# Landmarks-assisted Collaborative Deep Framework for Automatic 4D Facial Expression Recognition

Muzammil Behzad, Nhat Vo, Xiaobai Li and Guoying Zhao

Center for Machine Vision and Signal Analysis (CMVS), University of Oulu, Finland

Email: {muzammil.behzad, nhat.vo, xiaobai.li, guoying.zhao}@oulu.fi

*Abstract*— We propose a novel landmarks-assisted collaborative end-to-end deep framework for automatic 4D FER. Using 4D face scan data, we calculate its various geometrical images, and afterwards use rank pooling to generate their dynamic images encapsulating important facial muscle movements over time. As well, the given 3D landmarks are projected on a 2D plane as binary images and convolutional layers are used to extract sequences of feature vectors for every landmark video. During the training stage, the dynamic images are used to train an end-to-end deep network, while the feature vectors of landmark images are used train a long short-term memory (LSTM) network. The finally improved set of expression predictions are obtained when the dynamic and landmark images collaborate over multi-views using the proposed deep framework. Performance results obtained from extensive experimentation on the widely-adopted BU-4DFE database under globally used settings prove that our proposed collaborative framework outperforms the state-of-the-art 4D FER methods and reach a promising classification accuracy of 96.7% demonstrating its effectiveness.

## I. INTRODUCTION

Facial expressions (FEs) are important cues in understanding human emotions during their social communication. To better facilitate the understanding and analysis of such FEs, many researchers proposed facial expression recognition (FER) systems using the state-of-the-art computer vision based human-machine interaction methods. As a result, the research community has witnessed tremendous surge in FER systems towards potentially significant application areas like psychology, security, education, bio-medical and computing technology. The study concluded by Ekman and Friesen [1] serves as a pioneer contribution in this field dating back to 1970s. This work presented the six globally-adopted human facial expressions which are happiness, anger, sadness, fear, disgust and surprise.

In the past few years, several machine learning methods were presented for recognizing facial expressions with the help of static or dynamic 2D images. Despite promising contributions, however, emotion recognition still remains a challenging problem due to the sensitivity of 2D images towards lighting conditions, pose variations and occlusions [2]. This is why 2D based methods are not fully stable and could not potentially contribute to real-world applications beyond

a certain point. As a rival, 3D point clouds rescued and motivated novel FER directions via trending high-resolution and high-speed 3D data acquisition equipment. Although the data processing becomes complex, the significantly increased amount of data in terms of the facial deformation patterns over the depth axis considerably help the deep learning models to learn patterns effectively for automatic FER.

Importantly, although each facial expression is in fact a combination of different muscle movements ordered in a particular way, consequently triggering facial deformations [3], such compact cues are better captured in geometrical domain [4]. This is why the 3D face scans are quite convenient in representing such deformations, and therefore, predicting the emotions. The collection of different large-size and complex databases with various terabytes of data has supported such research on FER using 3D face scans. In this regard, the release of commonly known BU-3DFE [5] and Bosphorus [6] served as one of the pioneering datasets for investigating FER via static 3D data. The dynamic 3D data (referred as 4D), such as the BU-4DFE dataset [7], allows to perform 4D FER by fetching facial deformations both in the depth geometry and over time.

Contrary to the 3D FER methods which only rely on the static data at hand [8]–[11], 4D face scans enable deep networks to learn effectively for better analyzing and predicting facial expressions. In this regard, Sun *et al*. [12] and Yin *et al*. [7] proposed to work around Hidden Markov Models (HMM) to learn the facial muscle patterns over time. Similarly to benefit from local facial patterns, Drira and Amor [13], [14] introduced a deformation vector field mainly based on Riemannian analysis and combined with random forest. In another attempt, Sandbach *et al*. [15] represented 3D frames and its neighbors as Free-Form Deformation (FFD) and subsequently used HMM and GentleBoost as classifiers. Using the traditional Support Vector Machine (SVM), the authors represented geometrical coordinates and its normal as feature vectors [16], and as dynamic Local Binary Patterns (LBP) in an extended work [17]. Likewise, a spatio-temporal LBP-based feature was proposed in [18] to extract features from polar angles and curvatures.

