# 3D Dense Face Alignment via Graph Convolution Networks

Huawei Wei, Shuang Liang, Yichen Wei
Megvii Technology, Tongji University
{weihuawei, weiyichen}@megvii.com
shuangliang@tongji.edu.cn

## Abstract

*Recently, 3D face reconstruction and face alignment tasks are gradually combined into one task: 3D dense face alignment. Its goal is to reconstruct the 3D geometric structure of face with pose information. In this paper, we propose a graph convolution network to regress 3D face coordinates. Our method directly performs feature learning on the 3D face mesh, where the geometric structure and details are well preserved. Extensive experiments show that our approach gains a superior performance over state-of-the-art methods on several challenging datasets.*

## 1. Introduction

3D face reconstruction and face alignment have been widely used in film and animation production. The two tasks are highly related and well studied. Originally, 3D face reconstruction [32, 13] aims to recover the 3D face geometry and face alignment aims to locate a number of fiducial facial landmarks [30, 24]. Recently, a number of works have integrated them into a single task, so called *3D dense face alignment*. It aims to recover the 3D face geometry as well as its pose.

Previous 3D dense face alignment methods fall into two categories. The *first* category fits a parametric model, specifically, the coefficients of a 3D Morphable Model (3DMM) [4] and the projection matrix from a 2D face image [31, 20]. The recovered geometry accuracy is up to the capacity of the 3D Morphable Model. Usually, it is hard to recover the details that are missing in 3DMM.

The *second* category directly estimates the 3D face coordinates in grid-like data structure (e.g., volumetric representation [15] and uv position map [10]), so that convolutional neutral networks (CNNs) can be used for the regression. Such methods are capable of recovering the geometric details. However, their grid-like data structure is misaligned with the 3D geometry data format (usually a mesh). This causes errors by *data representation*. For example, volumetric representation [15] introduces quantization errors
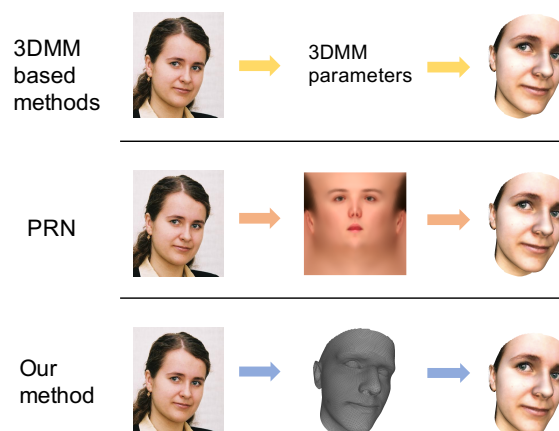


Figure 1. Illustration of three different 3D dense face alignment methods. Top row: 3DMM based methods [31, 20]. Middle row: PRN [10] uses uv position map to store 3D face coordinates. Bottom row: our graph convolution method directly regresses 3D face coordinates on the mesh.

during voxelization of the mesh. The generated 3D face is usually grainy. uv position map [10] introduces distortion, especially in the exterior area of the position map. This is illustrated in Fig. 1 (middle).

To address the problem in the second category, we propose to directly regress the 3D coordinates on the face mesh, as shown in Fig. 1. Because 3D face mesh is represented by a graph, ordinary CNNs are not suitable to process it. Although one can use fully connected network for such regression, this is brute force and computationally unaffordable. Instead, we propose to use graph convolution networks [9, 18] for the task, which are born to work with the graph data. Our method directly performs feature learning on the 3D face mesh, where the geometric structure is well preserved. There is no loss in data representation as in [10, 15]. 3D geometric details are recovered.

The contribution of this work is two fold. First, *for the first time*, we propose to use graph convolution for the 3D dense face alignment task and develop a novel implementation. We integrate well developed techniques in standard

CNNs, such as encoder-decoder architecture, residual learning [12] and Instance Normalization [26] with graph convolution. A coarse to fine strategy is developed using mesh sampling techniques [25, 11]. The resulting network architecture converges well in training and is fast in inference.

Second, comprehensive experiments are performed on several challenging datasets. Extensive quantitative and qualitative results show that our approach gains better results than state-of-the-art methods. This verifies the effect of our approach.

## 2. Related work

**3D dense face alignment** An early representative work is 3DDFA [31]. It uses a cascaded CNN to fit the 3DMM [4] parameters and the projection matrix from a single 2D face image. It shows promising alignment result across large poses, and defines the new problem called 3D dense face alignment. [20] utilizes multi-constraints such as face contours and SIFT feature points to estimate the 3D face shape. It provides a very dense 3D alignment. Although these 3DMM based methods achieve good performance, their accuracy is up to the capacity of the 3DMM. PRN [10] uses a simple encoder-decoder architecture to regress the 3D face coordinates, which are stored in uv position map. Although it achieves state-of-the-art performance, there are many stripes on the reconstructed 3D face. This is attributed to the geometry distortion caused by uv mapping.

