# A Multi-Face Challenging Dataset for Robust Face Recognition

Shiv Ram Dubey and Snehasis Mukherjee

*Abstract*— Face recognition in images is an active area of interest among the computer vision researchers. However, recognizing human face in an unconstrained environment, is a relatively less-explored area of research. Multiple face recognition in unconstrained environment is a challenging task, due to the variation of view-point, scale, pose, illumination and expression of the face images. Partial occlusion of faces makes the recognition task even more challenging. The contribution of this paper is two-folds: introducing a challenging multi-face dataset (i.e., IIITS_MFace Dataset) for face recognition in unconstrained environment and evaluating the performance of state-of-the-art hand-designed and deep learning based face descriptors on the dataset. The proposed IIITS_MFace dataset contains faces with challenges like pose variation, occlusion, mask, spectacle, expressions, change of illumination, etc. We experiment with several state-of-the-art face descriptors, including recent deep learning based face descriptors like VGGFace, and compare with the existing benchmark face datasets. Results of the experiments clearly show that the difficulty level of the proposed dataset is much higher compared to the benchmark datasets.

*Index Terms*— IIITS_MFace Dataset, Face detection, Face recognition, Challenging face dataset, Local binary pattern, Image descriptors.

## I. INTRODUCTION

Face detection and recognition from still images is an active research area in computer vision [1]. Most of the state-of-the-art face recognition approaches were restricted to the controlled environments such as frontal pose [2]. Detailed survey on face recognition tasks have been conducted many a time by the researchers [3], [4], [5].

Recently, recognizing faces in the wild images, has become an emerging area of research in computer vision [6]. Face recognition in unconstrained environment is still an unsolved problem due to the various levels of challenges like part or full occlusion of faces, varying illumination, multiple posture of faces, expressions on faces, etc. The face recognition task becomes even harder when multiple such challenges are present simultaneously. In order to facilitate the face detection and recognition research, we propose a multi-face challenging dataset including all such challenges discussed above. This dataset will be publicly available to the research community.

A few publicly available datasets exist in the literature for face recognition and detection, involving challenges like different side poses, occluded faces, varying light intensities,etc. For instance, the AT & T face database [7] has only grayscale and frontal face images. The AR face

The authors are with the Computer Vision Group, Indian Institute of Information Technology, Sri City, Andhra Pradesh-517646, India. Email: {srdubey, snehasis.mujherjee}@iiits.in.



Fig. 1: Sample images from original gallery set of the proposed IIITS_MFace dataset. The various challenges like pose variation, occlusion, illumination changes, orientations, etc. can be observed.

database [8] contains faces with different facial expressions, varying illumination, and occlusions in the face images. This database is having only single face images with uniform background. The CroppedYale dataset contains faces only with the illumination variations [9]. The LFW face dataset is challenging and captured under unconstrained environment [10] with single face images. For the comparison purpose, we have used CroppedLFW version of this dataset [11]. The PaSC face dataset consists of pose, illumination and blur effects [12]. Total 8718 faces from 293 subjects are present after applying the Viola Jones object detection method [13] for face localization in PaSC dataset. The PubFig dataset is another challenging dataset consisting of images in unconstrained environment [14]. Variations in lighting, expression and pose effects are present in the PubFig dataset with total 6472 images from 60 individuals. The dead urls are removed while downloading the PubFig images. However, none of the existing face datasets offer multiple faces in the images. The proposed IIITS_MFace is a new dataset for face recognition in images containing multiple faces. Moreover, in addition to the various challenges involved in the state-of-the-art datasets, the images of the proposed dataset are captured in uneven and varying background, which was missing in state-of-the-art datasets. Some sample images from the proposed dataset are shown in Figure 1.

