

---

# Mini-Project 4 - Paper Reproducibility Challenge

---

**Mario Francisco Munoz, Wassim Wazzi, Cesar Arnouk**  
COMP 551 - Applied Machine Learning

1 **Summary of the paper**

- 2 Blind face restoration is a task that normally relies on the geometric information in a face to then reconstruct that face.
- 3 The problem is that the information present in a degraded image is noisy and not very clear. This paper proposes a way
- 4 to use different "geometric priors" from a wide set of faces that are learned by a GAN to perform blind face restoration.
- 5 It achieves this by adding a spatial transform layer to the model.

6 **Scope of Reproducibility**

- 7 We attempt to reproduce similarity scores detailed in the paper and then proceed to push the model to its breaking point
- 8 by adding noise and compression to the input images until there is a radical change in the similarity of the images
- 9 coming out of the model. We also tested the model on different databases to see if occlusion, lighting, age, pose and
- 10 other factors influenced the model's ability to recreate faces.

11 **Methodology**

- 12 We used the same code as the authors. This was particularly easy as the clean, documented code was available on
- 13 GitHub. However, the training of the model was particularly difficult because of its complex structure, but also because
- 14 of the sheer amount of data needed to train it. Because of this, we used the pre-trained facial restoration model provided
- 15 by the authors of the paper for all the experiments. Using image libraries, we added noise and compressed the original
- 16 images from multiple databases to simulate natural degradation. We selected a sample of 100 images from each dataset
- 17 and proceeded to feed these examples through our pipeline. We evaluated the performance of the model using the same
- 18 similarity metrics that the authors detail in the paper. A qualitative analysis of the resulting images was also performed,
- 19 hypothesizing the different artifacts imagined by the generative model and discussing possible sources of failure given
- 20 the current model architecture.

21 **1 Model Analysis**

22 **1.1 Architecture**

- 23 The model begins by implementing a degradation removal solution by using a U-net (Ronneberger, Fischer, and Brox
- 24 2015) architecture to encode and decode the input image. Additionally, this U-net extracts latent features as well as
- 25 multi-resolution spatial features. The latent features discovered by the U-net are then passed through a multi-layer
- 26 perceptron that presents its output to a spatial feature transform algorithm. This is the code that bridges the U-Net to the
- 27 GAN at every stage of the GAN's generation. The spatial feature algorithm's job is to generate two affine transform
- 28 parameters ( $\alpha$  and  $\beta$ ) by taking the U-net's spatial features and applying convolutions. Additional care is provided to
- 29 ensure a balance of fidelity and realness by using channel-splitting and concatenating the original GAN features. These
- 30 parameters will modulate the GAN's features. This complex framework can be observed in Figure 3.

31 **1.2 Loss Functions**

- 32 The model employs various loss functions (see Equation 5) to calculate a total loss function. The first of these,
- 33 reconstruction loss, is meant to keep the models prediction  $\hat{y}$  close to the original  $y$  ground truth image (Equation
- 34 1). To achieve this, we add two terms. The first is the L1 reconstruction loss distance that can be calculated with  $\hat{y}$
- 35 and  $y$  directly. The second term is calculated using  $\phi(\hat{y}), \phi(y)$  denoting a pretrained VGG-19 network. This can be
- 36 perceived as a perceptual loss. Equation 2 shows the adversarial loss employed by the GAN. This loss function is meant
- 37 to incite the GAN to generate natural-looking pictures by looking at the discriminator prediction  $D(\hat{y})$ . Equation 3
- 38 shows a facial component loss from a region of interest (i.e. mouth, left eye). Its first term demonstrating the same
- 39 discriminative loss described in (Goodfellow et al. 2014). The second term in this equation denotes a feature style loss.
- 40 Equation 4 presents an identity loss, inspired by (Huang et al. 2017). This loss is defined by comparing the identity

41 of the ground truth image with the prediction image. This is obtained by using a pre-trained model. The authors use  
42 ArcFace (Deng et al. 2019). The total loss is presented in Equation 5 as the sum of all losses.

