

To Fit or Not to Fit: Model-based Face Reconstruction and Occlusion Segmentation from Weak Supervision

CHUNLU LI, Dept. of Automation, Southeast University, China Dept. of Mathematics and Informatics, University of Basel, Switzerland

ANDREAS MOREL-FORSTER, Dept. of Mathematics and Informatics, University of Basel, Switzerland

THOMAS VETTER, Dept. of Mathematics and Informatics, University of Basel, Switzerland

BERNHARD EGGER*, Chair of Visual Computing, Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany Dept. of Brain and Cognitive Sciences, Massachusetts Institute of Technology, USA

ADAM KORTYLEWSKI*, Max Planck Institute for Informatics, Germany Dept. of Computer Science, Johns Hopkins University, USA

3D face reconstruction under occlusions is highly challenging due to the large variability of occluders. Currently, the most successful methods fit a 3D face model through inverse rendering and assume a given segmentation of the occluder to avoid fitting the occluder. However, training an occlusion segmentation model requires large amounts of annotated data. In this work, we introduce a model-based approach for 3D face reconstruction that is highly robust to occlusions but does not require any occlusion annotations for training. In our approach, we exploit the fact that generative face models can only synthesize human faces, but not the occluders. We use this property to guide the decision-making process of an occlusion segmentation network and resulting in unsupervised training. The main challenge is that the model fitting and the occlusion segmentation are mutually dependent on each other, and need to be inferred jointly. We resolve this chicken-and-egg problem with an EM-type training strategy. This leads to a synergistic effect, in which the segmentation network prevents the face encoder from fitting to the occlusion, enhancing the reconstruction quality. The improved 3D face reconstruction, in turn, enables the segmentation network to better predict the occlusion. Qualitative and quantitative experiments on the CelebA-HQ, the AR databases, and the NoW challenge demonstrate that the proposed pipeline achieves the state-of-the-art 3D face reconstruction under occlusion. Moreover, the segmentation network localizes occlusions accurately despite being trained without any occlusion annotation. The code is available at <https://FakeLinkforDoubleBlind.com>.

CCS Concepts: • **Computing methodologies** → **Reconstruction**.

Additional Key Words and Phrases: Wireless sensor networks, media access control, multi-channel, radio interference, time synchronization

1 INTRODUCTION

Monocular 3D face reconstruction aims at estimating the pose, shape, and albedo of a face, as well as the illumination conditions and camera parameters of the scene. Solving for all these factors from a single image is an ill-posed problem. Model-based face autoencoders

* Denotes same contribution of Bernhard Egger and Adam Kortylewski.
 Authors' addresses: Chunlu Li, Dept. of Automation, Southeast University, Sipailou Street 2, Nanjing, Jiangsu, 210096, China and Dept. of Mathematics and Informatics, University of Basel, Switzerland, chunlu.li@unibas.ch; Andreas Morel-Forster, Dept. of Mathematics and Informatics, University of Basel, Spiegelgasse 1, Basel, CH4056, Switzerland, andreas.forster@unibas.ch; Thomas Vetter, Dept. of Mathematics and Informatics, University of Basel, Spiegelgasse 1, Basel, CH4056, Switzerland, thomas.vetter@unibas.ch; Bernhard Egger*, Chair of Visual Computing, Friedrich-Alexander-Universität Erlangen-Nürnberg, Cauerstraße 11, Erlangen, Bavarian, 91058, Germany and Dept. of Brain and Cognitive Sciences, Massachusetts Institute of Technology, USA, bernhard.egger@fau.de; Adam Kortylewski*, Max Planck Institute for Informatics, Saarbrücken, Saarland, 66123, Germany and Dept. of Computer Science, Johns Hopkins University, USA, akortyle@mpi-inf.mpg.de.



Fig. 1. Our proposed method conducts face reconstruction and occlusion segmentation jointly, and is operated under weak supervision. From top to bottom: target images, our reconstruction images, and the estimated occlusion masks.

[Tewari et al. 2017] overcome this problem by performing 3D reconstruction through fitting a 3D Morphable Model (3DMM) [Banz and Vetter 2003; Egger et al. 2020] to a target image. The 3DMM provides prior knowledge about the face albedo and geometry such that 3D face reconstruction from a single image becomes feasible, enabling face autoencoders to set the current state-of-the-art in 3D face reconstruction [Deng et al. 2019b]. The network architectures in the face autoencoders are devised to enable end-to-end reconstruction and to enhance reconstruction speed compared to optimization-based alternatives [Kortylewski et al. 2018b; Zhu et al. 2015], and sophisticated losses are designed to stabilize the training and to get better performance [Deng et al. 2019b].

A major remaining challenge for face autoencoders is that their performance in in-the-wild environments is still limited by nuisance factors such as occlusion, extreme illumination, and poses. Among those nuisances, occlusions are ubiquitous and inherently difficult to handle because of their wide variety in shape, appearance, and locations. A core problem caused by occlusions is that the face model adapts to occluded face regions and as a result, the reconstructed face will be distorted (as seen in our experiments). Therefore an important open question for occlusion-robust 3D face reconstruction is to decide which pixels to fit and which not to fit in a target image.

Existing solutions to face reconstruction under occlusions often follow a bottom-up approach. For example, a multi-view shape consistency loss is used as prior to regularize the shape variation of the same face in different images [Deng et al. 2019b; Feng et al. 2020;

Tiwari et al. 2022], or the face symmetry is used to detect occluders [Tran et al. 2018]. Most existing methods apply segmentation methods to locate the face region [Saito et al. 2016], or to detect skin [Deng et al. 2019b] and subsequently exclude the occluded image regions during reconstruction. These segmentation methods operate in a supervised manner, which is infeasible in practice due to the high cost and efforts for acquiring a great variety of occlusion annotations from in-the-wild images.

In this work, we introduce an approach for model-based face reconstruction that is highly occlusion-robust, without requiring any human occlusion annotation. In particular, we propose to train a face autoencoder and a segmentation network in a cooperative manner. The segmentation network answers the question of whether the face model should 'fit or not to fit' certain pixels so that the face reconstruction is not affected by the occlusion. To train the segmentation network in an unsupervised manner, we exploit the fact that generative face models can only synthesize human faces, but not the occluders. We use this property to guide the decision-making process of an occlusion segmentation network and resulting in unsupervised training. We find that the discrepancy between the target image and the rendered face image (Fig.1 1st and 2nd rows) can serve as a supervision signal to guide the training of the segmentation network. The face reconstruction network, in turn, becomes robust to occlusions by using the prediction from the segmentation network to mask out the occluded pixels during fitting. This leads to a synergistic effect, in which the occlusion segmentation first guides the face autoencoder to fit image regions that are easy to classify as face regions. The improved face fitting, in turn, enables the segmentation model to refine its prediction.

The training process follows the core idea of the Expectation-Maximization (EM) algorithm, by alternating between training the face autoencoder given the current estimate of the segmentation mask, and subsequently training the segmentation network based on the current 3D face reconstruction. The EM-like training strategy resolves the problem that the estimated occlusion segmentation depends on the estimated face model parameters and vice-versa. Importantly, the unsupervised training of the segmentation network is reached by regularizing and preserving the similarities among the target image and the reconstructed image under the estimated occlusion mask, and we introduce several losses to achieve this.

We demonstrate the effectiveness of our method by conducting experiments on the CelebA-HQ dataset [Liu et al. 2015], the AR database [Martinez and Benavente 1998] and the NoW challenge [Sanyal et al. 2019b], where we achieve state-of-the-art performance in 3D face reconstruction. Remarkably, our method is able to predict accurate occlusion masks without requiring any supervision during training.

In summary, we make the following contributions in this paper:

- (1) We introduce an approach for model-based 3D face reconstruction that is highly robust occlusion, without requiring any human occlusion annotation.
- (2) Our model achieves state-of-the-art performance at 3D face reconstruction under occlusions and provides accurate estimates of the facial occlusion masks on in-the-wild images.

2 RELATED WORK

Model-based face autoencoders [Tewari et al. 2017] solve the 3D face reconstruction task by fitting a face model to the target image with an encoder and a renderer, as well as the 3DMM, as the decoder. Typically, the encoder first estimates parameters from a target image, including the shape, texture, and pose of the target, and the illumination and camera settings from the scene. Then the renderer synthesizes a 2D image using the estimated parameters with an illumination model and a projection function. The face is reconstructed by retrieving the parameters which result in a synthesized image most similar to the target image. The 3DMM [Banz and Vetter 2003] plays a paramount role in the face autoencoders, because it parameterizes the latent distribution space of faces, and therefore can connect the encoder with the renderer and enable end-to-end training. The model-based face autoencoders have been proven effective in improving the reconstruction. They simplify the optimization step and enhance the reconstruction speed [Tewari et al. 2018], improve the details of shape and texture [Feng et al. 2020; Gecer et al. 2019; Richardson et al. 2017; Tran et al. 2018, 2019], and can also reconstruct more discriminative features [Deng et al. 2019b; Genova et al. 2018].