**Graph convolution networks** CNNs are suitable to process grid-like data such as images. By contrast, graph convolution networks (GCNs) are suitable to process graph data such as 3D mesh. A comprehensive overview of GCNs is provided by [6]. [21] define the first graph convolution operator on meshes by parameterizing the surface around each point using geodesic polar coordinates and performing graph convolution on the angular bins. Then, different parameterization methods are proposed by [5, 23], but the manners of graph convolution are similar. These methods only present generalization of convolutions to meshes. They do not design sampling operation on meshes. Therefore, coarse-to-fine features cannot be captured.

[7] develops the spectral graph convolution in Fourier space. However, it is computationally expensive and unable to obtain the local features on the graph. To address these problems, ChebyNet [9] formulates spectral convolution as a recursive Chebyshev polynomial, which avoids computing the Fourier basis. Recently, CoMA [25] extends ChebyNet to process 3D meshes. It constructs an autoencoder to learn a latent representation of 3D face and introduces a mesh pooling operator. By utilizing the spectral graph convolution and mesh sampling operations, CoMA [25] obtains state-of-the-art results in 3D face modeling. Motivated by CoMA [25], for the first time, this work uses graph convolution for 3D dense face alignment.
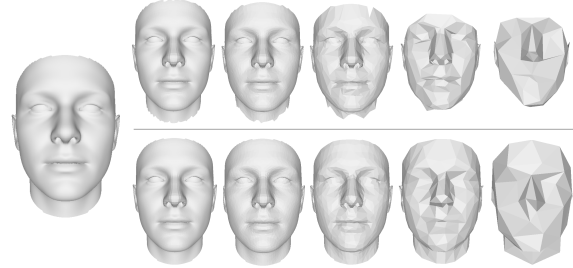


Figure 2. On the leftmost is the 3DMM [4] mesh used in this work. It has 105,954 triangle faces and 53,215 vertices. On the right are down sampled meshes, with 26356, 6528, 1599, 382, 84 triangle faces and 13304, 3326, 832, 208, 52 vertices, respectively. Top/bottom row: without/with boundary preservation [11].

## 3. Method

### 3.1. Graph convolution on face mesh

We briefly review spectral graph convolution, applied on face mesh. More details can be referred to [9].

A 3D face mesh is defined as $\mathcal{M} = (\mathcal{V}, W)$. $\mathcal{V}$ has $N$ 3D vertices on the 3D face surface, $\mathcal{V} \in \mathbb{R}^{N \times 3}$. $W$ is a sparse adjacency matrix of $\mathcal{V}$, $W \in \{0, 1\}^{N \times N}$. $W_{ij} = 1$ denotes there is an edge between vertices $i$ and $j$, and $W_{ij} = 0$ otherwise. The normalized Laplacian matrix is $L = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$, where $I$ is identity matrix and $D$ is a diagonal matrix with $D_{i,i} = \sum_{j=1}^{n} W_{i,j}$. The spectral graph convolution of $x$ and $y$ is defined as a Hadamard product in the Fourier space, $x*y = U \left( \left( U^T x \right) \odot \left( U^T y \right) \right)$. $U$ is the eigenvectors of Laplacian matrix [22]. Since $U$ is not sparse, this operation is computationally expensive. To reduce computation, [9] formulates spectral convolution with a kernel $g_\theta$ using a recursive Chebyshev polynomial, denoted as:

$$g_\theta(L) = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{L}), \tag{1}$$

where $\tilde{L} = 2L/\lambda_{\max} - I$ is the scaled Laplacian matrix, $\lambda_{\max}$ is the maximum eigenvalue of the Laplacian matrix, $\theta \in \mathbb{R}^K$ is the Chebyshev coefficients, and $T_k$ is the Chebyshev polynomial of order $k$, which is computed recursively as $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ with $T_0 = 1$ and $T_1 = x$.

For each layer, the spectral graph convolution is

$$h_j = \sum_{i=1}^{F_{in}} g_{\theta_{i,j}}(L) x_i, \tag{2}$$

where $x_i$ is the $i$-th feature of input $x \in \mathbb{R}^{N \times F_{in}}$, $h_j$ is the $j$-th feature of output $h \in \mathbb{R}^{N \times F_{out}}$. There are $F_{in} \times F_{out}$ vectors of Chebyshev coefficients in a convolution layer.

## 3.2. Face mesh sampling

We use the 3DMM [4] face mesh. It is shown in Fig. 2. Directly working on the high resolution in the original mesh is computationally prohibitive. As in standard CNNs, we develop a coarse-to-fine feature representation. The face mesh is down sampled to several different resolutions, using quadric error metrics [11]. Results are illustrated in Fig. 2. The idea is to minimize the quadric error between the simplified and the original meshes. Since 3D face mesh has open boundaries, performing simplification operation directly as did in [25] causes too much distortion. Therefore, we add boundary preservation constraints as in [11] to preserve the geometry fidelity. As shown in Fig. 2, adding boundary preservation constraints effectively preserves the face geometry. This improves the feature learning in the subsequent graph convolution networks.

