac{2\sigma_x \sigma_y}{(\sigma_x^2 + \sigma_y^2)} \quad (1)$$

where,

$$\begin{aligned} \mu_x &= \frac{1}{N} \sum_{i=1}^N x_i, \mu_y = \frac{1}{N} \sum_{i=1}^N y_i \\ \sigma_x^2 &= \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2, \sigma_y^2 \\ &= \frac{1}{N-1} \sum_{i=1}^N (y_i - \mu_y)^2 \\ \sigma_{xy}^2 &= \frac{1}{N-1} \sum_{k=1}^N (x_k - \mu_x)(y_k - \mu_y) \end{aligned}$$

In the Eq.1, the first, second and third components measure the linear correlation (correlation coefficient), luminance comparison and contrast comparisons between X and Y respectively. The dynamic range of correlation coefficient is [-1, 1] and for luminance and contrast comparisons it is [0, 1]. So the dynamic range of Q will be [-1, 1]. When applying this technique to images, the image quality index can be measured from the local regions of the image using a sliding window approach. The sliding window of a particular size is moved horizontally and vertically through the image; starting from the left corner of the image. The image quality index Q will be,

$$Q = \frac{1}{M} \sum_{j=1}^M Q_j \quad (2)$$

where M is the total number of steps and  $Q_j$  is the local quality index at  $j^{th}$  step.

### B. Structural Similarity Index (SSIM)

Structural Similarity Index (SSIM) is a modified version of Universal Quality Index. Like Universal Quality Index, the SSIM value will be high for similar frames/images. SSIM index collectively measure the changes in the

luminance, contrast and structure between two signals or images.

$$S(x, y) = f(l(x, y), c(x, y), s(x, y)) \quad (3)$$

where  $l(x, y)$  is the luminance comparison,  $c(x, y)$  is the contrast comparison and  $s(x, y)$  is the structural comparison. From the Eq.3, it is clear that the SSIM index makes three comparisons (luminance, contrast and structural) for similarity measurements. They are,

$$l(x, y) = \frac{2\mu_x \mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1} \quad (4)$$

$$c(x, y) = \frac{2\sigma_x \sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2} \quad (5)$$

$$s(x, y) = \frac{2\sigma_{xy} + C_3}{\sigma_x \sigma_y + C_3} \quad (6)$$

where  $C_1 = (k_1 L)^2$ ,  $C_2 = (k_2 L)^2$  and  $C_3 = \frac{C_2}{2}$  are small constants; L is the dynamic range of the pixel values, and  $K_1 \ll 1$  and  $K_2 \ll 1$  are scalar constants.

Finally by combining all the above mentioned comparisons, the Structural Similarity Index is given by,

$$S(x, y) = [l(x, y)]^\alpha [c(x, y)]^\beta [s(x, y)]^\gamma \quad (7)$$

where  $\alpha > 0$ ,  $\beta > 0$  and  $\gamma > 0$  are parameters used to adjust the relative importance of the three components. The simplified SSIM measure is obtained by substituting  $\alpha = \beta = \gamma = 1$  and  $C_3 = \frac{C_2}{2}$ , then

$$SSIM(x, y) = \frac{(2\mu_x \mu_y + C_1)(2\sigma_x \sigma_y + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (8)$$

Eq.8 is similar to the Eq.1 when  $C_1 = C_2 = 0$ . i.e. the Universal Quality Index is a special case of SSIM when  $C_1 = C_2 = 0$ . But, the Universal Quality Index will produce unstable results when either  $(\mu_x^2 + \mu_y^2)$  or  $(\sigma_x^2 + \sigma_y^2)$  is zero.

### III. FACE POSE CLASSIFICATION

Face pose identification is considered as a classification problem based on the SSIM value between a test image and reference ground truth images from multiple face-pose classes. Face pose is the relative orientation or angle of the face with the camera. The three degrees of freedom of a human face can be described as pitch, roll and yaw as shown in Fig.2.

