

Introduction

Face recognition and 3D face reconstruction seem to **contradict** each other. On one hand, face recognition This is formulated by prefers identity-sensitive features, but not every detail on faces; on the other hand, 3D reconstruction attempts to recover as much facial detail as possible, regardless whether the detail benefits or distracts facial identity recognition.

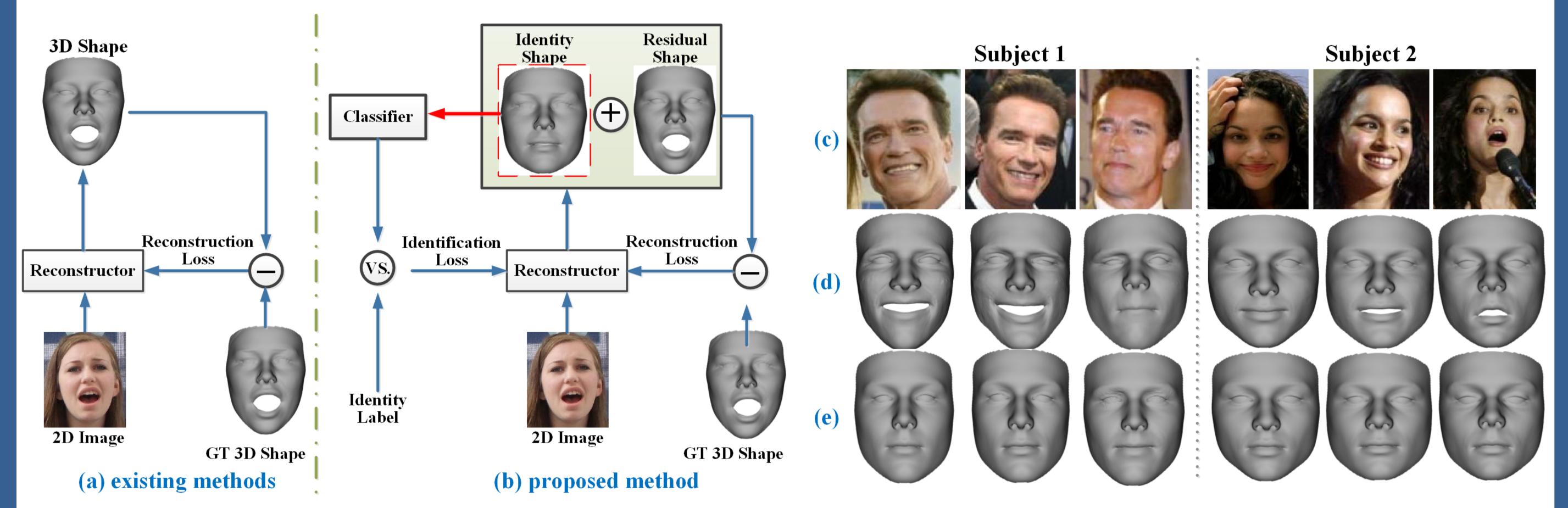


Figure: Comparison between the learning process of (a) existing methods and (b) our proposed method. GT denotes Ground Truth. (d) and (e) are **3D** face shapes and disentangled identity shapes reconstructed by our method for the images in (c) from

- . We propose a method which for the first time explicitly optimizes face recognition and 3D face reconstruction simultaneously. The method achieves state-of-the-art **3D** face reconstruction accuracy via joint discriminative feature learning and **3D** face reconstruction.
- •We devise an effective training process for the proposed network that can disentangle identity and nonidentity features in reconstructed 3D face shapes. The network, while being pre-trained by 3DMMgenerated data, can surmount the limited 3D shape space determined by the 3DMM bases, in the sense that it better captures identity-sensitive and identity-irrelevant features in **3D** face shapes.
- . We leverage the effectiveness of disentangled identity features in reconstructed 3D face shapes for improving face recognition accuracy, as being demonstrated by our experimental results. This further expands the application scope of **3D** face reconstruction.

Proposed Method

Based on the assumption that 3D face shapes are composed by identity-sensitive and identity-irrelevant parts, the **3**D face shape **s** of a subject is represented as

$$s = \bar{s} + \Delta s_{Id} + \Delta s_{Res}$$

where \bar{s} is the mean **3**D face shape (computed across all training samples with neutral expression), Δs_{ld} is the identity-sensitive difference between s and \bar{s} , and Δs_{Res} denotes the residual difference. A variety of sources could lead to the residual difference, for example, expression-induced deformations and temporary detail

References:

- [1] SphereFace: Liu et al. Sphereface: Deep hypersphere embedding for face recognition, CVPR 2017.
- [3] 3DDFA: Zhu et al. Face Alignment Across Large Poses: A 3D Solution, CVPR 2016.
- [5] 3DMM-CNN: Tran et al. Regressing robust and discriminative 3D morphable models with a very deep neural network, CVPR 2017.

