

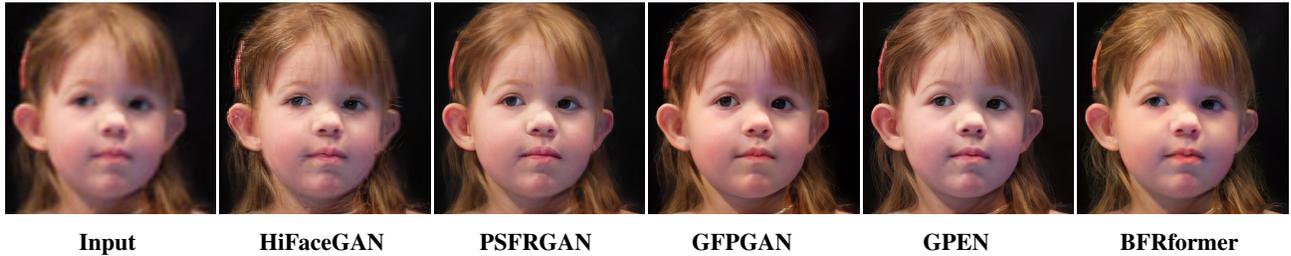
BFRFORMER: TRANSFORMER-BASED GENERATOR FOR REAL-WORLD BLIND FACE RESTORATION

Guojing Ge^{1,2,*}, Qi Song^{3,*}, Guibo Zhu^{1,2,5,6,†}, Yuting Zhang⁴, Jinglu Chen¹,
Miao Xin¹, Ming Tang¹, and Jinqiao Wang^{1,2,6}

¹Institute of Automation, Chinese Academy of Sciences ²Wuhan AI Research

³Hong Kong Baptist University ⁴China Telecom Corporation Ltd

⁵Shanghai Artificial Intelligence Laboratory ⁶University of Chinese Academy of Sciences



Input

HiFaceGAN

PSFRGAN

GFPGAN

GPEN

BFRformer

ABSTRACT

Blind face restoration is a challenging task due to the unknown and complex degradation. Although face prior-based methods and reference-based methods have recently demonstrated high-quality results, the restored images tend to contain over-smoothed results and lose identity-preserved details when the degradation is severe. It is observed that this is attributed to short-range dependencies, the intrinsic limitation of convolutional neural networks. To model long-range dependencies, we propose a Transformer-based blind face restoration method, named BFRFormer, to reconstruct images with more identity-preserved details in an end-to-end manner. In BFRFormer, to remove blocking artifacts, the wavelet discriminator and aggregated attention module are developed, and spectral normalization and balanced consistency regulation are adaptively applied to address the training instability and over-fitting problem, respectively. Extensive experiments show that our method outperforms state-of-the-art methods on a synthetic dataset and four real-world datasets. The source code, Casia-Test dataset, and pre-trained models is released at <https://github.com/s8Znk/BFRFormer>.

Index Terms— Blind face restoration, transformer, wavelet discriminator

* Equal contribution. † Corresponding author.

This work was supported in part by the National Key RD Program of China (No.2022ZD0160601), National Natural Science Foundation of China (No.62276260,62076235,61906195), Beijing Municipal Science and Technology Project (Z231100007423004), sponsored by Zhejiang Lab (No.2021KH0AB07).

1. INTRODUCTION

Blind face restoration aims at recovering high-quality face images from low-quality counterparts that suffer from unknown degradation. There are many applications such as old photo restoration and low-quality face recognition. The authentic, real-world low-quality images usually contain complex and diverse distributions that are impractical to imitate. For example, in the old photo restoration tasks, spatial uniformity and color fading are the major difficulties.

The current mainstream approaches are geometric prior methods [1][2][3][4][5][6][7], reference-based methods [8][9][10][11][12][13][14][15] and GAN prior-based methods [16][17]. Geometric priors can be facial landmarks [2][3], facial parsing maps [1][4], or facial component heatmaps [5]. However, those priors are estimated from degraded images and cannot offer accurate inputs when the degradation is severe. Besides, the geometric structures cannot provide sufficient information to recover facial details. The limitation of reference-based methods is that the learned codebook usually contains general characteristics of the training data set, losing identity-preserved details. GAN prior-based methods, instead of the previous facial priors such as geometric priors, have recently demonstrated high-quality results in this task. GFP-GAN [17] and GPEN [16] extract local detail information and global identity information from low-quality input and then leverage the pre-trained GAN as a generator, achieving a good balance between visual quality and fidelity. Although existing methods have recently demonstrated high-quality results for blind face restoration, the restoration images tend to generate

