

MoFRR: Mixture of Diffusion Models for Face Retouching Restoration

Jiaxin Liu¹, Qichao Ying¹, Zhenxing Qian¹, Sheng Li^{1†}, Runqi Zhang¹, Jian Liu², Xinpeng Zhang¹
¹ Fudan University ²Ant Group

Abstract

The widespread use of face retouching on social media platforms raises concerns about the authenticity of face images. While existing methods focus on detecting face retouching, how to accurately recover the original faces from the retouched ones has yet to be answered. This paper introduces Face Retouching Restoration (FRR), a novel computer vision task aimed at restoring original faces from their retouched counterparts. FRR differs from traditional image restoration tasks by addressing the complex retouching operations with various types and degrees, which focuses more on the restoration of the low-frequency information of the faces. To tackle this challenge, we propose MoFRR, Mixture of Diffusion Models for FRR. Inspired by DeepSeek’s expert isolation strategy, the MoFRR uses sparse activation of specialized experts handling distinct retouching types and the engagement of a shared expert dealing with universal retouching traces. Each specialized expert follows a dual-branch structure with a DDIM-based low-frequency branch guided by an Iterative Distortion Evaluation Module (IDEM) and a Cross-Attention-based High-Frequency branch (HFCAM) for detail refinement. Extensive experiments on a newly constructed face retouching dataset, RetouchingFFHQ++, demonstrate the effectiveness of MoFRR for FRR.

1. Introduction

Retouched face images, produced through various techniques such as face lifting, eye enlargement, whitening, and smoothing, are pervasive on social platforms. Most users acquire these images through straightforward and convenient methods to improve their appearance. Despite amusement and fun, the abuse of face retouching causes a lot of serious problems, including aesthetic degradation, commercial deception [2] and identity fraud [3]. To address these problems, many retouching detection schemes and regulatory countermeasures have been proposed. For instance, Norway [14] has enacted strict disclosure requirements for

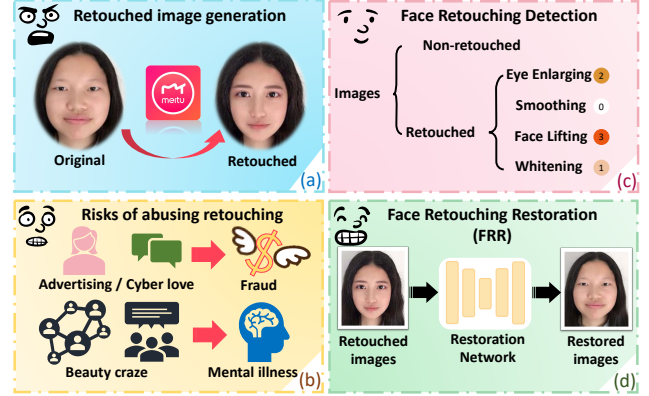


Figure 1. Application scenario of the proposed scheme (MoFRR). (a) People use face retouching in various applications, (b) face retouching poses risks such as fraud, societal security, and cultural psychological issues, (c) existing work for face retouching detection, (d) our proposed MoFRR method recovers the original image from the retouched version, thereby offering an additional layer of protection against face retouching fraud.

edited promotional content. US and Israel also have labeling requirements for retouched faces [9, 30]. Researchers have developed dedicated datasets and neural networks to auto-detect the presence of retouching within a digital image [7, 11, 42]. In application, once exaggerated retouching is detected, the corresponding images can be automatically removed, or the livestreams closed.

In the literature, researchers have focused mainly on the accurate detection of face retouching [20, 31, 42]. It has yet to be answered regarding how we could near-faithfully recover the original face image from the retouched one, which is important to trace back to the real identity of the severely retouched faces.

In this paper, we consider a new computer vision task, Face Retouching Restoration (FRR), to restore the original faces from the retouched ones, which serves as a further step for fighting against the misuse of face retouching. In literature, researchers have dedicated their efforts to dealing with similar tasks, including makeup removal [37] and image restoration (IR) [23, 27, 41]. The former transfers the style of a reference face to a makeup face for face restora-

†Corresponding author: Sheng Li (lisheng@fudan.edu.cn)

tion. The latter tries to restore high-quality textural details from low-quality images. These schemes work well in recovering textural features. However, in face retouching, it may not be sufficient to change only the textural features of the face image. Usually, the face structure has to be modified in order to have distinct differences in appearance. Consequently, noticeable structural changes can be made during face retouching, such as enlarged eyes, more prominent noses, and thinner face shapes. Thus, directly applying existing makeup removal or image restoration schemes may not be appropriate for FRR.

We believe that a good and appropriate FRR scheme should focus more on the restoration of low-frequency information of the faces, which significantly differs from the existing approaches for the IR and makeup removal tasks. FRR is a challenging task as the faces might be operated in a complex setting of face retouching with various types and degrees, including whitening, smoothing, eye enlargement, and face lifting. The logic among different touching types varies and the goals are also different. Thus, it might not be appropriate to perform the FRR using a single model. We also notice that the impact of different face retouching operations is independent in the face images.

Inspired by the Mixture of Experts (MoE) [4, 21], we propose MoFRR that restores the face images in a divide-and-conquer manner using an MoE framework with specialized and shared experts. A specialized expert is trained to explicitly restore the faces for a certain type of face retouching. The design of the shared expert is inspired by the recent success of DeepSeek [28], which aims to deal with universal retouching traces in different types of face retouching. We propose a novel wavelet DDIM [34] model, denoted as WaveFRR, for the design of specialized experts. Each WaveFRR contains a dual-branch structure including a DDIM-based low-frequency branch and a cross-attention-based high-frequency branch. The former contains an Iterative Distortion Evaluation Module (IDEM) that restores the wavelet low-frequency sub-band of the faces. The latter contains a High-Frequency Cross-Attention Module (HFCAM) that recovers the wavelet high-frequency sub-bands. For the shared expert, we adopt the ordinary DDIM architecture to make it general for various retouching conditions.

For training and evaluation, we newly construct a dataset named RetouchingFFHQ++, which contains over a million retouched face images from four commercial face retouching APIs. We have conducted comprehensive experiments including intra-API and cross-API tests, both demonstrate the effectiveness of our MoFRR for face retouching restoration. To summarize, our contributions are as follows.