Disentangling Features in 3D Face Shapes for Joint Face Reconstruction and Recognition Feng Liu¹, Ronghang Zhu¹, Dan Zeng¹, Qijun Zhao¹ and Xiaoming Liu²

¹College of Computer Science, Sichuan University

²Department of Computer Science and Engineering, Michigan State University

We further assume that Δs_{Id} and Δs_{Res} can be described by latent representations, c_{Id} and c_{Res} , respectively.

Here, f_{Id} (f_{Res}) is the mapping function that generates the corresponding shape component Δs_{Id} (Δs_{Res}) from the latent representation, with parameters θ_{Id} (θ_{Res}).

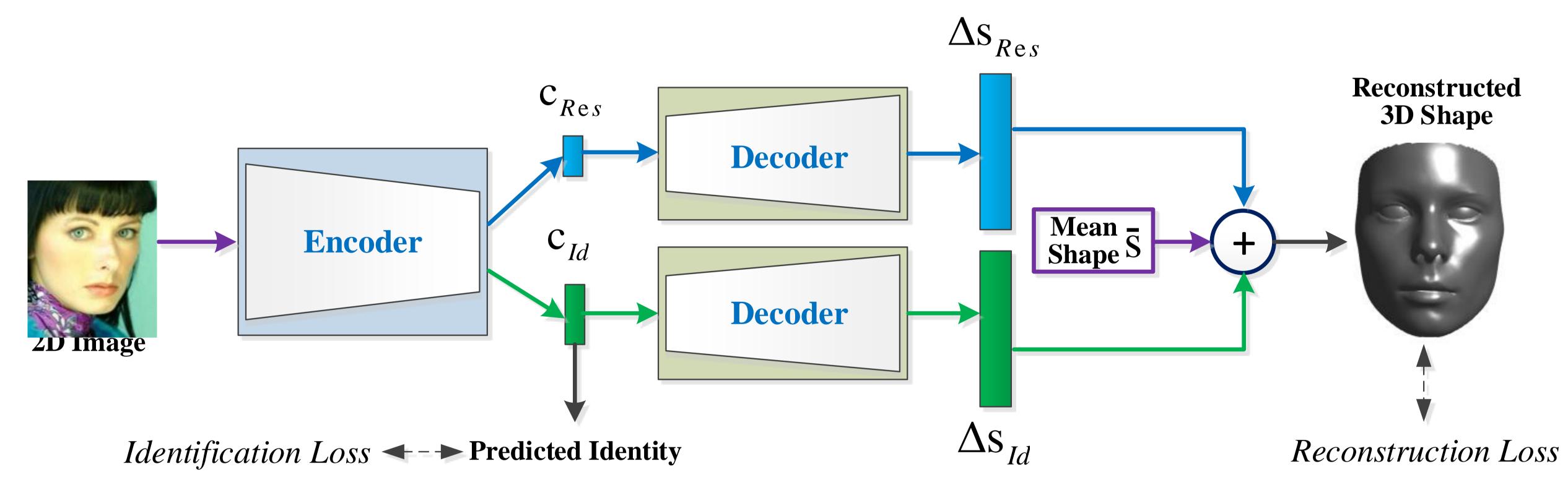


Figure: Overview of the proposed encoder-decoder based joint learning pipeline for face recognition and **3D** shape reconstruction.