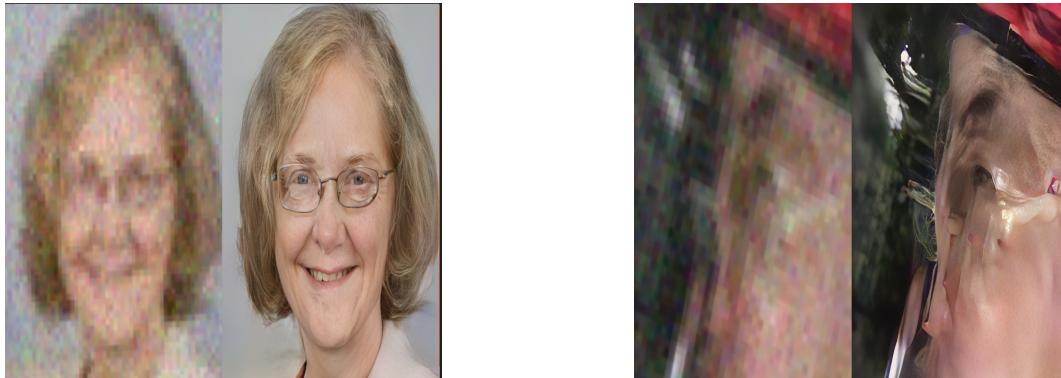
43 **2 Datasets**

44 Different databases and specific images were used for testing troubleshooting, however the main database used to  
45 provide comparable results is the "CelebA" database (Liu et al. 2015) as this one was used by the authors in the paper  
46 as well. We used two other databases to examine the difference in results. The first is UTK, a large face dataset with  
47 occlusions, variations in orientation and changes of lighting in a large span of age groups (UTKFace 2021). The second,  
48 an even more varied dataset with images of protests, faces with makeup and varying scale (Yang et al. 2016).

49 **2.1 Noise and Compression**

50 The authors add noise in the training (for data augmentation) and test (for natural image degradation simulation) datasets.  
51 We ran a data pipeline using OpenCV (Bradski 2000) that compressed images and added noise just as the authors  
52 detail in their paper and then tested other alternatives when the model experienced difficulties when inferring. This  
53 is a 3-part approach, starting with a gaussian blur kernel, a downsizing of the image, and finally the adding of white  
54 noise. This can be observed in Equation 6. The parameters used in these functions came directly from the paper, aiming  
55 to maximize the similarity in the degradation applied to the images during the testing phase. Because of our use of  
56 different varied datasets, the resolution in the faces varied greatly depending on the proportion of size that the face  
57 occupied in the image.

58 **3 Results**



(a) In this figure, we can observe a good reconstruction of a face with glasses. We can also note the lack of symmetry in the frames.  
(b) In this figure, we can observe the crude detection of features in an occluded face with a radical pose.

Figure 1: Two results from our testing from different databases

59 The model achieved very impressive results at times, especially when the faces are not obstructed, and degradation  
60 levels are not very high. However, when the input is too blurry, or the ratio of the size of the face compared to the image  
61 is small, the faces in the output become deformed, and some features are hallucinated on top of the faces (glasses for  
62 example).

63 **3.1 Similarity Metrics**

64 The metrics used in this paper are detailed in equations 7 to 11 in order to evaluate the performance of the model on all  
65 datasets. When compared to the original paper, we see that the results on the celebA dataset are not quite similar, this  
66 could due to many reasons. Firstly, we sub-sampled the dataset in order to maximize the GPU resources at hand, this  
67 might have influenced our findings if the batch was abnormally hard to infer. Another factor that could have influenced  
68 our findings was a difference in the noise and down scaling functions. Even though we carefully followed the article's  
69 guidelines, there could have been a difference in the Gaussian blur, down scaling or white noise functions. The upper  
70 bounds on the noise functions that we used were very high and this section was not thoroughly covered by the authors

	PSNR	SSIM	FID	LPIPS
<b>Input_celebA</b>	14.400068	0.531505	222.854079	0.517061
<b>Input_utk</b>	10.373555	0.459653	159.464849	0.566313
<b>Input_wider</b>	10.404472	0.508763	263.366030	0.603343
<b>Output_celebA</b>	13.727466	0.467375	239.039520	0.537694
<b>Output_utk</b>	9.931122	0.414266	148.868457	0.546444
<b>Output_wider</b>	9.595780	0.413394	313.560935	0.740526
<b>GT_celebA</b>	inf	1.000000	-0.000064	0.000000
<b>GT_utk</b>	inf	1.000000	-0.000086	0.000000
<b>GT_wider</b>	inf	1.000000	-0.000059	0.000000

Figure 2: Metric results we achieved with the 3 Datasets

Methods	LPIPS↓	FID↓	NIQE ↓	Deg.↓	PSNR↑	SSIM↑
<b>GFP-GAN (ours)</b>	<b>0.3646</b>	<b>42.62</b>	<b>4.077</b>	<b>34.60</b>	25.08	0.6777
Input	0.4866	143.98	13.440	47.94	25.35	0.6848
GT	0	43.43	4.292	0	$\infty$	1

Figure 3: Original results presented in the paper on the celebA dataset.