We show the complicacy in the proposed IIITS_MFace dataset, by applying state-of-the-art hand-designed as well as deep learning based face recognition techniques on the

TABLE I: A summary of gallery set in terms of the variations like Frontal/Non-frontal pose and Masked/Unmasked

| Subject ID | #Frontal Masked | #Frontal Un-masked | #Non-frontal Masked | #Non-frontal Unmasked | #Total Faces |
|---|---|---|---|---|---|
| 1 | 0 | 27 | 0 | 67 | 94 |
| 2 | 0 | 59 | 0 | 121 | 180 |
| 3 | 27 | 31 | 47 | 54 | 159 |
| 4 | 0 | 4 | 0 | 13 | 17 |
| 5 | 7 | 17 | 20 | 39 | 83 |
| 6 | 0 | 18 | 0 | 68 | 86 |
| 7 | 1 | 29 | 6 | 33 | 69 |
| Total | 35 | 185 | 73 | 395 | 688 |

TABLE II: A summary of gallery set in terms of the #faces with and without spectacles

| Subject ID | #With Spectacles | #Without Spectacles | #Total Faces |
|---|---|---|---|
| 1 | 0 | 94 | 94 |
| 2 | 0 | 180 | 180 |
| 3 | 159 | 0 | 159 |
| 4 | 17 | 0 | 17 |
| 5 | 0 | 83 | 83 |
| 6 | 0 | 86 | 86 |
| 7 | 0 | 69 | 69 |
| Total | 176 | 512 | 688 |

TABLE III: A summary of probe set in terms of the variations like Frontal/Non-frontal pose and Masked/Unmasked

| Subject ID | #Frontal Masked | #Frontal Un-masked | #Non-frontal Masked | #Non-frontal Unmasked | #Total Faces |
|---|---|---|---|---|---|
| 1 | 1 | 43 | 2 | 14 | 60 |
| 2 | 2 | 6 | 19 | 25 | 52 |
| 3 | 12 | 8 | 25 | 6 | 51 |
| 4 | 3 | 9 | 10 | 28 | 50 |
| 5 | 8 | 18 | 14 | 12 | 52 |
| 6 | 3 | 5 | 10 | 32 | 50 |
| 7 | 0 | 27 | 0 | 23 | 50 |
| Total | 29 | 116 | 80 | 140 | 365 |

TABLE IV: A summary of probe set in terms of the #faces with and without spectacles

| Subject ID | #With Spectacles | #Without Spectacles | #Total Faces |
|---|---|---|---|
| 1 | 24 | 36 | 60 |
| 2 | 18 | 34 | 52 |
| 3 | 32 | 19 | 51 |
| 4 | 12 | 38 | 50 |
| 5 | 12 | 40 | 52 |
| 6 | 12 | 38 | 50 |
| 7 | 0 | 50 | 50 |
| Total | 110 | 255 | 365 |

dataset. We emphasize on local image descriptors, considering the recent success of local image descriptors on face recognition task. Several efforts have been made to apply local image descriptors for face recognition. Local Binary Pattern (LBP) is proposed by Ahonen et al. for the face representation [15]. LBP is computed by finding a binary pattern of 1 and 0 for each neighbor of a center pixel. The bit is coded as 1 if the intensity value of neighbor is greater than or equal to the intensity value of center pixel; otherwise it is coded as 0. Local Ternary Pattern (LTP) is the extension of LBP by introducing two thresholds for uniform illumination robust face recognition [16]. The LBP over four derivative images corresponding to four directions are computed and concatenated to form the Local Derivative Pattern (LDP) [17]. The concept of high order directional gradient is used to find the Local Directional Gradient Pattern (LDGP) to extract the local information of the image [18]. In the recent advancements, Semi-structure Local Binary Pattern (SLBP) [20], Local Vector Pattern (LVP) [21], and Local Gradient Hexa Pattern (LGHP) [22] descriptors are proposed for the unconstrained face recognition. The VGGFace CNN descriptor [23] is very discriminative and based on the deep learning technique. We experimented with all these descriptors on the proposed dataset. Next we provide a detailed description of the proposed dataset.

## II. PROPOSED IIITS_MFACE DATASET

The images in the proposed IIITS_MFace dataset are captured by cameras of multiple mobile phones to make it more realistic with respect to the real world face recognition problem. A lot of variations in terms of pose, masked, spectacles, number of subjects, illumination, occlusion, etc. are present in the dataset to make it as unconstrained as possible. The proposed dataset is divided into two sections with seven subjects including six male and one female. The two sections of the proposed dataset are named as Gallery Set and Probe Set. The IIITS_MFace dataset is publicly available for research purpose only[1].

