# Few-Data Guided Learning Upon End-to-End Point Cloud Network for 3D Face Recognition

Yi Yu, Feipeng Da, and Ziyu Zhang

*Abstract*—3D face recognition has shown its potential in many application scenarios. Among numerous 3D face recognition methods, deep-learning-based methods have developed vigorously in recent years. In this paper, an end-to-end deep learning network entitled Sur3dNet-Face for point-cloud-based 3D face recognition is proposed. The network uses PointNet as the backbone, which is a successful point cloud classification solution but does not work properly in face recognition. Supplemented with modifications in network architecture and a few-data guided learning framework based on Gaussian process morphable model, the backbone is successfully modified for 3D face recognition. Different from existing methods training with a large amount of data in multiple datasets, our method uses Spring2003 subset of FRGC v2.0 for training which contains only 943 facial scans, and the network is well trained with the guidance of such a small amount of real data. Without fine-tuning on the test set, the Rank-1 Recognition Rate (RR1) is achieved as follows: 98.85% on FRGC v2.0 dataset and 99.33% on Bosphorus dataset, which proves the effectiveness and the potentiality of our method.

*Index Terms*—Deep learning, face recognition, point cloud.

## I. INTRODUCTION

WITH the development of deep learning and 3D measurement technology, 3D face recognition has shown its potential in many application scenarios. Several solutions to 3D face recognition have been proposed including feature-based, model-based, matching-based, and learning-based methods, among which learning-based methods are flourishing in recent years and have shown remarkable performance [1], [2].

However, whether deep learning can achieve good results, to a great extent, depends on training data. As is known, many articles use more and more data to train the network. Although these methods compare the recognition rate on the same dataset, they do not use the same training set, which makes the comparison unfair. As a result, we can only see that the reported recognition rate improves again and again, but cannot distinguish whether it is the effect of the method itself or caused by the increase of training data. Take the state-of-the-art method proposed in [3] as an example, six datasets are used in the training process including about $22K$ scans, and the real data used for training are ten times more than those for testing, which makes it easy to obtain satisfactory

recognition rate. However, in large-scale practical applications, the training set is usually smaller than the test set, and thus the recognition rate could be far below expectation, revealing that 3D face recognition is still an unsolved problem.

In this work, an end-to-end deep learning network entitled Sur3dNet-Face for point-cloud-based 3D face recognition is proposed. Taking advantage of a novel training framework upon Gaussian process morphable models (GPMM), and supplemented with a small amount of real data, the network is well trained.

The main contributions of this paper are as follows:

1) An end-to-end deep learning network for 3D face recognition is proposed to get face representations directly from 3D point clouds.

2) A few-data guided learning framework based on GPMM is established, upon that our network can be well trained with a small amount of real data.

3) The ablation study is analyzed to determine the parameters of the proposed network and the comparisons with other methods are conducted to prove the effectiveness and the potentiality of our method.

## II. RELATED WORK

### A. 3D Face Recognition

Numerous 3D face recognition methods have been developed, as reviewed in [1], [2], [4], including feature-based, model-based, matching-based, and learning-based methods.

The first three categories were widely studied in the early years. For example, Mian et al. [5] proposed a multimodal face recognition system, where the 3D and 2D faces are matched through modified ICP and SIFT descriptors respectively. Mohammadzade and Hatzinakos proposed Iterative Closest Normal Point (ICNP) [6] to match face surfaces. Liu et al. [7] developed the harmonic feature-based approach using energies in spherical harmonics at different frequencies. Elaiwat et al. [8] calculated Curvelet transform to extract features from semi-rigid regions. Lei et al. [9] proposed a keypoint-based Multiple Triangle Statistics (KMTS) method to handle pose variations.

As a latecomer, learning-based methods have gradually raised its head nowadays. Lei et al. [10] trained the Kernel Principal Component Analysis (KPCA) to extract feature representations. Song et al. [11] built 3D models to generate 2D images to improve the accuracy of 2D methods. Gilani and Mian [12] and Kim et al. [13] employed existing 2D deep neural networks to solve the 3D face recognition problem by projecting 3D surface into 2D space as depth map, azimuth map, and elevation map. Cai et al. [3] proposed a deep learning

technique based on facial component patches using depth map of the 3D face as the input of the traditional 2D network.

Some of these methods achieve decent accuracy, but the majority are still 2D-based networks, where 3D data are firstly projected into 2D images and then the traditional 2D networks are utilized to solve the problem.

### B. Point Cloud Network

There are many representations of 3D data, among which point clouds are widely used. Different from traditional deep learning networks upon 2D images, Charles et al. proposed PointNet [14] to directly handle point clouds, and the enhanced version PointNet++ [15] is established upon PointNet along with grouping and sampling techniques to synthesize both global and local features of point clouds.

Similarly, Li et al. proposed PointCNN [14] to learn $\mathcal{X}$-transform so that the convolution operator can work on un-ordered point clouds, and Komarichev et al. proposed A-CNN computing convolution directly on point clouds through annular convolution.

Graph Neural Network (GNN) is also introduced into point cloud networks. For example, DGCNN [16] uses EdgeConv, a graph-based operation, to carry out convolution in feature space, and DPAM [17] uses a graph network to take the place of sampling and grouping step in PointNet++.

These methods are designed for object classification and segmentation, though can also be applied to face recognition, the performance is much lower than face recognition methods mentioned in the previous subsection.

### C. Facial Data Generation

Whether deep learning can achieve good results, to a great extent, depends on training data. Different from 2D face datasets, the 3D face data are relatively inadequate for training a network.

Numerous 3D face generation methods have been developed. For example, Blanz and Vetter [18] proposed 3D face morphable model (3DMM) to model 3D faces and upon that Dou et al. [19] reconstructed 3DMM parameters with a deep neural network. Also, Lüthi et al. [20] proposed GPMM face model, a generalization of point distribution models.

Different from those model-based methods, GANFIT [21] and MMFace [22] utilize 2D faces to reconstruct 3D faces.

Gilani and Mian [12] and Kim et al. [13] also proposed data augmentation methods along with their face recognition networks to generate millions of 2D projected images specially for their networks from 3D faces on the ground that their networks use 2D images as the input.