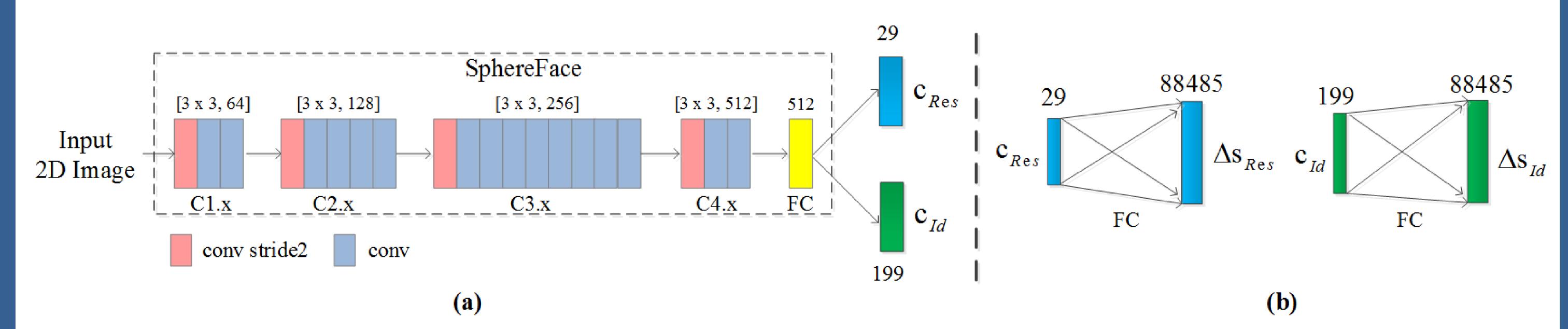


Figure: (a) Encoder in the proposed method is implemented based on SphereFace. It converts the input 2D image to latent identity and residual shape feature representations. (b) Decoders in the proposed method are implemented as a fully connected (FC) layer. They convert the latent representations to corresponding shape components.

Implementation Detail

Loss Functions. The overall loss to the proposed encoder-decoder network is defined by

 $\boldsymbol{L} = \lambda_{\boldsymbol{B}} \boldsymbol{L}_{\boldsymbol{B}} + \boldsymbol{L}_{\boldsymbol{C}},$

where λ_{R} is the Euclidean-based weight for the reconstruction loss L_{R} , L_{C} is the softmax loss.

- Training Data. To construct training data we develop a method for fitting **3D** morphable model (**3DMM**) to multiple 2D images of a subject, and apply it to CASIA-WebFace database, resulting in 488,848 images / 3D faces of **10**, **575** subjects.
- Training Process. We train our encoder-decoder network in three phases. In Phase I, we train the encoder by setting the target latent representations as 3DMM coefficients and using Euclidean loss. In Phase II, we train the decoder for the identity and residual components separately. In Phase III, the end-to-end joint training is conducted based on the pre-trained encoder and decoder.

 $\Delta S_{Id} = f_{Id}(C_{Id};\theta_{Id}), \ \Delta S_{Res} = f_{Res}(C_{Res};\theta_{Res}).$

Experimental Results

Two sets of experiments have been done to evaluate the effectiveness of the proposed method in 3D face reconstruction and face recognition. The MICC and BU3DFE databases are used for experiments of 3D face reconstruction, and the LFW, YTF and IJB-A databases are used in face recognition experiments. **3D** Shape Reconstruction Accuracy.

 Table: 3D face reconstruction acc
 Method VRN 3DDFA 3DM RMSE 5.34 2.73

• Face Recognition Accuracy.