Despite the advantages of the face autoencoders, their performance under occlusions is still limited. To solve this issue, some early methods [Romdhani and Vetter 2003] resort to robust fitting losses, but they are not robust to illumination variations and appearance variations in eye and mouth regions. In recent years, shape consistency losses have been used as prior to constrain the face shape across images of the same subject [Deng et al. 2019b; Feng et al. 2020; Sanyal et al. 2019b; Tiwari et al. 2022]. The variation of identity features of the 3D shape is restricted so that the shape reconstruction remains robust even in unconstrained environments. However, such methods usually need identity labels and do not promise robust texture reconstruction. Besides, many methods conduct face segmentation before reconstruction to lead the model to better fit the unoccluded face region. For example, a random forest detector for hair is proposed [Morel-Forster 2016] so that the face model does not fit the hair region, and a semantic segmentation network is trained to better locate the face region [Saito et al. 2016]. A skin detector is employed to impose different weights on the pixels during reconstruction to guide the network to put more attention on the skin-colored regions and prevent it from fitting the occlusions [Deng et al. 2019b]. However, the skin-colored occlusion, such as hair, hands, and so on, can not be distinguished correctly and the skin detector is sensitive to illumination. Yildirim et al. propose to explicitly model the 3D shape of certain types of occlusions and the shadow cast by them, in order to decompose the target into face regions and occlusion, and therefore the occlusions can be excluded during training [Yildirim et al. 2017]. However, the types of occlusions are limited. Generally, these off-the-shelf segmentation models require labeled data for training. Although using synthesized images can be used for training, there is a domain gap between the real images and the synthesized ones [Kortylewski et al. 2018a]. Unlike these methods, we merge the segmentation procedure into a model-based face autoencoder, which exploits the face model prior, and consequently does not require additional supervision.

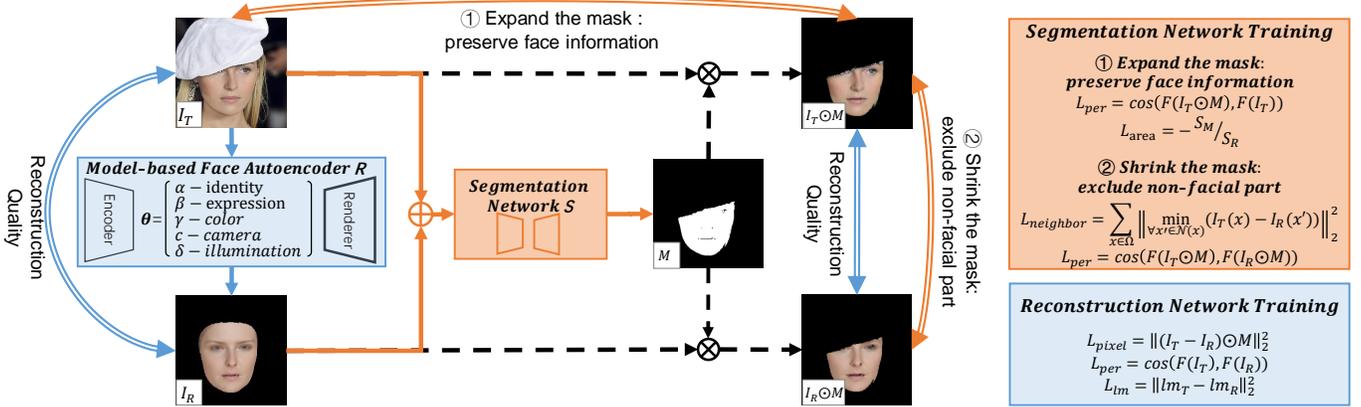


Fig. 2. Proposed pipeline. The solid single lines show the forward path: given a target image I_T , the reconstruction network, R , estimates the latent parameters and subsequently renders an image I_R , containing only the face. Then, I_T and I_R are stacked and fed into the segmentation network, S , which predicts the segmentation mask M . The dashed lines show that M is used to mask out the estimated occlusions in I_T and I_R to get assembly occlusion-free images, namely $I_T \odot M$ and $I_R \odot M$. The double-lined arrows indicate the compared image pairs in the losses for S (orange) and R (blue). For optimizing S , two groups of losses (shown in the orange rectangle on the right) are proposed: 1) losses expanding the mask so that the parts indicating face are preserved and 2) losses aiming at shrinking the mask so that the non-facial regions are excluded. S aims to minimize and reach the balance between the two groups of losses. The losses for R (seen in the blue box) compare $I_T \odot M$ and $I_R \odot M$ in the pixel level so that the reconstruction is not affected by the occlusion, and, since the perceptual features are sensitive to noises in the mask, I_T and I_R in the perceptual level. By training alternatively the two networks and exploiting the synergy between the segmentation and the reconstruction tasks, the proposed pipeline is capable of both reconstructing faces under occlusions robustly and conducting face segmentation.

The most relevant method to ours is proposed by Egger et al. [Egger et al. 2018]. They jointly adapt a face model to a target image and segment the target image into face, beard, and occlusion, and the segmentation models are trained with an EM-like algorithm, where different models for beard, foreground, and background are optimized in alternating steps. Compared to our method, their method does independent per-image optimization so the whole fitting-segmentation process is repeated given a new image, while ours solves both reconstruction and segmentation together as learning problems on a larger set of training data and only one forward run is enough for an input image after the model is trained. Therefore our method is much faster, more robust, and more effective. Besides, their method requires specific statistical occlusion models for beard, foreground, and background, while in our proposed method, the segmentation network is guided by several quality-controlling losses comparing a target image with the reconstructed face, which are intuitive and much easier to implement. De Smet et al. also propose to conduct face model fitting and occlusion segmentation jointly [De Smet et al. 2006], but they estimate the occlusions based on an appearance distribution model, which is sensitive to illumination variation and many other subtle changes in appearance. Maninchedda et al. propose to solve face reconstruction and segmentation in a joint manner [Maninchedda et al. 2016], but depth maps are required to provide supervision. In comparison, our pipeline learns from only weak supervision and does not need specific models for different types of occlusion. The face autoencoder also enables us to adapt the face model more efficiently. In addition, we integrate perceptual losses which enable the segmentation network to reason over semantic features instead of only over independent pixels, which increases the robustness to illumination and other factors.

3 APPROACH

We introduce a neural network-based pipeline that conducts 3D face reconstruction and occlusion segmentation jointly. In the following, we first discuss our proposed pipeline architecture (3.1) and then discuss how the model can be trained in an EM-type manner without any supervision regarding the occlusions (3.2). Finally, we discuss how the EM training is initialized in an unsupervised manner (3.3).

3.1 Network Architecture

Our goal is to robustly reconstruct the 3D face from a single target image I_T , even under severe occlusion. To solve this challenging problem, we integrate a model-based face autoencoder, R , with a segmentation network, S , and create synergy between them, as demonstrated in Fig. 2. For face reconstruction, the segmentation mask cuts the estimated occlusions out during model fitting, making the reconstruction network robust to occlusion. For segmentation, the reconstructed result provides reference, enhancing the segmentation accuracy. In this section, we explain how the two networks are connected together and how they benefit each other.

The model-based face autoencoder, R , is expected to reconstruct the complete face appearance from the visible face regions in the target image, I_T . It consists of an encoder and a computer graphics renderer as its decoder. The encoder estimates the latent parameters $\theta = [\alpha, \gamma, \phi, c] \in \mathbb{R}^{257}$, i.e. the 3D shape $\alpha \in \mathbb{R}^{144}$ and texture $\gamma \in \mathbb{R}^{80}$ of a 3DMM, as well as the illumination $\phi \in \mathbb{R}^{27}$ and camera parameters $c \in \mathbb{R}^6$ of the scene. Given the latent parameters, the decoder renders a face image $I_R = \mathbb{R}(\theta)$ of the target face.

Standard face autoencoders [Tewari et al. 2017] fit the face model parameters, regardless of whether the underlying pixels depict a



Fig. 3. In the presence of occlusion, the proposed method can reconstruct faces more faithfully than previous model-based face autoencoders. The images from top to bottom are: target images, results of the MoFA network [Tewari et al. 2017], and the results of ours.

face or occlusion. Consequently, the face model is distorted by the occluded face regions, as shown in the second row in Figure 3, it is obvious that the illumination, appearance, and shape are estimated incorrectly. To resolve this fundamental problem of face autoencoders, we introduce an unsupervised segmentation network, whose output can be used to mask the occlusions out during model fitting and therefore make the autoencoder robust to occlusion.

The segmentation network, S , takes the stacked target image I_T and the synthesized image I_R as input and predicts a binary mask, $M = S(I_T, I_R)$, to describe whether a pixel depicts the face (1) or not (0). Since I_R contains the estimated intact face, it provides the segmentation network with prior knowledge and helps the estimation.

The face autoencoder and the segmentation network are coupled together during training to induce a synergistic effect which makes the segmentation more accurate and reconstruction more robust under occlusion, as shown in the last row in Figure 3. In 3.2, we describe how the pipeline can be trained end-to-end, despite the entanglement between the two networks, and how the high-level losses work that relieve our pipeline of any annotation for occlusion.