### III. METHODOLOGY

In this section, we firstly introduce the architecture of our network, and then the training details are addressed. The overall procedure of the proposed method is shown in Fig. 1.

### A. Network Architecture

Different from the traditional learning-based 3D face recognition methods [12], [13], our method designs an end-to-end network directly inputting the coordinate of point cloud, so as to maintain the advantages of 3D data such as rotation invariance and transformation invariance. The output of our network is a feature vector, and the cosine distance between two feature vectors is calculated to reflect the probability that the two input faces are grabbed from the same subject.

Specifically, the forward process of our network can be represented as:

$$f = Sur3dNet\left(\Gamma\right) \tag{1}$$

where $\Gamma = \{x_1, x_2, \cdots, x_{N_0}\} \in \mathbb{R}^{N_0 \times 3}$ is the unordered input point cloud, $N_0$ denotes the number of points, $f \in \mathbb{R}^{256}$ is the output features.

The architecture of our proposed network is shown in Fig. 2 and the submodules in the figure will be introduced in the subsequent subsections.

### B. Normal Estimation

As is known, the normal vector is one of the most important attributes of point clouds. We calculate the normal vectors $n \in \mathbb{R}^{N_0 \times 3}$ of the input point cloud $\Gamma$ through Principal Component Analysis (PCA), during which the corresponding eigen values $e \in \mathbb{R}^{N_0 \times 1}$ can also be derived, which reflect the curvature of the surface. Therefore, the output of the normal estimation submodule is $\left[\begin{array}{ccc}\Gamma & n & e\end{array}\right] \in \mathbb{R}^{N_0 \times 7}$.

### C. Modified PointNet

PointNet [14] learns a function that maps a set of points to a feature vector, where multi-layer perceptrons (MLPs) are applied to every point individually before a max-pooling layer that aggregates features of all points to a global vector.

The backbone of our architecture is similar to PointNet, as is shown in Fig. 2(b), and the modifications are as follows.

*1) Ball Query With Physical Size:* The coordinates of all objects are normalized to the range of $(-1, 1)$ in PointNet. However, size is an important attribute of faces, which will be lost in the normalization process. To avoid this, we discard the normalization process and directly input the point cloud with original physical size with the unit of millimeter, and therefore, the radius of ball query in our network should also be measured in millimeters.

*2) Dithering Farthest Point Sampling (DFPS):* The traditional farthest point sampling (FPS) algorithm is an iterative process. Given input points $\Gamma = \{x_1, x_2, \cdots, x_{N_A}\}$, in order to obtain the output set $S = \{x_1, x_2, \cdots, x_{N_B}\}$ with $N_B$ points, $N_B$ iterations are required, and the formula for each iteration is as follows:

$$x_j = \arg\max\left(\min d\left(x_i, x_j\right)\right) \tag{2}$$

where $x_i \in S$ is the points already taken out before this iteration, $x_j \in \Gamma$ is the point to be taken out in this iteration, $d\left(x_i, x_j\right)$ is the Euclidean distance between $x_i$ and $x_j$, so that $x_j \in \Gamma$ is the most distant point relative to the set $S$.
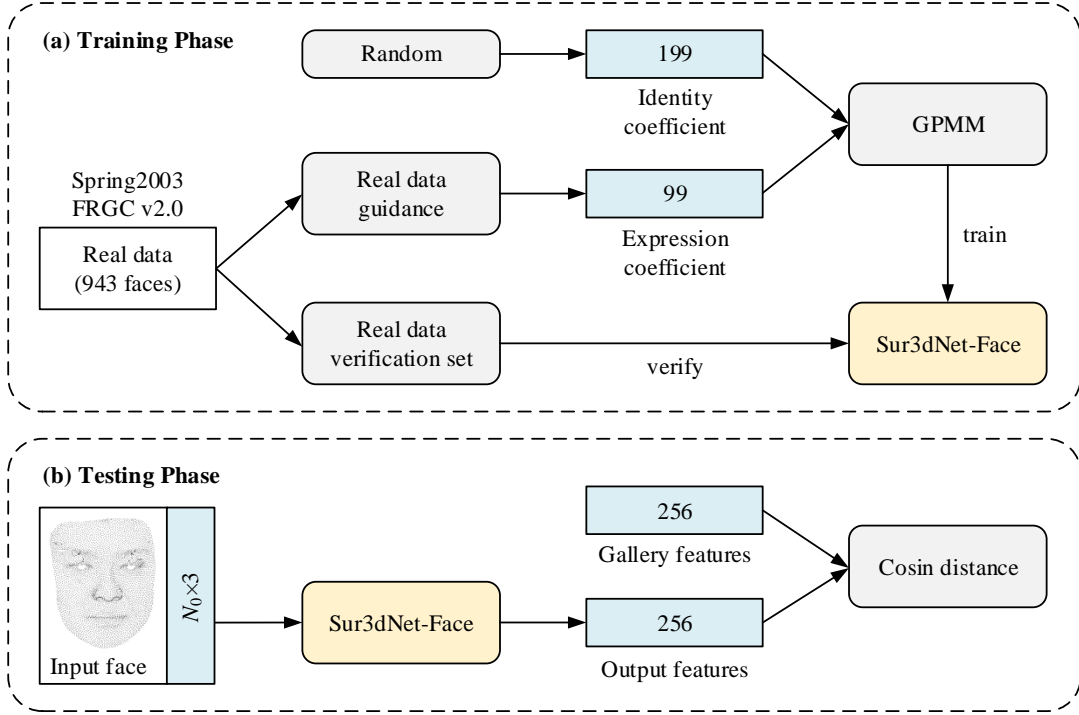
Fig. 1. An overview of the proposed few-data guided learning framework for face recognition.