- We are the first to consider the task of face retouching restoration, and we propose MoFRR, an MoE framework with both specialized and shared experts, to address this challenging problem.

- We propose a wavelet-based DDIM model for face restoration from a specific retouching type, where an IDEM module and an HFCAM module are designed to recover the low-frequency and high-frequency sub-bands of the face, respectively.
- We extend the large-scale face retouching dataset, RetouchingFFHQ [42], into RetouchingFFHQ++ to make it more suitable for training and evaluation of FRR schemes.

2. Related Works

2.1. Makeup Transfer and Removal

Over the past year, makeup removal [8, 37] has garnered significant attention in research, often explored in conjunction with makeup transfer [15]. PairedCycleGAN [6] utilizes an asymmetric framework that includes an additional sub-network specifically for makeup removal. PS-GAN++ [29] employs a makeup distillation network and an identity extraction network to facilitate makeup transfer and removal. SSAT [35] is proposed for simultaneous makeup transfer and removal, utilizing a Symmetric Semantic Corresponding Feature Transfer module for accurate semantic alignment. Recently, CSD-MT [36] introduces an unsupervised framework to achieve robust performance by adaptively combining high-frequency content details and low-frequency style features. Nevertheless, face retouching differs from making-up in that the former often involves changes to the facial structure while the latter does not.

2.2. Image Restoration

Many methods are proposed to utilize facial component dictionaries [26] or codebooks [16, 39] derived from high-quality images to guide the restoration process. DR2 [40] uses a DDPM to transform degraded images into coarse, degradation-invariant predictions, which are then enhanced to high-quality images. Restormer [43] uses an efficient Transformer to achieve superior performance on several image restoration tasks. ResDiff [33] proposes a DDPM to predict residuals, using frequency-domain loss and guided diffusion to improve generation speed and sample quality. However, image restoration primarily focuses on recovering high-frequency texture details for image quality improvement, while FRR often aims to restore a face image more semantically.

3. Method

3.1. FRR: Goals, Evaluations and Our Dataset

Goals and Evaluations. The goal of FRR is defined to blindly, i.e., with no other reference of template, reconstruct the original face images given the retouched ones via computer vision algorithms or deep learning methods.

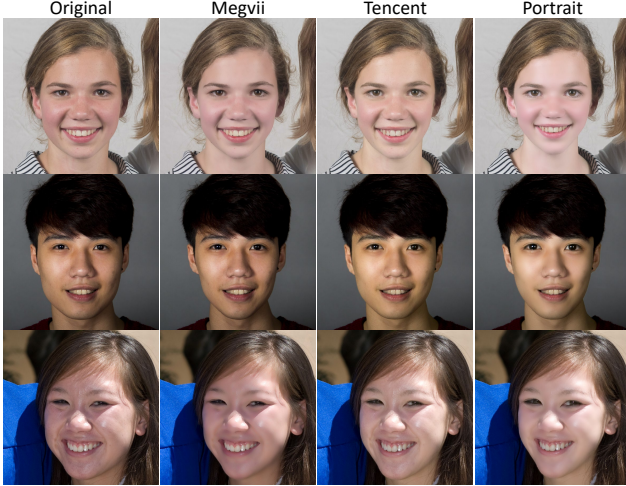


Figure 2. Examples from the RetouchingFFHQ++ dataset, where original faces are retouched by different retouching APIs. In each row, the applied retouching types and degrees are kept same.

| | Portrait | Megvii | Tencent | Alibaba | Total |
|--------|----------|--------|---------|---------|---------|
| Ori | - | - | - | - | 57910 |
| Single | 108267 | 66692 | 200100 | 29483 | 404542 |
| Multi | 206655 | 232641 | 173590 | - | 612886 |
| Total | 314922 | 299333 | 373690 | 29483 | 1075338 |

Table 1. Summary of RetouchingFFHQ++. The “Multi” column includes a mixture of 2+ face retouching operations.

Our evaluation framework prioritizes forensic admissibility over perceptual quality, explicitly discouraging methods that synthesize plausible but inauthentic facial features. We adopt a dual-metric approach: pixel-level fidelity is quantified through standard PSNR and SSIM [38] to measure spatial and structural reconstruction accuracy, while biometric veracity is assessed via feature-space cosine similarity using state-of-the-art face recognition models like AdaFace [24] and ArcFace [12], etc. This dual evaluation mechanism ensures FRR methods meet both visual reconstruction benchmarks and evidentiary standards required for judicial applications, where synthetic enhancements could render biometric evidence legally inadmissible.

The RetouchingFFHQ++ Dataset. In order to train networks for FRR, we prepare our first dedicated dataset for FRR based on RetouchingFFHQ [42], a recently proposed large-scale face retouching detection dataset that includes high-quality, fine-grained annotated retouched face images from commercial APIs such as Tencent, Megvii, and Alibaba. The original face images are from the famous FFHQ dataset [22], and the retouched subsets includes thousands of retouched images, covering whitening, smoothing, face lifting, eye enlarging, or combined operations of the above.

But the dataset is originally proposed for face retouching detection, so in order to better cater to our FRR task, we extend this dataset as RetouchingFFHQ++ in the following two aspects.

First, we expand the dataset using the widely-adopted PortraitPro 24 API [1], increasing the total number of images to over a million, as shown in Tab. 1. We show examples of the dataset in Fig. 2, and we include more details and statistics of the dataset in the supplement. From the exemplified images, we see that retouched images by the PortraitPro API can have noticeable visual differences compared to the existing images from other APIs, thus further diversifying the dataset.

Second, we propose a more reasonable criteria on the degree definition of each retouching. While RetouchingFFHQ directly defines the degrees according to the argument passed on generating the images, it could less reflect the actual degree due to API-level algorithmic discrepancy. To mitigate this, we redefine the degree of modification of each operation within images via statistically analyzing the PSNR distribution of all retouched face images. We re-categorize them into five groups with proportions of 15%, 25%, 25%, 25%, and 10%, based on ascending PSNR values, and labeled from 5 to 1 accordingly.