71 in the paper. The color jittering function they used was not described. This is emphasized by the difference in FID  
72 results. However, when looking at the relative difference between the PSNR in inputs and outputs, the results are more  
73 similar to those detailed in the paper. It is also interesting to note that our output under-performs the paper model in  
74 every metric.

### 75 3.2 Qualitative Evaluation

76 We present an array of different results, each labeled according to the important behavior witnessed in the image through  
77 Figures 3 until 14. Throughout this evaluation we can identify major factors that impede the proper reconstruction.  
78 These failure factors include partial face occlusion (see Figure 11, 13), radical poses different from a frontal, eye-level  
79 photograph (see Figure 7,9), improper or unbalanced lighting (see Figure 7, 12, 14), poor resolution (see Figure 5,6),  
80 lack of facial symmetry (see Figure 4,7,10) and lack of visible ethnic minorities in the dataset (see Figures 8,9). These  
81 factors all influence the clarity of the geometric priors that the network can extract from the original image and obviously  
82 gravely impact the reconstruction. It would be interesting to see an implementation of this model that allows the hard  
83 coding of certain latent facial features. For example, if I knew that my grandmother didn't wear glasses and had brown  
84 eyes, the sex, glasses and eye colour information could be passed on to the model to better the reconstruction of the  
85 original image. Mapping these latent codes could also contribute to a better understanding of the model's decisions  
86 and open a discussion about explicative artificial intelligence on the subject, as having an unexplained model could  
87 potentially limit its applications.

## 88 4 Discussion

89 All in all, the model tends to perform properly, but not in every case. It achieves proper face restoration for images than  
90 have gone through a reasonable amount of degradation. However, when we added noise that affected one or more of the  
91 aforementioned failure factors too much, the model started to show some inconsistencies (as can be seen in figures 4 to  
92 14). This is homogeneous with the results in the paper, as the authors stated that "when the degradation of real images is  
93 severe, the restored facial details by GFPGAN are twisted with artifacts."

94 **5 Appendix**

95 **5.1 Equations**

- Reconstruction loss

$$L_{reconstruction} = \lambda_{l1} \|\hat{y} - y\|_1 + \lambda_{per} \|\phi(\hat{y}) - \phi(y)\|_1 \quad (1)$$

- Adversarial loss

$$L_{adversarial} = -\lambda_{adversarial} E_{\hat{y}} \text{softplus}(D(\hat{y})) \quad (2)$$

- Component loss

$$L_{component} = \sum_{ROI} \lambda_{local} E_{\hat{y}ROI} [\log(1 - D_{ROI})] + \lambda_{fs} \|\text{Gram}(\psi(\hat{y}_{ROI})) - \text{Gram}(\psi(\hat{y}_{ROI}))\|_1 \quad (3)$$

- Identity preservation loss

$$L_{identitypreservation} = \lambda_{id} \|\eta(\hat{y}) - \eta(y)\|_1 \quad (4)$$

- Total loss

$$L_{total} = L_{rec} + L_{adv} + L_{comp} + L_{id} \quad (5)$$

- Noise and Compression equation:

$$x = [(y \cdot k_\sigma) \downarrow_r + n_\delta] \text{JPEG}_q \quad (6)$$

102 Values used by the authors:

103 the gaussian blur kernel was unspecified, so we used a range between 3 and 15

$$\downarrow_r = 1 : 8, \sigma = 0.2 : 10, n_\delta = 0 : 15 \quad (7)$$

- FID - Fréchet Inception Distance (Heusel et al. 2018)

$$FID = \|\mu - \mu_w\|_2^2 + \text{tr}(\Sigma + \Sigma_w - 2(\Sigma \cdot \Sigma_w)^{1/2}) \quad (8)$$