According to our analysis, the main reason for the poor performance of PointNet in face recognition is that FPS preferentially selects points at the edge region of the point cloud. For the closed point clouds generated from CAD model (e.g. ModelNet dataset), this strategy ensures that more corner points can be taken out. However, as compared in Fig. 3(a), in practical scenario, on the ground that the point cloud is grabbed from one direction, barely including one perspective of the object, the edge points are rather unstable. Therefore, the recognition rate of PointNet on the actual measured point cloud is far below expectation, especially on the face with different posture, where the boundary line can be quite different. Meanwhile, for the edge points, the neighbors clustered by ball query are aggregated on one side of the center, leading to high susceptibility of features to pose variation. To solve these problems, we propose dithering farthest point sampling (DFPS) as follows:

$$x_j = \arg\max\left(\min \lambda d\left(x_i, x_j\right)\right) \quad (3)$$

where

$$\lambda_j = \begin{cases} 0, d\left(x_j, x_{NT}\right) > R \\ e_j^p, else \end{cases} \quad (4)$$

where $x_{NT}$ is the coordinate of nose tip, $e_j$ is the eigen value of $x_j$ (as mentioned before in normal estimation subsection), $p$ is the weighting factor, $R$ is the valid radius beyond which the $\lambda_j$ will be set to 0 and the corresponding $x_j$ will never be selected by DFPS.

DFPS has different behavior in training phase and testing phase. During training, random variations are added to $R$ and $p$ to improve the adaptability of the network. Specifically, in the testing phase $R = 65$ and $p = 0$, while in the training

phase, $R$ and $p$ are random numbers in range of $(50, 80)$ and $(-0.2, 0.2)$ respectively.

The first row of (4) ensures that the points output by FPS are mostly selected in the central area that contains abundant facial features, avoiding the influence of unstable edge points, as is shown in Fig. 3(b).

Additionally, it is necessary to explain the second row of (4) to analyze its principle. As a rule, there is a lot of noise on the real captured point clouds. Owing to the tendency for FPS to select distant points, corner points and noise floating outside the surface of point cloud are prone to be selected. As the eigen value $e_j$ reflects the smoothness near the point $x_j$, by adjusting the value of $p$, DFPS will tend to select corner points (when $p$ increases) or points on the smooth surface (when $p$ decreases), as is shown in Fig. 3(c). For different measurement system, the smoothness of the output point cloud can be quite different, leading to difficulty in determining $p$, so we adopt a more concise strategy that randomly generates $p$ in training phase. With different $p$, the coordinates of the sampled points will dither slightly, so that the trained network can adapt to different types of noise. Our strategy actually adds more randomness to the traditional FPS, and the experiments show that when supplemented with the proposed dithering strategy, the network gets better results.

As is shown in Fig. 2(a), there are four modified PointNet layers in our network, and the parameters of each layer are shown in Table I.

These modifications are seemingly simple, but for the task of face recognition, they are of great significance. In the experiments section, the impacts of these modifications on the recognition rate are compared intuitively through the ablation
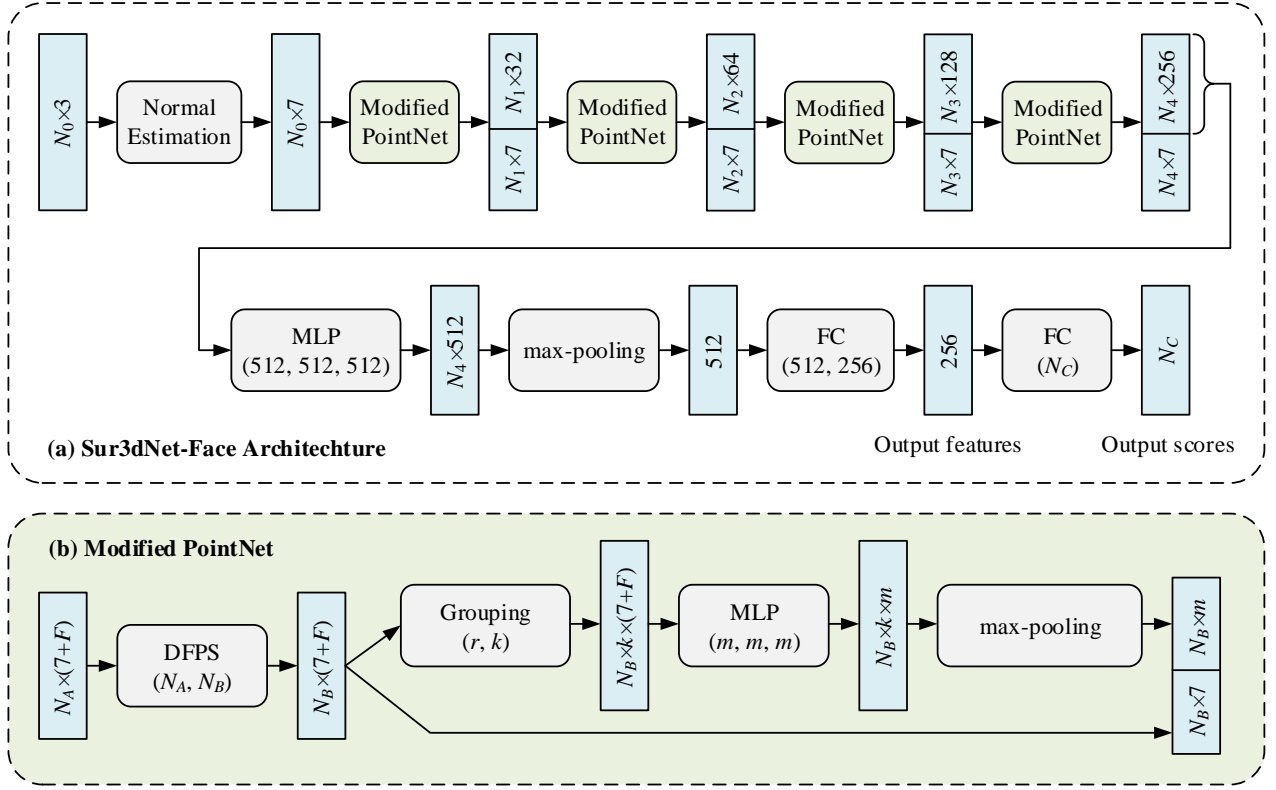
Fig. 2. Architecture of our proposed Sur3dNet-Face, an end-to-end point cloud network for 3D face recognition.