### A. Gallery Set

The images of the gallery set are captured from mobile phones with multiple people involved in some activities like talking, laughing, etc. A total of 180 such images are

[1]https://sites.google.com/a/iiits.in/snehasis-mukherjee/datasets-1

captured with minimum three and maximum five number of people in an image. Sample images of this set are shown in Figure 1. We have created a cropped version of the gallery set. All the visible faces in all the images are manually cropped and annotated with the subject labels. The cropped galley set comprises of 688 faces from 180 original multi-face gallery images. The co-ordinates of each face in each image is also provided to validate a face detection algorithm. The cropped version of gallery set can be used for the experiment purpose. The characteristics of gallery set such as frontal/non-frontal pose and masked/unmasked faces are summarized in Table I. In gallery set, a subject is either with spectacle or without spectacle. Table II highlights the subjects with/without spectacles. Only subjects 3 & 4 are with spectacles. Some cropped faces of gallery set are also shown in Figure 2.

### B. Probe Set

The probe set is created in the second section with same set of subjects used in galley set. For each subject, we provide a set of face images with differents poses, captured from mobile phones. Since these images are captured by

Fig. 2: Sample seven faces per subject from gallery set. Each row corresponds to a subject.



Fig. 3: Sample seven faces per subject from probe set. Each row corresponds to a subject.
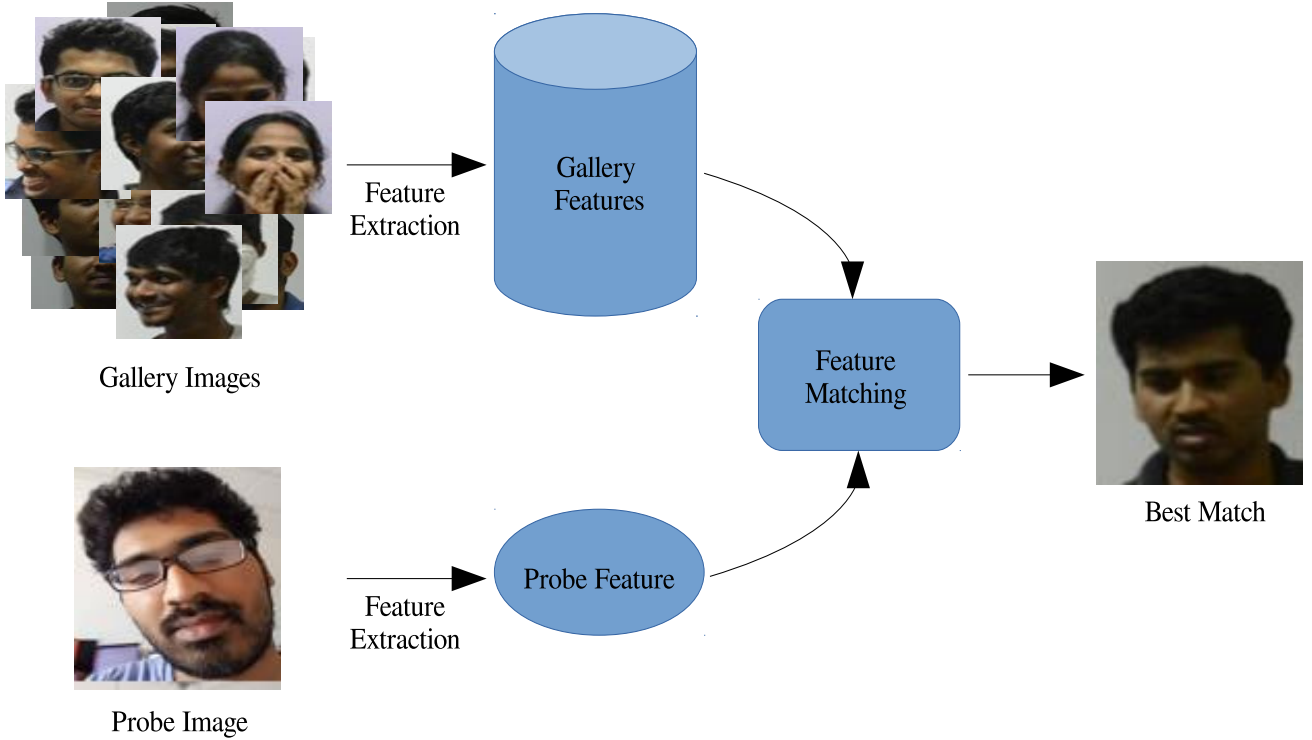
Fig. 4: The face recognition framework using local descriptors. The best matching face against a probe face is extracted based on the minimum distance between feature descriptors of probe face and gallery faces.

