

# Face Recognition from Sequential Sparse 3D Data via Deep Registration

Yang Tan<sup>1,2</sup>, Hongxin Lin<sup>1,2</sup>, Zelin Xiao<sup>1,2</sup>, Shengyong Ding<sup>\*</sup>, Hongyang Chao<sup>1</sup>

<sup>1</sup>School of Data and Computer Science, Sun Yat-sen University

<sup>2</sup>Pixtalks Tech

{tany36, linhx9, xiaozl}@mail2.sysu.edu.cn, marcding@163.com, isschhy@mail.sysu.edu.cn

## Abstract

Previous works have shown that face recognition with high accurate 3D data is more reliable and insensitive to pose and illumination variations. Recently, low-cost and portable 3D acquisition techniques like ToF(Time of Flight) and DoE based structured light systems enable us to access 3D data easily, e.g., via a mobile phone. However, such devices only provide sparse(limited speckles in structured light system) and noisy 3D data which can not support face recognition directly. In this paper, we aim at achieving high-performance face recognition for devices equipped with such modules which is very meaningful in practice as such devices will be very popular. We propose a framework to perform face recognition by fusing a sequence of low-quality 3D data. As 3D data are sparse and noisy which can not be well handled by conventional methods like the ICP algorithm, we design a PointNet-like Deep Registration Network(DRNet) which works with ordered 3D point coordinates while preserving the ability of mining local structures via convolution. Meanwhile we develop a novel loss function to optimize our DRNet based on the quaternion expression which obviously outperforms other widely used functions. For face recognition, we design a deep convolutional network which takes the fused 3D depth-map as input based on AMSoftmax model. Experiments show that our DRNet can achieve rotation error  $0.95^\circ$  and translation error  $0.28mm$  for registration. The face recognition on fused data also achieves rank-1 accuracy 99.2%, FAR-0.001 97.5% on Bosphorus dataset which is comparable with state-of-the-art high-quality data based recognition performance.

## 1. Introduction

Face recognition from 3D data has been attractive due to its inherent advantages of being insensitive to pose and

<sup>\*</sup>Corresponding author.

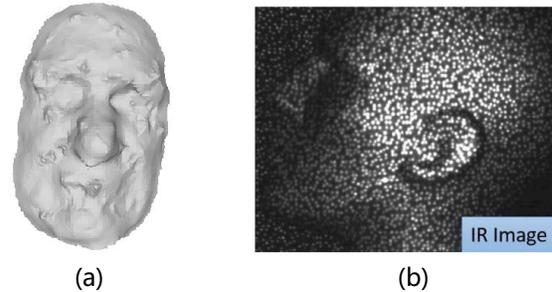


Figure 1. (a) The reconstructed face model from low-quality 3D data (around 1,000 points). (b) The IR image acquired by DoE based structured light system.

illumination variations. Previous literatures [30, 15] have shown that 3D face recognition from high-accuracy 3D data outperforms 2D recognition on a set of datasets. Recently, with the development of compact 3D acquisition techniques, e.g., time-of-flight(ToF) and DoE based structured light devices, people can easily access 3D data even with a mobile phone like iPhone X. The characteristic of such devices is that they can acquire 3D data at high frame speed while the quality of single frame is relatively poor, i.e., sparse and noisy. Figure 1(a) shows the reconstructed face model from low-quality data. For instance, structured light systems derive the depth information by observing the speckle disparity between the reference pattern and the projected pattern on a surface. Due to the limit number of emitters, usually there are only hundreds of speckles projected onto a face. Figure 1(b) [23] shows DoE based patterns projected on a face. A more accurate and computation effective approach is to derive the depth from speckles but it leads to a sparse point cloud.

In this paper, we aim at achieving high-performance face recognition for ToF and structured light based acquisition systems by using the fused sequential 3D data rather than a single frame. This is reasonable as the recognition process usually takes tens of milliseconds during which we can obtain several frames. Previous works [1, 17] also tried to fuse frames acquired by Kinect to obtain a super-resolved model

for face recognition. Note that the Kinect is not specifically designed for short distance scanning, so the face point cloud acquired by Kinect(dense but extremely noisy) is far from the new generation structured light devices(sparse but more accurate) designed for face recognition. In addition, they adopted common 3D registration algorithms like ICP [2] to align frames. The main shortcoming of such registration algorithms is incapable of handling large pose variations. Actually, in our experiments, we find that ICP algorithm almost can not handle difficult frame pairs when relative rotation angles  $\alpha > 60^\circ, \gamma > 40^\circ$ , where  $\alpha, \gamma$  represent the roll and yaw angles.