TABLE I
PARAMETERS OF THE FOUR MODIFIED POINTNET LAYERS

| Parameter | layer 1 | layer 2 | layer 3 | layer 4 | Comments |
|---|---|---|---|---|---|
| $N_A$ | 24576 | 4096 | 1024 | 256 | Input point number |
| $N_B$ | 4096 | 1024 | 256 | 64 | Output point number |
| $r$ | 4 | 8 | 16 | 32 | Radius of ball query (mm) |
| $k$ | 24 | 32 | 48 | 64 | Ball query point number |
| $F$ | 0 | 32 | 64 | 128 | Input feature number |
| $m$ | 32 | 64 | 128 | 256 | Output feature number |

study.

### D. Other Details

Our network has no limitation for the number of input points, but in the training phase, point clouds with different sizes cannot be stacked into a batch. In order to use a larger batch size, the input point clouds are firstly downsampled to $N_0 = 24576$ through random choice. Note that there is no such limitation in the test phase or in practical use, where the $N_0$ denotes the actual number of input points.

Batch normalization layer and dropout layer can significantly improve the performance of the network, so we add both on the FC layers (except the last one), and the dropout rate is 0.5. We use Adam [23] optimizer to train the network with the initial learning rate of 1e-3 multiplied by a factor of 0.1 every 10 epochs. Also, the weight decay is set to 1e-4, batch size is 32, and the total number of epochs is 35.

### E. Identification

Similar to traditional solutions, we take the 256-dimensional vector from the output of the penultimate FC layer as a face representation. After acquiring features of both gallery and probe, we calculate the cosine distance between them, and afterwards, identity of the probe is determined by the gallery with minimum distance.

### F. Training Data

*1) Gaussian Process Morphable Models:* In order to generate sufficient training data, we use the GPMM face model [20], a generalization of point distribution models, which assumes that any shape $\Gamma$ can be represented as a discrete set of points:

$$\Gamma = \{x_1, x_2, \cdots, x_N\} \tag{5}$$

where $N$ denotes the number of points.

Shape $\Gamma$ is represented as a vector $s \in \mathbb{R}^{3N}$:

$$s = \begin{bmatrix} x_{1x} & x_{1y} & x_{1z} & \cdots & x_{Nx} & x_{Ny} & x_{Nz} \end{bmatrix}^T \tag{6}$$
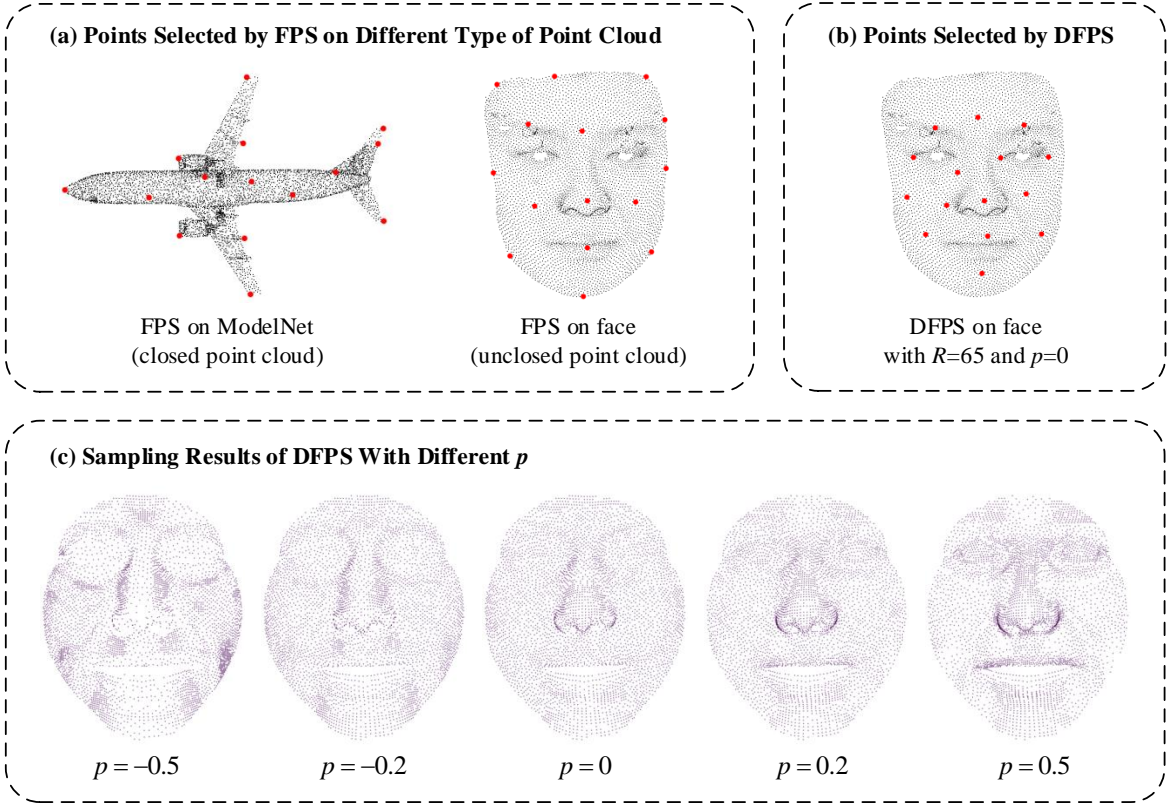
Fig. 3. Analysis on the dithering farthest point sampling (DFPS). In (a) and (b), the black points are the input and the first 15 points selected by the algorithm are marked in red. (c) shows the sampling results of DFPS with different $p$, where the corner points are more likely to be selected as $p$ increases.

GPMM model assumes that the shape variation can be modeled using a normal distribution:

$$s \sim \mathcal{N}(\mu, \Sigma) \tag{7}$$

According to the theory proposed by [20], any face $s \in \mathbb{R}^{3N}$ can be expressed as:

$$s = \bar{s} + B_S \sqrt{\Lambda_S}\alpha + B_E \sqrt{\Lambda_E}\beta \tag{8}$$

where the parameter $\alpha$ determines the shape of human face, $\beta$ determines the expression of human face, and $\bar{s}$ is the mean face model. $B_S$ and $B_E$ are the basis respectively for shape and expression, and similarly, $\Lambda_S$ and $\Lambda_E$ are the variance. For $\bar{s}$, $B_S$, $B_E$, $\Lambda_S$, and $\Lambda_E$, we directly adopt the values provided by [20].