3.2 EM-type Training

Due to the mutual dependencies between the face autoencoder and the segmentation network, we conduct an Expectation-Maximization (EM) like strategy, where we train the two networks in an alternating manner. This enables a stable convergence of the model training process. Similar to other EM-type training strategies, our training process starts from a rough initialization of the model parameters which is obtained in an unsupervised manner (as described at the end of this section). We then optimize the two networks in an alternating manner, as described in the following.

Training the segmentation network. When training the segmentation network, the parameters of the face autoencoder are fixed and only the segmentation network is optimized. Instead of hunting for labeled data, we propose four losses enforcing intrinsic similarities among the images. Each loss works to either include pixels indicating face or the opposite. The losses work either on the perceptual level or the pixel level, to fully exploit the visual clues. The perceptual-level losses compare the intermediate features of two images extracted by a pretrained face recognition model F

(we use Arcface [Deng et al. 2019a]). We use the cosine distance, $\cos(X, Y) = 1 - \frac{X \cdot Y}{\|X\| \|Y\|}$, to compute the distance between the features. Perceptual losses are common for training face autoencoders, which encourage encoding facial details that are important for face recognition [Feng et al. 2020]. We found that computing the perceptual losses is also very helpful to segmentation (see 4.4).

Since the proposed losses have overlapped or opposite functions to each other, only by reaching a balance among these losses can the network yield in good segmentation result. The proposed losses are as follows:

$$L_{neighbor} = \sum_{x \in \Omega} \left\| \min_{x' \in \mathcal{N}(x)} |I_T(x) - I_R(x')| \right\|_2^2 \quad (1)$$

$$L_{dist} = \cos(F(I_T \odot M), F(I_R \odot M)) \quad (2)$$

$$L_{area} = -S_M/S_R \quad (3)$$

$$L_{presv} = \cos(F(I_T \odot M), F(I_T)) \quad (4)$$

The pixel-level neighbour loss 1, $L_{neighbor}$, compares a pixel, $I_T(x)$, at location x on the target image, with the pixels on the rendered image in the neighbouring region, $\mathcal{N}(x)$ of this pixel, so that this loss is stable even if there are small misalignments. Note that it only accounts within the face region, Ω , predicted by the segmentation network, to encourage the mask to avoid those pixels with higher pixel-level reconstruction error.

Similarly, in Equation 2 we introduce a perceptual-level loss L_{dist} to encourage the mask to discard the parts with higher perceptual differences. These two losses, 1 and 2, aim at shrinking the mask where the pixel-level and perceptual differences are large. Without any other constraints, the segmentation network would output an all-zero mask to make them both 0. On the contrary, once there is a force to encourage the network to preserve some image parts, the segmentation network tends to preserve the parts with smaller losses, which in fact are the ones well-explained by the face model and therefore is much more likely to depict face.

Therefore, Equation 3 and 4 are proposed to counterwork 1 and 2. Equation 3 is an area loss, L_{area} that enlarges the ratio between the number of estimated facial pixels, S_M , and the number of pixels in the rendered face region, S_R . It prevents the segmentation network from discarding too many pixels. L_{presv} (Eq. 4), ensures that the perceptual face features remain similar when the occluders in the target image are masked out and encourages the model to preserve as much of the visible face region as possible. Likewise, the network would keep the most-likely face region to decrease Equation 3 and 4 in the presence of 1 and 2.

We use an additional regularization term, $L_{bin} = -\sum_x (M(x) - 0.5)^2$, to encourage the face mask M to be binary. The total loss for training the segmentation network is: $L_S = \eta_1 L_{neighbor} + \eta_2 L_{dist} + \eta_3 L_{area} + \eta_4 L_{presv} + \eta_5 L_{bin}$, with $\eta_1 = 15$, $\eta_2 = 3$, $\eta_3 = 0.5$, and $\eta_4 = 2.5$, and $\eta_5 = 10$. Analysis of the influence of the hyper-parameters is provided in the supplementary material.

During training, the segmentation network is guided seeking a balance between discarding pixels that cannot be explained well by the face autoencoder, while preserving pixels that are important to retain the perceptual representations of the target and rendered face images. Therefore no supervision for occlusions is required.

Training the face autoencoder. In the second step, we continue to optimize the parameters of the face autoencoder, while keeping the segmentation network fixed. The losses for training the face autoencoder include:

$$L_{pixel} = \left\| (I_T - I_R) \odot M \right\|_2^2 \quad (5)$$

$$L_{per} = \cos(F(I_T), F(I_R)) \quad (6)$$

$$L_{lm} = \left\| lm_T - lm_R \right\|_2^2 \quad (7)$$

Above are two reconstruction losses: L_{pixel} (5) at the image level and L_{per} (6) at the perceptual level, and a landmark loss (7) used to estimate the pose, where lm_T and lm_R stand for the 2D landmark coordinates on I_T and I_R , respectively [Deng et al. 2019b]. We set the weights for the landmarks on the nose ridge and inner lip as 20, and the rest as 1. A regularization term is also required for the 3DMM: $L_{reg} = \left\| \theta \right\|_2^2$. To sum up, the loss for training the face autoencoder can be represented as: $L_R = \lambda_1 L_{pixel} + \lambda_2 L_{per} + \lambda_3 L_{lm} + \lambda_4 L_{reg}$, where $\lambda_1 = 0.5$, $\lambda_2 = 0.25$, $\lambda_3 = 5e - 4$, and $\lambda_4 = 0.1$.

3.3 Unsupervised Initialization

As every other EM-type training strategy, our training needs to be roughly initialized. To achieve this, we generate preliminary masks using an occlusion robust loss [Egger et al. 2018] so that the initialization is unsupervised:

$$\log(P_{face}(x)) = -\frac{1}{2\sigma^2} (I_T(x) - I_R(x))^2 + N_c \quad (8)$$

$$M_{pre}(x) = \begin{cases} 1 & \text{if } (I_T(x) - I_R(x))^2 < \xi \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

We assume that the reconstruction error at pixel x in the face regions follows a zero-mean Gaussian distribution. Therefore, we can express the log-likelihood that a pixel belongs to the face regions as $\log(P_{face})$ (Eq. 8), where σ and N_c are constant. We also assume that the values of the non-face pixels follow a uniform distribution, i.e., $\log(P_{non-face})$ is a constant. Finally, a pixel at position x is classified as face or non-face by comparing the log-likelihoods. This reduces to thresholding of the reconstruction error with a constant parameter ξ (Equation 9). When ξ increases, the initialized masks allow the pixels on the target image to have a larger difference to the reconstructed pixels and encourage the reconstruction network to fit to these pixels. Empirically, we found that $\xi = 0.17$ leads to a good enough initialization.

To initialize the face autoencoder, the preliminary mask, $M_{pre}(x)$, is obtained in the forward pass using Eq. 8 and Equation 9, after the reconstructed face is rendered. Then $M_{pre}(x)$ is directly used to mask out the roughly-estimated occluded regions as in Eq. 5, preventing the face autoencoder from fitting to any possible occlusions. Subsequently, the segmentation network is pre-trained by using these preliminary masks as ground-truth labels.

4 EXPERIMENTS

In this section, results of systematic experiments show that our weakly-supervised method reaches the state-of-the-art face shape reconstruction accuracy and competitive occlusion segmentation

results compared to the state-of-the-art and methods that use full supervision in terms of occlusion labels. Our ablation study shows the effectiveness of the segmentation network and our proposed losses.

4.1 Experiment setting

Our face encoder shares the structure of the ResNet 50 [He et al. 2016] and uses the Basel Face Model (BFM) 2017 [Gerig et al. 2018] as the 3D face model, with the differentiable renderer proposed in [Koizumi and Smith 2020]. The segmentation network follows the UNet architecture [Ronneberger et al. 2015]. The proposed pipeline is trained on the CelebA-HQ trainset [Liu et al. 2015], following their protocol. Facial landmarks are detected using the method of [Bulat and Tzimiropoulos 2017], and images are pre-processed in the same way as [Deng et al. 2019b]. The perceptual features are extracted by the pre-trained ArcFace [Deng et al. 2019a]. The Adadelta optimizer is used, with an initial learning rate of 0.1, and a decay rate of 0.99 at every 5k iterations. The learning rate for the segmentation network is 0.06 times the one for the reconstruction network. In every 30k iterations, 25k iters are for the face autoencoder training, and the rest are for training the segmentation network. For initialization, the face autoencoder is trained for 300k iterations. Afterwards, the face autoencoder and segmentation network are trained jointly for 200k iterations. The speed is evaluated on an RTX 2080 Ti, with batch size 12. It takes about 120 hours for the initialization of the face autoencoder, and about 80 hours to train the complete pipeline. After the training, it takes 49 ms for reconstruction and 70 μ s for segmentation on average for one image. The reconstruction and the

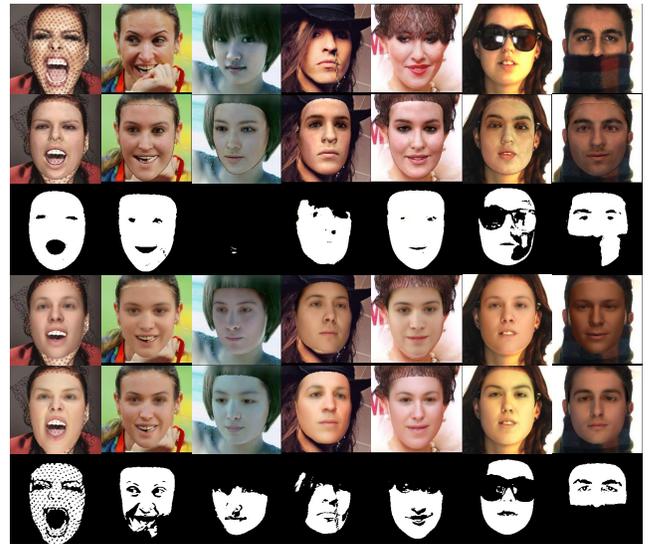


Fig. 4. Qualitative comparison on the reconstruction and segmentation results of the Deep3D [Deng et al. 2019b] network (the 2nd and 3rd rows), the MoFA network [Tewari et al. 2017] (the 4th row), and the proposed method (the last two rows) on occluded faces from the CelebA-HQ testset (the first 8 columns) and the AR database (the last 2 columns). Note that all the masks are binarized.