Yao *et al*. [19] applied Multiple Kernel Learning (MKL) by using the scattering operator [20] on 4D face scans to produce effective feature representations. In a similar attempt to recognize FEs, statistical shape model with global and local constraints were proposed in [21]. The authors claimed that local shape index and global face shape can together

help build a desirable FER system. A much effective deep network was proposed by Li *et al.* [22] to automate 4D FER using a dynamic geometrical image network. In this work, geometrical images were generated after the differential quantities were estimated. The final prediction step involved fusing the predicted scores from different geometrical images. Bejaoui *et al.* [23] recently proposed a sparse coding-based representation of LBP difference. They extracted a unified set of geometric and appearance features via Mesh-Local Binary Pattern Difference (mesh-LBPD), combined them into a compact representation via covariance matrices, and then applied sparse coding for effective 3D/4D FER.

Despite these attempts to automate FER via 4D data, we believe that the facial deformations should be appropriately extracted from the spatio-temporal 4D data for better network learning rather than simply tuning multiple parts of a deep network. Consequently, in this paper, we aim to fetch such deformations jointly from multi-views and some geometrical domains to propose a landmarks-assisted collaborative end-to-end deep framework for automatic 4D FER (LC4D). We project every 3D face scan to extract various geometrical images such as depth images and texture images. For a robust representation of the facial features, we aim to extract the features across various multi-views. To encode the muscle movements of different expressions from the temporal domain, we apply rank pooling to compute the dynamic images of the 4D data. As well, we represent the stored facial movements from the 3D landmarks by projecting them on a 2D plane. Afterwards, activations of these sequences of landmarks are computed using convolution layers, which are then trained on an LSTM network. Using decision-level fusion, the landmarks then collaborate with the dynamic images trained over another deep network for an improved FER by highlighting the correct classification probabilities. To our best knowledge, this is the first landmarks-assisted collaborative deep framework for 4D FER using multi-views.

The rest of the paper is sectioned as follows: Section II explains our proposed LC4D method for 4D FER in detail. In Section III, the experimental results of our framework are reported. Finally, Section IV concludes the paper.

## II. PROPOSED AUTOMATIC 4D FER METHOD

In this section, we discuss in detail the working mechanism of our proposed LC4D deep framework. First, we explain the filtering step to remove the unwanted and noisy mesh components in the given 3D point clouds. Then, we discuss the computation of different geometrical images in multi-views. Third, we elaborate the landmarks-assisted mechanism. Finally, the collaborative scheme for 4D FER is presented. An overview of our method is shown in Fig. 1.

### A. Pre-processing

The 3D point clouds from BU-4DFE contain noise and unwanted components like outliers, hair, and non-facial regions, which causes a deep model incapable of learning effectively, and therefore, should be filtered out.

Subsequently, we pre-process each 3D point cloud from $N$ 4D samples individually to combat such outliers. We define

$$I^{4D} = \{I_{nt}^{3D}\}, \ \forall t = \{1, 2, ..., T_n\} \text{ and } \forall n = \{1, 2, ..., N\}, \quad (1)$$

where $I^{4D}$ refers to 4D samples, and $I_{nt}^{3D}$ is the $n$th point cloud at $t$th frame. Note that (1) $\Rightarrow |I^{4D}| = N$, and $|I_{nt}^{3D}| = T_n$. The mesh with $M$ vertices can be denoted as

$$\mathbf{m} = [\mathbf{v}_1^T, ..., \mathbf{v}_M^T]^T = [x_1, y_1, z_1, ..., x_M, y_M, z_M]^T, \quad (2)$$

where $\mathbf{v}_j = [x_j, y_j, z_j]^T$ are coordinates of the $j$th vertex, and $\mathbf{m}_t$ is a mesh at $t$th frame such that $\mathbf{m}_n = \{\mathbf{m}_{nt}\}\forall n$. To crop the unwanted regions via given landmarks, we remove everything beyond the facial-border landmarks. For the head-hair, we trim regions above the forehead by using a threshold which is a fractional distance from eyebrows to the tip of nose. This is denoted as