Inspired by the success of Deep Learning in many vision tasks, we develop a PointNet-like [22] Deep Registration Network(DRNet) to regress the registration parameters between any point cloud pair. More precisely, we hope our neural network takes a pair of interpolated point clouds as input and outputs a vector from which we can derive the transformation parameters. As translation parameters do not have any particular constraints, we simply regress translation parameters by  $L_2$  loss. However, it is much more complex for rotation. The rotation can be expressed in several ways, e.g., a rotation matrix(nine parameters)  $\mathbf{R}$ , Euler angles( $\alpha, \beta, \gamma$ ), axis-angle( $\theta, v_x, v_y, v_z$ ) and quaternion( $\cos \frac{\theta}{2}, \sin \frac{\theta}{2} \cdot v_x, \sin \frac{\theta}{2} \cdot v_y, \sin \frac{\theta}{2} \cdot v_z$ ) where  $(v_x, v_y, v_z)$  is the rotation axis and  $\theta$  is the rotation angle. Considering the orthogonal constraints of rotation matrix and the non-unique parametrisation of Euler angles [13], we do not adopt these two expressions. Thus we design our loss function based on the unit quaternion system from the following facts. First, if a rotation is very small, then the rotation angle must also be very small, i.e.  $\cos \frac{\theta}{2} \rightarrow 1$  no matter what the rotation axis is. Second, if two rotations  $\mathbf{q}_1, \mathbf{q}_2$  are close, then the compositional rotation  $\mathbf{q}_3 = \mathbf{q}_1 \mathbf{q}_2^{-1}$  must be small, indicating the real part of  $\mathbf{q}_3$  approaching 1. Thus we define a loss function measuring the rotation angle between predicted and ground-truth pose for optimization. We find this loss function obviously outperforms other straight forward  $L_2$  loss function on axis-angle( $\theta, v_x, v_y, v_z$ ) expression.

For face recognition, we design a convolutional neural network FRNet to achieve high recognition performance on our fused sequential data. Unlike the models used by [15, 30], we adopt ResNet-18 [11] structure with AMSoftmax [27] loss function and find that simple augmentations for training data, e.g., pose variations and occlusions can surprisingly give good results without synthesizing any new identity. As there are no large scale existing sequential 3D datasets, we propose a method to generate desired data by sparse sampling and adding perturbations from existing high-quality face datasets. This approach enables us to obtain large amount of data at a very low cost while still producing meaningful results.

In summary, our contributions are follows:

- We raise a new challenging face recognition problem, i.e., face recognition from a sequence of sparse point clouds which will be common for structured light systems. Note the number of 3D points is only about 1,000 in our problem and much less than previous works [30, 15, 17].
- To handle large pose variations, we design a robust deep PointNet-like network DRNet for 3D point clouds rigid registration based on the unit quaternion expression with a carefully designed loss function. To the best of our knowledge, we are the first to align 3D facial point clouds via a neural network.
- We study how the 3D data quality impacts face recognition and demonstrate the possibility of achieving high recognition accuracy from very sparse point cloud sequence.

## 2. Related work

**Rigid Point Cloud Registration** The classical method for rigid point cloud registration is ICP [2] algorithm which aims at closing a pair of point clouds iteratively. However, the performance of ICP algorithm is heavily relied on the initial poses so it can not handle large pose variations. To address this problem, [20] adopted EGI to perform coarse registration and then refined the result by ICP, but the complexity was unacceptable for mobile devices.

At present, there is no specifically designed neural networks for point cloud registration. FacePoseNet [4] directly regressed 6DoF transform parameters between the generic 3D facial keypoints model and the keypoints on intensity images. [14] proposed a network to estimate camera poses from monocular images in quaternion form. These methods are intensity image based but our method adopts coordinate-maps with accurate 3D informations as input and it seems more reasonable to predict 3D poses.

**3D Face Recognition** For conventional methods, there are local and global descriptor based approaches. Local descriptors usually extract some sub-regions and local informations as features. [19] proposed a descriptor based on three face keypoints to describe local features. [10] performed 3D face recognition by comparing Euclidean and Geodesic distances between matched keypoints. Global descriptors treat face as an entity. [7] proposed using radial curves emanating from the nose tip to represent the facial surface. Some 3D Morphable Model [3] based methods used 3DMM parameters for face recognition but the fitting process required massive computations.

CNN based methods DeepFace [26] and FaceNet [25] brought remarkable improvements for 2D face recognition. DeepFace achieved an accuracy of 97.35% on LFW [23]