- NIQE - Natural Image Quality Evaluator (Mittal, Soundararajan, and Bovik 2013)

$$D(v_1, v_2, \Sigma_1, \Sigma_2) = \sqrt{\left( (v_1 - v_2)^T \left( \frac{\Sigma_1 + \Sigma_2}{2} \right)^{-1} (v_1 - v_2) \right)} \quad (9)$$

- LPIPS - Learned Perceptual Image Patch Similarity (Zhang et al. 2018)

$$d(x, x_0) = \sum_l \frac{1}{H_l W_l} \sum_{h,w} \|w_l \cdot (\hat{y}_{hw}^l - \hat{y}_{0hw}^l)\|_2^2 \quad (10)$$

- PSNR - Peak Signal to Noise Ratio

$$PSNR = 20 \cdot \log_{10} \left( \frac{\text{MAX}f}{\sqrt{MSE}} \right) \quad (11)$$

- SSIM - Structural Similarity

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma + c_2)}{(\mu_x^2\mu_y^2 + c_1)(\sigma_x^2\sigma_y^2 + c_2)} \quad (12)$$

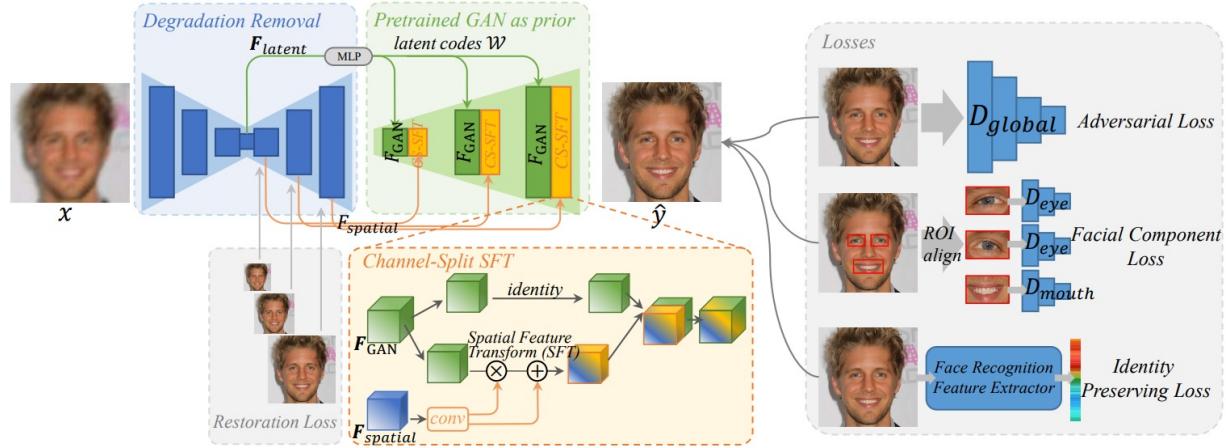


Figure 4: Overview of GFP-GAN framework



Figure 5: In this figure, we can observe the perfect case scenario for our model, a well-lit face without occlusion.



Figure 6: In this figure, we can observe a feature mismatch of a surgeon with two very different eyes and a deformation around the mouth.