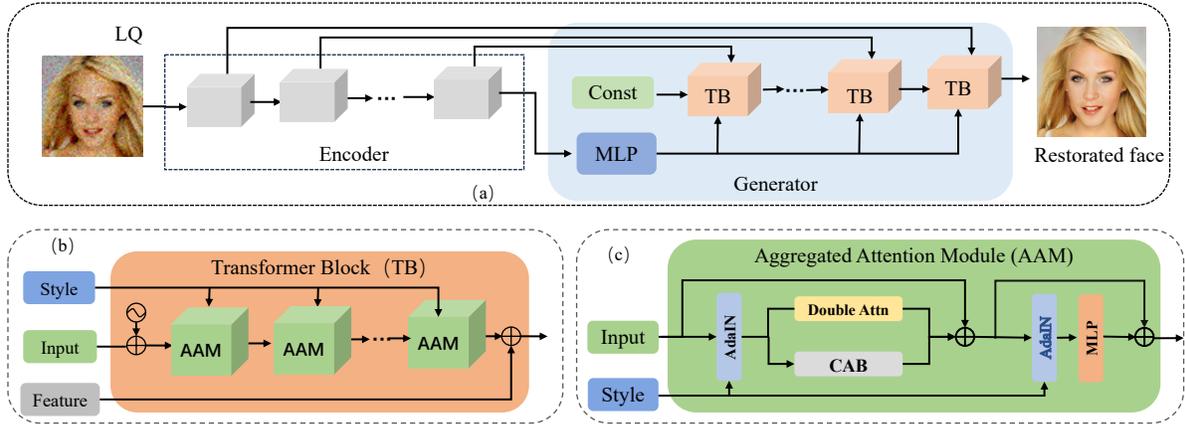


Fig. 1. Overview of BFRFormer framework. (a) It follows an encoder and generator architecture; (b) TB is a transform-based block of the generator; (c) Aggregated Attention Module (AAM) which combining channel attention extracting global information with double-attention extracting local information, to activate more input pixels for face restoration.

over-smoothed results and lose the identity-preserved details especially when the degradation is severe.

To address the above problems, we develop a Transformer-based blind face restoration method, BFRFormer, to take advantage of both GAN prior-based methods and Transformer, yet without geometric prior or reference prior. Compared with other embedded CNN in the GAN prior-based methods, the Transformer-based generator can effectively model long-range dependency, improving the over-smoothed problem and the identity-preserved details. In BFRFormer, to remove blocking artifacts and generate realistic face details, the wavelet discriminator and an aggregated attention module are developed, respectively. The wavelet discriminator suppresses the artifacts by the periodic artifact pattern which can be easily distinguished in the spectral domain. The aggregated attention module (AAM), which contains channel attention and double attention, aims to enlarge the effective receptive field of the low-quality input. To address training instability and over-fitting problems, spectral normalization (SN) and balanced Consistency Regulation (bCR) are adaptively applied. Existing test datasets, such as the CelebChild-Test containing 180 child faces of celebrities and the WebPhoto-Test consisting of 407 real-life faces, are not diverse or fair enough to represent real-world low-quality images. In order to fairly evaluate the generalization of the existing blind face restoration methods, we construct a new real-world test dataset. Additionally, we experimentally show that perception loss, facial ROI discriminator, and data augmentation can improve performance well.

In summary, this paper makes the following contributions:

- We develop a Transformer-based training method embedded in the GAN prior-based framework for blind face restoration in an end-to-end manner compared with Coderformer and VQFR.
- We propose a novel aggregated attention module, com-

binning channel attention extracting global information and double-attention extracting local information, to activate more effective pixels of the low-quality input.

- We construct one large test benchmark considering wider variation than the existing test datasets, including different ethnicities, different ages, different occlusion, and more than ten thousand persons for diversity.
- Extensive experiments show that our method outperforms state-of-the-art methods on a synthetic dataset and four real-world datasets.

2. METHOD

In this section, we will describe the overall pipeline, the details of the encoder, generator, and the loss function.