the subjects individually, a lot of variations are present in the image such as occlusion, spectacle, illumination, pose, viewpoint, blur, masked, etc. Total 365 images are present in the probe set consisting of nearly 50 images from each subject. A detailed description of the probe set is illustrated in Table III and IV along with the frontal/non-frontal/masked/unmasked/spectacles number of images. It can be noted that the subjects in gallery set have either used or not used the spectacles, whereas in the probe set, all the subjects except last one have mixed images with and without spectacles as depicted in Table IV. Some example faces of probe set are also shown in Fig. 3 in order to illustrate the complexity of the probe set. Next we illustrate the experiments made on proposed face dataset.

### III. FACE RECOGNITION USING LOCAL DESCRIPTORS

In this section, the nearest neighbour based face recognition framework using local descriptors is described as shown in Fig. 4. The features using a local descriptor is computed over gallery faces to create the gallery features database. The same descriptor is then used to extract the feature for any probe image. After computing the descriptors, the distance between probe feature and gallery features are computed. Finally, the class of probe face is recognized as the class of best matching gallery face based on the minimum distance between probe face and gallery faces.

Several state-of-the-art face descriptors including hand-crafted and deep learned like Local Binary Pattern (LBP) [15], Local Ternary Pattern (LTP) [16], Local Derivative Pat-

tern LDP [17], Local Directional Gradient Pattern (LDGP) [18], Semi-structure Local Binary Pattern (SLBP) [20], Local Vector Pattern (LVP) [21], Local Gradient Hexa Pattern (LGHP) [22] and VGGFace CNN descriptor [23] are tested over the proposed dataset to establish its complexity. Note that all these descriptors are proposed for face representation purpose and VGGFace CNN descriptor is very discriminative for face representation. The MatConvNet pre-trained model of VGGFace CNN descriptor is used in this paper[2]. Several distances such as Euclidean, L1, Cosine, Emd (Earth Mover Distance) and Chisq (Chi-square) [24] are also used in this paper to find the best performing distance measure for the proposed dataset.

### IV. EXPERIMENTAL RESULTS

The average recognition rate for the descriptors on the proposed IIITS_MFace dataset, is used as the evaluation criteria for the descriptors. The average recognition rate is computed by taking the mean of average accuracies obtained over all the subjects of the probe set. The average accuracy for a particular subject of probe set is computed by taking the mean of accuracies obtained by turning each image of that subject as the probe image. Until and otherwise not stated, L1 distance is used to compare the descriptors.

The average recognition rate over proposed dataset using different descriptors with L1 distance is summarized in Table V. The VGGFace descriptor is the best performing

[2]http://www.vlfeat.org/matconvnet/pretrained/

TABLE V: The average recognition rate using LBP, LTP, LDP, LDGP, SLBP, LVP, LGHP and VGGFace descriptors with L1 distance over proposed IIITS_MFace dataset.