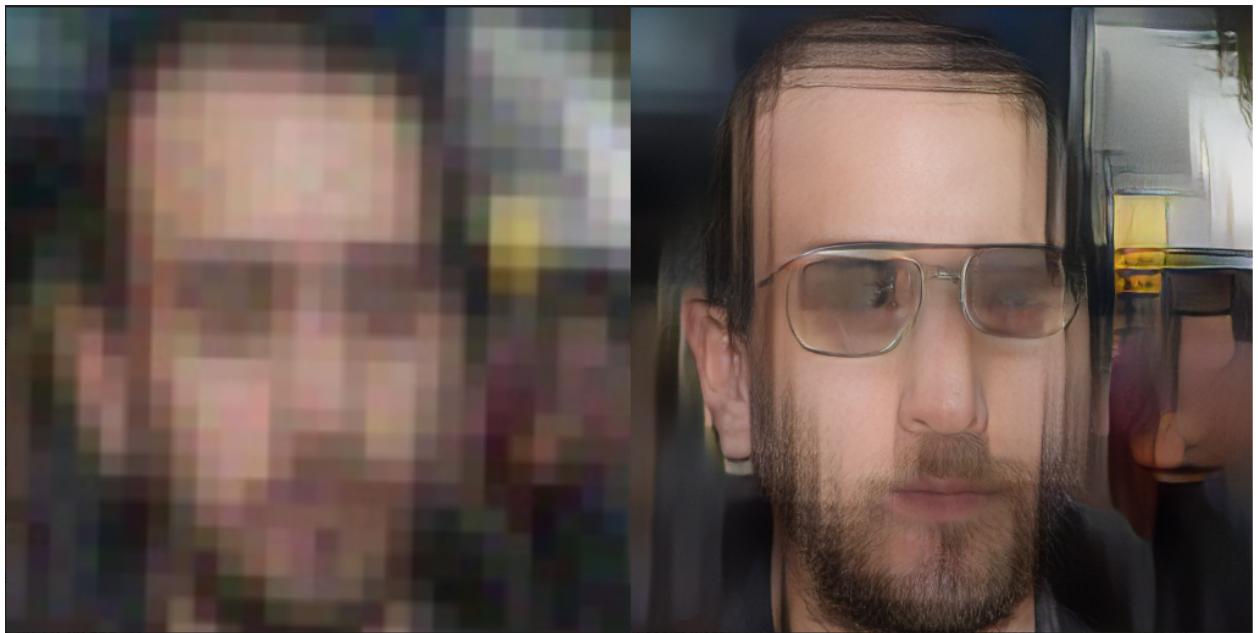


Figure 7: In this figure, we can observe the hallucination of glasses by the GAN model. Face alignment makes pixel distortion look like glasses.

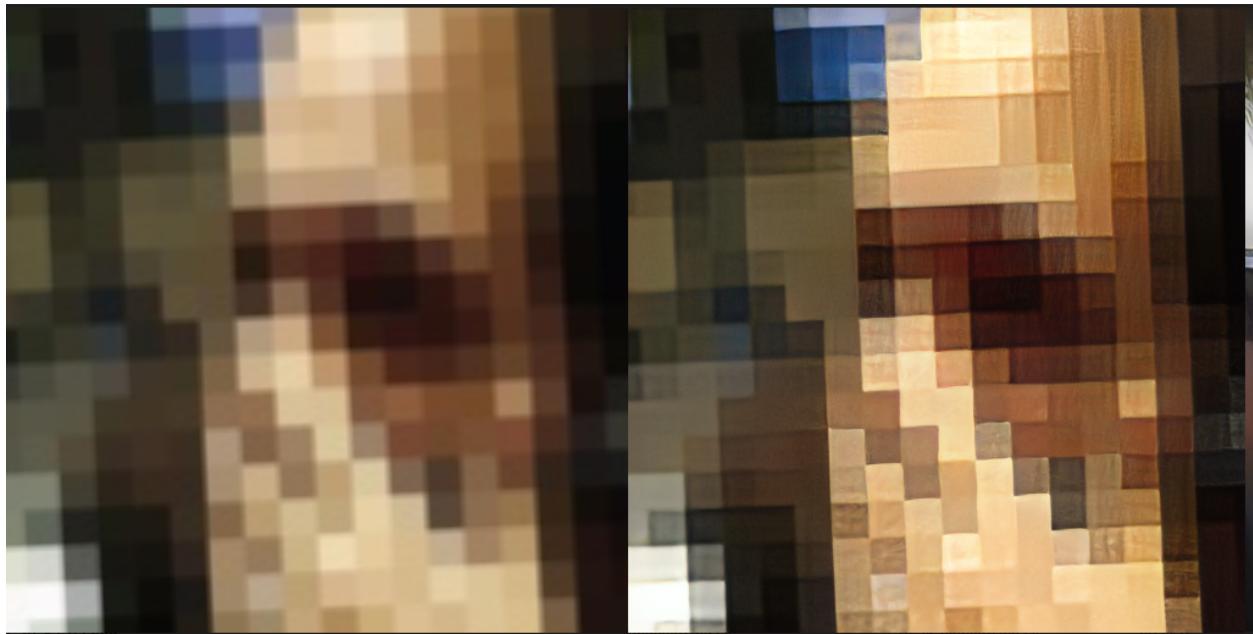


Figure 8: In this figure, we can observe a total failure at an attempt where the model recognized a face, but had extracted little geometric information



Figure 9: In this figure, we can observe another mismatch between features. The left and right eyes are not matched and the glare in the subject's left eye appears to be coming from a camera flash. Could the model be biased? After all, it was trained on a celebrity dataset.