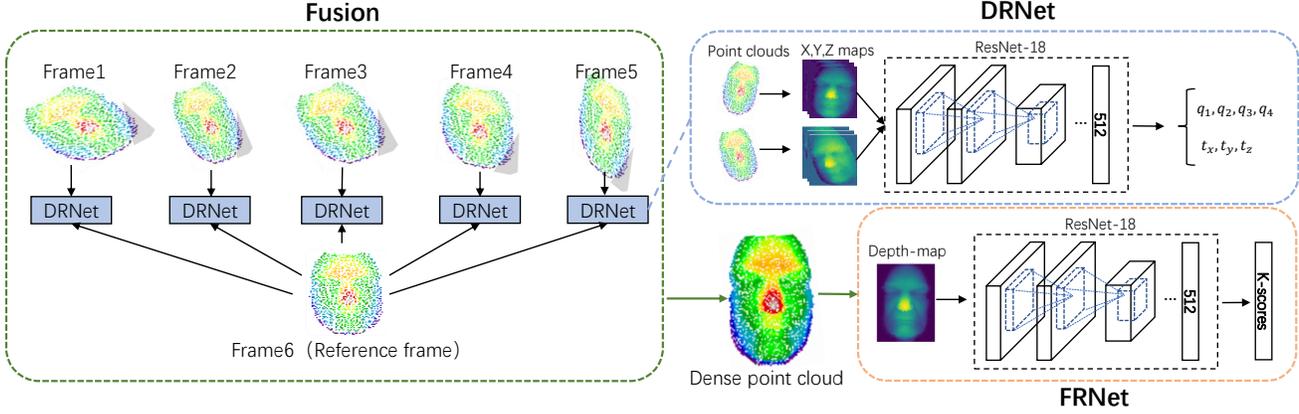


Figure 2. An overview of our Face Registration and Face Recognition framework.

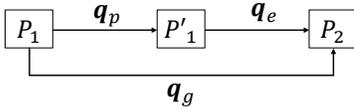


Figure 3.  $P_1, P'_1, P_2$  are source, predicted and target point clouds in different poses.  $\mathbf{q}_g$  is the ground-truth rotation quaternion from  $P_1$  to  $P_2$ .  $\mathbf{q}_p$  is the predicted quaternion from  $P_1$  to  $P'_1$ .  $\mathbf{q}_e$  is the quaternion from  $P'_1$  to  $P_2$ .

dataset outperforming the best conventional method 27% at that time. Later [15] proposed CNN based 3D face recognition pipeline and achieved comparable results. [30] proposed a synthesizing method to generate about 100 thousand identities for large scale training. These methods are trained on high-quality 3D data and can not be transferred to low-quality data directly.

**Sequential Methods** [29] used a sequence of temporal images to perform 2D face recognition. Recent works [1, 17, 12, 5] adopted a sequence of 3D data to perform depth fusion or morphology to reconstruct face models. They adopted ICP algorithm to perform registration for point clouds acquired by Kinect. As described in section 1, Kinect data is far from our desired data due to the high noise level. We study the sparser but more accurate data acquired by the new generation 3D face scanner.

### 3. Proposed framework

Our proposed framework is composed of two parts: Face Registration and Face Recognition. An overview of the framework is shown in Figure 2. We first use Deep Registration Network(DRNet) to reconstruct a dense 3D facial point cloud from 6 frames of sparse point clouds by registering and fusing. Then we feed the fused data to Face Recognition Network(FRNet).

### 3.1. Deep Registration Network(DRNet)

The most critical step to reconstruct the dense point cloud from a low-quality sequence is 3D registration, i.e., predict the rotation and translation transforms between two sparse point clouds. In order to handle large pose variations, we apply a deep neural network to regress the transformation parameters between two point clouds  $P_1, P_2$ . Inspired by PointNet [22], we introduce X,Y,Z coordinate-maps generated by projecting 3D points onto 2D image plane and then interpolating. Thus our network can be viewed as an implementation of 3D point cloud network which can still apply standard convolution and pooling operations to exploit hierarchical local structures. Our DRNet takes a pair of  $256 \times 256$  3-channel coordinate-maps as input. The network architecture is based on ResNet-18 [11] with a 512-dimensional fully connected layer and outputs a 7-dimensional vector  $\mathbf{p}$ , encoding a translation vector  $\mathbf{t} = (t_x, t_y, t_z)$  and a unit quaternion for rotation  $\mathbf{q} = (q_1, q_2, q_3, q_4)$  s.t.  $q_1^2 + q_2^2 + q_3^2 + q_4^2 = 1$ :

$$\mathbf{p} = (\mathbf{t}, \mathbf{q}) \quad (1)$$

#### 3.1.1 Interpretation of loss function

The reason we express rotation in quaternion form is the nice properties of quaternion system, e.g., uniqueness and interpretability. The interpretability lies in the fact that the unit quaternion for a rotation  $\mathbf{q}$  takes the following form:

$$\mathbf{q} = \left( \cos \frac{\theta}{2}, \sin \frac{\theta}{2} \cdot v_x, \sin \frac{\theta}{2} \cdot v_y, \sin \frac{\theta}{2} \cdot v_z \right) \quad (2)$$

with  $(v_x, v_y, v_z)$  representing the rotation axis and  $\theta$  being the rotation angle. Intuitively, we have two ways to construct the loss function, i.e., let the network directly regress quaternion  $(\cos \frac{\theta}{2}, \sin \frac{\theta}{2} \cdot v_x, \sin \frac{\theta}{2} \cdot v_y, \sin \frac{\theta}{2} \cdot v_z)$  or axis-angle  $(\theta, v_x, v_y, v_z)$  by L1/L2 loss functions. In this paper, we design our loss from another aspect which is more interpretable and demonstrate that optimizing the quaternion

form shows incomparable advantages. For ease of deduction, we first introduce some notations as shown in Figure 3 where  $P_1$ ,  $P'_1$  and  $P_2$  represent the point clouds of the same rigid object in different poses. We use the quaternion  $\mathbf{q}_p$  to encode the rotation between point cloud  $P_1$  and  $P'_1$ ,  $\mathbf{q}_g$  to encode the rotation between  $P_1$  and  $P_2$  and  $\mathbf{q}_e$  to encode the rotation between  $P'_1$  and  $P_2$ . In quaternion system, above definitions imply following equations:

$$\begin{aligned} \mathbf{P}'_1 &= \mathbf{q}_p \mathbf{P}_1 \mathbf{q}_p^{-1} \\ \mathbf{P}_2 &= \mathbf{q}_g \mathbf{P}_1 \mathbf{q}_g^{-1} \\ \mathbf{P}_2 &= \mathbf{q}_e \mathbf{P}'_1 \mathbf{q}_e^{-1} \\ \Rightarrow \mathbf{q}_e &= \mathbf{q}_g \mathbf{q}_p^{-1} \end{aligned} \quad (3)$$

According to the equation (2), it is not hard to see that if  $\mathbf{q}_p$  exactly equals  $\mathbf{q}_g$ , then  $\mathbf{q}_e = (1, 0, 0, 0)$  implying there is no pose difference between  $P'_1$  and  $P_2$ . So a well trained network should produce  $\mathbf{q}_p$  with  $\mathbf{q}_e = \mathbf{q}_g \mathbf{q}_p^{-1}$  approaching  $(1, 0, 0, 0)$ , thus the loss function can be defined as:

$$\begin{aligned} loss_1 &= \|\mathbf{q}_e - (1, 0, 0, 0)\|_2^2 \\ &= 2 - 2 \cos \frac{\theta_e}{2} \end{aligned} \quad (4)$$

It demonstrates that the network has a clear optimizing goal in quaternion form, i.e., minimizing the rotation angle  $\theta_e$  between predicted pose and target pose which has nothing to do with rotation axis. Actually, when the rotation angle approaches 0, we could say there is almost no rotation, no matter what the axis is. However, for another similar expression axis-angle  $(\theta, v_x, v_y, v_z)$  also encoding the rotation angle and axis, the network has to optimize four variables simultaneously and it is difficult to judge which variable is more important.

Interestingly, the equation (4)  $loss_1$  equals the common  $L_2$  loss function on  $\mathbf{q}_p$  defined as:

$$\begin{aligned} loss_2 &= \|\mathbf{q}_g - \mathbf{q}_p\|_2^2 \\ &= 2 - 2 \cos \frac{\theta_e}{2} \\ &= loss_1 \end{aligned} \quad (5)$$

For easy implementation, we define the complete loss function as:

$$loss = \|\mathbf{t}_g - \mathbf{t}_p\|_2 + \alpha \|\mathbf{q}_g - \mathbf{q}_p\|_2 \quad (6)$$

where  $\alpha$  is a scale factor initialized as 500 to balance the translation and rotation weights. It increases to  $1 \times 10^4$  when  $\|\mathbf{q}_g - \mathbf{q}_p\|_2^2 < 1 \times 10^{-4}$  during training. We adopt Adam [16] optimizer with weight decay  $5 \times 10^{-5}$  and the initial learning rate is 0.01. Experiments clearly show that our loss function is obviously better than other intuitive loss functions such as  $L_1$  loss on quaternion or axis-angle.

## 3.2. Face Recognition Network (FRNet)

We aim at achieving high-accuracy face recognition with the fused sequential data which will be described in section 4. We adopt the novel AMSoftmax [27] loss function to optimize our FRNet. The AMSoftmax loss function tries to separate different individuals on a sphere with a large margin. Our network takes the depth-map of size  $256 \times 256$  converted as input. Again we use the ResNet-18 architecture with a 512 dimensional fully connected layer as the output feature. The facial similarity is measured by the simple cosine distance on output features. We train FRNet by Adam optimizer with weight decay  $5 \times 10^{-5}$  and the initial learning rate is 0.01.

## 4. Data generation

As the new DoE based 3D scanners for face have not been mass produced, we are unable to construct a real large scale dataset to validate our method. Therefore, in this section, we will introduce the method to generate our simulated data from existing public 3D face datasets. To facilitate follow-up operations, we roughly align all face scans in datasets to a standard pose in advance.

### 4.1. Data for DRNet

Face scans in raw datasets are dense (more than 10,000 points) and clean, so we need to perform sparse sampling to generate our desired data. To simulate the distribution of sparse and random patterns on a moving face, we develop a sparse sampling strategy. First, we introduce noises and a random pose variation which is expressed in Euler angles  $\alpha \in [-45^\circ, 45^\circ]$ ,  $\beta \in [-20^\circ, 20^\circ]$ ,  $\gamma \in [-30^\circ, 30^\circ]$  representing roll, pitch and yaw angles respectively and the translation  $t_x, t_y, t_z \in [-8mm, 8mm]$  to a pre-aligned face scan. Noises come from a Gaussian distribution  $N(0, 4)$  and are randomly added to ten percent of points. Second, we project the point cloud onto a 2D plane divided into 1,000 grids of the same size and randomly select one point from each grid, thus we obtain a sparse point cloud containing about 1,000 points. For each face scan in raw datasets, we repeat the procedure above 6 times to obtain a sequence of sparse face data, while one of six frames is defined as the reference frame with  $\alpha = 0^\circ, \beta = 0^\circ, \gamma = 0^\circ$ . Note that the pose variation is relative to the standard pose and transform parameters will be recorded to calculate the relative transform between any pair of frames in the sequence.

