# A Generative Model for Depth-based Robust 3D Facial Pose Tracking (Supplementary Material)

Lu Sheng<sup>1</sup> Jianfei Cai<sup>2</sup> Tat-Jen Cham<sup>2</sup> Vladimir Pavlovic<sup>3</sup> King Ngi Ngan<sup>1</sup> <sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>Nanyang Technological University <sup>3</sup>Rutgers University {lsheng, knngan}@ee.cuhk.edu.hk, {asjfcai, astjcham}@ntu.edu.sg, vladimir@cs.rutgers.edu

# A. Probabilistic Face Parameterization

# A.1. Multilinear Representation

In this paper, we collect the  $N_{\rm id} \times N_{\rm exp} = 150 \times 47$ blendshapes to form the multilinear face representation, as suggested by Cao *et al.* [1]. Provided the shape dimension  $N_{\mathcal{M}} \times 3 = 11510 \times 3$ , the dataset forms a tensor  $\mathcal{D} \in \mathbb{R}^{3N_{\mathcal{M}} \times N_{\rm id} \times N_{\rm exp}}$ . By subtracting  $\mathcal{D}$  with the mean face mesh  $\bar{\mathbf{f}}$ , the applied data tensor is  $\mathcal{T} = \mathcal{D} - \bar{\mathbf{f}}$ .

The multilinear face representation requires the unary matrices  $\mathbf{U}_{id}$  and  $\mathbf{U}_{exp}$ , which can be derived from Highorder SVD (HOSVD) with respect to the second and third dimensions for identity and expression, respectively. The core tensor is thus  $\mathcal{C} = \mathcal{T} \times_2 \mathbf{U}_{id} \times_3 \mathbf{U}_{exp}$ , and inversely the data tensor can be derived as  $\mathcal{T} = \mathcal{C} \times_2 \mathbf{U}_{id}^\top \times_3 \mathbf{U}_{exp}^\top$ . Therefore, one face can be represented by linear combination of the data tensor with respect to the identity  $\mathbf{x}_{id} \in \mathbb{R}^{N_{id}}$  and expressions  $\mathbf{x}_{id} \in \mathbb{R}^{N_{exp}}$ , together with the mean face model

$$\begin{split} \mathbf{f} &= \bar{\mathbf{f}} + \mathcal{T} \times_2 \mathbf{x}_{id}^\top \times_3 \mathbf{x}_{exp}^\top \\ &= \bar{\mathbf{f}} + \mathcal{C} \times_2 (\mathbf{U}_{id}^\top \mathbf{x}_{id})^\top \times_3 (\mathbf{U}_{exp}^\top \mathbf{x}_{exp})^\top \\ &= \bar{\mathbf{f}} + \mathcal{C} \times_2 \mathbf{w}_{id}^\top \times_3 \mathbf{w}_{exp}^\top. \end{split}$$

 $\mathbf{w}_{id}$  and  $\mathbf{w}_{exp}$  are the parameters employed in this paper.

#### A.2. Learning the Identity and Expression Priors

Each face mesh in the training dataset is assigned a onehot label vector for identity  $\mathbf{x}_{id}$  and one for expression  $\mathbf{x}_{exp}$ . For example, the face mesh with the neutral expression from the first subject has  $\mathbf{x}_{id} = [1, 0, \dots, 0]^{\top} \in \mathbb{R}^{N_{id}}$  and  $\mathbf{x}_{exp}^{\top} = [1, 0, \dots, 0]^{\top} \in \mathbb{R}^{N_{exp}}$ .

Therefore, the mean face model  $\mathbf{\bar{f}}$  has  $\mathbf{\bar{x}}_{id} = \frac{1}{N_{id}}\mathbf{1}$  and  $\mathbf{\bar{x}}_{exp} = \frac{1}{N_{exp}}\mathbf{1}$ . The identity has  $\operatorname{Var}(\mathbf{x}_{id}) = \frac{1}{N_{id}}\mathbf{I} - \frac{1}{N_{id}^2}\mathbb{I} \simeq \frac{1}{N_{id}}\mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix and  $\mathbb{I}$  is the matrix filled by one. Similarly,  $\operatorname{Var}(\mathbf{x}_{exp}) \simeq \frac{1}{N_{exp}}\mathbf{I}$ . The mean and vari-

ance of  $\mathbf{w}_{id}$  and  $\mathbf{w}_{exp}$  are thus

$$oldsymbol{\mu}_{\mathrm{id}} = rac{1}{N_{\mathrm{id}}} \mathbf{U}_{\mathrm{id}}^{ op} \mathbf{1}, \, oldsymbol{\Sigma}_{\mathrm{id}} \simeq rac{1}{N_{\mathrm{id}}} \mathbf{U}_{\mathrm{id}}^{ op} \mathbf{U}_{\mathrm{id}} = rac{1}{N_{\mathrm{id}}} \mathbf{I}$$
 $oldsymbol{\mu}_{\mathrm{exp}} = rac{1}{N_{\mathrm{id}}} \mathbf{U}_{\mathrm{exp}}^{ op} \mathbf{1}, \, oldsymbol{\Sigma}_{\mathrm{exp}} \simeq rac{1}{N_{\mathrm{id}}} \mathbf{U}_{\mathrm{exp}}^{ op} \mathbf{U}_{\mathrm{exp}} = rac{1}{N_{\mathrm{exp}}} \mathbf{I}.$ 