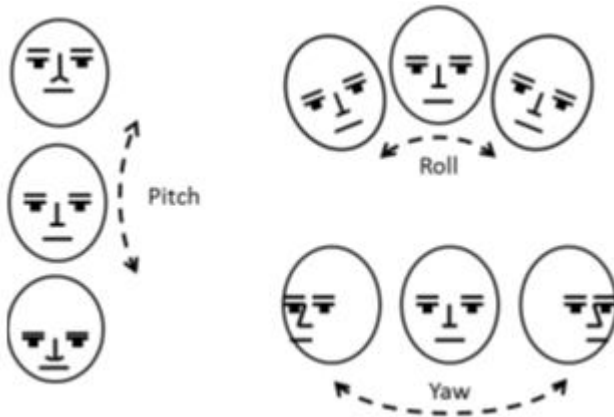


Fig.2. Face Pose Angles - pitch, roll and yaw [9]

SSIM value of the current frame/image is compared with all other face images in the database and is assigned to the class of the image which provides the maximum SSIM value. For calculating the SSIM value, the system contains an image database, which has reference face poses of this particular person. The input video will be subdivided into different frames and each frame's SSIM value will be calculated with all the images of the database. The current frame of the reference person is assigned to the pose class corresponding to the maximum SSIM value of the image from the database.

### IV. EXPERIMENTAL RESULTS

SSIM based facial pose classification algorithm is implemented in MATLAB and tested in the Pointing'04 [10] database. This database consists of 15 sets of images. Each set contains 2 series of 93 images of the same person at different poses. There are 15 people in the database, wearing glasses or not and having various skin color. The pose or head orientation is determined by the angles varying from -90 degrees to +90 degrees. Only luminance component is used here in order to reduce processing time.

The input video is sub-sampled into its frames. Each frame is compared with all the images of the database to calculate its SSIM value for all frames. High SSIM value will be obtained when the current frame and an image in the database will be similar to each other. The corresponding image in the database will give the current pose of the input frame.

In order to describe how the pose is assigned to a class, let us describe the face pose identification method for a person in the database. Let us consider frame number 50 of the input video which is shown in Fig.3 (a).

SSIM values which lie between frame 50 and all the images



a) Frame number 50 (b) Detected Pose from the Database

Fig.3. SSIM - Output for frame number 50

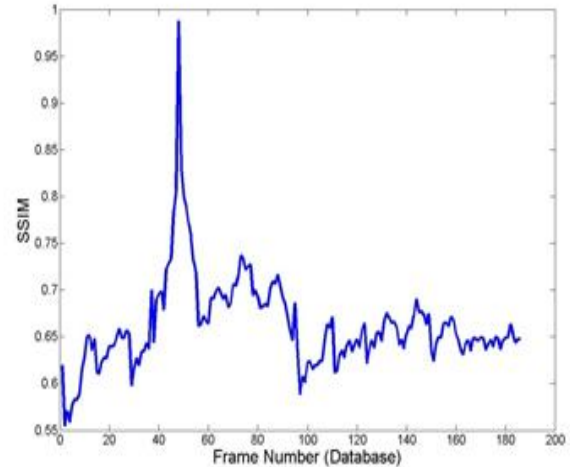


Fig. 4. SSIM Value of Frame number 50

in the database is calculated. Image representing a specified pose in the database with the highest SSIM value will be the pose of frame 50 as shown in Fig. 3(b).

As shown in Fig.4, SSIM value of frame number 50 with all other reference images in the database is varying. SSIM value is maximum for the 50th image in the database which is the same as the current frame under test. That image in the database represents the pose of the current frame 50 of our input video.

In this manner, we can identify any of the pose from the input video. The inputs and outputs for frame numbers 75, 100, 125 and 150 are as shown in the Fig. 5. SSIM values between the test frame and the images in the database are as shown in the Fig.7. From this graph it is clear that SSIM value is maximum for the image containing in the database which has same pose as in the input video and also that the SSIM value varies between 0 and 1.

#### A. Performance Analysis

If the number of sample images in the database is reduced, i.e. if all the poses are not available to us as ground truths, then the detected pose will be a pose nearest to the corresponding pose in the test input frame. For analyzing the performance experimentally, generate 'X' random numbers with uniform distribution. Corresponding numbered frames are left out from the training set and are used as frames for testing. This is repeated using different 'X' values. The 'X' values used here are 5, 10, 25, 50 and 100. SSIM values for this experiment are as shown in Fig.8, 9, 10, 11 and 12. From these figures, it is clear that the confidence for the detection is reducing, because by reducing the number of frames from the database, more frame's SSIM value will be close together.