	V	66.13 74.93 75.25	$\begin{array}{c}\pm \ \textbf{2.79}\\\pm \ \textbf{1.14}\end{array}$	65.70 <u>-</u> 74.50 <u>-</u>	± 2.81 ± 1.21	$\textbf{82.94} \pm$	2.75			$\begin{array}{r} \textbf{12.37} \pm \textbf{4} \\ \textbf{28.73} \pm \textbf{7} \end{array}$
		74.93 75.25	± 1.14	74.50 =	± 1.21	$\textbf{82.94} \pm$				
$ \begin{array}{c} \times \\ \checkmark \\ \hline \checkmark \\ \hline \checkmark \\ \hline \checkmark \end{array} $		75.25					1.14	$\textbf{60.40} \pm$	3.15	28.73 + 7
	V		+2.12	71 70					••••	
	Х			14.1J =	± 2.56	83.21 ±	1.93	$\textbf{59.40} \pm$	4.64	$\textbf{29.67} \pm \textbf{4}$
										10.00 ± 3
v	×									58.20 ± 12
×	\checkmark									52.60 ± 8
\checkmark	\checkmark									
\checkmark	×	94.43	\pm 1.47	94.40 =	± 1.52	98.12 ±	0.90	95.07 ±	2.39	74.54 ± 4
			You	Fube Fa	ices ()	(TF)				
\checkmark	×	73.26	\pm 2.51	73.08 -	± 2.65	80.41 ±	2.60	$\textbf{51.36} \pm$	5.11	$\textbf{24.04} \pm \textbf{4}$
×	\checkmark	77.34	$\pm \textbf{2.54}$	76.96 =	± 2.64	$\textbf{85.32} \pm$	2.63	$\textbf{63.16} \pm$	5.07	31.36 ± 5
\checkmark		79.56	$\pm\textbf{2.08}$	79.20 =	± 2.07	$\textbf{87.35} \pm$	1.92	$\textbf{69.08} \pm$	5.00	34.56 ± 6
\checkmark	×	68.10	\pm 2.93	67.96 <u>-</u>	± 3.12	$\textbf{74.95} \pm$	3.04	$\textbf{40.52} \pm$	3.65	$\textbf{12.20} \pm \textbf{2}$
\checkmark	×	88.28	\pm 1.84	88.32 -	± 2.16	95.95 ±	1.38	86.60 ±	3.95	51.12 ± 8
×	\checkmark	87.56	\pm 2.56	87.68 =	± 2.25	94.44 ±	1.38	$\textbf{84.80} \pm$	4.89	40.92 ± 8
\checkmark	V									
\checkmark	×	88.74	\pm 1.03	88.70 =	± 1.15	96.28 ±	0.63	$\textbf{89.00} \pm$	2.40	53.44 ± 4
	$ \begin{array}{c} \\ $	$ \begin{array}{c} \cdot \\ \checkmark \\ \times \\ \times \\ \times \\ \checkmark \\ \checkmark$	$$ \times 94.43 $$ \times 73.26 \times $$ 77.34 $$ $$ 79.56 $$ \times 68.10 $$ \times 88.28 \times $$ 88.28 \times $$ 88.28 \times $$ 88.28 $$ \times 88.28 $$ $$ 88.28 $$ \times 88.28 $$ $$ 88.28 $$ $$ 88.28 $$ $$ 88.29 $$ $$ 88.74	\checkmark ×94.43 \pm 1.47 \checkmark You \checkmark × \checkmark 73.26 \pm 2.51 \times \checkmark \checkmark 77.34 \pm 2.54 \checkmark \checkmark \checkmark 79.56 \pm 2.08 \checkmark × \checkmark 68.10 \pm 2.93 \checkmark × \checkmark 88.28 \pm 1.84 \times \checkmark \checkmark 88.26 \pm 2.56 \checkmark × \checkmark × \checkmark ×× \checkmark ×× \checkmark 88.80 \pm 2.21 \checkmark ××× </td <td>$\checkmark$$> 94.43 \pm 1.47 \ 94.40 =$YouTube Fa$\checkmark$$\times$$73.26 \pm 2.51 \ 73.08 =$$\times$$\checkmark$$77.34 \pm 2.54 \ 76.96 =$$\checkmark$$\checkmark$$79.56 \pm 2.08 \ 79.20 =$$\checkmark$$\checkmark$$68.10 \pm 2.93 \ 67.96 =$$\checkmark$$\times$$68.10 \pm 2.93 \ 67.96 =$$\checkmark$$\times$$88.28 \pm 1.84 \ 88.32 =$$\checkmark$$\checkmark$$88.28 \pm 1.84 \ 88.32 =$$\checkmark$$\checkmark$$88.80 \pm 2.21 \ 88.84 =$$\checkmark$$\checkmark$$88.74 \pm 1.03 \ 88.70 =$</td> <td>$\checkmark$94.43 \pm 1.4794.40 \pm 1.52YouTube Faces (Normalized Science Scie</td> <td>\checkmark94.43 \pm 1.4794.40 \pm 1.5298.12 \pmYouTube Faces (YTF)$\checkmark$$\times$73.26 \pm 2.5173.08 \pm 2.6580.41 $\pm$$\times$$\checkmark$77.34 \pm 2.5476.96 \pm 2.6485.32 $\pm$$\checkmark$$\checkmark$79.56 \pm 2.0879.20 \pm 2.0787.35 $\pm$$\checkmark$$\checkmark$68.10 \pm 2.9367.96 \pm 3.1274.95 $\pm$$\checkmark$$\times$88.28 \pm 1.8488.32 \pm 2.1695.95 $\pm$$\checkmark$$\checkmark$88.28 \pm 1.8488.32 \pm 2.2594.44 $\pm$$\checkmark$$\checkmark$88.80 \pm 2.2188.84 \pm 2.4095.37 $\pm$$\checkmark$$\checkmark$88.74 \pm 1.0388.70 \pm 1.1596.28 \pm</td> <td>$\begin{array}{c ccccccccccccccccccccccccccccccccccc$</td> <td>$\begin{array}{c ccccccccccccccccccccccccccccccccccc$</td> <td>$\begin{array}{c ccccccccccccccccccccccccccccccccccc$</td>	\checkmark $> 94.43 \pm 1.47 \ 94.40 =$ YouTube Fa \checkmark \times $73.26 \pm 2.51 \ 73.08 =$ \times \checkmark $77.34 \pm 2.54 \ 76.96 =$ \checkmark \checkmark $79.56 \pm 2.08 \ 79.20 =$ \checkmark \checkmark $68.10 \pm 2.93 \ 67.96 =$ \checkmark \times $68.10 \pm 2.93 \ 67.96 =$ \checkmark \times $88.28 \pm 1.84 \ 88.32 =$ \checkmark \checkmark $88.28 \pm 1.84 \ 88.32 =$ \checkmark \checkmark $88.80 \pm 2.21 \ 88.84 =$ \checkmark \checkmark $88.74 \pm 1.03 \ 88.70 =$	\checkmark 94.43 \pm 1.4794.40 \pm 1.52YouTube Faces (Normalized Science Scie	\checkmark 94.43 \pm 1.4794.40 \pm 1.5298.12 \pm YouTube Faces (YTF) \checkmark \times 73.26 \pm 2.5173.08 \pm 2.6580.41 \pm \times \checkmark 77.34 \pm 2.5476.96 \pm 2.6485.32 \pm \checkmark \checkmark 79.56 \pm 2.0879.20 \pm 2.0787.35 \pm \checkmark \checkmark 68.10 \pm 2.9367.96 \pm 3.1274.95 \pm \checkmark \times 88.28 \pm 1.8488.32 \pm 2.1695.95 \pm \checkmark \checkmark 88.28 \pm 1.8488.32 \pm 2.2594.44 \pm \checkmark \checkmark 88.80 \pm 2.2188.84 \pm 2.4095.37 \pm \checkmark \checkmark 88.74 \pm 1.0388.70 \pm 1.1596.28 \pm	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