$$\overline{I_{nt}^{3D}} = \eta_c(I_{nt}^{3D}), \quad (3)$$

where $\overline{I_{nt}^{3D}}$ is the cropped face and $\eta_c(.)$ is the cropping operation. This also updates the *mesh-space* in (2) as follows:

$$\overline{\mathbf{m}} = \eta_c(\mathbf{m}) = [\mathbf{v}_1^T, ..., \mathbf{v}_{\overline{M}}^T]^T = [x_1, y_1, z_1, ..., x_{\overline{M}}, y_{\overline{M}}, z_{\overline{M}}]^T, \quad (4)$$

where $\overline{\mathbf{m}}$ is the set of updated vertices such that $\overline{\mathbf{m}} \subseteq \mathbf{m}$ and $\overline{M} \leq M$. The filtering steps in (3) and (4) remove all the outliers that would potentially disturb the training process.

### B. Geometrical Images over Multi-views

Geometrical images provide efficient feature mapping from 3D to 2D [24]. Therefore, after pre-processing, we compute the depth images (DPI) as $f_D : \overline{I^{3D}} \rightarrow I_D$, and the texture images as $f_T : \overline{I^{3D}} \rightarrow I_T$ from the filtered meshes via 3D to 2D rendering, where $f_D$ and $f_T$ denote the function mapping to depth image $I_D \in \mathbb{R}^{K^2}$, and texture image $I_T \in \mathbb{R}^{K^2}$, respectively, where $K^2$ is the number of pixels. For sharp facial details, the contrast-limited adaptive histogram equalization is applied on DPIs to get enhanced-depth images (E-DPIs) as $I_{ED} = \eta_s(I_D)$, where $\eta_s(.)$ refers to the sharpening operator. Importantly, we also generate these images in alignment profiles from right-to-frontal-to-left for an effective collaboration at a later stage in the network.

Once we compute the pre-processed images from different domains (*e.g.*, texture and depth) over multi-views, as depicted in Fig. 1, the next step is to fetch and represent the variations of the facial deformations from the temporal domain. One optimal choice is to perform rank pooling over the projected 2D sequences of geometrical images to obtain their dynamic images [25]. Consequently, we perform rank pooling to represent the temporal dynamics of entire videos into a single RGB image. We compute the dynamic images for all 4D samples using the projected sequences over multi-views. All the extracted cross-domain dynamic images, shown in Fig. 1, incorporate the spatio-temporal patterns effectively which is favorable for training a deep network. It is worth mentioning that this idea is different from the one proposed in [22]. This is because they used geometrical images independently, while we use the dynamic
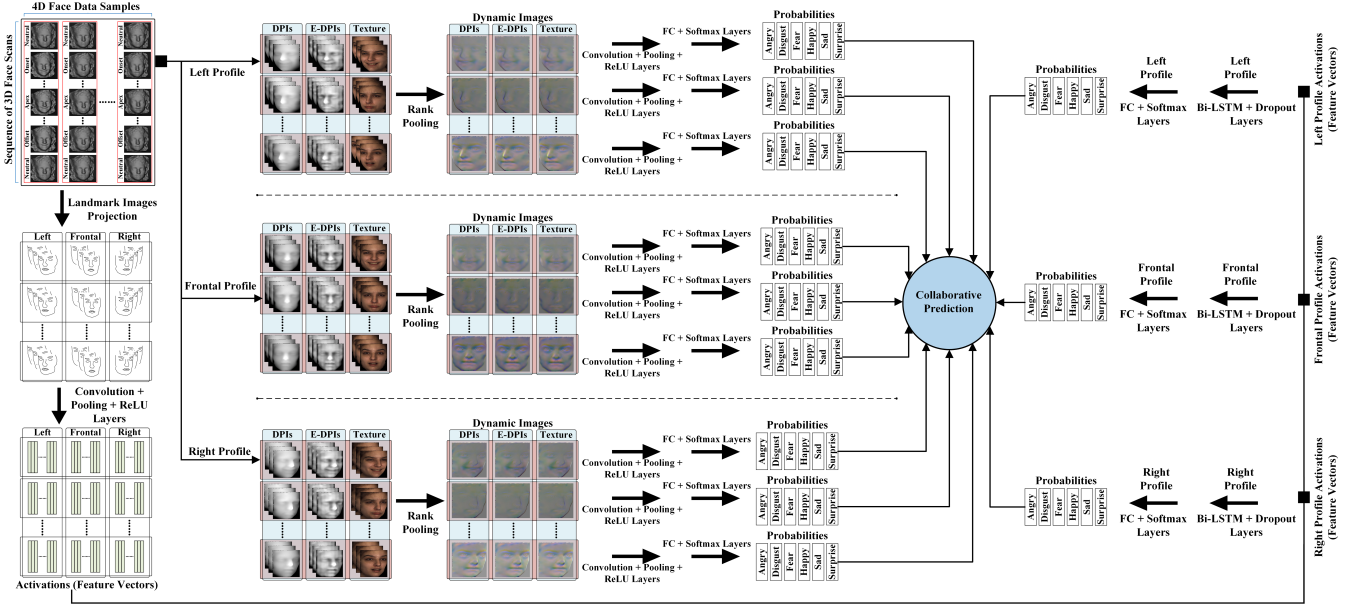
Fig. 1. The proposed LC4D method for 4D FER.