We also observe that the mean face  $\overline{f}$  corresponds to  $\overline{x}_{id}$ and  $\overline{x}_{exp}$ , thus it results in

$$\begin{split} \|\mathcal{T} \times_2 \bar{\mathbf{x}}_{\mathrm{id}}^\top \times_3 \bar{\mathbf{x}}_{\mathrm{exp}}^\top \|_2 &= \|\mathcal{C} \times_2 \boldsymbol{\mu}_{\mathrm{id}}^\top \times_3 \boldsymbol{\mu}_{\mathrm{exp}}^\top \|_2 \\ &\simeq \|\bar{\mathbf{f}} - \bar{\mathbf{f}}\|_2 = 0 \\ &\Longrightarrow \boldsymbol{\mu}_{\mathcal{M}} \simeq \bar{\mathbf{f}}, \end{split}$$

which means these priors will not bias the face model representations.

We need to mention that  $\mu_{id}$  and  $\mu_{exp}$  should not be zero vectors. If we assume that  $\mu_{id} = 0$ , it will result in the independence of the face model with respect to the varying of the expression parameters, *i.e.*,  $C \times_2 \mathbf{w}_{id}^\top \times_3 \mathbf{w}_{exp}^\top \simeq \mathbf{0}$  in which the  $\mathbf{w}_{id}$  is usually near to the zero vector. The same problem happens when  $\mu_{exp} = \mathbf{0}$ . What's worse, such priors lead to instable solutions, *i.e.*,  $C \times_3 \mathbf{w}_{exp}^\top$  suffers such a small magnitude that the online adaptation system (derived from section 4.2) with respect to  $\mathbf{w}_{id}$  will be unstable.

#### **B.** Probabilistic Facial Pose Tracking

#### **B.1. Ray Visibility Score**

The ray visibility score is the Kullback-Leibler divergence between  $p_{\mathcal{P}}(\mathbf{y})$  and  $p_{\mathcal{Q} \to \mathcal{P}}(\mathbf{y}; \boldsymbol{\theta})$ , written as:

$$S(Q, \mathcal{P}; \boldsymbol{\theta}) = D_{KL}[p_{Q \to \mathcal{P}}(\mathbf{y}; \boldsymbol{\theta}) || p_{\mathcal{P}}(\mathbf{y})]$$
$$= \sum_{n=1}^{N_{\mathcal{M}}} D_{KL}[p_{Q \to \mathcal{P}}(y_n; \boldsymbol{\theta}) || p_{\mathcal{P}}(y_n)]$$

where  $D_{KL}[p_{Q\to \mathcal{P}}(y_n; \boldsymbol{\theta}) || p_{\mathcal{P}}(y_n)]$  is further derived as

$$\frac{\gamma_n}{2\sigma_o^2} \left\{ \sigma_o^2 + e^{2\alpha} \mathbf{n}_n^\top \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})} \mathbf{n}_n + \Delta(\mathbf{T}(\boldsymbol{\theta}) \circ \boldsymbol{\mu}_{\mathcal{M},[n]}; \mathbf{p}_n)^2 \right\} \\ + \frac{\gamma_n}{2} \left\{ \log 2\pi + \log \sigma_o^2 \right\} - (1 - \gamma_n) \log U_{\mathcal{O}} \\ - \frac{1}{2} \left\{ 1 + \log 2\pi + \log(\sigma_o^2 + e^{2\alpha} \mathbf{n}_n^\top \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})} \mathbf{n}_n) \right\}.$$

# **B.2.** Optimization of the Ray Visibility Score

We apply the quasi-Newton update  $\theta^{(t)} = \theta^{(t-1)} + \Delta \theta$ using the trust region approach for  $S(Q, P; \theta^{(t-1)})$  given the previous  $\gamma^{(t-1)}$ , in which the key ingredient is the gradient and Hessian calculations with respect to  $\theta$ . Straightforward calculation is not trivial, especially for rotation vector  $\omega$ , but  $S(Q, P; \theta)$  can be written in an approximated form to reduce the complexity.