The confidence level for the pose classification reduces by leaving out frames from the database. Confidence level is

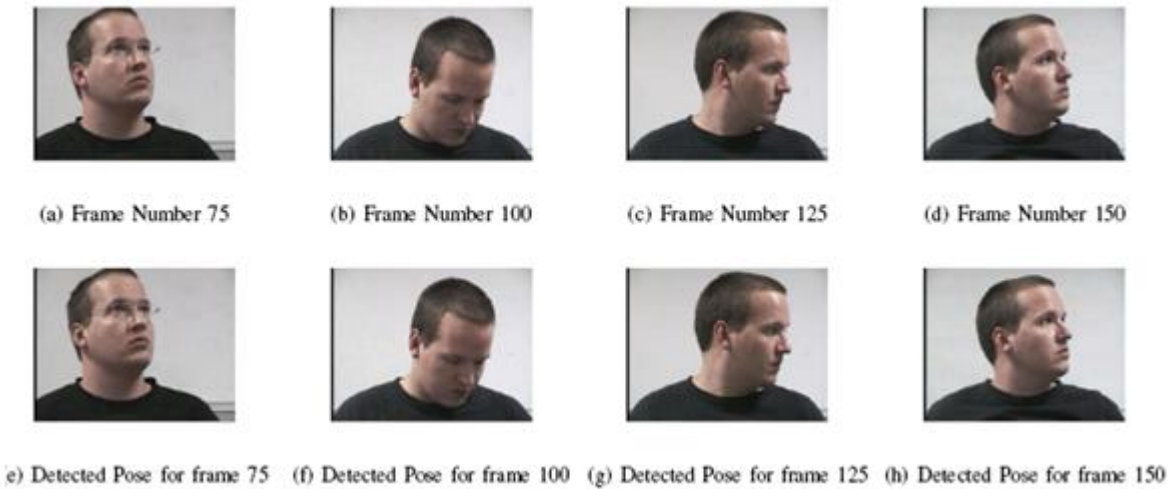


Fig.5.Different frames and their detected poses

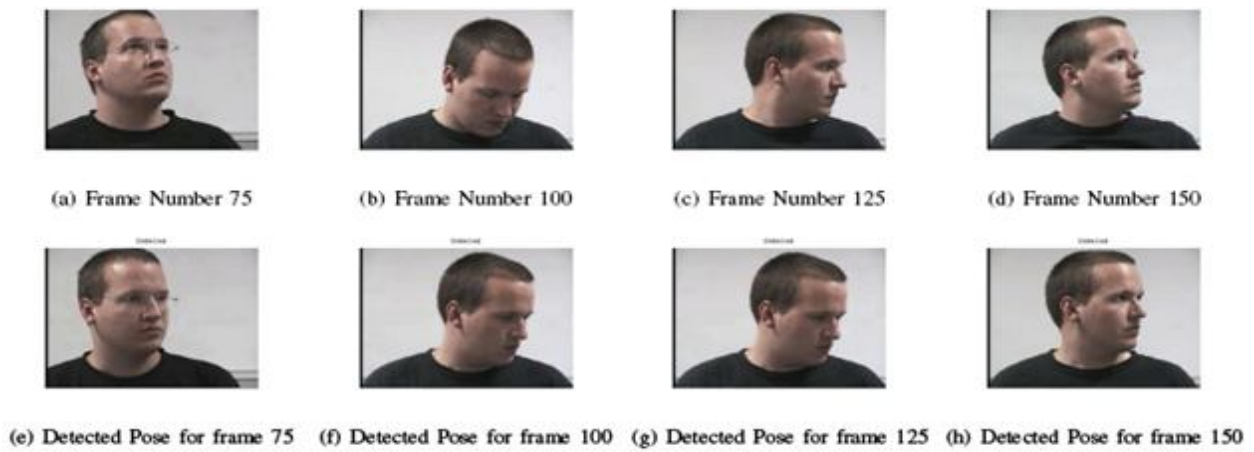


Fig. 6. Different frames and their detected poses in a reduced database (Leaving 10 frames from the database)

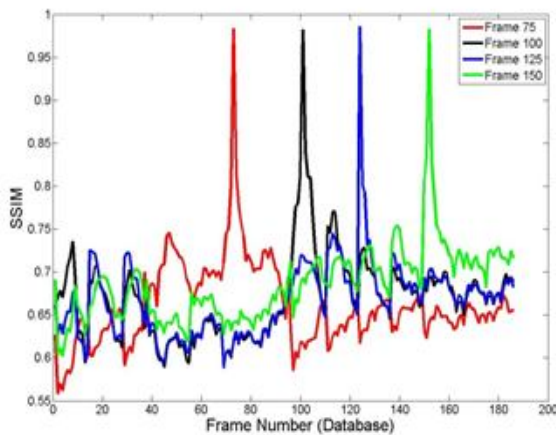


Fig.7. SSIM Value of Frame number 75,100,125 and 150

measured by taking the average of the ratio between peak values to the non-peak values. The confidence level measurement for different databases is as shown in the Table. I. From this, it is clear that, the confidence level for detection is decreasing by reducing the number of images in the database.