3.2. FRR Model with Multi Experts

The overall pipeline of our proposed MoFRR, as depicted in Fig. 3, is a Mixture of Diffusion Models for FRR. We explicitly decompose the task into several steps. First, given a targeted image, we first predict which types of face retouching operations are performed. Second, we employ several expert networks, where apart from a reserved shared expert as inspired by DeepSeek, the rest of experts (denoted as specialized experts) are respectively specialized in removing one certain type of face retouching in the given image. As a result, the experts provide several versions of intermediate recovered images. Third, with the intermediate images from all the activated experts, we merge these images with the original image through the Combine Module to obtain the final restored image.

Concretely, we train a router to identify which face retouching operations have been applied to the images and to activate the following expert networks. We implement it by the ResNet-MAM [42], which is adapted for multi-label classification and outputs an N -dimensional binary label vector. Here in our paper, we focus on the most typical types of face retouching operations, which are *whitening*, *smoothing*, *face lifting*, and *eye enlarging*, so $N = 4$.

$$[b_w, b_s, b_f, b_e] = \mathbb{1}(\text{Router}(X) > \text{th}), \quad (1)$$

where $X \in \mathbb{R}^{H \times W \times C}$ represents the input retouched face image. b_M for $M \in \{w, s, f, e\}$ denotes the prediction result of whether the input image contains the correspond-

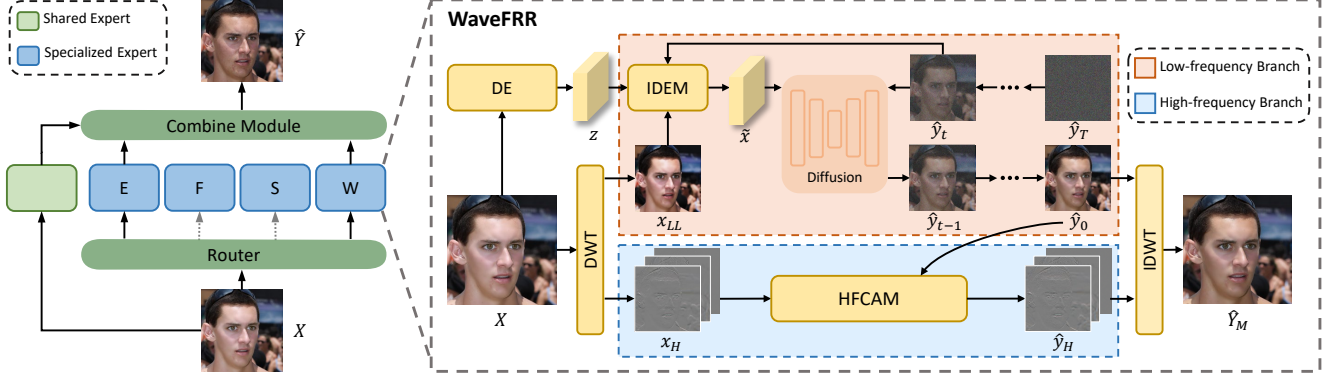


Figure 3. Left: Overview of MoFRR, where input image is processed by the router and selectively sent into the specialized experts for different retouching types. A shared expert is continuously engaged inspired by Deepseek, and the result images from the experts are jointly sent to a lightweight Combine Module to produce the final restored image. Right: the WaveFRR design, where the image is decomposed using DWT, and the DDIM process works on low-frequency subband restoration guided by the transformed degree estimation result (by IDEM). Then HFCAM module restores the high-frequency sub-bands via cross-attending the restored lower sub-bands.

ing face retouching operation M . $\mathbb{1}(\cdot > \text{th})$ represents if the prediction exceeds a given threshold (here 0.5 by default).

Next, four (N) independent experts, named WaveFRR models, process the image once activated by the router. The design of a WaveFRR model is illustrated in Fig. 3. Each expert is specialized in estimating and restoring a particular face retouching operation. It begins with predicting the degree of the specific face retouching operation, which is trained ahead and guides the subsequent restoration process. Then, considering that face retouching primarily emphasizes content modifications, such as cheeks and eyes, rather than global high-frequency details, we first perform the discrete wavelet transformation (DWT) that decomposes the input image into lower and higher frequencies. For lower frequencies, we propose sampling a retouching-free version via conditional DDIM [18, 34], guided by the predicted retouching degree and the original lower frequencies. For higher frequencies, we modify them by cross-attending [13] the generated lower frequencies for more refined retouching removal. The modified frequencies are then combined via Inverse Discrete Wavelet Transform (IDWT) to produce the restored image \hat{Y}_M , which removes only the specific face retouching operation M .

For the shared expert, we adopt a baseline DDIM architecture without wavelet decomposition, where we hope the architectural divergence from WaveFRR explicitly could promote complementary functional specialization. While WaveFRR experts focus on frequency-aware sub-band recovery, the shared expert captures global retouching patterns across spectral domains, denoted as \hat{Y}_{se} .

Finally, we gather all restored outputs of the activated experts and design a lightweight Combine Module to give the final retouching-free image \hat{Y} . The Combine Module employs a UNet [32], which leverages all images and the re-

touched image to generate results while preserving fidelity to the input image, as follows:

$$\hat{Y} = \text{Combine_Module}([X, \hat{Y}_{se}, \hat{Y}_{M_0}, \hat{Y}_{M_1}, \dots]). \quad (2)$$

3.3. WaveFRR Model

Degree Estimator and DWT. Each expert begins with estimating the degree of its specialized face retouching operation. Afterward, we train a ResNet50 [17] as Degree Estimators, using our newly defined degree labels as supervision. It predicts the global retouching degree $z = \text{DE}_M(X)$ for the face retouching operation M . Besides the estimator, we apply the DWT to decompose it into low-frequency and high-frequency sub-bands, denoted as $\{x_{LL}, x_H\}$, where $x_{LL} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$ and $x_H \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times 3C}$ represent one low-frequency sub-band, i.e., LL, and three high-frequency sub-bands, i.e., HH, LH and HL.