*2) GPMM-Based Data Generation:* According to formula (8), the face generated by GPMM depends on two parameters, $\alpha$ and $\beta$, where $\alpha$ determines the identity and $\beta$ determines the expression.

As GPMM is not specially designed for training neural networks, when using faces generated by formula (8) for training, the trained network can recognize the frontal faces. However, in practical application with posture and noise, in order to improve the variety of the data, the formula should be modified as:

$$s = f\left(\bar{s} + B_S \sqrt{\Lambda_S}\alpha + B_E \sqrt{\Lambda_E}\beta + \delta\right) \tag{9}$$

where $f(\bullet)$ is the random rotation function, $\delta \sim \mathcal{N}(0, \sigma_\delta^2)$ is the Gaussian noise, $\alpha \sim \mathcal{N}(0, \sigma_\alpha^2)$ is the shape coefficient, $\beta \sim \mathcal{N}(0, \sigma_\beta^2)$ is the expression coefficient.

*3) Real Data Guided GPMM-Based Data Generation:* The randomly generating strategy appears difficult to cover the common expressions in the real faces. Taking the disgust expression as an example, we observe that few of the faces with random coefficients look like the disgust. To make our training data further closer to the real data, we propose an enhanced strategy matching the real face with GPMM model, so as to use the real data as a guidance for the expression coefficient to generate training data.

We denote a pair of real faces in the dataset as $(\Gamma_N, \Gamma_E)$, where $\Gamma_N$ and $\Gamma_E$ are respectively the neutral face and the expressive face of a certain subject. According to [24], the registration problem can be solved as:

$$\begin{aligned} \alpha_N &= \arg\max_{\alpha} p\left(\alpha\right) p\left(\Gamma_N \,|\alpha, \bar{S}\right) \\ \beta_E &= \arg\max_{\beta} p\left(\beta\right) p\left(\Gamma_E \,|\beta, \Gamma_N\right) \end{aligned} \tag{10}$$

where $\alpha_N$ is the shape coefficient of both $\Gamma_N$ and $\Gamma_E$, and $\beta_E$ is the expression coefficient of $\Gamma_E$.

We calculate $\beta_E$ of each expressive face in the Spring2003 subset of FRGC v2.0 in advance. When generating a face, one of $\beta_E$ is randomly selected to obtain $\beta_F$ as:

$$\beta_F = \lambda\beta_E + (1 - \lambda)\beta \tag{11}$$

where $\lambda$ is a random value between $(0, 1)$, $\beta \sim \mathcal{N}(0, \sigma_\beta^2)$.

Afterwards, $\beta_F$ is substituted into formula (9) to generate the face, so that the faces generated can cover more expressions similar to real faces in datasets. Experiments show that this strategy can considerably improve the recognition rate of faces with expressions.

### G. Verification Set

To avoid overfitting, a verification set separated from the training set is usually used in deep learning to verify the performance of the network. In our method, all the training samples are generated from GPMM model upon the same formula, leading to similar distribution. So if we directly take a part of the training samples as the verification set, the loss of the verification set will have almost the same trend as that of the training set, and thus cannot play the role of verifying whether the network is overfitted. Existing public datasets contain too few samples for training, but they are quite sufficient as the verification set. Therefore, we put forward a strategy training the network with the generated data and verifying upon the real data.

Since real data in the verification set do not share the identity with our training data, feature vectors extracted by the network from the real data are used to calculate the cosine distance. Through cosine distance, we can evaluate the recognition performance on the real data verification set through Rank-1 Recognition Rate (RR1), Verification Rate (VR) under $FAR = 1e-3$, and Area Under ROC Curve (AUC). Taking these three commonly used indicators into account, we define the loss of the verification set as:

$$loss = 1 - VR \times RR1 \times AUC \qquad (12)$$

In the experiments, we observe that the loss initially decreases with the training epoch, and then it starts to increase, indicating that the network starts to overfit. The network parameters with the lowest verification loss are saved as the final result of the training process.

## IV. EXPERIMENTS

In the experiments, the proposed network is implemented on PyTorch [25] and is training with i7-8700K CPU and two GTX1080TI GPU.

We use Face Recognition Grand Challenge (FRGC) v2.0 dataset [26] and Bosphorus dataset [27] to evaluate the performance of the proposed face recognition method.

### A. Ablation Study

It takes about 50 hours on our hardware to train the network with the training set of 3000 identities, each with 200 different expressions. In order to make more comparisons in a shorter time, we use training set of 3000 identities for ablation study unless otherwise specified, and the rank-1 recognition rate on Bosphorus dataset excluding occlusion subset and posture subset is the main indicator to evaluate the performance.

TABLE II
RANK-1 RECOGNITION RATE (RR1) UNDER DIFFERENT COMBINATION OF $r$ AND $k$ ON BOSPHORUS DATASET

| $r$ | $k$ | RR1 |
|---|---|---|
| 3,6,12,24 | 16,24,32,48 | 86.97% |
| 3,6,12,24 | 24,32,48,64 | 80.53% |
| 4,8,16,32 | 16,24,32,48 | 95.17% |
| 4,8,16,32 | 24,32,48,64 | 97.85% |
| 4,8,16,32 | 40,48,56,64 | 96.66% |
| 5,10,20,40 | 24,32,48,64 | 94.50% |
| 5,10,20,40 | 40,48,56,64 | 95.15% |

*1) Parameters in DFPS:* There are two new parameters introduced in DFPS, namely $R$ and $p$. Fig. 4 (a) demonstrates the rank-1 recognition rate on Bosphorus dataset by training the network with different $R$ values under $p = 0 \pm 0.2$. It can be seen from the blue line that when $45 \leqslant R \leqslant 75$, the recognition rate has little difference. Meanwhile, when the random $R$ strategy is adopted, where $R$ is added by a random variation between -15 and 15 during training phase, as shown in orange line, the recognition rate is further improved.