2.1. Overall Pipeline

Given a low-quality face image, the purpose of blind face restoration is to reconstruct the high-quality image with realistic facial details while reserving identity information of the low-quality face image. The overall framework of BFRFormer is depicted in Figure 1. It comprises three modules: encoder, generator, and loss function.

2.2. Encoder network

The encoder network is a simple convolution network, compared with RRDBNet [18] used in GLEN [19] and VQFR [14]. The encoder extracts the multi-resolution features containing shallow features and deep features of the input image.

2.3. Generator network

As shown in Figure 1, the generator adapts a style-based architecture [20][21][22] that takes shallow features and style vectors as inputs where style vectors are injected into each Transformer Block (TB). The multi-level shallow features can

add or concatenate to the output TB. Furthermore, the shallow features are concatenated rather than added to the Aggregated Attention Module (AAM) output.

2.3.1. Transformer Block (TB)

The Transformer Block (TB) is the basic building block of the generator. The input features of TB are the output of the previous TB, shallow features achieved from the encoder, and style vector achieved from MLP. The input feature X is first processed by the AAM model with the style vector, and then the features obtained from the encoder are concatenated with the feature from AAM. Specifically, we have

$$\begin{aligned} X &= \text{AAM}(X, \text{Style}), & (1) \\ X &= \text{Concat}(\text{Feature}, X), & (2) \\ Y &= \text{Upsample}(X), & (3) \end{aligned}$$

where the concatenated are further Up-sampled.

2.3.2. Aggregated Attention Module (AAM)

A channel attention block in parallel with the double-attention followed by AdaIN [23]. To avoid the possible conflict of channel attention block and double attention on optimization and visual representation, a small constant α is multiplied by the output of CAB. The whole process of Aggregated Attention Module (AAM) is computed as:

$$\begin{aligned} X_N &= \text{AdaIN}(X_N, \text{Style}), & (4) \\ X_M &= \text{DoubleA}(X_N) + \alpha \text{CAB}(X_N) + X, & (5) \\ X &= \text{MLP}(\text{AdaIN}(X_M, \text{Style})) + X_M, & (6) \end{aligned}$$

where X_N and X_M denote the intermediate features. Y represents the output of AAM. MLP denotes a multi-layer perception.

2.3.3. Channel Attention Block

A CAB consists of two standard convolution layers with a GELU activation function between them and a channel attention (CA) module. The details is shown in Figure 1.

2.3.4. Double Attention

Double attention [21] aims to achieve an enlarged receptive field and allows a single Transformer block to simultaneously achieve the context of the local and shifted windows.

2.4. Loss Functions

In the Transformer-based GAN training process, the loss function is extremely important for the stability of training. Our losses involve several aspects, including pixel-level, component-level, and image-level. In the training process of baseline, perceptual loss, and component-level loss are not used. The detailed discussions are as following.

Pixel-level loss. L1 loss and perceptual loss are used in this paper. L2 loss is verified that it can not deal with high-frequency content in the image, resulting in overly smooth

images [1]. The perceptual loss further improves the image quality as shown in Table 2.

Component-level loss. Following [17], we only focus on regions $r \in \{\text{eyes, mouth}\}$. Specifically, the loss functions are formulated with discriminative loss

$$\mathcal{L}_{roi} = \sum_{r \in ROI} E_{\hat{y}_r} [\log(1 - D_r(\hat{y}_r))]. \quad (7)$$

Image-level loss Adversarial loss and identity preserving loss are used as image-level loss. We employ the wavelet discriminator in SWAGAN [24] to solve the blocking artifacts.

$$\mathcal{L}_{adv} = -\mathbb{E}_x [\text{softplus}(D_g(x))] \quad (8)$$

3. EXPERIMENTS

3.1. Datasets and Evaluation Metric

To evaluate our model, we use one synthetic dataset [25] and four different real-world datasets to compare the proposed method with other blind face restoration methods. CelebA-Test is the synthetic dataset with 5,000 CelebA-HQ images. The generation way is the same as that during training.

Casia-Test. CASIA dataset [26] contains low-quality images in the wild, containing 10,575 subjects and 494,414 images. We randomly select three photos for each subject and use a quality assessment algorithm to compute the score of each face image. The lowest two images are selected as our target images. Among these, too many large poses are selected, and four experts are selected to balance the dataset to make it more impartial. At last, we get 16,683 images containing 10,575 identities to guarantee diversity with different races, ages, extreme poses, and light conditions.