Low-Frequency Branch. The LL sub-band of DWT contains major information of X . We propose to sample a retouching-free LL sub-band of X via conditional DDIM, guided by the predicted retouching degree and the original LL sub-band. For denoised image \hat{y}_t at timestamp t , our IDEM module utilizes the multi-scale channel attention network [10], denoted as $\text{MCA}(\cdot)$, which produces the pixel-wise condition $\tilde{x} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times 2C}$ given \hat{y}_t, x_{LL} and z as:

$$F = \text{MCA}(z + x_{LL}) \otimes x_{LL} + (1 - \text{MCA}(z + x_{LL})) \otimes z, \quad (3)$$

$$\hat{R} = \text{MCA}(F + \hat{y}_t) \otimes \hat{y}_t + (1 - \text{MCA}(F + \hat{y}_t)) \otimes F, \quad (4)$$

$$\tilde{x} = \text{IDEM}(\hat{y}_t, x_{LL}, z) = \text{Concat}(x_{LL}, \hat{R}), \quad (5)$$

where $\hat{R} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$ is the pixel-wise distortion map.

The DDPM [18] uses a Markov chain to restore the image step by step from Gaussian noise. The forward diffusion process progressively transforms the ground truth

y_0 into noise data in T steps with a variance schedule $\{\beta_1, \beta_2, \dots, \beta_T\}$, which can be formulated as:

$$q(y_t | y_{t-1}) = \mathcal{N}(y_t; \sqrt{1 - \beta_t} y_{t-1}, \beta_t \mathbf{I}), \quad (6)$$

where \mathcal{N} denotes the Gaussian distribution. The denoising process gradually removes noise from a randomly sampled Gaussian noise $\hat{y}_T \sim \mathcal{N}(0, \mathbf{I})$. We use conditional denoising process with the condition \tilde{x} , i.e.,

$$p_\theta(\hat{y}_{0:T} | \tilde{x}) = p(\hat{y}_T) \prod_{t=1}^T p_\theta(\hat{y}_{t-1} | \hat{y}_t, \tilde{x}). \quad (7)$$

The sampled lower-frequency subband at the last timestamp is denoted as \hat{y}_0 , which is regarded as the retouching-free lower frequency of X . To accelerate inference, we take DDIM [34] as the sampling scheme.

High-Frequency Branch and IDWT. We further introduce the HFCAM module to adjust the high-frequency sub-bands for refined retouching removal. The famous visual patch-wise Cross-Attention [13] mechanism is efficient in providing large receptive field only with several layers' effort. Thus, we propose to patchify \hat{y}_0 , x_H , and apply a lightweight one-layer CA and a three-layered convolution network to modify the higher-frequency bands.

$$\hat{y}_H = x_H + \text{Conv}(\text{CA}(\hat{y}_0, x_H)). \quad (8)$$

Finally, we combine the restored frequencies via the inverse wavelet transformation, i.e., $\hat{Y}_M = \text{IDWT}(\hat{y}_0, \hat{y}_H)$.

3.4. Objective Functions and Training Details

To train the proposed method, we first separately train the router and the Degree Estimators ahead of the rest modules for restoration. Next, these models are fixed and we train the low- and high-frequency branches within WaveFRR. Once the experts are prepared, we train the Combine Module that gets the intermediate images and produces the final result. Note that while we train MoFRR based on hybrid face images that have come through 0-4 face retouching operations, we only train WaveFRR models using single-operated images since they only remove their selected type, and train the shared expert on subset B for capturing global retouching patterns.

The objective functions of the proposed method include three categories, namely, frequency loss \mathcal{L}_{freq} , spacial loss \mathcal{L}_{space} , and classification loss \mathcal{L}_{class} . The total loss for the WaveFRR is given by $\mathcal{L} = \mathcal{L}_{freq} + \mathcal{L}_{space} + \mathcal{L}_{class}$.

The frequency loss \mathcal{L}_{freq} includes three parts, namely, IDEM loss \mathcal{L}_{IDEM} , diffusion loss \mathcal{L}_{simple} , and high-frequency loss \mathcal{L}_{high} , i.e., $\mathcal{L}_{freq} = \mathcal{L}_{IDEM} + \mathcal{L}_{simple} + \mathcal{L}_{high}$. The IDEM loss uses the ℓ_2 loss to quantify the difference between the pixel-wise distortion and the true residual of the low-frequency sub-band, which is given by $\mathcal{L}_{IDEM} = \| \hat{R} -$

$(y_{LL} - x_{LL}) \|^2$, where y_{LL} is the LL sub-band of the ground truth image. The training objective for diffusion models is $\mathcal{L}_{simple} = \mathbb{E}_{y_0, t, \epsilon_t \sim \mathcal{N}(0, \mathbf{I})} [\| \epsilon_t - \epsilon_\theta(y_t, \tilde{x}, t) \|^2]$. We utilize the ℓ_2 loss and Total Variation (TV) [5] loss to optimize the HFCAM, defined as: $\mathcal{L}_{high} = \lambda_1 \| \hat{y}_H - y_H \|^2 + \lambda_2 \text{TV}(\hat{y}_H)$, where λ_1 and λ_2 are set as 0.1 and 0.01, respectively.

The spatial loss \mathcal{L}_{space} is to maximize the PSNR between restored and original image. We use a hybrid ℓ_1 & SSIM [38] loss, defined as $\mathcal{L}_{hyb}(\cdot, \hat{\cdot}) = \| \cdot - \hat{\cdot} \|_1 + (1 - \text{SSIM}(\cdot, \hat{\cdot}))$, where the symbol \cdot denotes the ground-truth value of the symbol $\hat{\cdot}$, to minimize the difference between the restored face image and the ground truth image for both WaveFRR and Combine Module, i.e., $\mathcal{L}_{space} = \mathcal{L}_{hyb}(Y, \hat{Y}) + \sum_M \mathcal{L}_{hyb}(Y, \hat{Y}_M)$.

As for the router and Degree Estimator, we respectively train them based on cross-entropy loss \mathcal{L}_{class} , i.e., $\mathcal{L}_{class} = -\sum_{i=1}^5 c_i \log(\hat{c}_i)$, where c and \hat{c} respectively denote the true label and the predicted probability for class i .