Fig. 4 (b) depicts the rank-1 recognition rate on Bosphorus dataset by training the network with different $p$ values under $R = 65 \pm 15$. In fact, even if we use the same parameters for training, the recognition rate will occasionally show a variation of 0.5% in each reproduction. Owing to the relatively small impact of parameter $p$ on the results, it appears difficult to determine which is the best value from these experimental results. However, as can be seen from the figure, the random strategy plays a positive role.

Through this ablation study, the parameters of DFPS is determined as follows: in the testing phase $R = 65$ and $p = 0$, while in the training phase, $R$ and $p$ are random numbers in range of $(50, 80)$ and $(-0.2, 0.2)$ respectively for each mini-batch.

*2) Ball Query Radius:* Ball query radius has a significant impact on the recognition rate. We observed that when $r$ and $k$ match a certain proportion, the recognition rate goes higher. Specifically, a smaller $r$ is more suitable for a smaller $k$, and vice versa. Among all the tested combinations, as shown in Table II, the combination of $r = 4, 8, 16, 32$ and $k = 24, 32, 48, 64$ gets the best result, and these parameters are also given in Table I. Note that the recognition rate in Table II is obtained without using real data guided generation and verification.

*3) Real Data Guided Generation and Verification:* The previous experiments are aimed at network parameters, among which the highest rank-1 recognition rate we achieve is 97.85%. Afterwards, we make further experiments verifying the effect of training with a small amount of real data.

There are two real-data-based techniques proposed in this paper, namely, real data guided generation and real data verification set. In the experiments, we randomly select about half of the faces in Spring2003 subset of FRGC v2.0 for real data guided generation and the other half for real data verification set. The results show that our network achieves rank-1 recognition rate of 98.29% with the former technique, real data guided generation, and when both techniques are applied, the
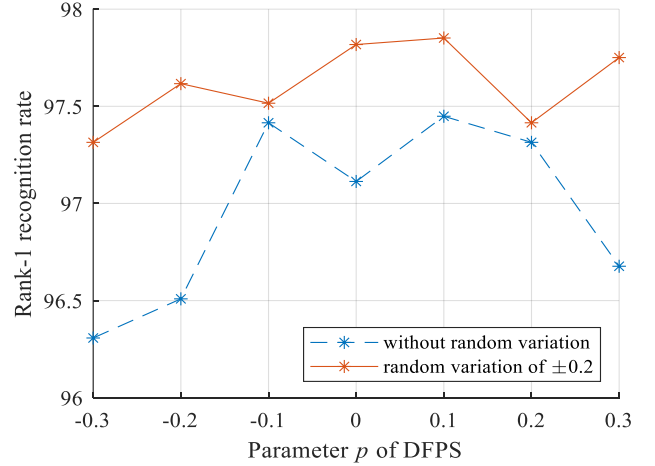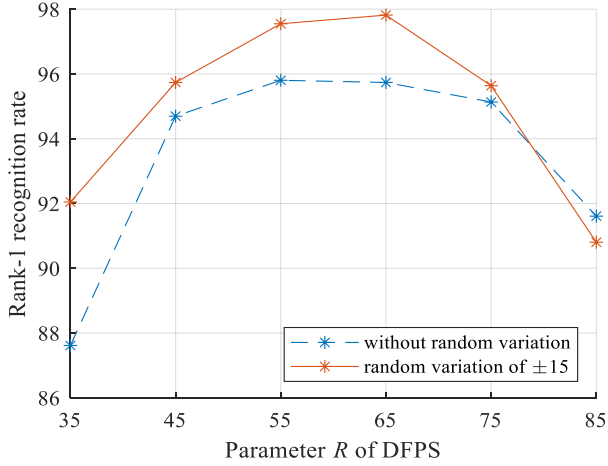
Fig. 4. Comparisons on different $R$ and $p$. The results in this figure are obtained without using real data guided generation and verification.
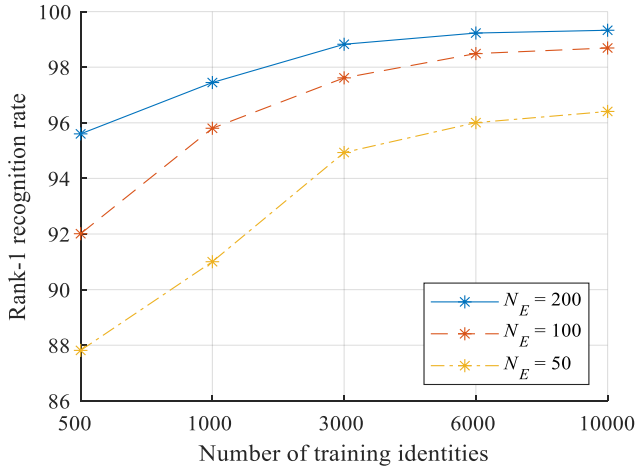


Fig. 5. Comparisons on different volume of training data. $N_E$ denotes the number of generated expressions for each identity.

rank-1 recognition rate reaches 98.83% on Bosphorus dataset.

*4) Training Data Volume:* The training data can be arbitrarily generated, but due to the limitation of hardware performance, we test 10000 identities at most, each containing 200 expressions, and the recognition rate on Bosphorus dataset is displayed in Fig. 5.

From the experimental results, the rank-1 recognition rate reaches 99.33% of training with 10000 identities each with 200 expressions. Also, the increasing trend of recognition rate in the figure indicates that if a larger data volume is used, a slightly higher recognition rate can possibly be obtained.

### B. Comparisons With Other Methods

*1) Results on FRGC v2.0:* FRGC v2.0 dataset [26] containing 4007 scans of 466 subjects in total is divided into three partitions as Spring2003 subset, Fall2003 subset, and Spring2004 subset. Protocol [26] uses the Spring2003 for training and the remaining for testing. We follow this protocol taking Spring2003 as real data to train the network. Specifically, we randomly select about half of the faces in

Spring2003 for real data guidance and the other half for real data verification set.

The comparisons between our proposed method and the state-of-the-art methods on FRGC v2.0 dataset are shown in Table III, where the most representative face recognition methods are compared. Some methods use the corresponding 2D photos of the 3D faces, which we denote as (2D+3D). Also, there are some methods using fine-tuning to further improve the recognition rate, which are marked with (FT).