During the down-sampling procedure, we follow the method in [25] to compute the up-sampling matrix. More details can be found in [25]. It is used to restore the mesh resolution in the coarse-to-fine feature learning.

## 3.3. Our proposed network

Our proposed network consists of an encoder and a decoder. The encoder is Resnet-50 [12]. It encodes the input 2D face image into a feature vector. The decoder is shown in Fig. 3. It restores the 3D face mesh vertices from the encoded feature in a hierarchical, coarse-to-fine manner. It consists of 6 graph convolution residual blocks, called ResGCN Block in this work, and two graph convolution layers. Each ResGCN block corresponds to a mesh resolution in Fig. 2 and performs feature learning on that level. It is followed by an up-sampling operation, except for the last one.

The ResGCN Block is illustrated in Fig. 4. The identity path in the block helps convergence in the training. Instead of the commonly used Batch Normalization (BN) [14], Instance Normalization (IN) [26] is used in the block. This is because there is no strong statistical connection between face meshes under different poses. On the contrary, there is strong correlation among coordinates of a single face mesh since their combination determines the orientation of the face. Therefore, IN is preferred over BN. Leaky Relu [28] is the activation function. Each block contains two graph convolution layers. The channel numbers of the six ResGCN Block are 128, 64, 32, 32, 16 and 16 respectively. The output features from the last ResGCN block are fed in the last two graph convolution layers to generate the 3D face vertices' coordinates.

In experiments, we found the network converges well in training. If the identity path in the block is removed, or BN is used instead of IN, the training does not converge.

**Loss function** Similar to [25], we adopt $\mathcal{L}_1$ loss for predicted 3D face vertices' coordinates $\tilde{Y}$. For each 2D face image, its ground truth 3D mesh $\mathcal{M}$ is obtained by Basel Face Model(BFM) [4].

We also use a smooth loss $\mathcal{L}_{smooth}$, denoted as:

$$\mathcal{L}_{smooth} = ||(D - W)\tilde{Y}||_2. \tag{3}$$

Note that $(D - W)$ is the unnormalized Laplacian matrix [22] of mesh $\mathcal{M}$. Thus, $(D - W)\tilde{Y}$ denotes the difference between each vertex and its surrounding vertices. Adding the smoothness loss regularizes the training. A weight $\alpha$ is used to balance the smoothness and $\mathcal{L}_1$.

## 4. Experiments

### 4.1. Training

We use 300W-LP [31] as our training dataset, the same as [10]. 300W-LP contains more than 60K unconstrained images with fitted 3DMM parameters and pose coefficients. The 3DMM parameters are obtained by Basel Face Model (BFM) [4]. For each 2D face image in 300W-LP [31], we first use BFM to generate the corresponding 3D mesh. The pose coefficients are then utilized to transform it as the ground truth. We compute the bounding box using the labeled 68 facial keypoints provided by 300W-LP, and then use the the bounding box to crop the image and resize the cropped image to size $256 \times 256$. Following the same setting as in [10], we randomly rotate the input image by -45 to 45 degrees and perturb it with a random translation of 10% of the input size. In addition, a random scale from 0.85 to 1.15 is added. The order number $K$ of all graph convolution layers is set as 3. The dimension of the encoding vector is 256. We train our network for 80 epochs with batch size 50. We choose Adam optimizer [17], where the initial learning rate is 0.001, decayed by half every 20 epochs. The weight $\alpha$ of smooth loss is 0.1. All experiments are conducted on the Geforce GTX 1080 Ti GPU using Tensorflow [1].

### 4.2. Evaluation datasets and metrics

To verify the performance on both face alignment and face reconstruction tasks, we choose the following three datasets as our evaluation benchmarks.

**AFLW2000-3D** [31] contains the first 2000 images from AFLW [19] whose annotations include 68 3D facial landmarks and the fitted 3DMM parameters. Performance of our method on face alignment and face reconstruction tasks is evaluated on this dataset.

**AFLW-LFPA** is constructed by [16]. It contains 1299 images with a balanced distribution of face postures. For each image, 34 facial landmarks are provided. This database is used for evaluation on face alignment task.

