

Supplementary Material for: Learning to Deblur and Rotate Motion-Blurred Faces

Givi Meishvili¹
 givi.meishvili@inf.unibe.ch
 Attila Szabó³
 attila.szabo@huawei.com
 Simon Jenni²
 jenni@adobe.com
 Paolo Favaro¹
 paolo.favaro@inf.unibe.ch

¹ University of Bern, Switzerland
² Adobe Research
³ Huawei, Noah's Ark Lab
 (work was done before joining)

1 Additional Implementation Details

This section clarifies some details regarding the losses in the training objective of viewpoint encoder E_{3D} (section 3.4).

1.1 Training Objectives

Reconstruction Losses for E_{3D} . As mentioned in section 3.4 of the paper, the training objective for the viewpoint encoder E_{3D} is given by

$$\min_{E_{3D}} \sum_{i=1}^n \sum_{v=1}^{V_i} \ell_{im}(\phi(E_{3D}(x_i^v)), y_i^v) + \ell_{3D}(E_{3D}(x_i^v), y_i^v) + \lambda_c(|\alpha_i|^2 + |\beta_i|^2 + |\delta_i|^2), \quad (1)$$

where ℓ_{im} and ℓ_{3D} represents the following combination of different reconstruction losses:

$$\begin{aligned} \ell_{im}(x, y) &= \lambda_{id} \mathcal{L}_{id}(x, y) + \lambda_{edge} \mathcal{L}_{edge}(x, y) + \lambda_{data} |x - y| \\ \ell_{3D}(x, y) &= \lambda_{lan} \mathcal{L}_{lan}(x, y) + \lambda_{mview} \mathcal{L}_{mview}(x) \\ \mathcal{L}_{lan}(x, y) &= |\mathcal{Q}_{basel}(x) - \mathcal{Q}_{image}(y)|_2^2 \\ \mathcal{L}_{mview}(x) &= \sum_v |\bar{\alpha} - \alpha^v|^2 + |\bar{\beta} - \beta^v|^2 + |\bar{\delta} - \delta^v|^2. \end{aligned}$$

where, $\lambda_{data} = 5$, $\lambda_{id} = 0.5$, $\lambda_{edge} = 30$, $\lambda_{mview} = 0.25$ and $\lambda_{lan} = 1$. We use the perspective camera model in the renderer ϕ , with an empirically selected focal length for the 3D-2D projection. The term \mathcal{L}_{lan} is a MSE between 2D projections of facial landmarks of the predicted mesh and pre-computed landmarks in sharp images. \mathcal{Q}_{basel} projects the 3D landmark vertices of the reconstructed mesh onto the image (obtaining 68 facial landmarks), and \mathcal{Q}_{image}



Figure 1: **Qualitative sample on real-world motion blurred face.** The first column corresponds to the blurry input image. All the other columns are output sequences rotated by a different amount. Rows from 1 to 5 correspond to the appropriate frame in the output sequence. The last column is the copy of the previous one with rectangles on top of different facial regions. Rectangles are at a fixed location with respect to the image in all frames. Note how the both eyes and the nose move upwards as we go from the top to the bottom.

extracts landmarks using the method of [10] from the ground-truth targets. \mathcal{L}_{mview} ensures that the identity, texture, and expression parameters of the BFM are consistent across views for samples of our multi-view dataset.

2 Additional Qualitative Results

A real-world deblurring example is presented in Figure 1. Some more qualitative examples of our multi-view reconstructions on VIDTIMIT[10] can be found in Figures 2 and 3.



Figure 2: **Qualitative samples on VIDTIMIT.** The first column corresponds to the blurry input image. All the other columns are output sequences rotated by a different amount. Rows from 1 to 5 correspond to the appropriate frame in the output sequence.