Assume we have the mean and variance for the current warped face model point  $\bar{\mathbf{q}}_n = \mathbf{T}(\theta) \circ \boldsymbol{\mu}_{\mathcal{M},[n]}$  and  $\boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})}$ , we are interested in the incremental pose  $\Delta \theta$  in the ray visibility score. The incremental rotation matrix is approximated a skew-symmetric representation, *i.e.*,  $\mathbf{R}(\Delta \omega) = (\mathbf{I} + [\Delta \omega]_{\times})$ . Therefore, we have

$$\begin{split} &\Delta(\mathbf{T}(\boldsymbol{\theta} + \Delta \boldsymbol{\theta}) \circ \boldsymbol{\mu}_{\mathcal{M},[n]}; \mathbf{p}_n) \\ &= \mathbf{n}_n^\top \left\{ e^{\Delta \alpha} \mathbf{R}(\Delta \omega) \bar{\mathbf{q}}_n + \Delta \mathbf{t} - \mathbf{p}_n \right\} \\ &= \mathbf{n}_n^\top \left\{ (\mathbf{I} + [\Delta \omega]_\times) e^{\Delta \alpha} \bar{\mathbf{q}}_n + \Delta \mathbf{t} - \mathbf{p}_n \right\} \\ &= \mathbf{n}_n^\top \left\{ -[e^{\Delta \alpha} \bar{\mathbf{q}}_n]_\times \Delta \omega + \Delta \mathbf{t} + e^{\Delta \alpha} \bar{\mathbf{q}}_n - \mathbf{p}_n \right\} \end{split}$$

and the term about the warped variance is

$$e^{2\Delta\alpha} \mathbf{n}_{n}^{\top} \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega}+\Delta\boldsymbol{\omega})} \mathbf{n}_{n}$$

$$= e^{2\Delta\alpha} \mathbf{n}_{n}^{\top} \mathbf{R}(\Delta\boldsymbol{\omega}) \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})} \mathbf{R}(\Delta\boldsymbol{\omega})^{\top} \mathbf{n}_{n}$$

$$= e^{2\Delta\alpha} \mathbf{n}_{n}^{\top} (\mathbf{I} + [\Delta\boldsymbol{\omega}]_{\times}) \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})} (\mathbf{I} + [\Delta\boldsymbol{\omega}]_{\times})^{\top} \mathbf{n}_{n}$$

$$= (\mathbf{n}_{n} + [\mathbf{n}_{n}]_{\times} \Delta\boldsymbol{\omega})^{\top} e^{2\Delta\alpha} \boldsymbol{\Sigma}_{\mathcal{M},[n]}^{(\boldsymbol{\omega})} (\mathbf{n}_{n} + [\mathbf{n}_{n}]_{\times} \Delta\boldsymbol{\omega})$$

 $\Delta(\mathbf{T}(\boldsymbol{\theta} + \Delta \boldsymbol{\theta}) \circ \boldsymbol{\mu}_{\mathcal{M},[n]}; \mathbf{p}_n)^2$  and  $e^{2\Delta\alpha} \mathbf{n}_n^\top \Sigma_{\mathcal{M},[n]}^{(\boldsymbol{\omega} + \Delta \boldsymbol{\omega})} \mathbf{n}_n$  are approximated quadratic forms with respect to  $\Delta \boldsymbol{\omega}$ , thus the gradient and Hessian calculations about the rotation vector become much simpler. The related calculations with respect to the scale factor  $\Delta \alpha$  and translation vector  $\Delta \mathbf{t}$  involve the previous approximations and thus become simpler as well.

Moreover, to reduce the computational complexity even further, we can force the variance  $\Sigma_{\mathcal{M},[n]} = \sigma_{[n]}^2 \mathbf{I}$ , with  $\sigma_{[n]}$  derived from the canonical face point representation. Thus  $\Sigma_{\mathcal{M},[n]}^{(\omega)} = \Sigma_{\mathcal{M},[n]} = \sigma_{[n]}^2 \mathbf{I}$  and  $\mathbf{n}_n^\top \Sigma_{\mathcal{M},[n]}^{(\omega)} \mathbf{n}_n = \sigma_{[n]}^2$ are constant with varying  $\omega$ . In this case, only  $\Delta(\mathbf{T}(\boldsymbol{\theta}) \circ \boldsymbol{\mu}_{\mathcal{M},[n]}; \mathbf{p}_n)^2$  makes contribution to the update of  $\omega$ .

### **B.3.** Online Identity Adaptation

The online identity adaptation can be progressively personalized to the test subject, as visualized in Figure 2. With



Figure 1. Examples of facial pose results with the visibility detection. The third column shows the visibility masks. The last column shows the personalized face models warped on the point clouds. Best viewed in color.

different poses are used to update the identity distribution, the face model is continuously adapted to the test subject.

# C. Facial Poses with Visibility Detection

In addition to the facial pose results visualized in the paper, we also illustrate some examples with the personalized face models, as shown in Figure 1. The proposed method can effectively check the visibility of a face model with respect to the input point cloud, and its pose estimation is robust to severe occlusions, *e.g.*, profiled faces, accessories and hands, *etc*.

#### References

 C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou. Facewarehouse: A 3D facial expression database for visual computing.



Figure 2. We continuously adapt the identities of the face model to different users. (a)-(c) are three examples showing that the face model can be gradually personalized when the facial depth data from different poses are captured during the tracking process. The face model is initialized with the generic face model.

IEEE Trans. Vis. Comput. Graphics, 20(3):413-425, 2014. 1