**Florence** [2] is a common used benchmark for 3D face reconstruction. It consists of high-resolution 3D scans of 53
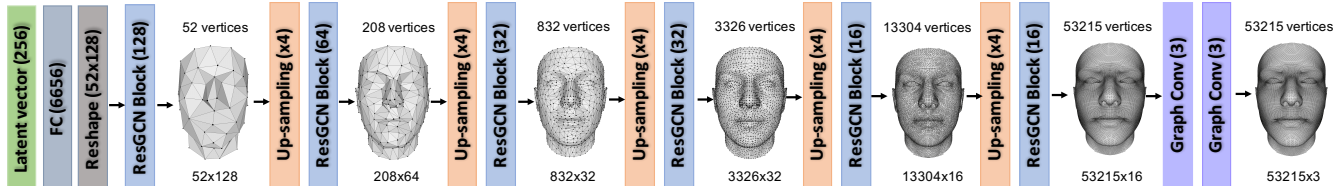
Figure 3. The structure of the decoder. It consists of 6 graph convolution residual blocks, corresponding to the 6 resolution levels of the face mesh in Fig. 2. Each block is followed by an up-sampling operation except for the last one. The last graph convolution layer generates the 3D face vertices' coordinates.
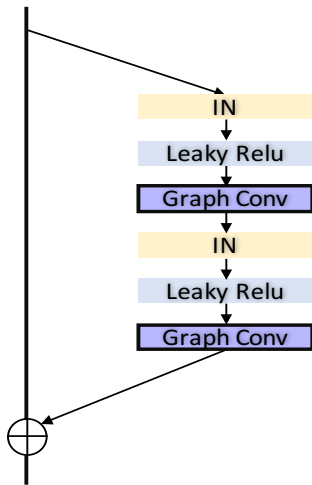


Figure 4. Our proposed graph convolution residual block.

subjects. We follow the protocol of [15] to generate renderings with different poses, using a pitch of -15, 20 or 25 degrees and each of the five evenly spaced rotations between -80 and 80 degrees.

We employ the Normalized Mean Error(NME) as the evaluation metric, which is:

$$\text{NME} = \frac{1}{N} \sum_{i=1}^{N} \frac{\left\| Y_i - \tilde{Y}_i \right\|_2}{d}, \qquad (4)$$

where $Y_i$ is the $i$-th coordinate of the ground truth 3D face and $\tilde{Y}_i$ denotes the $i$-th coordinate of the predicted 3D face. $d$ is the normalization factor.

## 4.3. Results

In this part, we first present the 2D and 3D face alignment performance compared with several state-of-the-art methods on AFLW2000-3D and AFLW-LFPA datasets. Then the results of 3D face reconstruction on Florence are shown. At last, ablation study about our different experimental setting is demonstrated. The qualitative results are shown in Fig. 5, notice that our method can guarantee good face alignment and 3D face reconstruction performance even in cases of large pose, occlusion and weak illumination.

### 4.3.1 Face alignment

We first conduct face alignment experiments on 68 *sparse* facial landmarks on AFLW2000-3D. We follow [31] to use bounding box size as the normalization factor. Several recent state-of-the-art methods are selected for comparison, including 3DDFA [31], DeFA [20], 3D-FAN [8] and PRN [10]. As shown in Fig. 6, our approach outperforms other methods both on 2D and 3D face alignment tasks by a large margin. Specifically, more than 10% higher performance is achieved compared with the best 3D face alignment method. It shows that our approach can locate landmarks more accurately.

In addition, we also present our *dense* face alignment results compared with other recent methods including 3DDFA [31], DeFA [20] and PRN [10] on AFLW2000-3D. Follow the setting of PRN, we select 45K points from the largest common face region of all compared methods. The quantitative results is illustrated in Fig. 7. Our method is superior to the best state-of-the-art method. In addition, our generated 3d faces have better visual quality. Examples and comparison are illustrated in Fig. 8. For better display, we rotate the 3d mesh into a front face and zoom in the nose area. Note that there are many stripes on face surface of PRN. By contrast, our results are smoother. Besides, our generated face has finer details and better correspondence to the ground truth. By comparison, PRN adopts uv position map as the regression target, there is geometric distortion between it and the 3D coordinates. Therefore, it produces less precise results.

To investigate the performance of our method across different poses and datasets, we conduct face alignment experiments with different yaw angles on AFLW2000-3D and AFLW-LPFA datasets. Following the protocol of [31], we randomly select 696 images from AFLW2000-3D, whose absolute yaw angles with small, medium and large values are 1/3 each. 68 points of AFLW2000-3D and 34 points of AFLW-LPFA are selected for evaluation. As illustrated in Table 1, our method surpasses the previous methods by a large margin. Specifically, our performance outperforms PRN more than 9% both on AFLW2000-3D and AFLW-