It can be seen from the table that some deep-learning-based methods have achieved satisfactory recognition rate in recent years, especially those with fine-tuning that can considerably improve the recognition rate. However, in the actual application scenario, due to the dynamic changes of the face gallery, the effect of fine-tuning is much lower than expectation. Among those methods without using fine-tuning, our method is quite competitive, which achieves rank-1 recognition rate of 98.85% and verification rate of 96.75% on FRGC v2.0 dataset.

As is known, whether deep learning can achieve good results, to a great extent, depends on training data. Different from those methods training with more and more data to get a slightly higher recognition rate, the only real data we use in the training process are the 943 faces in Spring2003 subset. With such a small amount of data, we achieve a higher recognition rate than some of the latest methods, which proves the effectiveness and the potentiality of our method.

*2) Results on Bosphorus:* Bosphorus dataset [27] includes totally 4666 scans collected from 105 subjects (60 men and 45 woman aged between 25 and 35) with poses changes, expression variations, and typical occlusions. The yaw rotation of faces is from 10 degrees to 90 degrees in Bosphorus dataset. This paper does not involve occlusion and large posture, so the occlusion subset and the posture subset are excluded. The neutral scans with file name containing N_N_0 are used to form the gallery features.

The comparisons between our method and the state-of-the-art methods on Bosphorus dataset are shown in Table IV. In order to make a fair comparison, the recognition rate under

TABLE III
RANK-1 RECOGNITION RATE (RR1) AND VERIFICATION RATE (VR)
UNDER FAR = $1e-3$ ON FRGC V2.0 DATASET

| Method | RR1 | VR |
|---|---|---|
| Mian et al. [5] (2008) (2D+3D) | 96.10% | 98.60% |
| Huang et al. [28] (2012) | 97.60% | 98.40% |
| Liu et al. [7] (2013) | 96.94% | 90.00% |
| Elaiwat et al. [8] (2015) | 97.10% | 99.20% |
| Lei et al. [9] (2016) | 96.30% | 98.30% |
| Al-Osaimi [29] (2016) | 96.49% | 98.69% |
| Ouamane et al. [30] (2017) | - | 96.65% |
| Ouamane et al. [30] (2017) (2D+3D) | - | 98.32% |
| Gilani et al. [31] (2018) | 98.50% | 98.70% |
| Gilani and Mian [12] (2018) | 97.06% | - |
| Gilani and Mian [12] (2018) (FT) | 99.88% | - |
| Cai et al. [3] (2019) (FT) | 100% | 100% |
| Ours | 98.85% | 96.75% |

TABLE IV
RANK-1 RECOGNITION RATE (RR1) AND VERIFICATION RATE (VR)
UNDER FAR = $1e-3$ ON BOSPHORUS DATASET

| Method | RR1 | VR |
|---|---|---|
| Mian et al. [5] (2008) (2D+3D) | 96.40% | - |
| Huang et al. [28] (2012) | 97.00% | - |
| Liu et al. [7] (2013) | 95.63% | 81.40% |
| Berretti et al. [32] (2013) | 95.67% | - |
| Elaiwat et al. [8] (2015) | - | 91.10% |
| Lei et al. [9] (2016) | 98.90% | - |
| Al-Osaimi [29] (2016) | 92.41% | 93.5% |
| Ouamane et al. [30] (2017) (2D+3D) | - | 96.17% |
| Gilani et al. [31] (2018) | 98.5% | - |
| Gilani and Mian [12] (2018) | 96.18% | - |
| Gilani and Mian [12] (2018) (FT) | 100% | - |
| Cai et al. [3] (2019) (FT) | 99.75% | 98.39% |
| Ours | 99.33% | 97.70% |

the same subsets of Bosphorus is displayed. For the methods that do not provide the result of each subset, the recognition rate of the complete dataset is shown.

## V. DISCUSSIONS

Although some existing methods have achieved satisfactory recognition rate, our method has some notable advantages, upon that we make further discussions.

### A. Computational Complexity

Our method designs an end-to-end network directly inputting the point clouds, which is faster than those methods requiring a complex preprocessing. For example, [9] reports that they need 3.16 s for preprocessing, including detection and alignment, region segmentation, model registration, and other steps; and [13] reports that they need 6.08 s to generate the 2D images including depth map, azimuth map, and elevation map from 3D point clouds.

On the hardware platform described in experiments section, in the training phase with the faces generated in advance, the forward time and backward time of our network for each mini-batch (with batch size of 32) is 0.208 s and 0.095 s, respectively. When being applied to actual applications, it takes about 0.035 s to detect the nose tip and calculate the normal vectors,

and 0.105 s to extract the feature of a single face using a single GPU, and there is no additional preprocessing required.

### B. Data Generation and Overfitting

There is a common problem in existing deep-learning-based face recognition methods, that is, owing to the lack of data, more faces need to be generated by interpolating existing datasets with each other [12], [13]. Although the generated faces used for training are considered different, they have implicit overlaps with the test set, meaning that the results may be obtained to some extent by overfitting.

Comparably, the training data in our method are absolutely generated from GPMM model, which have little intersection with real datasets. Therefore, what we listed are actually cross-dataset results without any fine-tuning, which are closer to the results in actual application than results of those methods that take part of the datasets to generate training data and use the other part to test.

## VI. CONCLUSION

An end-to-end deep learning network entitled Sur3dNet-Face for point-cloud-based 3D face recognition is presented in this paper, along with the concrete approach for training. Coupled with real data guided generation and real data verification set, a few-data guided learning framework based on Gaussian process morphable model is proposed, upon that the common problem in 3D face deep learning of lacking training data is overcome.

Different from existing methods training with more and more data to get a slightly higher recognition rate, the only real data we use in the training process are the 943 faces in Spring2003 subset of FRGC v2.0. With such a small amount of data, we achieve a higher recognition rate than some of the latest methods, which proves the effectiveness and the potentiality of our method.

Furthermore, the ablation study of our method have been analyzed, and the validity has been proved by experiments. Without any fine-tuning on test set, the Rank-1 Recognition Rate (RR1) and Verification Rate (VR) are achieved as follows: 98.85% (RR1) and 96.75% (VR) on FRGC v2.0 dataset, and 99.33% (RR1) and 97.70% (VR) on Bosphorus dataset.