| Descriptor | Subject1 | Subject2 | Subject3 | Subject4 | Subject5 | Subject6 | Subject7 | Mean |
|---|---|---|---|---|---|---|---|---|
| LBP | 16.67 | 19.23 | 58.82 | 6 | 19.23 | 20 | 52 | 27.42 |
| LTP | 16.67 | 17.31 | 54.90 | 8 | 5.77 | 18 | 54 | 24.95 |
| LDP | 21.67 | 30.77 | 64.71 | 0 | 5.77 | 4 | 2 | 18.42 |
| LDGP | 8.33 | 15.38 | 43.14 | 10 | 1.92 | 8 | 72 | 22.68 |
| SLBP | 11.67 | 23.08 | 45.10 | 4 | 34.62 | 6 | 88 | 30.35 |
| LVP | 18.33 | 26.92 | 64.71 | 18 | 7.69 | 12 | 56 | 29.09 |
| LGHP | 25 | 30.77 | 58.82 | 10 | 26.92 | 4 | 100 | 36.50 |
| VGGFace | 83.33 | 51.92 | 68.63 | 32 | 92.31 | 50 | 100 | 68.31 |

TABLE VI: Confusion matrix of average recognition rate using VGGFace descriptor with L1 distance over proposed IIITS_MFace dataset. The True Positive Values are highlighted in bold.

| Subjects | Subject1 | Subject2 | Subject3 | Subject4 | Subject5 | Subject6 | Subject7 |
|---|---|---|---|---|---|---|---|
| Sub1 | **50** | 3 | 4 | 1 | 2 | 0 | 0 |
| Sub2 | 8 | **27** | 2 | 1 | 12 | 2 | 0 |
| Sub3 | 0 | 1 | **35** | 4 | 4 | 2 | 5 |
| Sub4 | 9 | 4 | 20 | **16** | 1 | 0 | 0 |
| Sub5 | 0 | 0 | 1 | 0 | **48** | 3 | 0 |
| Sub6 | 6 | 1 | 6 | 0 | 12 | **25** | 0 |
| Sub7 | 0 | 0 | 0 | 0 | 0 | 0 | **50** |

TABLE VII: The average recognition rate using each descriptor with different distances over proposed IIITS_MFace dataset. The top value in a row is highlighted in bold face.

| Descriptor | Euclidean Distance | L1 Distance | Cosine Distance | Emd Distance | Chi-square Distance |
|---|---|---|---|---|---|
| LBP | 25.94 | **27.42** | 27.38 | 18.78 | 26.66 |
| LTP | 22.73 | 24.95 | 22.73 | 22.39 | **25.57** |
| LDP | 19.01 | 18.42 | 20.36 | **22.87** | 22.85 |
| LDGP | 23.88 | 22.68 | **25.48** | 16.08 | 21.68 |
| SLBP | 30.27 | 30.35 | 29.99 | 25.96 | **32.42** |
| LVP | 23.67 | 29.09 | 26.27 | 24.07 | **30.22** |
| LGHP | 29.09 | **36.50** | 33.31 | 23.84 | 35.38 |
| VGGFace | 62.58 | 68.31 | 68.11 | 36.55 | **69.39** |

TABLE VIII: A comparison of proposed IIITS_MFace dataset with AT&T, AR, Yale, LFW, PaSC and PubFig datasets. Here, 'Y', 'N' and 'P' represent the presence, absence and partial presence of effects like Non-frontal (NoFront), Masked, Occlusion (Occl), Mixed-spectacle (MixSpec), Illumination variation (IllVar), Extreme-Illumination (ExtIll), Background Variation (BackVar), and MutiFace VGGFace. The last row presents the accuracy in % using VGGFace CNN descriptor over each database using L1 distance measure.