Table 1. RMSE on the CelebA-HQ testsets and the AR testset.

Testset	MoFA [Tewari et al. 2017]	Supervised MoFA	Supervised MoFA-cutmix	Supervised MoFA-cutout	Deep3D [Deng et al. 2019b]	Proposed
CelebA-Unoccluded	8.77 ± 0.40	8.71 ± 0.38	8.75 ± 0.39	8.72 ± 0.40	8.49 ± 0.39	8.38 ± 0.42
CelebA-Occluded	9.20 ± 0.45	9.01 ± 0.45	9.04 ± 0.44	8.99 ± 0.45	8.79 ± 0.45	8.71 ± 0.48
CelebA-Overall	8.99 ± 0.47	8.86 ± 0.44	8.90 ± 0.44	8.85 ± 0.45	8.64 ± 0.44	8.55 ± 0.48
AR-Overall	9.53 ± 0.33	9.34 ± 0.33	9.33 ± 0.32	9.28 ± 0.31	9.11 ± 0.37	8.93 ± 0.35

segmentation networks have 25.6M and 34.5M parameters, respectively. There is no fine-tuning with 3D data on any of the testsets in our experiments.

Baselines. We compare our method with two state-of-the-art model-based face autoencoders, i.e. the MoFA [Tewari et al. 2017] and the Deep3D [Deng et al. 2019b]. Additionally, to achieve fair comparison between our proposed weakly-supervised method and supervised methods, we train our reconstruction network in supervised settings, in which the supervised pipelines use the ground truth (GT) masks provided by the CelebA-HQ database to exclude occlusions during training. The GT masks of the CelebA-HQ database stand for the merge of their manually labeled masks for skin, hair, accessories, and so on. Two data augmentation methods for occlusion handling, i.e. the cutmix [Yun et al. 2019] and cutout [DeVries and Taylor 2017], are also implemented to enhance the performance of the supervised pipelines. We refer to the three baselines as Supervised MoFA, MoFA cutmix, and MoFA cutout, respectively.

Databases The CelebA-HQ testset [Liu et al. 2015] and the AR database [Martinez and Benavente 1998] are used for evaluating the effectiveness of fitting and face segmentation. For the AR database, 120 manually-segmented results in [Egger et al. 2018] are used as GT masks. We also evaluate the shape reconstruction accuracy on the subsets of the NoW database [Sanyal et al. 2019a]. The standard deviation is provided after ‘±’. The CelebA-HQ database and the AR database are publicly available and the occlusion labels for the AR dataset are publicly shared by the authors of [Egger et al. 2018].

4.2 Reconstruction Quality

Fig. 4 shows some results of the Deep3D network, the MoFA network, and our proposed method for qualitative comparison. Note that all the masks are binarized by rounding the pixels. The segmentation masks provided by the Deep3D result from a skin detector which assumes that skin color follows the simple multivariate Gaussian distribution. It shows that in our segmentation results, some skin-colored occlusions are better detected, and some small occlusions are also located well. Furthermore, our segmentation is more robust to illumination variations. It can also be observed from the reconstructed images that the illumination and texture of the faces are better estimated. Visually speaking, our method reaches competitive fitting results and improved segmentation masks. Please refer to the supplementary materials for more quantitative results.

Image fitting accuracy shows how much the fitting results get misled by occlusions. We evaluate the Root Mean Square Error (RMSE) between the input image and the reconstructed image inside visible face regions, with provided GT segmentation masks. We compare different methods on the AR database, CelebA-HQ testset (referred to as ‘CelebA-Overall’), and two randomly-selected occluded (750 random images) and unoccluded subsets (558 random images), referred to as ‘CelebA-Occluded’ and ‘CelebA-Unoccluded’,

Table 2. Reconstruction error (mm) on the NoW Challenge [Sanyal et al. 2019a].

Method	median	mean	std
Deep3D [Deng et al. 2019b]	1.23	1.54	1.29
DECA [Feng et al. 2020]	1.09	1.38	1.18
PRNet [Feng et al. 2018]	1.50	1.98	1.88
RingNet [Sanyal et al. 2019a]	1.21	1.53	1.31
3DMM-CNN [Tuan Tran et al. 2017]	1.84	2.33	2.05
Proposed	1.04	1.30	1.10

Table 3. Reconstruction error (mm) on the non-occluded and occluded data in the NoW validation subset.

Method	Unoccluded Subset			Occluded Subset		
	median	mean	std	median	mean	std
Deep3D [Deng et al. 2019b]	1.33	1.67	1.41	1.40	1.73	1.41
DECA [Feng et al. 2020]	1.18	1.47	1.24	1.29	1.56	1.29
MoFA [Tewari et al. 2017]	1.35	1.69	1.42	1.36	1.69	1.41
Supervised MoFA	1.02	1.25	1.04	1.05	1.29	1.09
Supervised MoFA-cutmix	1.05	1.28	1.04	1.08	1.33	1.11
Supervised MoFA-cutOUT	1.03	1.28	1.06	1.09	1.34	1.10
Proposed	1.03	1.25	1.03	1.07	1.34	1.19

respectively. As shown in Tab. 1, our fitting accuracy is competitive even to the fully supervised MoFA with data augmentation (i.e. the cutout and cutmix).

Shape reconstruction accuracy is evaluated on the NoW Dataset. The cumulative errors on the testset shown in Tab. 2 indicate that our results reach the state-of-the-art even with a considerably smaller amount of training data and without constraints on identity consistency. To further evaluate the occlusion robustness, 62 pairs of images in the evaluation set are selected with comparable poses with or without occlusions in its publicly-available evaluation set. Tab. 3 shows that the shape reconstruction accuracy of our pipeline is barely affected by occlusions, and reaches a similar level as the fully-supervised pipelines. Please refer to the supplementary for a more detailed analysis.

4.3 Occlusion Segmentation

The accuracy of occlusion segmentation is indicated by four indices: accuracy (ACC), precision (Positive Predictive Value, PPV), recall rate (True Positive Rate, TPR), and F1 score (F1). These indices are only calculated inside the rendered regions for all pipelines. Tab. 4 shows the results on the AR database, where the results of [Egger et al. 2018] are provided for comparison. We separate the AR dataset into three subsets, which include faces without occlusions (neutral), faces with glasses (glasses), and faces with scarves (scarf). According to Tab. 4, the masks predicted by our method show a higher accuracy, recall rate, and F1 score, and competitive precision compared to the skin detector used in [Deng et al. 2019b] and the segmentation method proposed in [Egger et al. 2018].

Table 4. Evaluation of occlusions segmentation accuracy on the AR testsets.

Method	Unoccluded				Glasses				Scarf			
	ACC	PPV	TPR	F1	ACC	PPV	TPR	F1	ACC	PPV	TPR	F1
Deep3D [Deng et al. 2019b]	0.88	0.93	0.94	0.93 ± 0.04	0.88	0.92	0.92	0.92 ± 0.04	0.80	0.80	0.93	0.86 ± 0.05
Egger et al. [Egger et al. 2018]	-	-	-	0.90	-	-	-	0.87	-	-	-	0.86
Proposed	0.88	0.96	0.91	0.93 ± 0.03	0.88	0.98	0.85	0.91 ± 0.04	0.86	0.97	0.81	0.88 ± 0.05

Table 5. Ablation study on the AR testsets and the NoW evaluation subset.