images of different domains to collaborate over multi-views. More importantly, as explained in the subsequent section, we also deploy a deep landmarks-based network ultimately improving the classification scores for 4D FER, which, to the best of our knowledge, has never been reported in the literature before.

### C. Landmarks-assisted Learning

Since landmarks encode key facial points that potentially represent an expressed emotion, our landmarks-assisted approach significantly benefits the proposed framework for 4D FER. To do so, we first similarly pre-process all the given 83 landmarks for each frame. Then, we project them on a 2D plane across various multi-views for representing the facial deformations stored in landmarks over time as binary images. With the sequences of such projected landmark images as input, we extract its activations by using convolution and pooling layers. For this purpose specifically, we used a pre-trained GoogLeNet [26] to convert the videos of landmark images into sequences of feature vectors containing appropriate feature presentation of each frame as a vector.

Finally, we create a long short-term memory (LSTM) network with a sequence input layer, Bi-LSTM layer with 2000 hidden units, 50% dropout layer followed by FC, Softmax and classification layer. This is first trained on the sequences of these feature vectors, and is then used to predict the expressions. Note that the parameterized model in [27] just use landmarks to train and then predict expressions in a straight forward manner. Conversely, we propose a deep framework where the extracted sequences of feature vectors of the projected landmark images across multi-views are used to train an LSTM network first, and are then used in collaboration with the rest of the framework for an improved expression prediction as outlined in the next section.

### D. Collaborative Prediction

For improved predictions, a final collaborative step is performed to tailor the voting scores of different expressions using various collaborative elements. Specifically, since our network is trained after multiple Convolution+Pooling+ReLU layers, we jointly utilize the expression probabilities from different geometrical domains over multi-views. While doing so, the contributions from landmarks-assisted network are equally respected to have updated and much improved predictions. This is because even though the patterns from dynamic images already discriminate different expressions, the landmarks-assisted approach via multi-view further enhances the likelihood of correct predictions. For six expressions, the predicted and collaboratively-updated probabilities, respectively, are as follows:

$$C = [\mathbf{c}_1^T, ..., \mathbf{c}_N^T]^T, \text{ and} \tag{5}$$

$$C(n,l) = \frac{1}{|\Theta|} \sum_{\forall \theta \in \Theta} [C_{DI}(n,l_\theta) + C_{LI}(n,l_\theta)], \ \forall n,l. \tag{6}$$

Here, $\mathbf{c} = \{\rho_l\}$ for $l = \{1,...,6\}$, is the predicted probabilities of six expressions for $n$th sample, $\Theta$ is all view angles over which multi-views are collected, while the subscript $DI$ and $LI$ refers to dynamic and landmark images, respectively. The finally updated predictions $F(n)$ are computed as maximum of all the expression probabilities

$$F(n) = \max\{C(n)\}, \ \forall n = \{1,2,...,N\}. \tag{7}$$

### III. EXPERIMENTAL RESULTS AND ANALYSIS

With the extensive training, we evaluate our experimental results to analyze the improvement in prediction performance based on our proposed LC4D method for FER. We opted the generally used BU-4DFE dataset for our experiments. This dataset contains 58 females and 43 males (a total of 101

subjects) each having all six human facial expressions, i.e., happiness, anger, sadness, fear, disgust and surprise. Every expression of a subject contains dynamic 3D data of raw face scans with a frame rate of 25 frames per second (fps) lasting approximately 3 to 4 seconds.