Figure 5. The qualitative results of our 3D dense face alignment method. The first row of each person is the alignment results(68 landmarks are plotted), the second row is the face rendered by the corresponding depth image, the last row is the 3D reconstruction results.
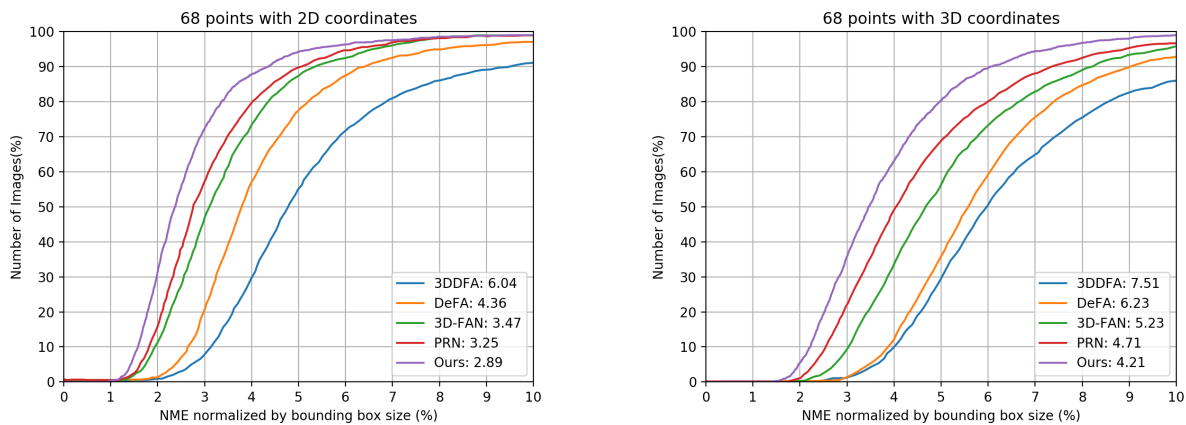


Figure 6. Cumulative Errors Distribution (CED) curves of 68 points face alignment on AFLW2000-3D. The left is result of 2D face alignment. The right is the result of 3D face alignment.

LPFA. This verifies that our method has good face alignment performance even in the case of large face pose.

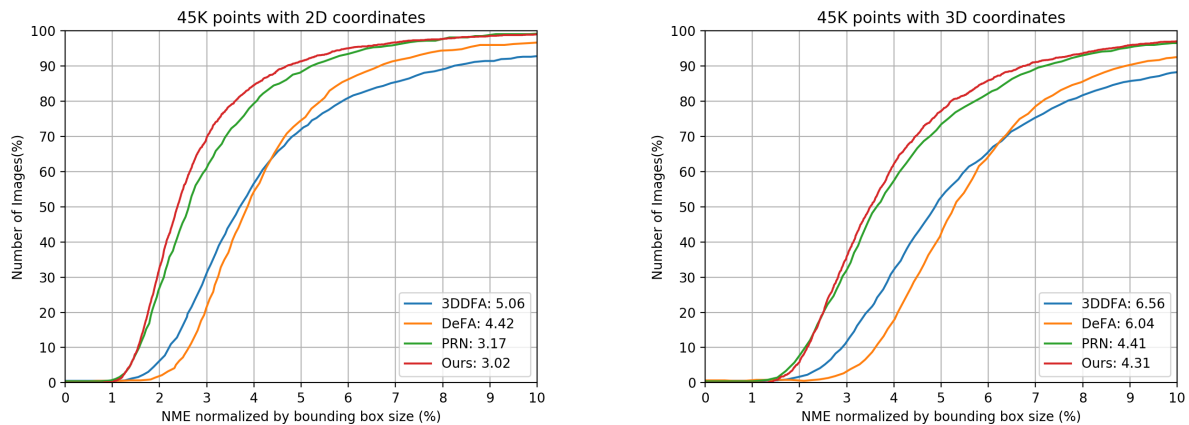Some face alignment examples of AFLW2000-3D are presented in Fig. 9, we find that in some cases, our predicted