| Traits | AT&T | AR | Yale | LFW | PaSC | PubFig | IIITS_MFace (Ours) |
|---|---|---|---|---|---|---|---|
| NoFront | N | N | N | Y | Y | Y | Y |
| Masked | N | Y | N | N | N | N | Y |
| Occl. | N | N | N | Y | N | N | Y |
| MixSpec | Y | Y | N | Y | N | N | Y |
| IllVar | N | N | Y | Y | Y | Y | Y |
| ExtIll | N | N | Y | N | N | N | N |
| BackVar | N | N | N | N | Y | Y | Y |
| MultiFace | N | N | N | N | N | P | Y |
| VGGFace Result (%) | 100 | 89.98 | 76.56 | 88.37 | 85.45 | 86.73 | 68.31 |

one with 68.31% average recognition rate among all the descriptors. Among hand-crafted descriptors, the LGHP is the best performing descriptor. Whereas, the LDP is the least performing descriptor because it is more suited to the frontal faces. The performance of most of the descriptor is better for Subject 7 because it is the only female subject. All the descriptors are failed to perform well in case of Subject 4 due to the following reasons: a) the number of faces in gallery set corresponding to subject 4 is just 17, b) all faces of subject 4 in the gallery set are unmasked and with spectacles, and c) the faces of subject 4 in probe set are mixed with

huge amount of pose variations, with/without spectacles and masked/unmasked. Overall, despite of being recent, well-known and highly discriminative, these face descriptors are failed to perform well over the proposed face dataset. Table VI illustrates the confusion matrix over proposed dataset obtained using the VGGFace CNN descriptor. It can be noted that most of the Subject 4 and 6 probe faces are recognized as the Subject 1, 3 ad 5 due to the huge amount of illumination change in the probe faces of Subject 4 and 6 as compared to the gallery faces. Subjects 1, 3 and 5 are also facing the problems like illumination, background, occlusion and

masking.

In order to find out which distance is better suited for the proposed IIITS_MFace dataset, we have conducted an experiment with different distance measures such as Euclidean (Eucld), L1, Cosine, Earth Movers Distance (Emd) and Chisq (Chi-square). The average recognition rate using all the descriptors are presented in Table VII. It can be noted that the Chi-square distance is performing well with LTP, SLBP, LVP and VGGFace descriptors. The Euclidean distance is not recommended to be used for the proposed dataset. Though, we have used L1 distance in other experiments, the best result (i.e., 69.39% accuracy) is obtained using VGGFace descriptor using Chi-square distance.

There are challenging datasets available in the literature with challenges like different side poses, occluded faces, varying light intensities, etc. These datasets are discussed in the Introduction section. However, the proposed IIITS_MFace dataset is much more challenging compared to the other existing face datasets such as AT&T, AR, Yale, LFW, PaSC and PubFig, as depicted in Table VIII. The result of VGGFace descriptor is lowest over the proposed IIITS_MFace dataset, which shows its difficulty and robustness.

From the experimental results, we can deduce that the proposed IIITS_MFace dataset is more challenging compared to the existing face datasets even for the deep learned VGGface descriptor, which makes it more realistic for the experiments to meet the real world scenario.

## V. CONCLUSION

A multi-face challenging IIITS_MFace dataset is proposed in this paper to validate the performance of hand-crafted local descriptors as well as deep learned CNN descriptor against the different kind of variations. The difficulties like pose, illumination, occlusion, masking, spectacle, background etc. are present in the dataset. The recent state-of-the-art face image descriptors such as LBP, LGHP, VGGFace etc. are used to test the complexity of the IIITS_MFace dataset. The results in terms of the average recognition rate support the challenges present in the dataset as the best performing VGGFace CNN descriptor achieved only 69.39% of accuracy in best setting. In general, the VGGFace CNN descriptor is very discriminative and performs reasonably good for face recognition. Several distance measures are also tested and found that the Chi-square distance is better suited for this dataset. In future, the number of subjects and number of samples in the dataset may be increased to facilitate applying some deeper neural network architecture for more robust training.

## ACKNOWLEDGMENT

The authors would like to thank all the individuals who have been involved in the process of data collection. Special thanks to Kanv Kumar and Naveen Thella for capturing the images.