Method	AR-unoccluded				AR-glasses				AR-scarf				NoW Evaluation Set		
	ACC	PPV	TPR	F1	ACC	PPV	TPR	F1	ACC	PPV	TPR	F1	median	mean	std
Pretrained	0.75	0.95	0.77	0.85 ± 0.05	0.78	0.97	0.72	0.82 ± 0.05	0.70	0.89	0.62	0.73 ± 0.07	1.06	1.32	1.14
Baseline	0.81	0.96	0.83	0.89 ± 0.04	0.81	0.97	0.76	0.85 ± 0.05	0.79	0.96	0.71	0.82 ± 0.07	1.06	1.32	1.15
Neighbour	0.85	0.95	0.88	0.91 ± 0.04	0.84	0.95	0.81	0.87 ± 0.04	0.83	0.94	0.79	0.85 ± 0.06	1.06	1.32	1.15
Perceptual	0.89	0.96	0.92	0.94 ± 0.03	0.89	0.98	0.87	0.92 ± 0.04	0.87	0.97	0.84	0.90 ± 0.05	1.06	1.32	1.14
Proposed	0.88	0.96	0.91	0.93 ± 0.03	0.88	0.98	0.85	0.91 ± 0.04	0.86	0.97	0.81	0.88 ± 0.05	1.05	1.31	1.14

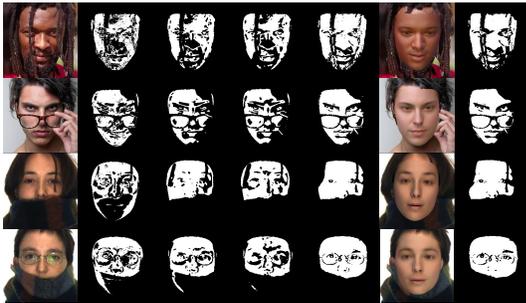


Fig. 5. Qualitative comparison for ablation study. From left to right: target images, masks estimated by the 'Pretrained', 'Baseline', 'Neighbour', and 'Perceptual' pipelines, and the reconstruction results and predicted masks of the proposed pipeline.

4.4 Ablation Study

In this section, we verify the usefulness of the segmentation network and the proposed neighbour loss, $L_{neighbour}$ (Eq. 1), and the coupled perceptual losses, L_{dist} (Eq. 2) and L_{presv} (Eq. 4), for its training. We compare the segmentation performances of ablated pipelines on the AR testset, since the samples are with heavier occlusion, and test the shape reconstruction quality on the NoW evaluation subset. The pre-trained model using the occlusion robust function is referred to as 'Pretrained'. We refer to the segmentation network trained without the neighbour loss or perceptual losses as 'Baseline' and in order to compensate for the lack of such losses we use the pixel-wise reconstruction loss L_{pixel} . The 'Neighbour' pipeline refers to the proposed segmentation network trained only with $L_{neighbour}$ and without the two perceptual losses, and the 'Perceptual' pipeline stands for our proposed network trained only with two perceptual losses, L_{dist} and L_{presv} and without $L_{neighbour}$.

The results in Tab. 5 indicate that with the segmentation network, the segmentation results excel the pretrained model in almost all the indices, which verifies the usefulness of the segmentation network. Comparison among the segmentation results of the 'Baseline', 'Neighbour' and the 'Perceptual' pipelines shows that both losses contribute significantly to the segmentation accuracy.

Fig. 5 provides a visual comparison among the ablated pipelines. It highlights that the occlusion robust function is not robust to illumination variations, and the segmentation network brings great

benefit to the robustness to illumination. The neighbour loss encourages the network to produce smoother results, and the perceptual losses help to locate the occlusions more accurately. Generally, the reconstruction performance of our proposed method are the best one and the segmentation accuracies are also competitive.

4.5 Limitations and Societal Impact

Despite the accurate reconstruction and segmentation proven in the experiments, there are several limitations.

The main issue is that the segmentation performance relies largely on the generative ability of the face model. If the model is too expressive, and can even fit the occlusion, then the segmentation network will not expel the occlusions properly. Likewise, if the model under-fits the target images, the segmentation network tends to regard more pixels as occlusions. For example, in the eye region, the eye gaze usually cannot be properly reconstructed, so some pixels are regarded as occlusions. Although it is solved partially by the neighbour loss, which enhances the robustness to small misalignments, there is still space for improvement. We assume that this problem can be alleviated by an enhanced face appearance model or shape model, such as the detail model proposed by [Feng et al. 2020]. Additionally, we only predict occlusions inside the rendered face region. We assume that a full-face model or a head model can solve this problem.

As for the societal impact, in general, our proposed pipeline has the potential to bring face reconstruction to the real world and to save costs of occlusion labeling, which is generally required in many existing deep-learning-based methods. The model-based reconstruction methods improved by our method could contribute to many applications, including Augmented Reality (AR), Virtual Reality (VR), surveillance, 3D design, and so on. Each of these applications may bring societal and economic benefits, and risks at the same time: The application of AR or VR could bring profits to the entertainment industry and also may result in unethical practices such as demonizing the image of others, identity fraud, and so on. The application of surveillance could help arrest criminals, yet might also invade the privacy and safety of others. The application in 3D design enables the quick capture of the 3D shape of an existing face but might also cause problems in portrait rights.

5 CONCLUSION

In this paper, we have shown how to solve face reconstruction and occlusion segmentation jointly in a weakly-supervised way, so as to enhance the robustness to occlusions for model-based face autoencoders in unconstrained environments. Comprehensive experiments have shown that our method reaches state-of-the-art reconstruction accuracy on the NoW Challenge and provides better segmentation masks as well.

Theoretically, our proposed method can be integrated with the existing face autoencoders. More importantly, we believe that the fundamental concepts of our approach can go beyond the context of face reconstruction and will inspire future work. More specifically, we expect that our pipeline will be extended to many other implementations, such as human body reconstruction, or object reconstruction, with a reliable generative model. We also expect that the masks will be useful for other tasks, e.g. image completion, recognition, or more.

REFERENCES

- Volker Blanz and Thomas Vetter. 2003. Face recognition based on fitting a 3d morphable model. *IEEE Transactions on pattern analysis and machine intelligence* 25, 9 (2003), 1063–1074.
- Adrian Bulat and Georgios Tzimiropoulos. 2017. 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*.
- M. De Smet, R. Fransens, and L. Van Gool. 2006. A Generalized EM Approach for 3D Model Based Face Recognition under Occlusions. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, Vol. 2. 1423–1430. <https://doi.org/10.1109/CVPR.2006.26>
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019a. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4690–4699.
- Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. 2019b. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.
- Terrance DeVries and Graham W Taylor. 2017. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552* (2017).
- Abdallah Dib, Cedric Theobalt, Junghyun Ahn, Philippe-Henri Gosselin, Christian Theobalt, and Louis Chevallier. 2021. Towards High Fidelity Monocular Face Reconstruction with Rich Reflectance using Self-supervised Learning and Ray Tracing. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Bernhard Egger, Sandro Schönborn, Andreas Schneider, Adam Kortylewski, Andreas Morel-Forster, Clemens Blumer, and Thomas Vetter. 2018. Occlusion-aware 3d morphable models and an illumination prior for face image analysis. *International Journal of Computer Vision* 126, 12 (2018), 1269–1287.
- Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhler, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 2020. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)* 39, 5 (2020), 1–38.
- Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. 2020. Learning an Animatable Detailed 3D Face Model from In-The-Wild Images. *arXiv preprint arXiv:2012.04012* (2020).
- Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. 2018. Joint 3D Face Reconstruction and Dense Alignment with Position Map Regression Network. In *ECCV*.
- Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. 2019. GANFIT: Generative Adversarial Network Fitting for High Fidelity 3D Face Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kyle Genova, Forrester Cole, Aaron Maschinot, Aaron Sarna, Daniel Vlasic, and William T Freeman. 2018. Unsupervised training for 3d morphable model regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8377–8386.
- Thomas Gerig, Andreas Morel-Forster, Clemens Blumer, Bernhard Egger, Marcel Luthi, Sandro Schönborn, and Thomas Vetter. 2018. Morphable face models—an open framework. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 75–82.
- Jianzhu Guo, Xiangyu Zhu, Yang Yang, Fan Yang, Zhen Lei, and Stan Z Li. 2020. Towards Fast, Accurate and Stable 3D Dense Face Alignment. In *European Conference on Computer Vision (ECCV)*. 152–168.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- Tatsuro Koizumi and William AP Smith. 2020. “Look Ma, no landmarks!”—Unsupervised, model-based dense face alignment. In *European Conference on Computer Vision*. Springer, 690–706.
- Adam Kortylewski, Andreas Schneider, Thomas Gerig, Bernhard Egger, Andreas Morel-Forster, and Thomas Vetter. 2018a. Training deep face recognition systems with synthetic data. *arXiv preprint arXiv:1802.05891* (2018).
- Adam Kortylewski, Mario Wieser, Andreas Morel-Forster, Aleksander Wiczeorek, Sonali Parbhoo, Volker Roth, and Thomas Vetter. 2018b. Informed MCMC with Bayesian neural networks for facial image analysis. *arXiv preprint arXiv:1811.07969* (2018).
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Fabio Maninchedda, Christian Häne, Bastien Jacquet, Amaël Delaunoy, and Marc Pollefeys. 2016. Semantic 3d reconstruction of heads. In *European conference on computer vision*. Springer, 667–683.
- A. Martinez and Robert Benavente. 1998. The AR face database. *Tech. Rep. 24 CVC Technical Report* (01 1998).
- Andreas Morel-Forster. 2016. *Generative shape and image analysis by combining Gaussian processes and MCMC sampling*. Ph. D. Dissertation. University_of_Basel.
- Elad Richardson, Matan Sela, Roy Or-El, and Ron Kimmel. 2017. Learning Detailed Face Reconstruction From a Single Image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sami Romdhani and Thomas Vetter. 2003. Efficient, Robust and Accurate Fitting of a 3D Morphable Model. In *ICCV*, Vol. 3. 59–66.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 234–241.
- Shunsuke Saito, Tianye Li, and Hao Li. 2016. Real-time facial segmentation and performance capture from rgb input. In *European conference on computer vision*. Springer, 244–261.
- Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael Black. 2019a. Learning to Regress 3D Face Shape and Expression from an Image without 3D Supervision. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael J Black. 2019b. Learning to regress 3D face shape and expression from an image without 3D supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7763–7772.
- Jiaxiang Shang, Tianwei Shen, Shiwei Li, Lei Zhou, Mingmin Zhen, Tian Fang, and Long Quan. 2020. Self-Supervised Monocular 3D Face Reconstruction by Occlusion-Aware Multi-view Geometry Consistency. In *European Conference on Computer Vision (ECCV)*, Vol. 12360. 53–70.
- Ayush Tewari, Michael Zollhofer, Hyeonwoo Kim, Pablo Garrido, Florian Bernard, Patrick Perez, and Christian Theobalt. 2017. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 1274–1283.
- Ayush Tewari, Michael Zollhöfer, Pablo Garrido, Florian Bernard, Hyeonwoo Kim, Patrick Pérez, and Christian Theobalt. 2018. Self-Supervised Multi-Level Face Model Learning for Monocular Reconstruction at Over 250 Hz. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hitika Tiwari, Vinod K Kurmi, KS Venkatesh, and Yong-Sheng Chen. 2022. Occlusion Resistant Network for 3D Face Reconstruction. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 813–822.
- Anh Tuan Tran, Tal Hassner, Iacopo Masi, and Gerard Medioni. 2017. Regressing Robust and Discriminative 3D Morphable Models With a Very Deep Neural Network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1599–1608.
- Anh Tuan Tran, Tal Hassner, Iacopo Masi, Eran Paz, Yuval Nirkin, and Gérard Medioni. 2018. Extreme 3D Face Reconstruction: Seeing Through Occlusions. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3935–3944. <https://doi.org/10.1109/CVPR.2018.00414>
- Luan Tran, Feng Liu, and Xiaoming Liu. 2019. Towards high-fidelity nonlinear 3D face morphable model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1126–1135.
- Anh Tuan Tran, Tal Hassner, Iacopo Masi, and Gérard Medioni. 2017. Regressing robust and discriminative 3D morphable models with a very deep neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5163–5172.
- Ilker Yildirim, Michael Janner, Mario Belledonne, Christian Wallraven, Winrich Freiwald, and Josh Tenenbaum. 2017. Causal and compositional generative models in online perception.. In *CogSci*.

- Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 6023–6032.
- Xiangyu Zhu, Junjie Yan, Dong Yi, Zhen Lei, and Stan Z Li. 2015. Discriminative 3D morphable model fitting. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, Vol. 1. IEEE, 1–8.

A OVERVIEW

In this section we reveal more implementation details and provide more analytical and visual comparisons on the shape reconstruction and segmentation performances.

B IMPLEMENTATION DETAILS

B.1 The neighbor loss

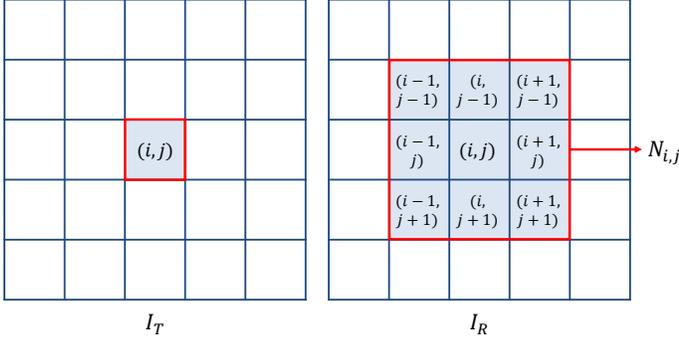


Fig. 1. Visual explanation of the neighbor loss (Eq. 1). For an pixel at (i, j) on the target image I_T (left), we search for the most similar pixel in intensity in its neighboring region, $N_{i,j}$, on the reconstructed image I_R (right).

In this paper, we introduce a new image-level neighbor loss, $L_{neighbor}$, that compares one pixel in the target image to a small region in the reconstructed image:

$$L_{neighbor} = \sum_{x \in \Omega} \left\| \min_{x' \in N(x)} |I_T(x) - I_R(x')| \right\|_2^2 \quad (1)$$

As shown in Fig. 1, for every pixel $I_{T(i,j)}$ in the target image, we search in a 3×3 neighborhood $N(i, j)$ in the reconstructed image $I_{R_{N(i,j)}}$ for the pixel that is most similar to $I_T(i, j)$ in intensity. This neighbour loss accounts for small misalignments of the face model during segmentation.

C QUANTITATIVE ANALYSIS

C.1 Reconstruction Performance on the NoW Challenge.

Fig. 2 shows the cumulative error curves of the proposed method and the state-of-the-arts regarding the NoW Challenge testset. With a higher percentage of sampling points with lower errors, our proposed method performs the best on the NoW Challenge.

We further compare analytically the distributions of reconstruction errors of DECA [Feng et al. 2020] and our proposed method on the NoW validation set, as shown in Fig. 3. To further disentangle the influence of occlusions from other factors, we categorize the samples according to the yaw angles (the angles are rounded off to the nearest 10), and use the error bars under different poses to reflect the distribution of the reconstruction errors. It is obvious from the plots that our proposed method exceeds DECA in mean errors and also yields in much lower variations, even without identity supervision (which emphasize the shape consistency of a same identity) and with significantly less training data. Besides, the lower

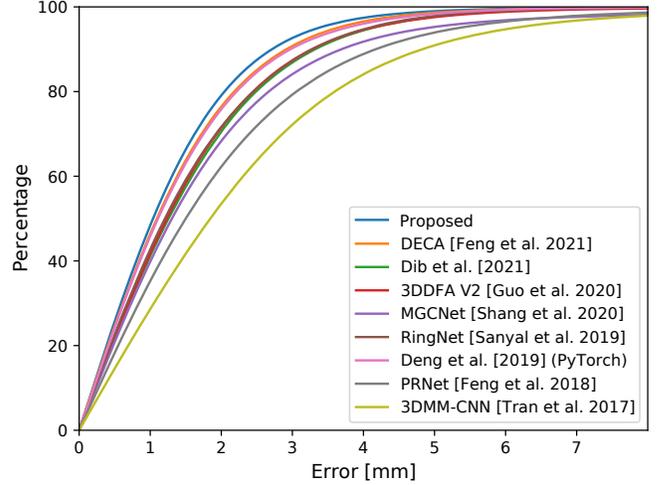


Fig. 2. Quantitative comparison of the 3D reconstruction accuracy on the NoW [Sanyal et al. 2019a] challenge. The methods shown include: DECA [Feng et al. 2020], the work of Dib et al. [Dib et al. 2021], 3DDFA V2 [Guo et al. 2020], MGCNet [Shang et al. 2020], RingNet [Sanyal et al. 2019a], Deep3D (pytorch version) [Deng et al. 2019b], PRNet [Feng et al. 2018], and 3DMM-CNN [Tran et al. 2017].

Table 1. Evaluation of occlusion segmentation accuracy on the CelebA-HQ testsets.

Dataset	Index	Deep3D [Deng et al. 2019b]	Proposed
CelebA-Unoccluded	ACC	0.95	0.92
	PPV	0.98	0.99
	TPR	0.97	0.93
	F1	0.97 ± 0.06	0.96 ± 0.02
CelebA-Occluded	ACC	0.84	0.86
	PPV	0.86	0.95
	TPR	0.96	0.87
	F1	0.90 ± 0.08	0.91 ± 0.06
CelebA-Overall	ACC	0.89	0.89
	PPV	0.92	0.97
	TPR	0.96	0.90
	F1	0.93 ± 0.07	0.93 ± 0.05

or comparable means and standard deviations of the errors under occlusions than under unoccluded conditions proves the effectiveness of our method in improving the occlusion robustness.

C.2 Segmentation Accuracy on the Celeb A HQ testset.

Tab. 1 indicates that the masks predicted by our method show a competitive accuracy, precision, and F1 score, compared to the skin detector used in [Deng et al. 2019b].