During the training stage, we randomly feed pairs of frames from the training sequential dataset to regress rotation and translation parameters.

### 4.2. Data for FRNet

Our goal is to demonstrate that the fused data from a sparse sequence of 6 frames can also achieve comparable

results to the high-quality data. To study the performances of our FRNet under different data qualities, we generate 4 types of data as follows:

- **Fused data** For each sparse sequence, we align other 5 frames to the reference frame by our DRNet, then the union of these 6 aligned point clouds is defined as fused data. As the fused data contains around 6,000 points, it is possible to conduct denoising. Specifically, for each 3D point  $(x, y, z)$ , we calculate the mean z-coordinate  $z_m$  of neighbor points within  $radius = 3mm$ . When  $|z - z_m| > 2$ , update  $z = z_m$ . This denoising strategy can effectively remove outliers but does not smooth the point cloud excessively. During the training stage, we first augment the fused data with pose variations  $\alpha, \beta, \gamma \in [-10^\circ, 10^\circ]$ . Then we project point clouds onto image plane and interpolate it to generate depth-maps. We randomly occlude depth-maps with 1-6 patches of size in  $[0, 20]$  for augmentation.
- **High-quality data** We select raw point clouds (more than 10,000 accurate points) from 3D face datasets. We also adopt the same augmentation strategy described above to generate depth-maps.
- **Low-quality data** We select reference frames (around 1,000 points) from sparse sequences. We also adopt the same augmentation strategy described above to generate depth-maps.
- **Sequential data** We directly feed sparse sequences to FRNet without registration and fusion. We also adopt the same augmentation strategy described above for each frame to generate depth-maps. Note that we need to expand the input channels of FRNet from one to six.

## 5. Experiments

In this section, we first evaluate the performance of face registration by DRNet and then evaluate face recognition performances on different types of face data.

### 5.1. Datasets

Table 1 shows the most popular 3D face datasets including ND-2006 [8], Bosphorus [24], CASIA [28] and UMBDB [6]. The complete ND-2006 dataset is used to construct our training set. We select 2,900 scans (gallery 105, probes 2,795) from Bosphorus, 3,555 scans (gallery 123, probes 3,432) from CASIA, 749 scans (gallery 122, probes 627) from UMBDB except side faces and extremely occluded samples to construct testing sets.

Specially, we construct two testing sets for face registration, i.e., standard set and difficult set. Both of them contain 4,000 pairs of point clouds generated from Bosphorus dataset. The difference is that samples in

standard set uniformly pose in  $\alpha \in [-45^\circ, 45^\circ], \beta \in [-20^\circ, 20^\circ], \gamma \in [-30^\circ, 30^\circ]$ , but samples in difficult set pose in  $\alpha \in [-45^\circ, -30^\circ] \cup [30^\circ, 45^\circ], \beta \in [-20^\circ, 20^\circ], \gamma \in [-30^\circ, -20^\circ] \cup [20^\circ, 30^\circ]$ , where  $\alpha, \beta, \gamma$  are Euler angles described in section 4.1

Table 1. Details of datasets

Name	IDs	Scans	Expressions	Pose	Occlusion
ND-2006	888	13,450	Multiple	$\pm 15^\circ$	None
Bosphorus	105	4,666	7	$\pm 90^\circ$	4 types
CASIA	123	4,674	6	Frontal	None
UMBDB	143	1,473	4	Frontal	7 types

### 5.2. Evaluation of Face Registration

We quantitatively evaluate the registration result by rotation error and translation error. The rotation error  $\theta_e$  is defined as how many degrees still need to be rotated from the predicted pose to the ground-truth pose, which is derived from the real part of  $\mathbf{q}_e = \mathbf{q}_g \mathbf{q}_p^{-1}$  shown in Figure 3. The translation error is measured by  $t_e = \|\mathbf{t}_g - \mathbf{t}_p\|_2$  where  $\mathbf{t}_g$  and  $\mathbf{t}_p$  are the ground-truth translation and predicted translation.

Table 2 shows that our DRNet achieves an impressive registration performance, especially on difficult testing set. We can see that both the ICP algorithm and our DRNet work well on standard testing set, i.e., the relative rotation angle between two point clouds satisfies  $\Delta\alpha < 60^\circ$  and  $\Delta\gamma < 40^\circ$ . However, for large pose variations, ICP algorithm usually falls into a local optimum and produces large errors. Actually, on difficult testing set, we find that about 11.8% registrations fail with ICP algorithm, i.e., meaningless alignment as shown in Figure 4, while none registration fails with DRNet. Note that on standard testing set, ICP algorithm performs slightly better than DRNet since ICP can recursively refine the registration results. We observe the same effect with DRNet, i.e., a second registration by DRNet will produce better results as shown in Table 2. In addition, Figure 5 shows fused faces from sequences of sparse data.