## REFERENCES

[1] M. Pietikäinen, A. Hadid, G. Zhao, and T. Ahonen, "Face analysis using still images," in *Computer Vision Using Local Binary Patterns*. Springer, 2011, pp. 151–168.
[2] M. Bereta, W. Pedrycz, and M. Reformat, "Local descriptors and similarity measures for frontal face recognition: a comparative analysis," *Journal of Visual Communication and Image Representation*, vol. 24, no. 8, pp. 1213–1231, 2013.
[3] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 12, pp. 1349–1380, 2000.
[4] D. Huang, C. Shan, M. Ardabilian, Y. Wang, and L. Chen, "Local binary patterns and its application to facial image analysis: a survey," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 41, no. 6, pp. 765–781, 2011.
[5] B. Yang and S. Chen, "A comparative study on local binary pattern (lbp) based face recognition: Lbp histogram versus lbp image," *Neurocomputing*, vol. 120, pp. 365–379, 2013.
[6] E. Learned-Miller, G. B. Huang, A. RoyChowdhury, H. Li, and G. Hua, "Labeled faces in the wild: A survey," in *Advances in face detection and facial image analysis*. Springer, 2016, pp. 189–248.
[7] "AT&T face database," http://www.cl.cam.ac.uk/research/dtg/attarchive/-facedatabase.html.
[8] "AR face database," http://www2.ece.ohio-state.edu/~aleix/ARdatabase.html.
[9] K. Lee, J. Ho, and D. Kriegman, "Acquiring linear subspaces for face recognition under variable lighting," *IEEE Trans. Pattern Anal. Mach. Intelligence*, vol. 27, no. 5, pp. 684–698, 2005.
[10] "LFW face database," http://vis-www.cs.umass.edu/lfw/.
[11] "Croppedlfw face database," http://conradsanderson.id.au/lfwcrop/.
[12] J. R. Beveridge, P. J. Phillips, D. S. Bolme, B. A. Draper, G. H. Givens, Y. M. Lui, M. N. Teli, H. Zhang, W. T. Scruggs, K. W. Bowyer *et al.*, "The challenge of face recognition from digital point-and-shoot cameras," in *Biometrics: Theory, Applications and Systems (BTAS), 2013 IEEE Sixth International Conference on*. IEEE, 2013, pp. 1–8.
[13] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, vol. 1. IEEE, 2001, pp. I–I.
[14] N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar, "Attribute and simile classifiers for face verification," in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 365–372.
[15] T. Ahonen, A. Hadid, and M. Pietikainen, "Face description with local binary patterns: Application to face recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 12, pp. 2037–2041, 2006.
[16] X. Tan and B. Triggs, "Enhanced local texture feature sets for face recognition under difficult lighting conditions," *IEEE transactions on image processing*, vol. 19, no. 6, pp. 1635–1650, 2010.
[17] B. Zhang, Y. Gao, S. Zhao, and J. Liu, "Local derivative pattern versus local binary pattern: face recognition with high-order local pattern descriptor," *IEEE transactions on image processing*, vol. 19, no. 2, pp. 533–544, 2010.
[18] S. Chakraborty, S. K. Singh, and P. Chakraborty, "Local directional gradient pattern: a local descriptor for face recognition," *Multimedia Tools and Applications*, vol. 76, no. 1, pp. 1201–1216, 2017.
[19] S. R. Dubey and S. Mukherjee, "Ldop: Local directional order pattern for robust face retrieval," *arXiv preprint arXiv:1803.07441*, 2018.
[20] K. Jeong, J. Choi, and G.-J. Jang, "Semi-local structure patterns for robust face detection," *IEEE Signal Processing Letters*, vol. 22, no. 9, pp. 1400–1403, 2015.
[21] K.-C. Fan and T.-Y. Hung, "A novel local pattern descriptorlocal vector pattern in high-order derivative space for face recognition," *IEEE transactions on image processing*, vol. 23, no. 7, pp. 2877–2891, 2014.
[22] S. Chakraborty, S. Singh, and P. Chakraborty, "Local gradient hexa pattern: A descriptor for face recognition and retrieval," *IEEE Transactions on Circuits and Systems for Video Technology*, 2016.
[23] O. M. Parkhi, A. Vedaldi, A. Zisserman *et al.*, "Deep face recognition." in *BMVC*, vol. 1, no. 3, 2015, p. 6.
[24] "Pairwise distance between two sets of observations," http://in.mathworks.com/help/stats/pdist2.html.