C.3 Hyper-parameter Analysis

In this section we systematically evaluate the influence of the hyper-parameters, η_1 to η_5 , used for segmentation. Recall that the total

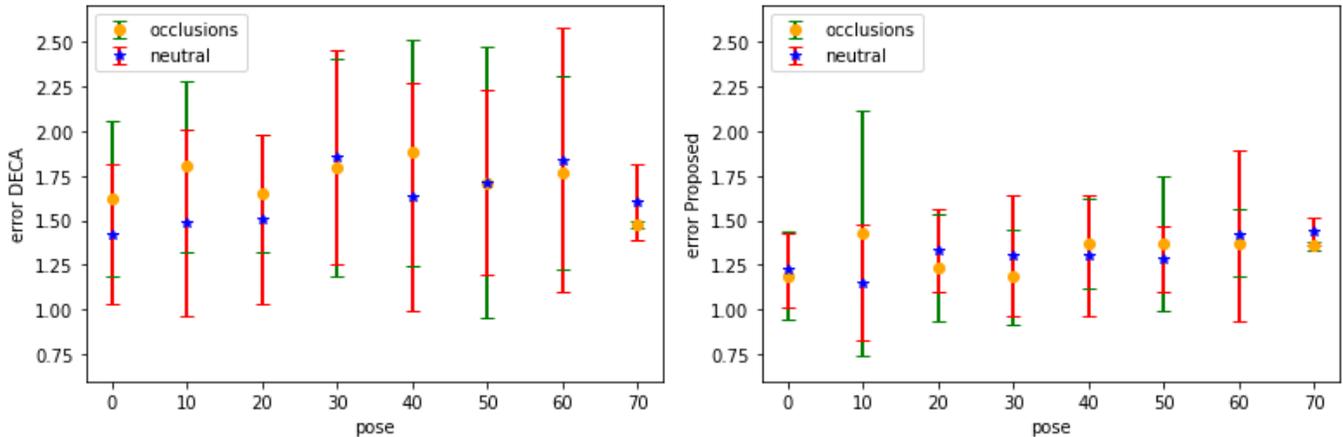


Fig. 3. The distribution of the reconstruction errors under on the neutral and occluded subsets of the NoW validation set. The results of DECA [Feng et al. 2020] are on the left, and ours on the right. The x axis indicates the approximated poses of the samples, and the y axis denotes the reconstruction error.

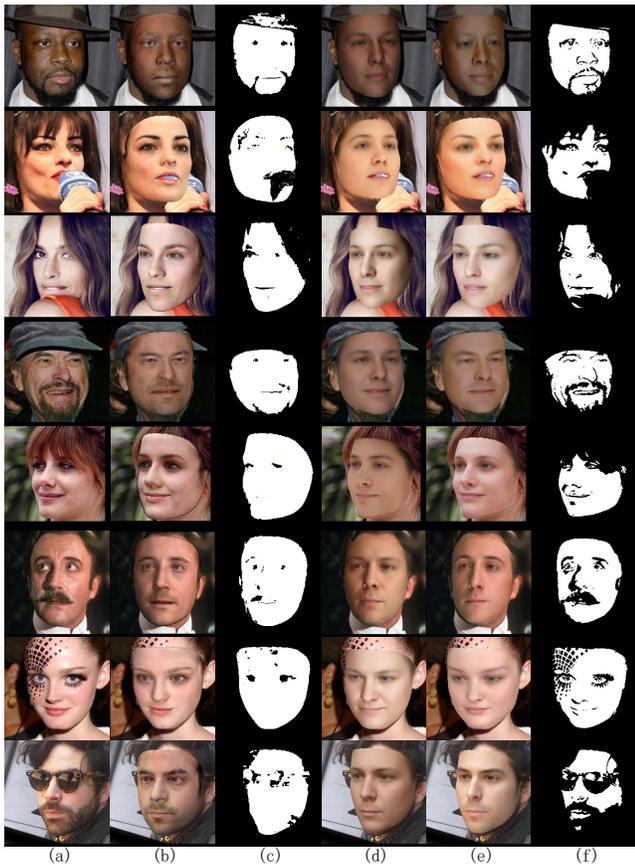


Fig. 4. Comparison on **random samples** in the Celeb A HQ [Liu et al. 2015] testset. (a) Target image. (b) and (c) Reconstruction and segmentation results of the Deep3D network [Deng et al. 2019b]. (d) Reconstructed result of the MoFA network [Tewari et al. 2017]. (e) and (f) Reconstruction and segmentation results of ours.

loss for training the segmentation network is:

$$L_S = \eta_1 L_{neighbor} + \eta_2 L_{dist} + \eta_3 L_{area} + \eta_4 L_{presv} + \eta_5 L_{bin} \quad (2)$$

, with $\eta_1 = 15$, $\eta_2 = 3$, $\eta_3 = 0.5$, and $\eta_4 = 2.5$, and $\eta_5 = 10$. We call this set of parameters as 'standard parameters'. We use the control variates method, namely changing one of the 5 parameters while fixing the others, to evaluate the influence of each hyper-parameters. The accuracy (ACC), precision (Positive Predictive Value, PPV), recall rate (True Positive Rate, TPR), and F1 score (F1) are taken to indicate the segmentation performance. We use the AR dataset [Martinez and Benavente 1998] because the segmentation labels are more accurate.

As shown in Fig. 5, it is clear that with the increase of the neighbour loss $L_{neighbor}$ and the perceptual-level loss L_{dist} , the segmentation network tends to regard more pixels as non-facial. On the contrary, with the increase of the area loss L_{area} and the pixel-wise preserve loss L_{presv} , the segmentation network takes more pixels as face. This observation is consistent with our theory in section 3.2. Fig. 5 also indicates that the indices are positively related to the area loss L_{area} and preserving loss L_{presv} , and are negatively related to the neighbour loss $L_{neighbor}$ and the perceptual-level loss L_{dist} . The binary loss, L_{bin} , barely affects the segmentation.

D QUALITATIVE COMPARISON AND ABLATION STUDY

D.1 Qualitative Comparison on Face Reconstruction and Segmentation

In this section, we provide more visual results of our method on the Celeb A HQ testset [Liu et al. 2015], the AR dataset [Martinez and Benavente 1998], and the NoW Challenge [Sanyal et al. 2019a]. Fig. 4 shows the results on faces with general occlusions. Fig. 6 shows the performance under extreme lighting. Fig. 7 shows the performance of segmenting skin-colored occlusions. Fig. 8 shows the robustness of our method to large poses. Fig. 9 and 10 shows more samples for ablation study.

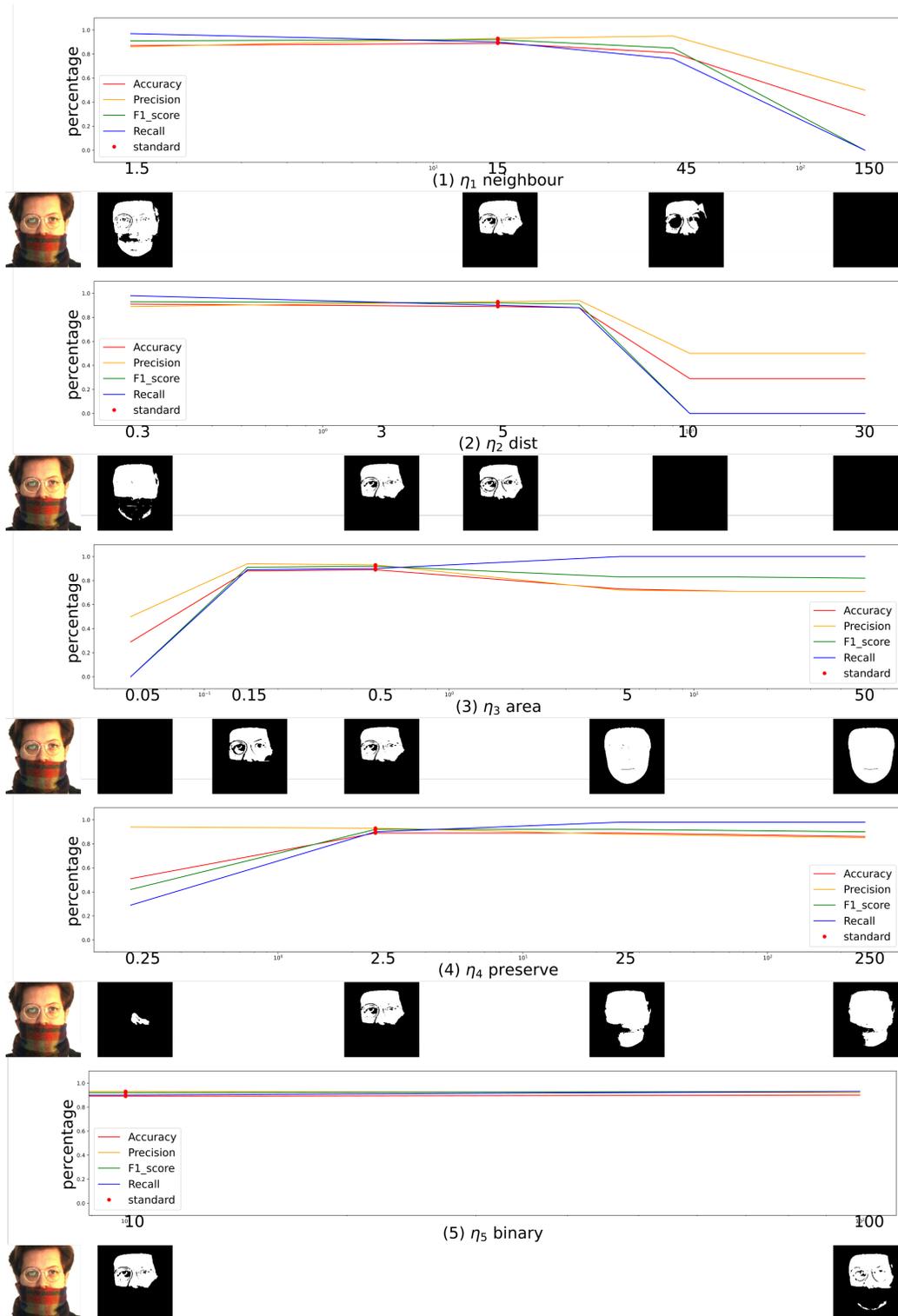


Fig. 5. Analysis of hyper-parameters. The subplots show the change of for indices, namely accuracy, precision, F1 score, and recall rate, with the change of the hyper-parameters. The corresponding segmentation results are shown below each subplot. In each subplot, to evaluate the effect of each hyper parameter η_i , the other hyper-parameters $\eta_j (j \neq i)$ are fixed. The red dots denote the 'standard' positions where the set of parameters in the paper is used.