For having a fair comparison of results with other state-of-the-art methods, we follow [22] and choose similar experimental settings. In particular, we employ a 10-fold subject-independent cross-validation (10-CV), and use a 60-20-20 split of the data for training, validation and testing, respectively over five iterations. Instead of using key-frames [19] or employing sliding windows [15], we use entire 3D sequences. For the training using dynamic images, we use the pre-trained VGGNet [29] as the deep network, and therefore resize our images to $224 \times 224$. For the landmarks, we use the pre-trained GoogLeNet to extract feature vectors from activations, and then train an LSTM network from scratch on the extracted activations. All of our experiments are carried out on a GP100GL GPU (Tesla P100-PCIE), and the overall training time takes approximately one day.

We show in Table I the extensive comparisons of the classification accuracies calculated at various collaboration stages of our proposed method. As shown, the multi-views significantly help in revealing potential patterns and assist both landmarks as well as the dynamic images for better learning. It can be seen that promising results are achieved when all collaborators help in prediction. Similarly, we also show the confusion matrix of our experiments in Fig. 2. Despite angry and disgust being some error cases due to their similarities, Fig. 2 indicates that our method correctly predicts emotions most of the time showing its effectiveness.

Finally, in Table II, we present accuracies achieved on the BU-4DFE dataset by our proposed framework and sev-

TABLE II
ACCURACY (%) COMPARISON WITH THE STATE-OF-THE-ART ON THE
BU-4DFE DATASET.

| Method | Experimental Settings | Accuracy |
|---|---|---|
| 2012 - Sandbach et al. [15] | 6-CV, Sliding window | 64.60 |
| 2011 - Fang et al. [17] | 10-CV, Full sequence | 75.82 |
| 2015 - Xue et al. [30] | 10-CV, Full sequence | 78.80 |
| 2010 - Sun et al. [12] | 10-CV, - | 83.70 |
| 2016 - Zhen et al. [4] | 10-CV, Full sequence | 87.06 |
| 2018 - Yao et al. [19] | 10-CV, Key-frame | 87.61 |
| 2012 - Fang et al. [16] | 10-CV, - | 91.00 |
| 2018 - Li et al. [22] | 10-CV, Full sequence | 92.22 |
| 2014 - Ben Amor et al. [14] | 10-CV, Full sequence | 93.21 |
| 2016 - Zhen et al. [28] | 10-CV, Full sequence | 94.18 |
| 2019 - Bejaoui et al. [23] | 10-CV, Full sequence | 94.20 |
| 2018 - Zhen et al. [28] | 10-CV, Key-frame | 95.13 |
| **Ours** | 10-CV, Full sequence | **96.70** |

eral state-of-the-art methods [4], [12], [14]–[17], [19], [22], [28], [30]. As shown in the table, our proposed method outperforms the existing methods while predicting the correct expression during classification. This is because of the extensively collaborative nature of our proposed framework in which the predictions are refined when the probability scores are updated from neighboring resources. The refinement in probabilities come from the fact that facial deformation patterns are well-captured in the geometric domain and its temporal movements are caught in the dynamic images. Importantly, the significant amount of assistance received from landmarks images also helped in making a reliable classification decision. Consequently, our LC4D framework reached an accuracy of 96.7% for 4D FER.