4. Experiments

4.1. Experimental Setups

Implementation Details We use Adam [25] as the optimizer with a learning rate of 1×10^{-4} . The WaveFRR model is trained till convergence for about 2×10^5 iterations on eight NVIDIA A100 GPUs with a default batch size of 16, and the overall pipeline is trained for 50 epoch to converge. Input and output images have a resolution of 1024×1024 . More details can be found in the supplement.

Test Cases and Baseline Preparations. We employ two test cases, namely, intra-API tests, and cross-API tests. For intra-API tests, we train and test the models with a mixture of data from all APIs. For cross-API tests, we train our MoFRR and all comparing models on a mixture of subsets except Portrait and test on the Portrait subset, additionally validating the generalizability of models. In both test cases, the according dataset material is divided evenly into 8:1:1 as training/evaluation/test set.

There are no FRR methods available in the literature. We select a list of representative work from several relevant fields as baselines, including reference-free makeup removal method PairedCycleGAN [6] (PCGAN for short), as well as image restoration methods containing Pix2pix [19], Restormer [43], DR2 [40], and ResDiff [33]. The comparison methods cover a variety of mainstream architectures, including GAN, Transformer, and diffusion models. For fair comparison, we retrain them on our RetouchingFFHQ++ dataset for fair comparisons.

4.2. Intra-API Performance

In Fig. 4, we provide a visual comparison with other methods for FRR on the RetouchingFFHQ++ dataset. The images are randomly chosen from the subset that contains

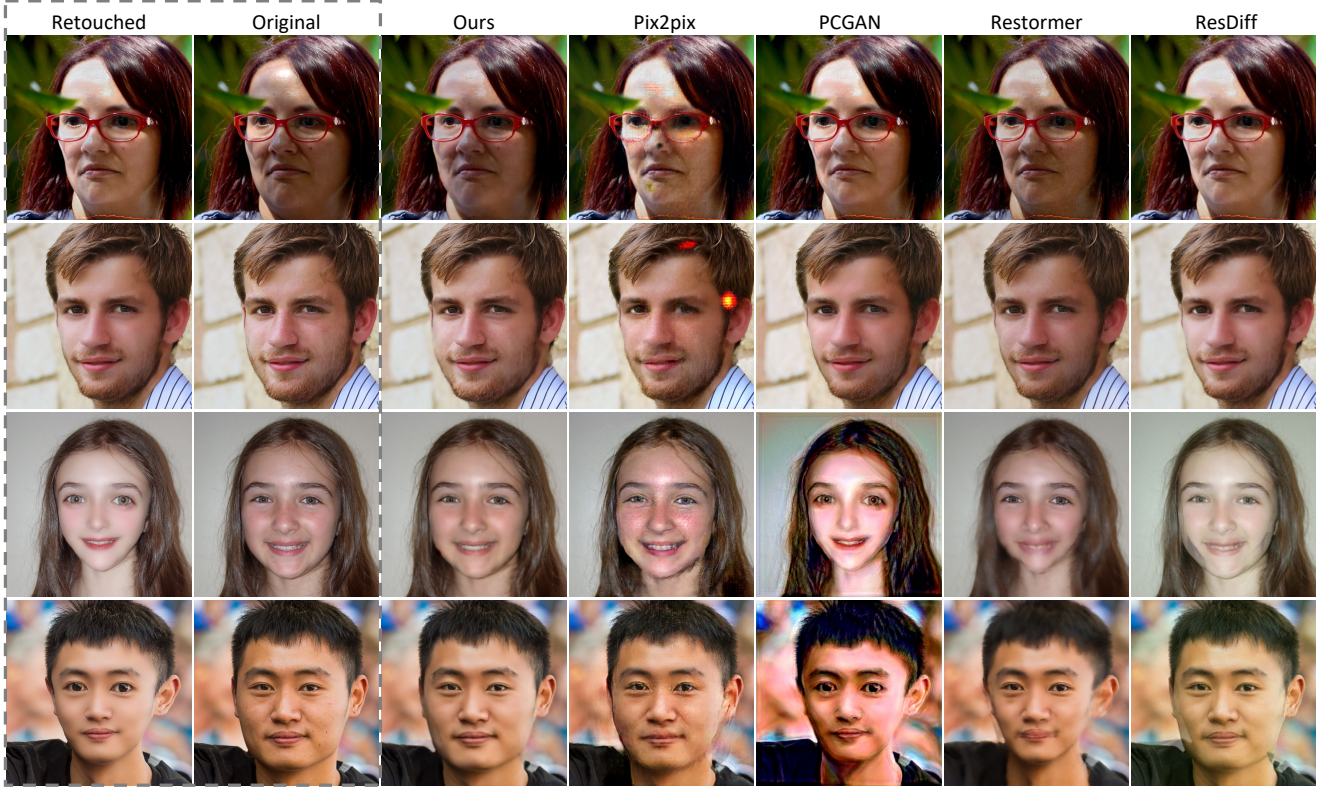


Figure 4. Qualitative comparison of the restored faces using different methods. The images are retouched through multiple operations, and our method achieves higher restoration quality with fewer artifacts. Zoom in for best view.

| Methods | White | | Smooth | | Face | | Eye | | Multi | |
|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| Input | 29.14 | 0.940 | 35.59 | 0.934 | 29.55 | 0.903 | 35.82 | 0.952 | 28.03 | 0.887 |
| Pix2pix (CVPR17) | 27.72 | 0.883 | 28.55 | 0.858 | 27.34 | 0.856 | 28.41 | 0.870 | 28.73 | 0.897 |
| PCGAN (CVPR18) | 22.47 | 0.684 | 23.53 | 0.639 | 23.04 | 0.627 | 23.63 | 0.634 | 22.17 | 0.581 |
| Restormer (CVPR22) | 29.89 | 0.918 | 37.59 | 0.943 | 29.83 | 0.903 | 33.75 | 0.942 | 28.19 | 0.888 |
| DR2 (CVPR23) | 24.33 | 0.723 | 25.97 | 0.719 | 24.39 | 0.712 | 25.71 | 0.719 | 23.87 | 0.711 |
| ResDiff (AAAI24) | 26.21 | 0.891 | 35.94 | 0.935 | 29.86 | 0.897 | 35.84 | 0.937 | 28.98 | 0.889 |
| MoFRR | 33.11 | 0.949 | 38.06 | 0.943 | 31.26 | 0.913 | 38.05 | 0.958 | 34.47 | 0.959 |