The outputs for different frames for this reduced database are as shown in Fig.6 and Fig.13. There are 10 frames removed from the database in both cases and SSIM value will

TABLE I: CONFIDENCE LEVEL MEASUREMENT

Number of Training images reduced from database	Confidence measure( Peak/ Non-peak Average)	Time (Sec)
5	1.4588	11.68
25	1.3336	10.30
50	1.3116	8.70
75	1.3005	7.04
100	1.2508	5.54

bemaximum for a pose that is closest to the input test posesince the original pose is not available in the database. Thedetected pose for frame number 50 is shown in Fig. 13. Here,the detected pose is not exactly similar with the input frame. The SSIM value is maximum for some other image and isshown as the detected pose of this frame. SSIM values offrame 50 for reduced database are as shown in the Fig. 14. The inputs and outputs for frame numbers 75, 100, 125 and 150are as shown in the Fig. 6 with a reduced number of imagesin the database.

## V. CONCLUSION

This paper proposed a novel approach for face pose



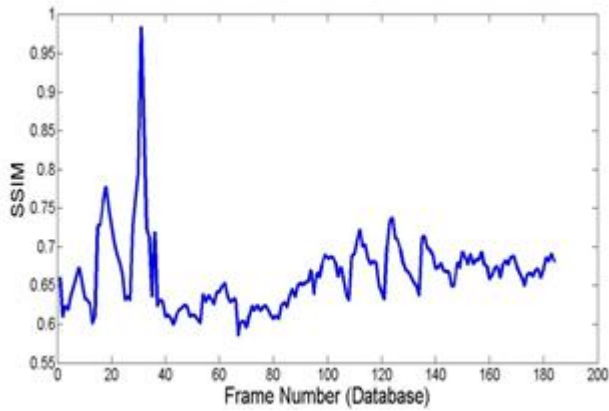


Fig. 8. SSIM Value for a reduced database (Leaving 5 frames from the database)

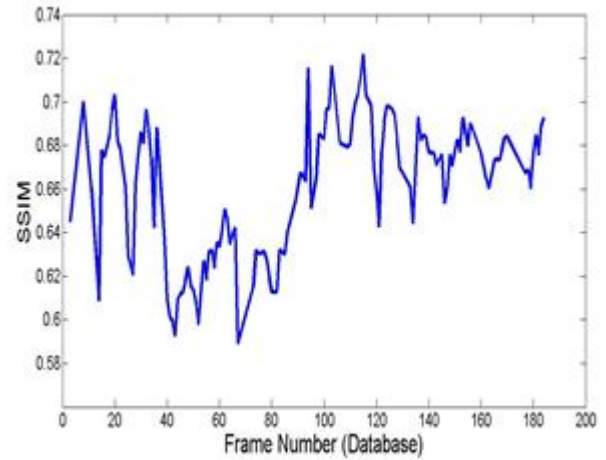


Fig.11. SSIM Value for a reduced database (Leaving 50 frames from the database)

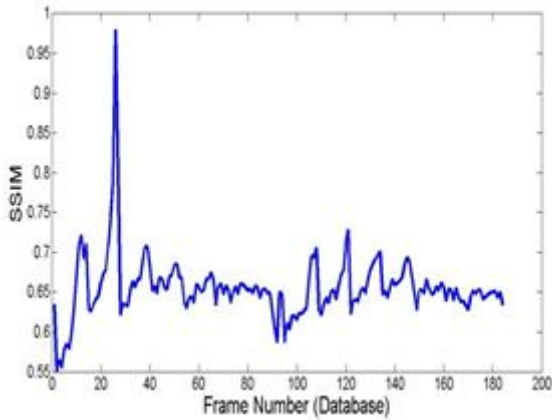


Fig. 9. SSIM Value for a reduced database (Leaving 10 frames from the database)

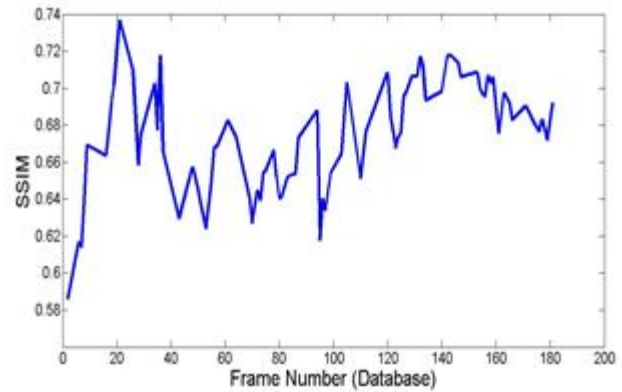


Fig.12. SSIM Value for a reduced database (Leaving 100 frames from the database)