In the future research, we will optimize the process of data generation to generate more occluded and angled faces, so as to further improve the applicability of our method.

## REFERENCES

[1] H. Patil, A. Kothari, and K. Bhurchandi, "3-d face recognition: features, databases, algorithms and challenges," *Artificial Intelligence Review*, vol. 44, no. 3, pp. 393–441, Oct 2015.
[2] S. Soltanpour, B. Boufama, and Q. M. J. Wu, "A survey of local feature methods for 3d face recognition," *Pattern Recognition*, vol. 72, 2017.
[3] Y. Cai, Y. Lei, M. Yang, Z. You, and S. Shan, "A fast and robust 3d face recognition approach based on deeply learned face representation," *Neurocomputing*, vol. 363, pp. 375–397, 2019.
[4] Y. Sun, X. Wang, and X. Tang, "Deep learning face representation from predicting 10,000 classes," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, June 2014, pp. 1891–1898.
[5] A. S. Mian, M. Bennamoun, and R. Owens, "Keypoint detection and local feature matching for textured 3d face recognition," *International Journal of Computer Vision*, vol. 79, no. 1, pp. 1–12, Aug 2008.

[6] H. Mohammadzade and D. Hatzinakos, "Iterative closest normal point for 3d face recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 2, pp. 381–397, Feb 2013.

[7] P. Liu, Y. Wang, D. Huang, Z. Zhang, and L. Chen, "Learning the spherical harmonic features for 3-d face recognition," *IEEE Transactions on Image Processing*, vol. 22, no. 3, pp. 914–925, 2013.

[8] S. Elaiwat, M. Bennamoun, F. Boussaid, and A. El-Sallam, "A curvelet-based approach for textured 3d face recognition," *Pattern Recognition*, vol. 48, no. 4, pp. 1235–1246, Apr. 2015.

[9] Y. Lei, Y. Guo, M. Hayat, M. Bennamoun, and X. Zhou, "A two-phase weighted collaborative representation for 3d partial face recognition with single sample," *Pattern Recognition*, vol. 52, pp. 218–237, 2016.

[10] Y. Lei, M. Bennamoun, M. Hayat, and Y. Guo, "An efficient 3d face recognition approach using local geometrical signatures," *Pattern Recognition*, vol. 47, no. 2, pp. 509–524, Feb. 2014.

[11] X. Song, Z. H. Feng, G. Hu, J. Kittler, and X. J. Wu, "Dictionary integration using 3d morphable face models for pose-invariant collaborative-representation-based classification," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 11, pp. 2734–2745, Nov 2018.

[12] S. Zulqarnain Gilani and A. Mian, "Learning from millions of 3d scans for large-scale 3d face recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[13] D. Kim, M. Hernandez, J. Choi, and G. Medioni, "Deep 3d face identification," in *2017 IEEE International Joint Conference on Biometrics (IJCB)*, 2017, pp. 133–142.

[14] R. Q. Charles, H. Su, M. Kaichun, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 77–85.

[15] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 5099–5108.

[16] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *ACM Trans. Graph.*, vol. 38, no. 5, Oct. 2019.

[17] J. Liu, B. Ni, C. Li, J. Yang, and Q. Tian, "Dynamic points agglomeration for hierarchical point sets learning," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 7545–7554.

[18] V. Blanz and T. Vetter, "A morphable model for the synthesis of 3d faces," in *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, ser. SIGGRAPH '99. USA: ACM Press/Addison-Wesley Publishing Co., 1999, p. 187–194.

[19] P. Dou, S. K. Shah, and I. A. Kakadiaris, "End-to-end 3d face reconstruction with deep neural networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1503–1512.

[20] M. Lüthi, T. Gerig, C. Jud, and T. Vetter, "Gaussian process morphable models," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 8, pp. 1860–1873, 2018.

[21] B. Gecer, S. Ploumpis, I. Kotsia, and S. Zafeiriou, "Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1155–1164.

[22] H. Yi, C. Li, Q. Cao, X. Shen, S. Li, G. Wang, and Y. Tai, "Mmface: A multi-metric regression network for unconstrained face reconstruction," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7655–7664.

[23] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.

[24] T. Gerig, A. Morel-Forster, C. Blumer, B. Egger, M. Luthi, S. Schoenborn, and T. Vetter, "Morphable face models - an open framework," in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, 2018, pp. 75–82.

[25] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035.

[26] P. J. Phillips, P. J. Flynn, T. Scruggs, K. W. Bowyer, J. Chang, K. Hoffman, J. Marques, J. Min, and W. Worek, "Overview of the face recognition grand challenge," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2005, pp. 947–954.

[27] A. Savran, N. Alyüz, H. Dibeklioğlu, O. Çeliktutan, B. Gökberk, B. Sankur, and L. Akarun, "Bosphorus database for 3d face analysis," in *Biometrics and Identity Management*, B. Schouten, N. C. Juul, A. Drygajlo, and M. Tistarelli, Eds. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 47–56.

[28] D. Huang, M. Ardabilian, Y. Wang, and L. Chen, "3-d face recognition using elbp-based facial description and local feature hybrid matching," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 5, pp. 1551–1565, 2012.

[29] F. R. Al-Osaimi, "A novel multi-purpose matching representation of local 3d surfaces: A rotationally invariant, efficient, and highly discriminative approach with an adjustable sensitivity," *IEEE Transactions on Image Processing*, vol. 25, no. 2, pp. 658–672, 2016.

[30] A. Ouamane, A. Chouchane, E. Boutellaa, M. Belahcene, S. Bourennane, and A. Hadid, "Efficient tensor-based 2d+3d face verification," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 11, pp. 2751–2762, Nov 2017.

[31] S. Z. Gilani, A. Mian, F. Shafait, and I. Reid, "Dense 3d face correspondence," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 7, pp. 1584–1598, 2018.

[32] S. Berretti, N. Werghi, A. [del Bimbo], and P. Pala, "Matching 3d face scans using interest points and local histogram descriptors," *Computers & Graphics*, vol. 37, no. 5, pp. 509–525, 2013.