Table 2. Quantitative comparison of face retouching restoration. All compared models are retrained on our RetouchingFFHQ++ dataset for fair comparison. For each column, all the models are trained and tested on the same single-operated images, or all images for “Multi”.

face images retouched by multiple operations. As can be seen in the figures, methods like Pix2pix [19] and ResDiff [33] can recover the facial structure to some extent, but fail to generate high-quality face images without artifacts. Restormer [43] and PCGAN [6] perform better in restoring facial skin color, but struggle with structural restoration. In contrast, it can be observed that our MoFRR restores the original images more faithfully with fewer artifacts, particularly around the face and eyes. Besides, we also interestingly find that the baseline methods that apply the diffusion network generally provides better results compared to

the non-diffusion ones. The reason might be that diffusion models also accept the “divide and conquer” methodology in nature, and therefore the image qualities are generally better, as aligned with many recent findings. However, the fidelity of their results are noticeably worse compared to ours, especially from the view that ResDiff in many cases would apply wrong type of retouching.

Table 2 shows the quantitative comparison of FRR w.r.t. overall PSNR and SSIM. To better investigate the model performance against different retouching types as well as combined retouching, in the leading four groups we report

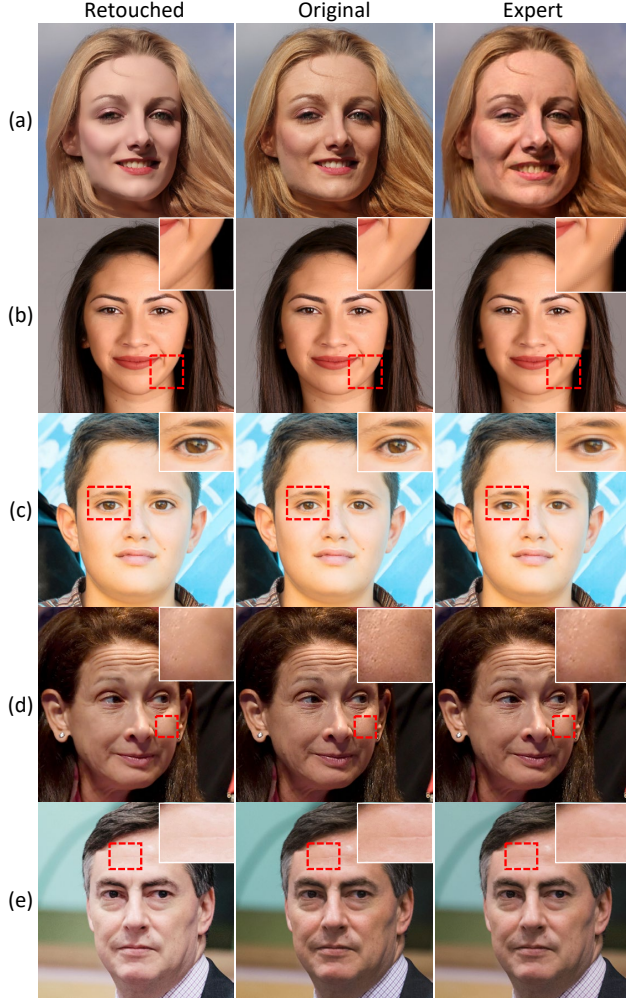


Figure 5. Visualization of images generated by five specific experts. (a) shared expert, (b) face lifting expert, (c) eye enlarging expert, (d) smoothing expert, and (e) whitening expert.

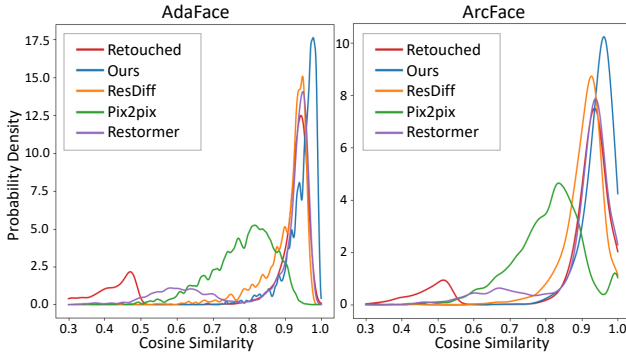


Figure 6. Probability density distribution of facial feature similarity among retouched/restored and original face images.

the performance on images that only has the specific retouching, and in the last group (named “Multi”), the per-

| Methods | Single | | Multi | |
|-----------|--------------|--------------|--------------|--------------|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| Input | 34.45 | 0.956 | 26.41 | 0.879 |
| Pix2pix | 28.18 | 0.881 | 26.68 | 0.886 |
| PCGAN | 23.18 | 0.626 | 21.30 | 0.569 |
| Restormer | 34.83 | 0.954 | 26.50 | 0.882 |
| DR2 | 24.72 | 0.709 | 22.63 | 0.676 |
| ResDiff | 34.93 | 0.961 | 26.79 | 0.885 |
| MoFRR | 36.65 | 0.971 | 31.28 | 0.938 |

Table 3. Cross-API performance. Models trained on all other datasets are tested on the Portrait subset to evaluate their generalization ability.

| Methods | Single | | Multi | |
|-------------------|--------------|--------------|--------------|--------------|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| Input | 32.73 | 0.934 | 28.03 | 0.887 |
| WaveFRR | 35.18 | 0.939 | 31.09 | 0.904 |
| w/o IDEM | 34.34 | 0.936 | 32.94 | 0.915 |
| w/o Degree | 34.43 | 0.937 | 32.95 | 0.921 |
| w/o HFCAM | 34.69 | 0.937 | 33.04 | 0.926 |
| w/o shared expert | 35.37 | 0.941 | 32.42 | 0.913 |
| w/o router | 35.63 | 0.942 | 34.19 | 0.947 |
| MoFRR | 36.12 | 0.943 | 34.47 | 0.959 |