Figure 7. Cumulative Errors Distribution (CED) curves of 45K points face aligment on AFLW2000-3D. The left is result of 2D face alignment. The right is the result of 3D face alignment.
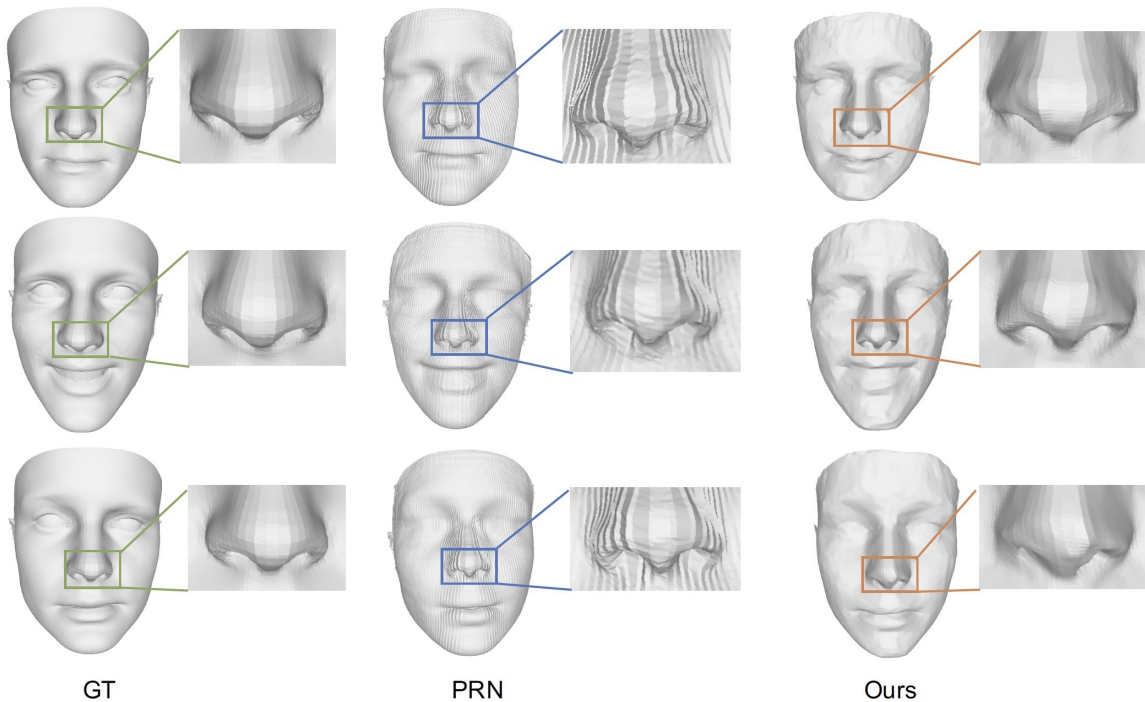


Figure 8. Examples of the reconstructed 3D faces. The left is the ground truth, the middle is result of PRN and the right is our result. The nose area is zoomed in for better view.

landmarks are more accurate than the ground truth. This is due to that the ground truth points of AFLW2000-3D are generated by the semi-automatic annotation pipeline of [31] rather than manual annotation. This phenomenon shows the high accuracy of our method in locating landmarks of faces.

### 4.3.2 3D face reconstruction

In addition to face alignment, we also conduct experiments on 3D face reconstruction task. In this subsection, we inves-

tigate the face reconstruction performance of our method on Florence dataset [2]. Several state-of-the-art methods are chosen for comparison, including 3DDFA [31], VRN [15] and PRN [10]. Following the experimental setting of [15], we render the testing images of Florence with different poses, the details have been introduced in Section 4.2. Before input to the network, the image are cropped using the bounding box computed from the ground truth point cloud. We select the most common 19K points of all compared methods to perform evaluation. As the outputs of different

| | AFlW2000-3D | | | | AFLW-LFPA |
|---|---|---|---|---|---|
| Method | o to 30 | 30 to 60 | 60 to 90 | Mean | Mean |
| SDM [27] | 3.67 | 4.94 | 9.67 | 6.12 | - |
| 3DDFA [31] | 3.78 | 4.54 | 7.93 | 5.42 | - |
| 3DDFA+SDM [31] | 3.43 | 4.24 | 7.17 | 4.94 | - |
| Yu *et al.* [29] | 3.62 | 6.06 | 9.56 | - | - |
| 3DSTN [3] | 3.15 | 4.33 | 5.98 | 4.49 | - |
| DeFA [20] | - | - | - | 4.50 | 3.86 |
| PRN [10] | 2.75 | 3.51 | 4.61 | 3.62 | 2.93 |
| Ours | **2.44** | **3.26** | **4.35** | **3.35** | **2.65** |

Table 1. Performance comparison (NME) between our method and other state-of-the-art methods on AFLW2000-3D and AFLW-LPFA benchmarks. The first best result in each category is highlighted in bold, the lower is the better.
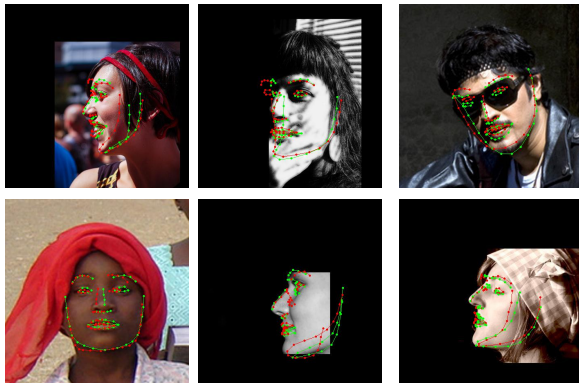


Figure 9. Examples from AFLW2000-3D. Our face alignment results are more accurate than the ground truth in some cases. Red is the ground truth, green is our predicted landmarks.
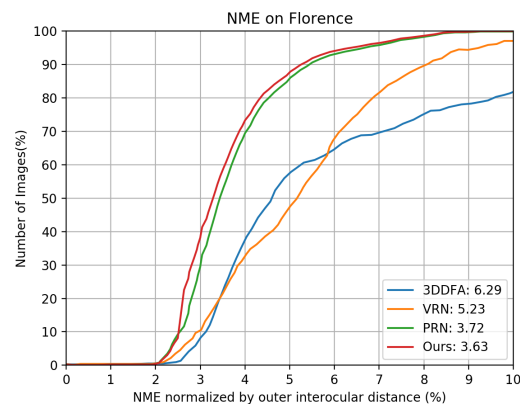


Figure 10. Cumulative Errors Distribution (CED) curves of 19K points face reconstruction on orence.



