

Hybrid Approach for 3D Head Reconstruction: Using Neural Networks and Visual Geometry

Oussema Bouafif ^{1,2}

Bogdan Khomutenko¹

Mohamed Daoudi²

¹MCQ-Scan ²IMT Lille Douai

Lille, France

()-Scan





Introduction and Motivation



[Thies et al. 2016]

Facial Animation



[FotoNation.com]

Face Recognition



Avatar Digitization



[Gilani et al. 2014]

Gender Classification

Related Work

Monocular 3D face reconstruction methods







Proposed Method



Synthetic Data Generation



 $X = X_0 + W y$

- •X:3D head,
- X₀ : Mean face shape
- W : Principal components,
- $y \sim N(0, 1)$: Coefficient vector shape.

The Liverpool-York Head Model

(LYHM)



Complete head model

- Textures
- Hair model (USC-HairSalon)
- Eyes model
- Glasses model

- Illumination conditions
- Shadows
- Poses
- Coefficients for PBM
- Background image

Synthetic Data Generation

r



Complete head model

- Local geometric properties,
- Completely independent prediction,
- Translation and scale invariance

[Klasing et al., 2009]



Normal surface (N)



Landmarks map

Network Structure



Inspired by [Su et al.,2018]

$$Loss = L_{\mathcal{N}} + L_{\mathcal{Z}}$$
$$L_{\mathcal{N}} = ||\mathcal{N}_{\text{GT}} - \mathcal{N}||_2^2 \quad , \quad L_{\mathcal{Z}} = ||\mathcal{Z}_{\text{GT}} - \mathcal{Z}||_2^2$$

3D Morphable Model fitting

41-14- 131 1	3DMM fitting	
3D Morphable Model	Morphable Model's Normals	Projection Model
LYHM $X = X_0 + W_y$	$\tilde{n}_{i} = \sum_{j=1}^{11510} \frac{(q_{i,j} - p_{i}) \times (q_{i,j+1} - p_{i})}{\ (q_{i,j} - p_{i}) \times (q_{i,j+1} - p_{i})\ }$	$a=\frac{1}{z}KR(p+t)$
$X \in \mathbb{R}^{3 \times 11510} \rightarrow \text{Generated 3D head}$	$n_i = \frac{\tilde{n}_i}{\ \tilde{n}_i\ }$	a ∈ \mathbb{R}^2 → projected point $\mathbb{R} \in SO(3)$ → rotation matrix z → normalization of a 3D point
$\begin{array}{l} X_0 \in \mathbb{R}^{3 \times 11510} \to \text{ Mean head} \\ W \in \mathbb{R}^{3 \times 11510 \times 100} \to \text{Principal components} \\ y \in \mathbb{R}^{100} \to \text{ Parameter vector} \end{array}$	$ \begin{array}{c} n_i \rightarrow \text{normal vector} \\ p_i \rightarrow \text{vertex location} \\ \{q_{i,1}, q_{i,2}, \dots, q_{i,11510}\} \rightarrow \text{adjacent vertices} \end{array} $	K → projection matrix $\begin{pmatrix} f & 0 & u_0 \\ 0 & f & v_0 \end{pmatrix}$ t ∈ R ³ → translation vector

3D Morphable Model fitting



Parametric Regression



Multi-view 3D Face Reconstruction



Experimental Results

Training Data set :

- 60,000 facial images (30,000 for males and also for females).
- **3**000 epochs, learning rate : $1e^{-5}$, 32 : batch size, *RMSpropoptimizer*
- random blur effect, gaussian noise
- $λ_N$ =1, $λ_Z$ =0.8, $λ_P$ =0.4,



Evaluation



Evaluation

Stereo vs Mono Fitting



- (Multi-view fitting): 3D head reconstruction using all images in the same fitting process.
- (Mono fitting) : 3D head exploiting only the frontal image in the fitting process (third row).
- (GT) : the ground-truth of the 3D head mesh.

Final Evaluation

Absolute denth error Point-to-plane Root Mean Square Error (RMSE)

RMSE
$$\rightarrow \epsilon = \sqrt{\frac{\sum_{i=1}^{N} ((p_i - q_i).n_i)^2}{N}}$$

N → number of vertices

 $P_i \rightarrow$ vertex of the reference model

 $q_i \rightarrow$ nearest neighbor from the reconstructed model

 $n \rightarrow normal vector of the reference model$

Method	Ours (Mono)	Ours (Multi)	RingNet [3]	PRN [4]	R-C-Nets [5]
RMSE	1.74±0.44	1.67±0.43	1.90±0.49	1.86±0.47	1.60±0.41
μ	2.21	2.17	3.42	1.83	1.64
σ	1.08	1.04	1.58	1.70	1.69
ñ	2.09	2.05	3.23	1.43	1.27
δ _{90%}	3.61	3.53	5.60	3.46	3.00

Quantitative comparison on the BU-3DFE [2] dataset. Lowers are better

[2] L. Yin, X. Wei, Y. Sun, J. Wang, and M. J. Rosato, "A 3D facial expression database for facial behavior research," in7th international conference on automatic face and gesture recognition (FGR06). IEEE,2006, pp. 211–216.

[3] S. Sanyal, T. Bolkart, H. Feng, and M. J. Black, "Learning to regress 3D face shape and expression from an image without 3D supervision" in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 7763–7772.

[4] 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), 2018, pp. 534–551.

[5]. Deng, J. Yang, S. Xu, D. Chen, Y. Jia, and X. Tong, "Accurate 3D face reconstruction with weakly-supervised learning: From single image to image set," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019

Conclusion

- Hybrid 3D face reconstruction approach (machine learning and geometric).
 Reconstruction from a single or multiple images
- Trained only on synthetic data but generalizes on real world data.
- Evaluation :
- Effectiveness of using the (Z) map.
 - Performance on *BU3D-FE* dataset.
 - New error calculation method (Point-to-Plane).

- Limitation :
 - No facial expressions and a limited age range (LYHM).
 - Synthetic data can have unrealistic features.

Thank you !