Table 4. Ablation study over different model modules w.r.t. the ground truth. We fine-tune the models in partial settings till convergence. Test: intra-API.

formance on images with two or more retouching operations. From the results, our method outperforms other methods on the four single-operated subdatasets. Although PCGAN [6] exhibits excellent performance on the makeup transfer task, it struggles on the novel task of FRR with PSNR of around 23 dB. DR2 [40] performs well on face restoration but struggles with FRR due to the addition of excessive detail noise, resulting in higher image quality but lower fidelity. Other baselines have their own advantages in different subsets but do not perform as well and as comprehensively as MoFRR, particularly on the whitening subset where MoFRR achieves 3.97 dB improvement in PSNR. Particularly, our proposed MoFRR demonstrates superior performance compared to all baselines on the multi-operated subset. The distortion effects of multi-operated images are more pronounced. MoFRR effectively restores retouched face images, which surpasses the second-best method by 5.49 dB, achieving the highest improvement of 6.44 dB in PSNR. Multi-operated images better represent actual application scenarios, indicating that our method has significant advantages in practical use. The significant improvement observed in the multi-operated subset further validates effectiveness of our MoE framework.

Figure 5 provides the visual results generated by specific experts. To facilitate a clearer visual observation, we pro-



Figure 7. Qualitative comparison of the ablation studies, which proves the necessity of employing shared expert like DeepSeek and the IDEM design.

vide a sub-figure in each row as a zoom-in on some region-of-interest areas. It can be observed that four operation-specific experts perform well in restoring the targeted retouching operation, where we indeed see that the corresponding retouching is successfully reverted while other types of retouching operations that the expert is not specialized in are kept unchanged. We also see that the shared expert is able to do multiple types of retouching removal in parallel, but the generated result is not close to the ground-truth images. This finding is generally in line with our hypothesis that though one expert is not enough yet could steadily contribute to the overall quality improvement via seeing more retouched samples compared to the specialized experts.

Impact on Face Recognition. Figure 6 illustrates the probability density distributions of facial feature similarity between retouched/restored and original images. Each curve represents the cosine similarity distribution of a specific method’s results compared to the original faces, where a value closer to 1.0 indicates near-perfect feature alignment. The blue curve of MoFRR exhibits a sharp peak near 1.0 across both AdaFace [24] and ArcFace [12], which means the faces recovered by using our proposed method are more similar to the original faces at the feature level. The narrow width of the blue curve’s peak further highlights its high confidence in restoring identity features. While existing methods ResDiff [33] and Restormer [43] improve over retouched images, their curves show flatter distribution and lower peak values than ours. The curve of Pix2pix [19] has the lowest peak value and the flattest distribution, implying feature misalignment which may be caused by the heavy artifacts. Besides, the long left tail of the red curve indicates some retouched images are extremely processed with severe identity degradation. All methods can restore the extreme cases to some extent, while our MoFRR suppress the left-tail density more effectively.

4.3. Cross-API Performance

In a real-world scenario, training data and testing data are likely to come from different sources, so cross-API performance analysis helps us probe into the estimated performance of models in real-world conditions. From the results shown in Tab. 3, we observe that while there is a modest degradation in overall PSNR, MoFRR consistently delivers commendable restoration outcomes. Compared to the baselines, which achieve poor performance on both single- and multi-operated images, MoFRR demonstrates a notable advantage with improvements of 1.2 dB and 3.87 dB. This substantial performance gain, particularly for multi-operation images, can be attributed to the specialized expertise of multiple expert networks within MoFRR.

4.4. Ablation Study

Architectural analysis. The results are shown in Tab. 4. It can be observed that removing any one of modules in WaveFRR degrades performance compared to the full model. Both the Degree Estimator and IDEM handle distortion predictions, with the former focusing on global predictions and the latter on pixel-level predictions. The removal of IDEM results in significant performance degradation due to its refined iterative predictions. The combination of the two modules makes the prediction more effective by indicating the location and extent of retouching, as well as HF-CAM demonstrates unique effectiveness in artifact removal. The performance degradation of removing the shared expert mainly occurs in multi-operated cases due to its effect on global structural restoration. In contrast, removing the router has minimal impact on performance, as the Degree Estimator can reduce the impact of false activation while the router contributes positively to enhancing efficiency.

Comparison with Single-Expert Mode. We also test the performance where we use one single WaveFRR as the whole MoFRR. Results show that it performs well on single-operated images, but limited when restoring all operations simultaneously. The multi-expert mechanism also offers flexibility that once novel types of face retouching appear, a plug-in of corresponding expert networks would mitigate the huge cost of fully retraining the model.

5. Conclusion

We propose face retouching restoration, a novel task of recovering the original faces from retouched ones. We propose a Mixture of Diffusion Models for FRR, including specialized WaveFRR models that uses a dual-branch structure, i.e., a low-frequency recovery using DDIM, and a high-frequency refinement guided by the restored lower sub-bands. Extensive experiments on our dedicated RetouchingFFHQ++ dataset prove the effectiveness of our model on FRR against varied combination of retouching operations.

Acknowledgment This work was supported by Ant Group through CCF-Ant Research Fund.