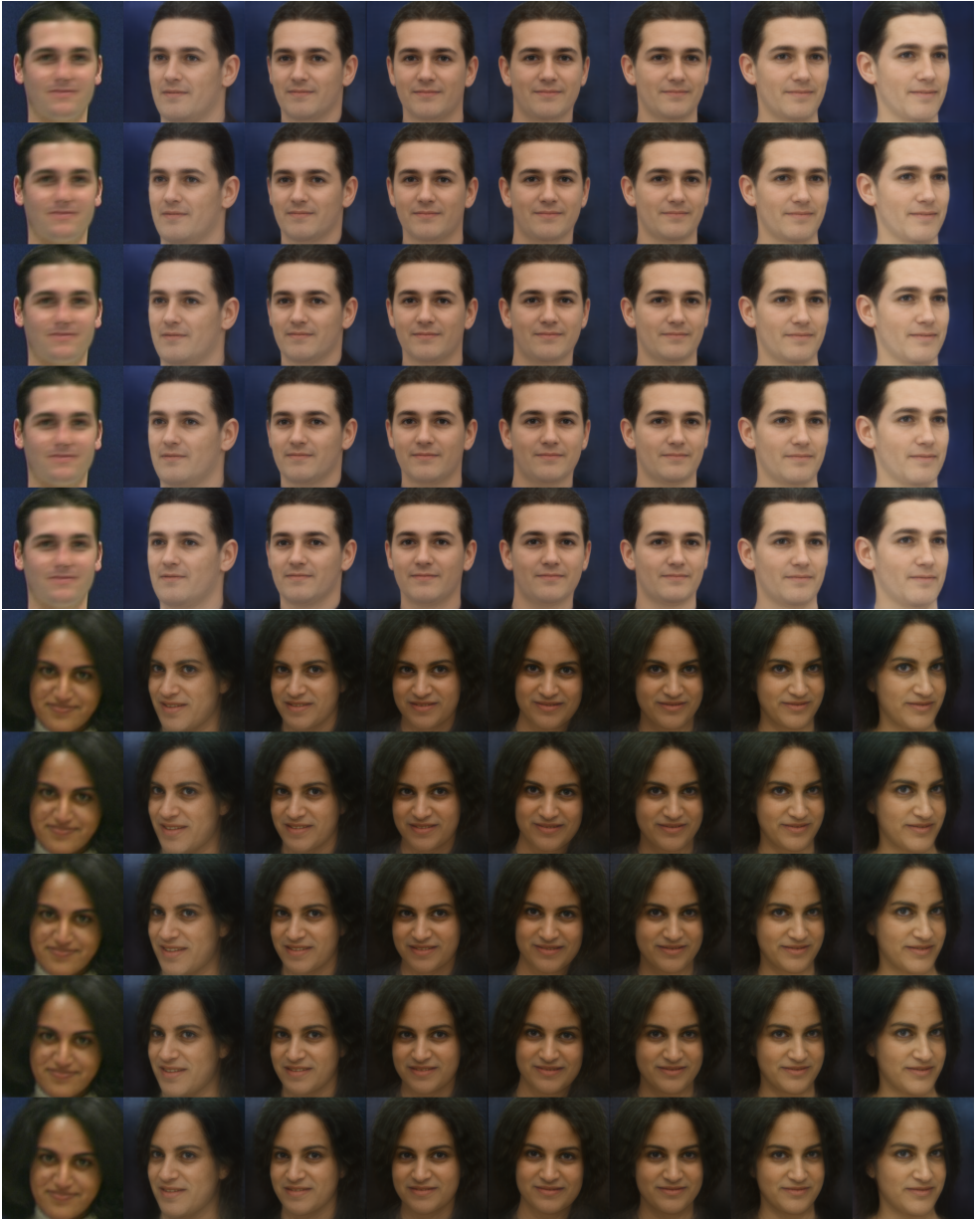


Figure 3: **Qualitative samples on VIDTIMIT.** The first column corresponds to the blurry input image. All the other columns are output sequences rotated by a different amount. Rows from 1 to 5 correspond to the appropriate frame in the output sequence.

References

- [1] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017.
- [2] Conrad Sanderson and Brian C. Lovell. Multi-region probabilistic histograms for robust and scalable identity inference. In Massimo Tistarelli and Mark S. Nixon, editors, *Advances in Biometrics*, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.