Figure 10: In this figure, we can observe an odd texture around the man’s hair, indicating maybe fewer examples of ethnic minorities were represented in the training data as the model struggles to find the right hair texture.

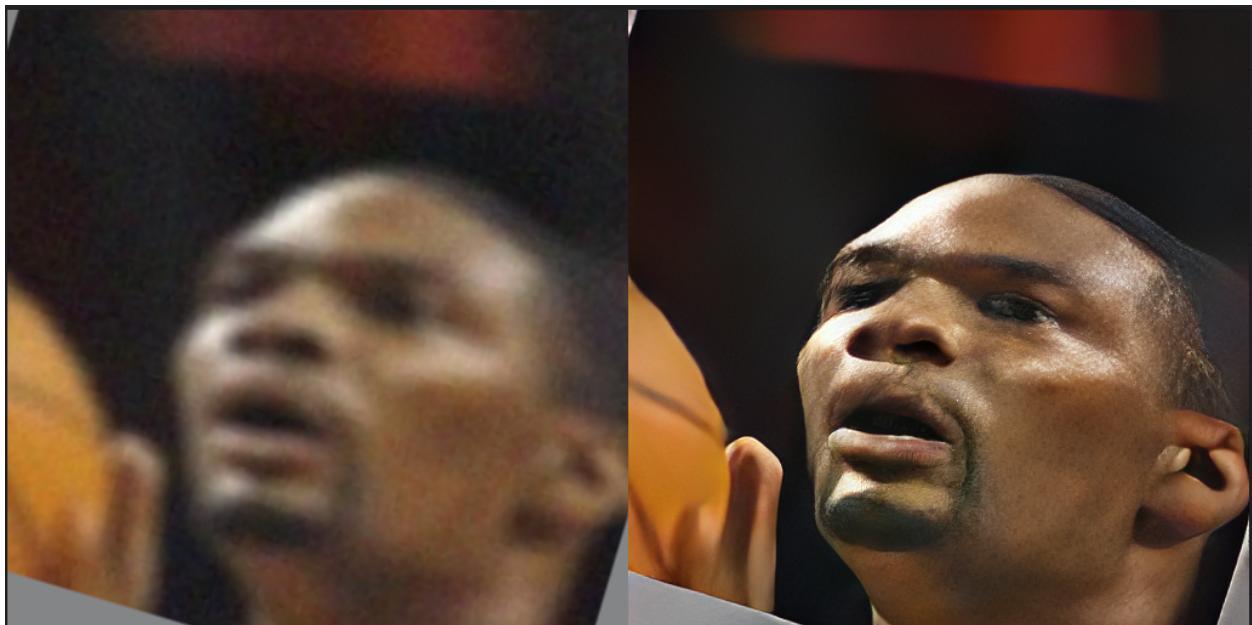


Figure 11: In this figure, we can observe a face with a difficult pose. The model does not perform as well.



Figure 12: In this figure, we can observe another eye feature mismatch with the same flash lighting phenomenon.



Figure 13: In this figure, we can observe the model's failure to reconstruct the glasses on the man's face.