References

- [1] Anthropic. Portraitpro 24. <https://www.anthropics.com/portraitpro/>, 2024. 3
- [2] Khadijah Ateq, Mohammed Alhajji, and Noara Alhusseini. The association between use of social media and the development of body dysmorphic disorder and attitudes toward cosmetic surgeries: a national survey. *Frontiers in Public Health*, 12:1324092, 2024. 1
- [3] Aparna Bharati, Richa Singh, Mayank Vatsa, and Kevin W Bowyer. Detecting facial retouching using supervised deep learning. *IEEE Transactions on Information Forensics and Security*, 11(9):1903–1913, 2016. 1
- [4] Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*, 2024. 2
- [5] Stanley H Chan, Ramsin Khoshabeh, Kristofor B Gibson, Philip E Gill, and Truong Q Nguyen. An augmented lagrangian method for total variation video restoration. *IEEE Transactions on Image Processing*, 20(11):3097–3111, 2011. 5
- [6] Huiwen Chang, Jingwan Lu, Fisher Yu, and Adam Finkelstein. Pairedcyclegan: Asymmetric style transfer for applying and removing makeup. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 40–48, 2018. 2, 5, 6, 7
- [7] Cunjian Chen, Antitza Dantcheva, and Arun Ross. Automatic facial makeup detection with application in face recognition. In *2013 international conference on biometrics (ICB)*, pages 1–8. IEEE, 2013. 1
- [8] Ying-Cong Chen, Xiaoyong Shen, and Jiaya Jia. Makeupgo: Blind reversion of portrait edit. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4501–4509, 2017. 2
- [9] United States Congress. Truth in advertising act of 2014, 2014. <https://www.congress.gov/bills/113th-congress/house-bill/4341>. 1
- [10] Yimian Dai, Fabian Gieseke, Stefan Oehmcke, Yiquan Wu, and Kobus Barnard. Attentional feature fusion. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3560–3569, 2021. 4
- [11] Antitza Dantcheva, Cunjian Chen, and Arun Ross. Can facial cosmetics affect the matching accuracy of face recognition systems? In *2012 IEEE Fifth international conference on biometrics: theory, applications and systems (BTAS)*, pages 391–398. IEEE, 2012. 1
- [12] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. 3, 8
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 4, 5
- [14] Forbrukertilsynet. The marketing control act. <https://www.forbrukertilsynet.no/english/the-marketing-control-act>, 2021. 1
- [15] Qiao Gu, Guanzhi Wang, Mang Tik Chiu, Yu-Wing Tai, and Chi-Keung Tang. Ladrn: Local adversarial disentangling network for facial makeup and de-makeup. In *Proceedings of the IEEE/CVF International conference on computer vision*, pages 10481–10490, 2019. 2
- [16] Yuchao Gu, Xintao Wang, Liangbin Xie, Chao Dong, Gen Li, Ying Shan, and Ming-Ming Cheng. Vqfr: Blind face restoration with vector-quantized dictionary and parallel decoder. In *European Conference on Computer Vision*, pages 126–143. Springer, 2022. 2
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 4
- [19] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 5, 6, 8
- [20] Anubhav Jain, Puspita Majumdar, Richa Singh, and Mayank Vatsa. Detecting gans and retouching based digital alterations via dad-hcnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 672–673, 2020. 1
- [21] Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 6(2): 181–214, 1994. 2
- [22] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 3
- [23] Mohammed Adib Khan, Morteza Noferesti, and Naser Ezzati-Jivan. Pasd: A performance analysis approach through the statistical debugging of kernel events. In *2023 IEEE 23rd International Working Conference on Source Code Analysis and Manipulation (SCAM)*, pages 151–161. IEEE, 2023. 1
- [24] Minchul Kim, Anil K Jain, and Xiaoming Liu. Adaface: Quality adaptive margin for face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18750–18759, 2022. 3, 8
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015. 5
- [26] Xiaoming Li, Chaofeng Chen, Shangchen Zhou, Xianhui Lin, Wangmeng Zuo, and Lei Zhang. Blind face restoration via deep multi-scale component dictionaries. In *European conference on computer vision*, pages 399–415. Springer, 2020. 2

- [27] Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Yu Qiao, Wanli Ouyang, and Chao Dong. Diffbir: Toward blind image restoration with generative diffusion prior. In *European conference on computer vision*, pages 430–448. Springer, 2024. 1
- [28] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 2
- [29] Si Liu, Wentao Jiang, Chen Gao, Ran He, Jiashi Feng, Bo Li, and Shuicheng Yan. Psgan++: Robust detail-preserving makeup transfer and removal. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8538–8551, 2021. 2
- [30] IBT News. Supermodels without photoshop: Israel’s ‘photoshop law’ puts focus on digitally altered images, 2015. <http://goo.gl/3PvaEf>. 1
- [31] Christian Rathgeb, Angelika Botaljov, Fabian Stockhardt, Sergey Isadskiy, Luca Debiasi, Andreas Uhl, and Christoph Busch. Prnu-based detection of facial retouching. *IET Biometrics*, 9(4):154–164, 2020. 1
- [32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 4
- [33] Shuyao Shang, Zhengyang Shan, Guangxing Liu, LunQian Wang, XingHua Wang, Zekai Zhang, and Jinglin Zhang. Resdiff: Combining cnn and diffusion model for image super-resolution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8975–8983, 2024. 2, 5, 6, 8
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 4, 5
- [35] Zhaoyang Sun, Yaxiong Chen, and Shengwu Xiong. Ssat: A symmetric semantic-aware transformer network for makeup transfer and removal. In *Proceedings of the AAAI Conference on artificial intelligence*, pages 2325–2334, 2022. 2
- [36] Zhaoyang Sun, Shengwu Xiong, Yaxiong Chen, and Yi Rong. Content-style decoupling for unsupervised makeup transfer without generating pseudo ground truth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7601–7610, 2024. 2
- [37] Shuyang Wang and Yun Fu. Face behind makeup. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016. 1, 2
- [38] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 3, 5
- [39] Zhouxia Wang, Jiawei Zhang, Tianshui Chen, Wenping Wang, and Ping Luo. Restoreformer++: Towards real-world blind face restoration from undegraded key-value pairs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2
- [40] Zhixin Wang, Ziyang Zhang, Xiaoyun Zhang, Huangjie Zheng, Mingyuan Zhou, Ya Zhang, and Yanfeng Wang. Dr2: Diffusion-based robust degradation remover for blind face restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1704–1713, 2023. 2, 5, 7
- [41] Rongyuan Wu, Lingchen Sun, Zhiyuan Ma, and Lei Zhang. One-step effective diffusion network for real-world image super-resolution. *Advances in Neural Information Processing Systems*, 37:92529–92553, 2024. 1
- [42] Qichao Ying, Jiaxin Liu, Sheng Li, Haisheng Xu, Zhenxing Qian, and Xinpeng Zhang. Retouchingffhq: A large-scale dataset for fine-grained face retouching detection. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 737–746, 2023. 1, 2, 3
- [43] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5728–5739, 2022. 2, 5, 6, 8