Figure 11. Examples of our reconstructed 3D faces from Florence.

methods are not aligned, we follow [10] to utilize Iterative Closest Points algorithm to find the nearest points between the network output and the ground truth point cloud. After aligning the generated point cloud, Mean Squared Error normalized by outer interocular distance of 3D coordinates is adopted as the evaluation metric. The quantitative result is shown in Fig. 10, our performance is slightly better than PRN. This is due to the fact that the labeled 3D face meshes of our training data 300W-LP are from 3DMM fitting by [31], while the 3d faces of Florence are acquired from a structured-light scanning system [2], there exists a large gap between the two kinds of mesh annotations. So our model trained on 300W-LP does not bring much performance improvement on Florence data. In spite of this, our method still achieve a good correspondence from 2D images to the 3D face meshes, some examples are shown in Fig. 11. As is illustrated, the face shape and expression details are well captured by our method.

### 4.3.3 Ablation study

In this part, we analyze the effectiveness of smooth loss and present quantitative and qualitative results of different weights $\alpha$. Fig. 12 shows the quantitative face alignment performance comparisons on AFLW2000-3D when $\alpha$ is set to 0, 0.01, 0.1 and 1.0. The corresponding qualitative results are presented in Fig. 13. As is illustrated, 0.1 is an appropriate weight which can well balance the the performance and smoothness of the generated 3D faces. This can be interpreted as the following fact, when the value of $\alpha$ is small (0 or 0.01), the generation of 3D face mainly depends on the guidance of $\mathcal{L}_1$, so that the generated surfaces are not smooth. When it becomes large ($\alpha$=1), the smooth
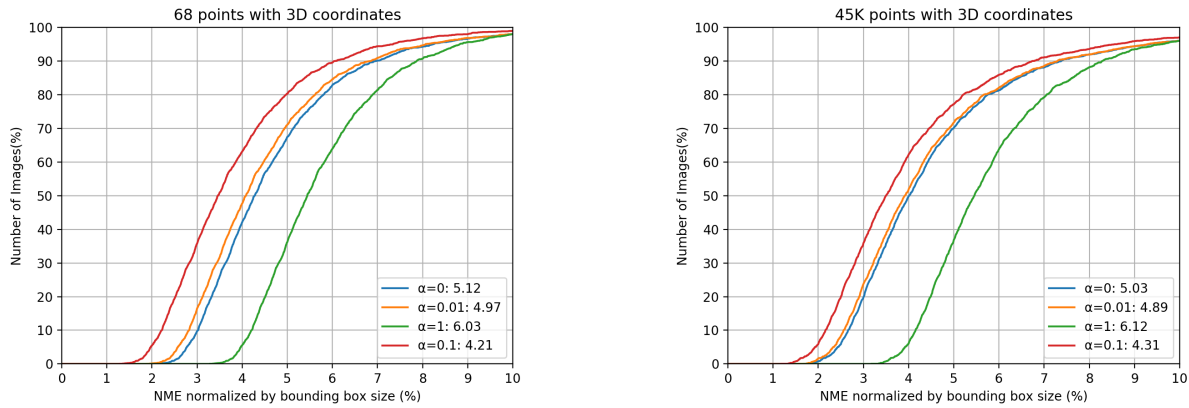
Figure 12. Influence of the smooth loss on 3D face alignment performance. The left is performance with 68 points and the right is performance with 45K points.
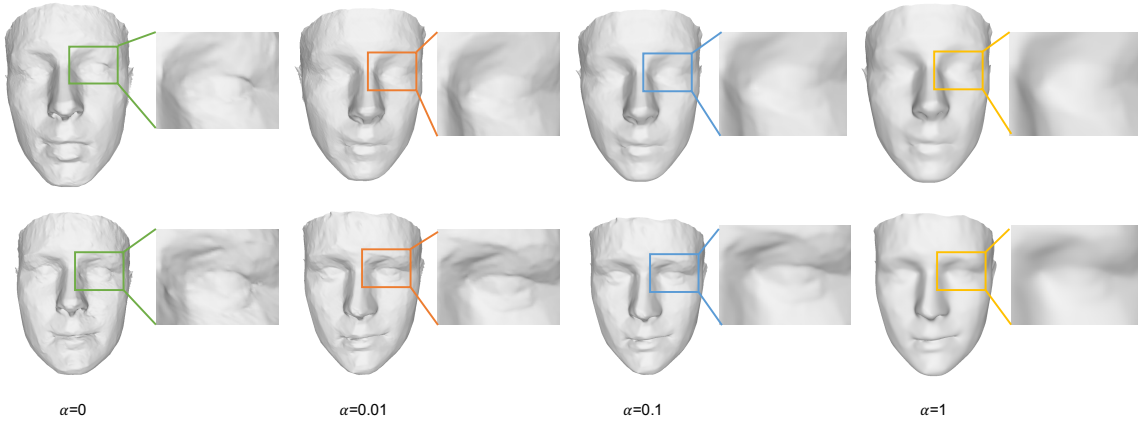