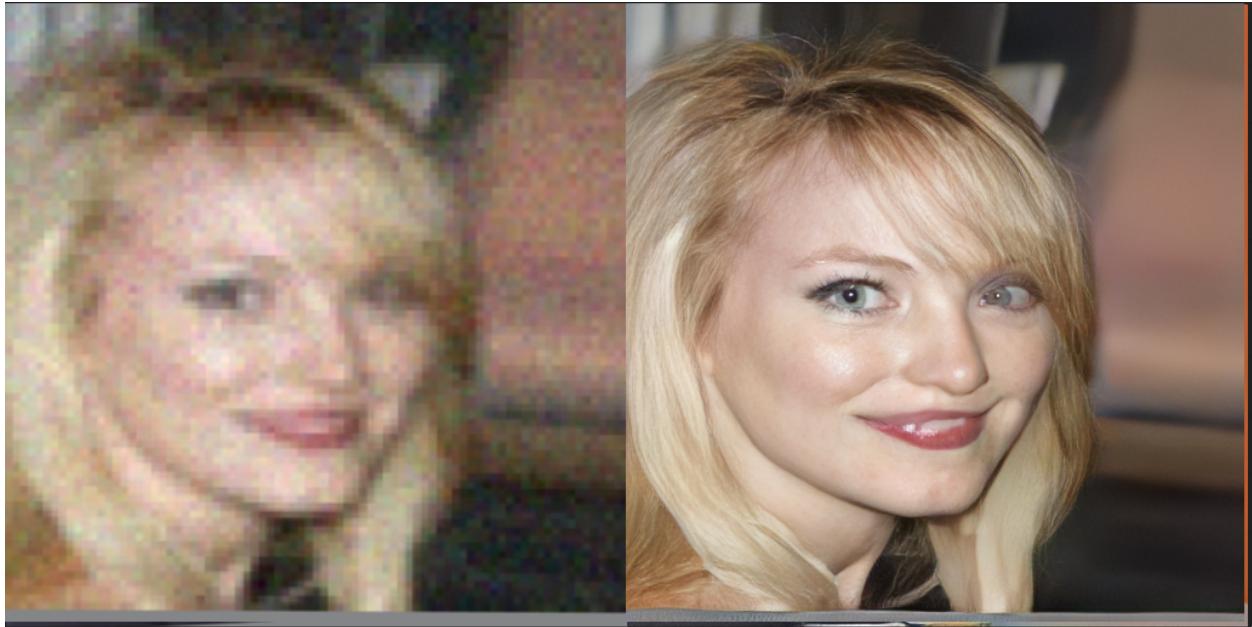


Figure 14: In this figure, we can observe an eye alignment mismatch.



Figure 15: In this figure, we can observe a more difficult reconstruction for the model, given that the face is occluded and the model has no geometric information about hands.



Figure 16: In this figure, we can observe another eye color and alignment mismatch.

110 **References**

- 111 Bradski, G. (2000). “The OpenCV Library”. In: *Dr. Dobb’s Journal of Software Tools*.
- 112 Mittal, Anish, Rajiv Soundararajan, and Alan C Bovik (2013). “Making a ‘Completely Blind’ Image Quality Analyzer”.
- 113 en. In: p. 4.
- 114 Goodfellow, Ian et al. (2014). “Generative Adversarial Nets”. en. In: p. 9.
- 115 Liu, Ziwei et al. (Dec. 2015). “Deep Learning Face Attributes in the Wild”. In: *Proceedings of International Conference*  
116 *on Computer Vision (ICCV)*.
- 117 Ronneberger, Olaf, Philipp Fischer, and Thomas Brox (May 2015). “U-Net: Convolutional Networks for Biomedical  
118 Image Segmentation”. en. In: *arXiv:1505.04597 [cs]*. arXiv: 1505.04597. URL: <http://arxiv.org/abs/1505.04597> (visited on 11/28/2021).
- 119 Yang, Shuo et al. (2016). “WIDER FACE: A Face Detection Benchmark”. In: *IEEE Conference on Computer Vision*  
120 *and Pattern Recognition (CVPR)*.
- 121 Huang, Rui et al. (Oct. 2017). “Beyond Face Rotation: Global and Local Perception GAN for Photorealistic and  
122 Identity Preserving Frontal View Synthesis”. en. In: *2017 IEEE International Conference on Computer Vision*  
123 (*ICCV*). Venice: IEEE, pp. 2458–2467. ISBN: 978-1-5386-1032-9. DOI: 10.1109/ICCV.2017.267. URL: <http://ieeexplore.ieee.org/document/8237529/> (visited on 12/02/2021).
- 124 Heusel, Martin et al. (Jan. 2018). “GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash  
125 Equilibrium”. en. In: *arXiv:1706.08500 [cs, stat]*. arXiv: 1706.08500. URL: <http://arxiv.org/abs/1706.08500>  
126 (visited on 12/07/2021).
- 127 Zhang, Richard et al. (2018). “The Unreasonable Effectiveness of Deep Features as a Perceptual Metric”. en. In: p. 10.
- 128 Deng, Jiankang et al. (2019). “ArcFace: Additive Angular Margin Loss for Deep Face Recognition”. en. In: p. 10.
- 129 UTKFace (2021). en-US. URL: <https://susanqq.github.io/UTKFace/> (visited on 12/13/2021).