As the discussion in section 3.1.1, a rotation can be expressed in axis-angle form  $(\theta, v_x, v_y, v_z)$  and quaternion form  $(\cos \frac{\theta}{2}, \sin \frac{\theta}{2} \cdot v_x, \sin \frac{\theta}{2} \cdot v_y, \sin \frac{\theta}{2} \cdot v_z)$ . We give the registration results on DRNets based on these two forms with L1 and L2 loss functions in Table 3. Experiments clearly show that Quaternion-L2 is superior to others as expected.

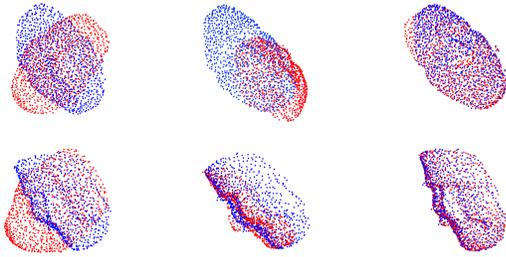
Table 2. Average registration errors on standard testing set and difficult testing set.

Method	Standard		Difficult	
	$\theta_e(^{\circ})$	$t_e(mm)$	$\theta_e(^{\circ})$	$t_e(mm)$
<i>ICP</i>	<b>0.77</b>	0.27	7.64	1.92
<i>Ours</i>	1.80	0.38	1.83	0.45
<i>Ours*</i>	1.08	<b>0.25</b>	<b>0.95</b>	<b>0.28</b>

\*represents performing registration twice.

Table 3. Average registration errors of different loss functions and expressions on standard testing set.

Method	$\theta_e(^{\circ})$	$t_e(mm)$
Axis-angle-L1	4.67	0.71
Axis-angle-L2	3.17	0.41
Quaternion-L1	2.79	0.60
Quaternion-L2	<b>1.80</b>	<b>0.38</b>



(a) Before registration (b) Registered by ICP (c) Registered by DRNet

Figure 4. Two failure samples of ICP algorithm. (a) column shows two pairs of sparse point clouds before registration. Blue represents the target point cloud. (b) column shows the registration result of ICP with significant errors. (c) column shows the well aligned result by our DRNet.

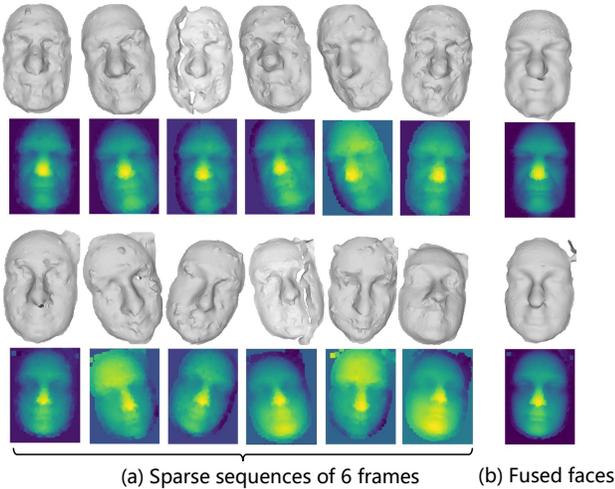


Figure 5. (a) shows two sequences of low-quality 3D data. (b) shows the reconstructed faces after registration and fusion, which contain more facial details.

## 5.3. Evaluation of Face Recognition

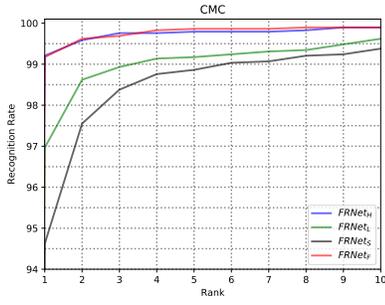
### 5.3.1 Settings

For comparison, we evaluate the performances of face recognition with data described in section 4.2. Below are experimental settings:

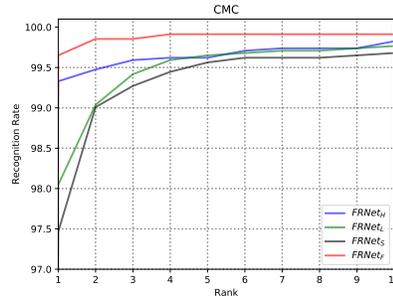
- $FRNet_F$ : Use the fused data for training and testing. It is designed to demonstrate whether our fusion strategy is effective.
- $FRNet_H$ : Use the high-quality data for training and testing. It is designed to compare with the state-of-the-art high-quality data based methods and demonstrate whether our FRNet is well designed.
- $FRNet_H^f$ : Use the high-quality data for training but fused data for testing to study whether the high-quality data based model is capable of adapting to our fused data.
- $FRNet_L$ : Use the low-quality data for training and testing. We want to know if it is possible to achieve high performance recognition just from one single frame of sparse face data.
- $FRNet_S$ : Use the sequential data for training and testing. It is designed to demonstrate whether the registration is necessary.