Figure 13. Examples of generated 3D faces with different smoothness. The eye area is zoomed in for better view.

loss $\mathcal{L}_{smooth}$ seriously affects the dominant role of $\mathcal{L}_1$ in training procedure. As a result, over-smoothing 3D faces are generated. In contrast, setting $\alpha$ as 0.1 is an appropriate option, which can not only achieve better performance, but also improve the visual quality of 3D face.

## 5. Conclusions

In this paper, we propose graph convolution networks to solve the problem of 3D dense face alignment. Our network captures coarse-to-fine features of face mesh in a hierarchical manner and generate 3D face coordinates directly. Extensive experiments show that our approach gains a superior performance both on face alignment and 3D face reconstruction tasks over other state-of-the-art methods.

## References

[1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.

[2] A. D. Bagdanov, A. Del Bimbo, and I. Masi. The florence 2d/3d hybrid face dataset. In *Proceedings of the 2011 joint ACM workshop on Human gesture and behavior understanding*, pages 79–80. ACM, 2011.

[3] C. Bhagavatula, C. Zhu, K. Luu, and M. Savvides. Faster than real-time facial alignment: A 3d spatial transformer network approach in unconstrained poses. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3980–3989, 2017.

[4] V. Blanz, T. Vetter, et al. A morphable model for the synthesis of 3d faces. In *Siggraph*, volume 99, pages 187–194, 1999.

[5] D. Boscaini, J. Masci, E. Rodolà, and M. Bronstein. Learning shape correspondence with anisotropic convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 3189–3197, 2016.

[6] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.

[7] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.

[8] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1021–1030, 2017.

[9] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems*, pages 3844–3852, 2016.

[10] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 534–551, 2018.

[11] M. Garland and P. S. Heckbert. Surface simplification using quadric error metrics. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 209–216. ACM Press/Addison-Wesley Publishing Co., 1997.

[12] 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.

[13] P. Huber, G. Hu, R. Tena, P. Mortazavian, P. Koppen, W. J. Christmas, M. Ratsch, and J. Kittler. A multiresolution 3d morphable face model and fitting framework. In *Proceedings of the 11th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2016.

[14] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

[15] A. S. Jackson, A. Bulat, V. Argyriou, and G. Tzimiropoulos. Large pose 3d face reconstruction from a single image via direct volumetric cnn regression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1031–1039, 2017.

[16] A. Jourabloo and X. Liu. Large-pose face alignment via cnn-based dense 3d model fitting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4188–4196, 2016.

[17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[18] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.

[19] M. Koestinger, P. Wohlhart, P. M. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *2011 IEEE international conference on computer vision workshops (ICCV workshops)*, pages 2144–2151. IEEE, 2011.

[20] Y. Liu, A. Jourabloo, W. Ren, and X. Liu. Dense face alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1619–1628, 2017.

[21] J. Masci, D. Boscaini, M. Bronstein, and P. Vandergheynst. Geodesic convolutional neural networks on riemannian manifolds. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 37–45, 2015.

[22] R. Merris. Laplacian matrices of graphs: a survey. *Linear algebra and its applications*, 197:143–176, 1994.

[23] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5115–5124, 2017.

[24] X. Peng, R. S. Feris, X. Wang, and D. N. Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *European conference on computer vision*, pages 38–56. Springer, 2016.

[25] A. Ranjan, T. Bolkart, S. Sanyal, and M. J. Black. Generating 3d faces using convolutional mesh autoencoders. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 704–720, 2018.

[26] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.

[27] X. Xiong and F. De la Torre. Global supervised descent method. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2664–2673, 2015.

[28] B. Xu, N. Wang, T. Chen, and M. Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.

[29] R. Yu, S. Saito, H. Li, D. Ceylan, and H. Li. Learning dense facial correspondences in unconstrained images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4723–4732, 2017.

[30] E. Zhou, H. Fan, Z. Cao, Y. Jiang, and Q. Yin. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 386–391, 2013.

[31] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 146–155, 2016.

[32] X. Zhu, Z. Lei, J. Yan, D. Yi, and S. Z. Li. High-fidelity pose and expression normalization for face recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 787–796, 2015.