Ablation Study.

Table: Reconstruction and recognition accuracy on different test data sets when identity disentangling and identification loss are used or not used. Refer to the paper for test data set details.

Training Dhace	Idontity Dicontonaling	Idantification Lago	Re	econstruction RM	Recognition Accuracy on			
ITAITIITY FITASE	e Identity Disentangling	Identification L055	MICC	BU3DFE (pose)	BU3DFE (exp.)	LFW	YTF	
	×	×	$\textbf{2.51} \pm \textbf{0.57}$	$\textbf{2.54} \pm \textbf{0.67}$	$\textbf{2.62} \pm \textbf{0.73}$			
	\checkmark	×	$\textbf{2.23} \pm \textbf{0.48}$	$\textbf{2.31} \pm \textbf{0.55}$	$\textbf{2.45} \pm \textbf{0.62}$	$\textbf{68.00} \pm \textbf{2.21}$	$\textbf{69.19} \pm \textbf{1.91}$	
	\checkmark	\checkmark	$\textbf{2.00} \pm \textbf{0.32}$	$\textbf{2.01} \pm \textbf{0.49}$	$\textbf{2.19} \pm \textbf{0.54}$	$\textbf{94.43} \pm \textbf{1.47}$	$\textbf{88.74} \pm \textbf{1.03}$	

Reconstruction errors are further reduced after incorporating identification loss in Phase III. Recognition accuracy is significantly improved from Phase II to Phase III. This reveals the limited discrimination power of 3DMM representations and the importance of CNN-based joint learning in expanding the representation and discrimination capacity of 3DMM-like bases.

[2] VRN: Jackson et al. Large pose 3D face reconstruction from a single image via direct volumetric CNN regression, ICCV 2017. [4] 3DMM-CNN: Tran et al. Regressing robust and discriminative 3D morphable models with a very deep neural network, CVPR 2017. [6] DRGAN: Tran et al. Disentangled Representation Learning GAN for Pose-Invariant Face Recognition, CVPR 2017.



Table: 3D face reconstruction accuracy (RMSE) under different yaw angles on the BU3DFE database.

ccuracy or	the M	ICC database.	Method	±90°	± 80 °	±70°	±60°	± 50 °	±40°	± 30 °	±20°	±10°	0 °	Avg.
MM-CNN	3DSR	Proposed	VRN	6.96	6.20	6.14	6.01	5.91	5.50	4.93	3.86	3.70	3.66	5.29
2.20	2.07	2.00	3DDFA 3DMM-CNN									2.55 2.19		
												2.10		
			Proposed	2.09	2.04	2.03	2.03	2.00	1.99	2.03	2.01	1.97	1.93	2.01

• Computational Efficiency. We run the methods on a PC (with an Intel Core i7-5930K @ 3.5GHz, 32GB RAM and an GeForce GTX 1080) for 700 images, and calculate the average runtime per image.

> Table: Efficiency comparison of different methods. Method VRN 3DDFA 3DMM-CNN 3DSR Proposed Time (ms) 55.68 39.17 30.12 29.80 4.79