TABLE I
FER ACCURACY REACHED ON THE BU-4DFE DATASET. [LP = LEFT
PROFILE, RP = RIGHT PROFILE, FP = FRONTAL PROFILE]

| Collaborator(s) | Multi-view Profile(s) | FER Accuracy (%) |
|---|---|---|
| Landmark Images | LP | 75.40 |
| | FP | 78.70 |
| | RP | 77.60 |
| | RP + FP | 85.50 |
| | LP + FP | 84.20 |
| | RP + LP | 83.80 |
| | RP + FP + LP | **88.80** |
| Dynamic Images | LP | 78.30 |
| | FP | 80.20 |
| | RP | 79.20 |
| | RP + FP | 83.20 |
| | LP + FP | 82.10 |
| | RP + LP | 81.30 |
| | RP + FP + LP | **84.70** |
| Landmark and Dynamic Images | LP | 83.40 |
| | FP | 91.40 |
| | RP | 87.70 |
| | RP + FP | 93.60 |
| | LP + FP | 92.10 |
| | RP + LP | 88.80 |
| | RP + FP + LP | **96.70** |



Fig. 2. Confusion matrix of predicting expressions on BU-4DFE database.

## IV. CONCLUSIONS

We proposed a 4D FER method via landmarks-assisted collaborative end-to-end deep framework. In this framework, different geometrical images were extracted first and were later used to capture facial movements over time in terms of compact dynamic images. Additionally, efficient feature representations were extracted from landmark images that were then used to train an LSTM network. An effective collaboration step performed over multi-views served as an added advantage of our deep framework. With a promising accuracy of 96.7%, our method outperformed the state-of-the-art 4D FER methods in terms of classification accuracy.

## REFERENCES

[1] P. Ekman and W. V. Friesen, "Constants Across Cultures in the Face and Emotion," *Journal of Personality and Social Psychology*, vol. 17, no. 2, pp. 124–129, 1971.

[2] B. Fasel and J. Luettin, "Automatic Facial Expression Analysis: A Survey," *Pattern Recognition*, vol. 36, no. 1, pp. 259–275, 2003.

[3] T. Fang, X. Zhao, O. Ocegueda, S. K. Shah, and I. A. Kakadiaris, "3D Facial Expression Recognition: A Perspective on Promises and Challenges," in *Face and Gesture*, pp. 603–610, 2011.

[4] Q. Zhen, D. Huang, Y. Wang, and L. Chen, "Muscular Movement Model-based Automatic 3D/4D Facial Expression Recognition," *IEEE Transactions on Multimedia*, vol. 18, no. 7, pp. 1438–1450, 2016.

[5] L. Yin, X. Wei, Y. Sun, J. Wang, and M. J. Rosato, "A 3D Facial Expression Database for Facial Behavior Research," in *International Conference on Automatic Face and Gesture Recognition*, pp. 211–216, 2006.

[6] A. Savran, N. Alyü, H. Dibeklioğlu, O. Çeliktutan, B. Gökberk, B. Sankur, and L. Akarun, "Bosphorus Database for 3D Face Analysis," in *European Workshop on Biometrics and Identity Management*, pp. 47–56, 2008.

[7] L. Yin, X. Chen, Y. Sun, T. Worm, and M. Reale, "A High-resolution 3D Dynamic Facial Expression Database," in *International Conference on Automatic Face and Gesture Recognition*, pp. 1–6, 2008.

[8] X. Zhao, D. Huang, E. Dellandrea, and L. Chen, "Automatic 3D Facial Expression Recognition based on a Bayesian Belief Net and a Statistical Facial Feature Model," in *International Conference on Pattern Recognition*, pp. 3724–3727, 2010.

[9] H. Li, L. Chen, D. Huang, Y. Wang, and J. Morvan, "3D Facial Expression Recognition via Multiple Kernel Learning of Multi-scale Local Normal Patterns," in *International Conference on Pattern Recognition*, pp. 2577–2580, 2012.

[10] Q. Zhen, D. Huang, Y. Wang, and L. Chen, "Muscular Movement Model based Automatic 3D Facial Expression Recognition," in *International Conference on MultiMedia Modeling*, pp. 522–533, 2015.

[11] H. Li, H. D. D., Huang, Y. Wang, X. Zhao, J. M. M. L., and Chen, "An Efficient Multimodal 2D+3D Feature-based Approach to Automatic Facial Expression Recognition," *Computer Vision and Image Understanding*, vol. 140, pp. 83–92, 2015.