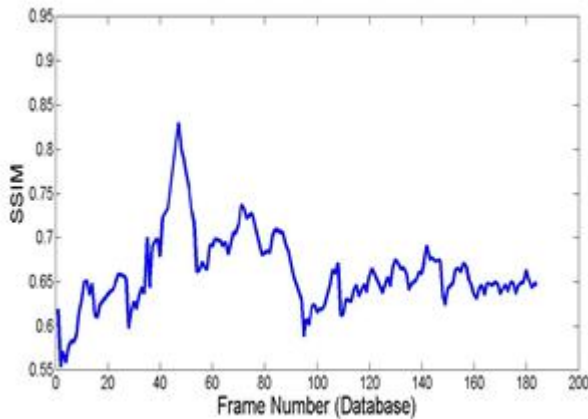


Fig.10. SSIM Value for a reduced database (Leaving 25 frames from the database)



(a) Frame Number (b) Detected Pose from the Reduced Database

Fig.13. Detected Pose in a Reduced Database (Leaving 10 frames from the database)

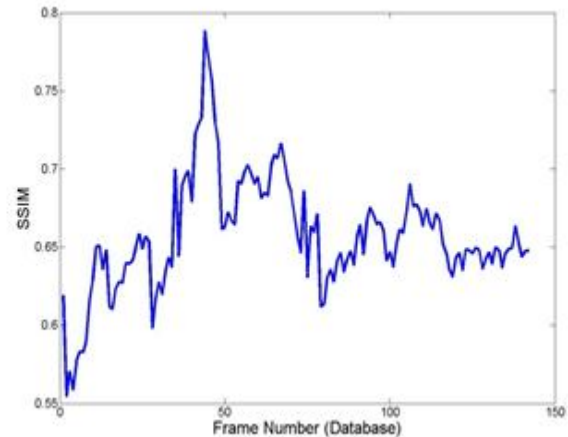


Fig.14. SSIM Values for the frame shown in Fig.13

identification based on structural similarity index. The method gave good results for cases with large number of training samples for a pre-determined user and could be a good tool for a pose authentication system. Current work focuses on applying pose identification algorithm for a vehicle driver distraction system. While it was seen that a reduced database resulted in a reduced confidence in the obtained result, an authentication problem would start with possibility to have a large database, under which the proposed method of pose identification gives extremely accurate output.

## REFERENCES

- [1] E. Murphy-Chutorian and M. Trivedi, "Head pose estimation in computervision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 4, pp. 607–626, 2009.
- [2] S. Niyogi and W. Freeman, "Example-based head tracking," in *Proceedings of the Second International Conference on Automatic Face and Gesture Recognition, 1996*, IEEE, 1996, pp. 374–378.
- [3] D. Beymer, "Face recognition under varying pose," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 1994. Proceedings CVPR'94*, 1994. IEEE, 1994, pp. 756–761.
- [4] J. Sherrah, S. Gong, and E. Ong, "Face distributions in similarity space under varying head pose," *Image and Vision Computing*, vol. 19, no. 12, pp. 807–819, 2001.
- [5] G. Goudelis, A. Tefas, and I. Pitas, "Automated facial pose extraction from video sequences based on mutual information," *Circuits and Systems for Video Technology, IEEE Transactions on*, vol. 18, no. 3, pp. 418–424, 2008.
- [6] T. Pappas, R. Safranek, and J. Chen, "Perceptual criteria for image quality evaluation," *Handbook of image and video processing*, pp. 669–684, 2000.
- [7] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *Image Processing, IEEE Transactions on*, vol. 13, no. 4, pp. 600–612, 2004.
- [8] Z. Wang and A. Bovik, "A universal image quality index," *Signal Processing Letters, IEEE*, vol. 9, no. 3, pp. 81–84, 2002.
- [9] "Head pose angles," 2012. [Online]. Available: <http://msdn.microsoft.com/en-us/library/jj130970.aspx>
- [10] N. Gourier, D. Hall, and J. Crowley, "Estimating face orientation from robust detection of salient facial structures," in *FG Net Workshop on Visual Observation of Deictic Gestures, 2004*, pp. 1–9.