### 5.3.2 Performances

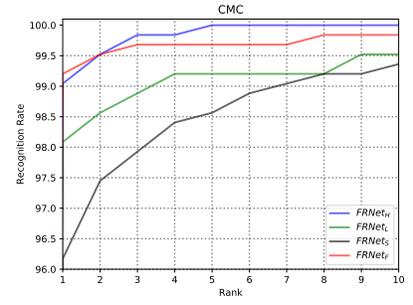
We follow the evaluation protocol described in [15] and adopt common criteria CMC and ROC. Figure 6 shows CMC and ROC curves on Bosphorus, CASIA and UMBDB testing sets. Table 4 shows the result of False Acceptance Rate at 0.001. Table 5 shows the rank-1 result compared with state-of-the-art methods. Firstly, our  $FRNet_H$  and  $FRNet_F$  achieve comparable and state-of-the-art results on testing sets, which shows that our fused data is capable of performing high-accuracy recognition as high-quality data. Note that the  $FRNet_F$  is a little bit better than  $FRNet_H$  on some criteria and we consider it may be caused by the noise of fused data which ease the overfitting. Secondly,  $FRNet_L$  just achieves Rank-1 97.0% and FAR-0.001 88.8% on Bosphorus which demonstrates that one single frame of sparse data is unable to achieve high recognition performance. Thirdly, the performances of  $FRNet_S$  are also not really good, only achieving Rank-1 94.6% and FAR-0.001 80.1% on Bosphorus. We consider that it is hard for the network to utilize the complementary information from sequential frames without registration and fusion.



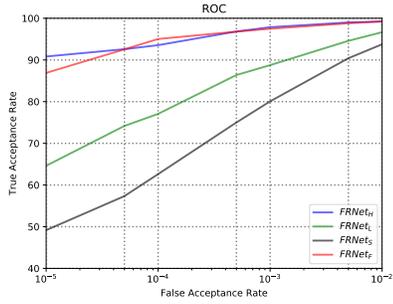
(a) CMC curve on Bosphorus



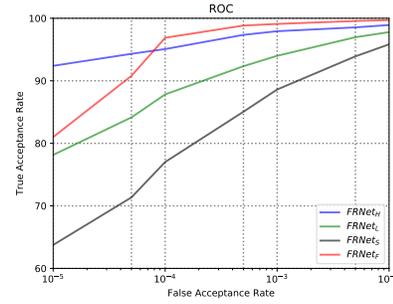
(b) CMC curve on CASIA



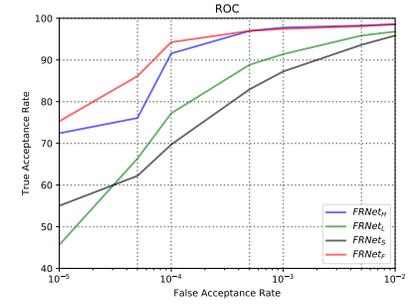
(c) CMC curve on UMBDB



(d) ROC curve on Bosphorus



(e) ROC curve on CASIA



(f) ROC curve on UMBDB

Figure 6. Face recognition performances on 3 testing sets

Table 4. Comparison of False Acceptance Rate (%) at 0.001 with different settings

Method	Bosphorus	CASIA	UMBDB
$FRNet_F$	97.5	<b>99.1</b>	97.4
$FRNet_H$	<b>97.9</b>	97.9	<b>97.8</b>
$FRNet_H^f$	96.8	96.9	97.6
$FRNet_L$	88.8	94.0	91.4
$FRNet_S$	80.1	88.6	87.2

Table 5. Comparison of Rank-1 recognition accuracy (%) with state-of-the-art methods

Method	Bosphorus	CASIA	UMBDB
Xu <i>et al.</i> [28](2006)	-	83.9	-
Mian <i>et al.</i> [21](2007)	96.4	82.5	69.3
Gilani <i>et al.</i> [9](2018)	98.6	85.4	78.6
Lei <i>et al.</i> [18](2016)	98.9	-	-
Kim <i>et al.</i> [15](2017)	99.2	-	-
Zulqarnain <i>et al.</i> [30](2018)	<b>100.0</b>	<b>99.7</b>	97.2
$FRNet_F$	99.2	<b>99.7</b>	<b>99.2</b>
$FRNet_H$	99.2	99.3	99.0
$FRNet_H^f$	99.3	98.9	99.0
$FRNet_L$	97.0	98.0	98.1
$FRNet_S$	94.6	97.5	96.2

## 6. Conclusion

In this paper, we propose a framework to achieve high-accuracy face recognition from sequential sparse and noisy 3D data. Unlike previous works relying on ICP algorithm for registration, we propose a deep convolutional network DRNet to regress the transformation parameters with a carefully designed loss function. Our DRNet is able to achieve rotation error  $0.95^\circ$  and translation error  $0.28mm$  even for large pose variations (difficult testing set). Using the fused data by DRNet for face recognition, we achieve rank-1 99.2%, 99.7%, 99.2% and FAR-0.001 97.5%, 99.1%, 97.4% on Bosphorus, CASIA and UMBDB datasets which is a comparable result with the performance tested on high-quality data.