[12] Y. Sun, X. Chen, M. Rosato, and L. Yin, "Tracking Vertex Flow and Model Adaptation for Three-dimensional Spatiotemporal Face Analysis," *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, vol. 40, no. 3, pp. 461–474, 2010.

[13] H. Drira, B. B. Amor, M. Daoudi, A. Srivastava, and S. Berretti, "3D Dynamic Expression Recognition based on a Novel Deformation Vector Field and Random Forest," in *International Conference on Pattern Recognition*, pp. 1104–1107, 2012.

[14] B. B. Amor, H. Drira, S. Berretti, M. Daoudi, and A. Srivastava, "4D Facial Expression Recognition by Learning Geometric Deformations," *IEEE Transactions on Cybernetics*, vol. 44, no. 12, pp. 2443–2457, 2014.

[15] G. Sandbach, S. Zafeiriou, M. Pantic, and D. Rueckert, "Recognition of 3D Facial Expression Dynamics," *Image and Vision Computing*, vol. 30, no. 10, pp. 762–773, 2012.

[16] T. Fang, X. Zhao, O. Ocegueda, S. K. Shah, and I. A. Kakadiaris, "3D/4D Facial Expression Analysis: An Advanced Annotated Face Model Approach," *Image and Vision Computing*, vol. 30, no. 10, pp. 738 – 749, 2012.

[17] T. Fang, X. Zhao, S. K. Shah, and I. A. Kakadiaris, "4D Facial Expression Recognition," in *International Conference on Computer Vision Workshops*, pp. 1594–1601, 2011.

[18] M. Reale, X. Zhang, and L. Yin, "Nebula Feature: A Space-time Feature for Posed and Spontaneous 4D Facial Behavior Analysis," in *International Conference and Workshops on Automatic Face and Gesture Recognition*, pp. 1–8, 2013.

[19] Y. Yao, D. Huang, X. Yang, Y. Wang, and L. Chen, "Texture and Geometry Scattering Representation-based Facial Expression Recognition in 2D+3D Videos," *ACM Transactions on Multimedia Computing, Communications, and Applications*, 2018.

[20] J. Bruna and S. Mallat, "Invariant Scattering Convolution Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1872–1886, 2013.

[21] D. Fabiano and S. Canavan, "Spontaneous and Non-spontaneous 3D Facial Expression Recognition using a Statistical Model with Global and Local Constraints," in *International Conference on Image Processing*, pp. 3089–3093, 2018.

[22] W. Li, D. Huang, H. Li, and Y. Wang, "Automatic 4D Facial Expression Recognition using Dynamic Geometrical Image Network," in *International Conference on Automatic Face Gesture Recognition*, pp. 24–30, 2018.

[23] H. Bejaoui, H. Ghazouani, and W. Barhoumi, "Sparse coding-based representation of lbp difference for 3d/4d facial expression recognition," *Multimedia Tools and Applications*, vol. 78, pp. 22773–22796, Aug 2019.

[24] X. Yang, D. Huang, Y. Wang, and L. Chen, "Automatic 3D Facial Expression Recognition using Geometric Scattering Representation," in *International Conference and Workshops on Automatic Face and Gesture Recognition*, pp. 1–6, 2015.

[25] H. Bilen, B. Fernando, E. Gavves, and A. Vedaldi, "Action Recognition with Dynamic Image Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2799–2813, 2018.

[26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.

[27] D. Fabiano and S. Canavan, "Deformable synthesis model for emotion recognition," in *2019 14th IEEE International Conference on Automatic Face Gesture Recognition (FG 2019)*, pp. 1–5, May 2019.

[28] Q. Zhen, D. Huang, H. Drira, B. B. Amor, Y. Wang, and M. Daoudi, "Magnifying Subtle Facial Motions for Effective 4D Expression Recognition," *IEEE Transactions on Affective Computing*, 2018.

[29] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," in *British Machine Vision Conference*, 2015.

[30] M. Xue, A. Mian, W. Liu, and L. Li, "Automatic 4D Facial Expression Recognition using DCT Features," in *Winter Conference on Applications of Computer Vision*, pp. 199–206, 2015.