Fig. 7. Comparison on samples with occlusions that the **skin detector in [Deng et al. 2019b] fails to locate** in the Celeb A HQ testset [Liu et al. 2015]. (a) Target image. (b) and (c) Reconstruction and segmentation results of the Deep3D network [Deng et al. 2019b]. (d) Reconstructed result of the MoFA network [Tewari et al. 2017]. (e) and (f) Reconstruction and segmentation results of ours.



Fig. 6. Comparison on samples with **extreme illumination** conditions in the Celeb A HQ [Liu et al. 2015] and the AR [Martinez and Benavente 1998] testsets. (a) Target image. (b) and (c) Reconstruction and segmentation results of the Deep3D network [Deng et al. 2019b]. (d) Reconstructed result of the MoFA network [Tewari et al. 2017]. (e) and (f) Reconstruction and segmentation results of ours.

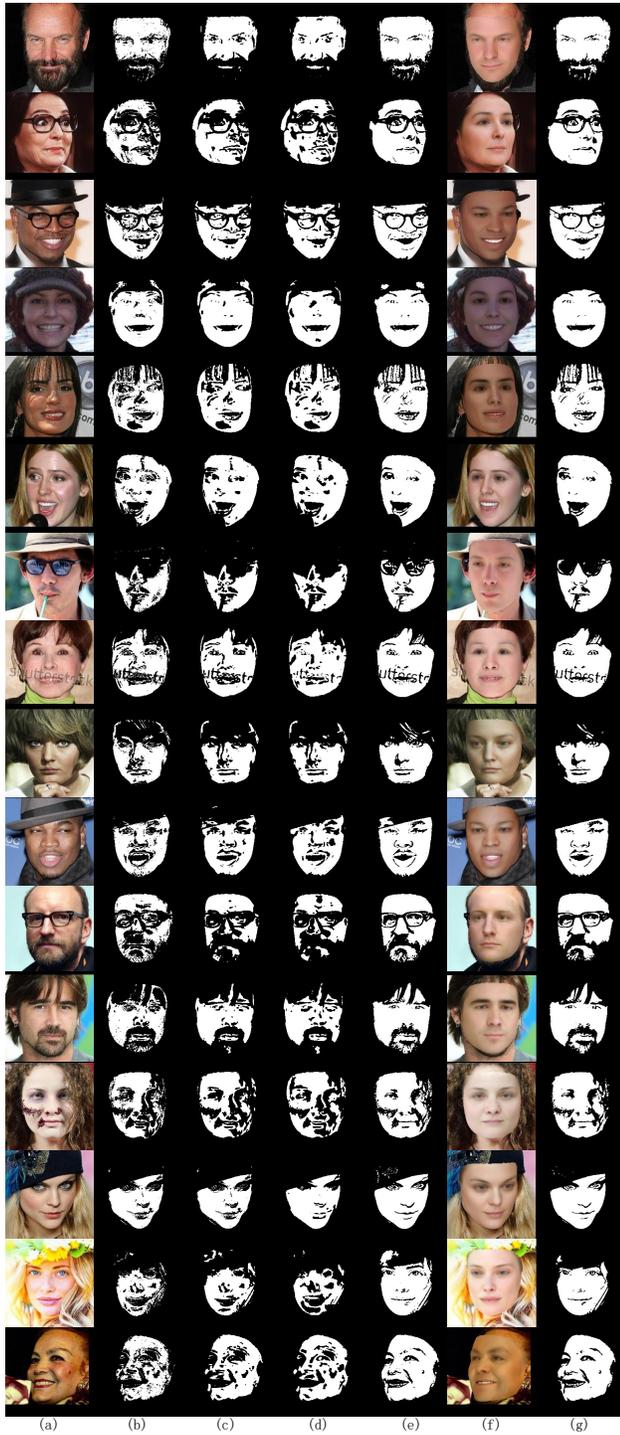


Fig. 9. Qualitative comparison for ablation study on the Celeb A HQ testset [Liu et al. 2015]. From left to right are (a) target images, masks estimated by the (b) 'Pretrained', (c) 'Baseline', (d) 'Neighbour', and (e) 'Perceptual' pipelines, and (f) the reconstruction results and (g) predicted masks of the proposed pipelines, respectively.

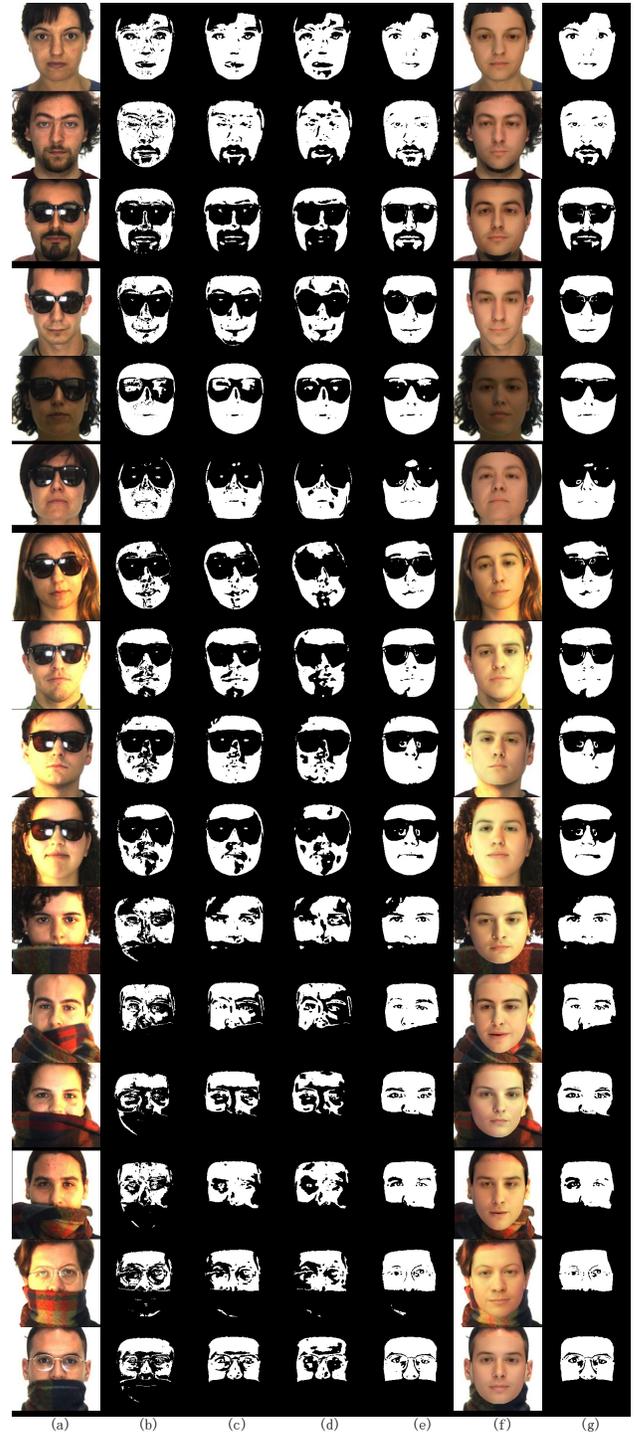


Fig. 10. Qualitative comparison for ablation study on the AR testset [Martinez and Benavente 1998]. From left to right are (a) target images, masks estimated by the (b) 'Pretrained', (c) 'Baseline', (d) 'Neighbour', and (e) 'Perceptual' pipelines, and (f) the reconstruction results and (g) predicted masks of the proposed pipelines, respectively.



Fig. 8. Comparison on samples with occlusions and **large poses** in the NoW Database [Sanyal et al. 2019a] shows that our method can effectively handle occlusions even when there are large poses. (a) Target image. (b) and (c) Reconstruction and segmentation results of the Deep3D network [Deng et al. 2019b]. (d) Reconstructed result of the MoFA network [Tewari et al. 2017]. (e) and (f) Reconstruction and segmentation results of ours.