## References

- [1] S. Berretti, P. Pala, and A. Del Bimbo. Face recognition by super-resolved 3d models from consumer depth cameras. *IEEE Transactions on Information Forensics and Security*, 9(9):1436–1449, 2014. **1, 3**
- [2] P. J. Besl and N. D. McKay. Method for registration of 3-d shapes. In *Sensor Fusion IV: Control Paradigms and Data Structures*, volume 1611, pages 586–607. International Society for Optics and Photonics, 1992. **2**
- [3] 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*, pages 187–194. ACM Press/Addison-Wesley Publishing Co., 1999. 2
- [4] F.-J. Chang, A. T. Tran, T. Hassner, I. Masi, R. Nevatia, and G. Medioni. Faceposenet: Making a case for landmark-free face alignment. In *Computer Vision Workshop (ICCVW), 2017 IEEE International Conference on*, pages 1599–1608. IEEE, 2017. 2
- [5] J. Choi, A. Sharma, and G. Medioni. Comparing strategies for 3d face recognition from a 3d sensor. In *RO-MAN, 2013 IEEE*, pages 19–24. IEEE, 2013. 3
- [6] A. Colombo, C. Cusano, and R. Schettini. Umb-db: A database of partially occluded 3d faces. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pages 2113–2119. IEEE, 2011. 5
- [7] H. Drira, B. B. Amor, A. Srivastava, M. Daoudi, and R. Slama. 3d face recognition under expressions, occlusions, and pose variations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(9):2270–2283, 2013. 2
- [8] T. C. Faltemier, K. W. Bowyer, and P. J. Flynn. Using a multi-instance enrollment representation to improve 3d face recognition. In *IEEE International Conference on Biometrics: Theory, Applications, and Systems*, pages 1–6, 2007. 5
- [9] S. Z. Gilani, A. Mian, F. Shafait, and I. Reid. Dense 3d face correspondence. *IEEE transactions on pattern analysis and machine intelligence*, 40(7):1584–1598, 2018. 7
- [10] S. Gupta, M. K. Markey, and A. C. Bovik. Anthropometric 3d face recognition. *International journal of computer vision*, 90(3):331–349, 2010. 2
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 3
- [12] G.-S. J. Hsu, Y.-L. Liu, H.-C. Peng, and P.-X. Wu. Rgb-d-based face reconstruction and recognition. *IEEE Transactions on Information Forensics and Security*, 9(12):2110–2118, 2014. 3
- [13] A. Kendall, R. Cipolla, et al. Geometric loss functions for camera pose regression with deep learning. In *CVPR*, volume 3, page 8, 2017. 2
- [14] A. Kendall, M. Grimes, and R. Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocation. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015. 2
- [15] D. Kim, M. Hernandez, J. Choi, and G. Medioni. Deep 3d face identification. *arXiv preprint arXiv:1703.10714*, 2017. 1, 2, 3, 6, 7
- [16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *Computer Science*, 2014. 4
- [17] Y.-C. Lee, J. Chen, C. W. Tseng, and S.-H. Lai. Accurate and robust face recognition from rgb-d images with a deep learning approach. In *BMVC*, 2016. 1, 2, 3
- [18] 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*, 52:218–237, 2016. 7
- [19] H. Li, D. Huang, J.-M. Morvan, Y. Wang, and L. Chen. Towards 3d face recognition in the real: a registration-free approach using fine-grained matching of 3d keypoint descriptors. *International Journal of Computer Vision*, 113(2):128–142, 2015. 2
- [20] A. Makadia, A. Patterson, and K. Daniilidis. Fully automatic registration of 3d point clouds. In *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, volume 1, pages 1297–1304. IEEE, 2006. 2
- [21] A. Mian, M. Bennamoun, and R. Owens. An efficient multi-modal 2d-3d hybrid approach to automatic face recognition. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (11):1927–1943, 2007. 7
- [22] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017. 2, 3
- [23] S. Ryan Fanello, C. Rhemann, V. Tankovich, A. Kowdle, S. Orts Escolano, D. Kim, and S. Izadi. Hyperdepth: Learning depth from structured light without matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5441–5450, 2016. 1, 2
- [24] A. Savran, B. Sankur, and M. T. Bilge. Comparative evaluation of 3d vs. 2d modality for automatic detection of facial action units. *Pattern recognition*, 45(2):767–782, 2012. 5
- [25] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 2
- [26] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708, 2014. 2
- [27] F. Wang, J. Cheng, W. Liu, and H. Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018. 2, 4
- [28] C. Xu, T. Tan, S. Li, Y. Wang, and C. Zhong. Learning effective intrinsic features to boost 3d-based face recognition. In *European Conference on Computer Vision*, pages 416–427. Springer, 2006. 5, 7
- [29] O. Yamaguchi, K. Fukui, and K.-i. Maeda. Face recognition using temporal image sequence. In *fg*, page 318. IEEE, 1998. 3
- [30] S. Zulqarnain Gilani and A. Mian. Learning from millions of 3d scans for large-scale 3d face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1896–1905, 2018. 1, 2, 3, 7