3.2. Ablation Study

To better understand the roles of different components of BFRFormer and the training strategy, in this section, we conduct an ablation study by introducing some variants of BFRFormer and comparing with their BFR performance.

Baseline. The BFRFormer-simple uses a basic Transformer as the generator, without using the CAB module, perception loss, and local ROI loss, only using the basic loss functions: Adv loss, ID loss, and L1 loss. As shown in Table 2, it can achieve high perceptual quality (low FID and NIQE) and produce surprisingly good results compared with GPEN [16]. Through attribution analysis, BFRFormer utilizes more information compared to GPEN.

Facial Component Discriminator. Facial Component discriminators [17], including right eyes, left eyes, and mouth, enhance the perceptually significant face components. Component discriminators with feature style loss could better capture the eye distribution and restore the plausible details [17]. The effectiveness is shown in Table 2 and Figure 2.

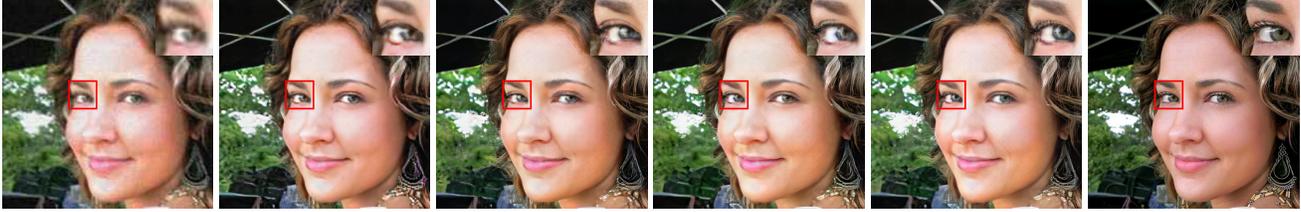


Fig. 2. Comparison of our variants of BFRFormer. (a) Low-quality input; (b) Perceptual Loss; (c) bCR; (d) CAB; (e) Facial Component; (f) Ground Truth.

Datasets Method	LFW-Test		CelebChild		WebPhoto		Casia-Test	
	FID↓	NIQE↓	FID ↓	NIQE ↓	FID↓	NIQE ↓	FID↓	NIQE ↓
PULSE	64.86	5.10	102.74	5.23	86.45	5.14	-	-
HiFaceGAN [27]	64.50	4.51	113.0	4.86	116.1	4.88	79.15	5.06
PSFRGAN [1]	51.89	5.09	107.40	4.80	88.45	5.58	60.13	4.90
GPEN [16]	54.85	3.90	106.90	4.08	80.70	4.34	58.78	4.30
GFPGAN [17]	49.96	3.88	111.78	4.35	87.35	4.14	61.35	4.35
CodeFormer [28]	52.02	4.01	116.23	4.98	83.19	4.70	56.16	4.01
VQFR [14]	50.64	3.59	105.18	3.94	75.38	3.61	54.01	3.99
Ours-Simple	49.81	3.95	111.18	4.08	77.89	3.75	57.13	4.14
Ours	48.35	3.81	103.89	3.93	79.53	3.65	55.56	3.97

Table 1. Quantitative comparison on the *real-world* LFW-Test, CelebChild, WebPhoto and Casia-Test.

\mathcal{L}_{PER}	\mathcal{L}_{ROI}	DA	CAB	FID↓
×	×	×	×	27.94
✓	×	×	×	23.01
✓	✓	×	×	22.46
✓	✓	✓	×	22.86
✓	✓	✓	✓	20.51

Table 2. FID between the different variants of BFRFormer

Method	PSNR↑	FID↓	LPIPS ↓
DFDNet [9]	21.10	59.09	0.45
HiFaceGAN [27]	20.17	66.09	0.38
PSFRGAN [1]	19.86	62.05	0.32
GFPGAN [17]	21.17	58.36	0.28
GPEN [16]	19.85	59.70	0.29
VQFR [14]	20.51	58.01	0.28
CodeFormer [28]	22.18	60.62	0.30
Ours	22.83	57.37	0.27

Table 3. Quantitative comparison (PNSR, FID, and LPIPS) of different BFR methods on Synthetic Datasets

Data Augmentation. Recent successes in Generative Adversarial Networks have affirmed the importance of using more data in GAN training. Augment data can improve the performance of the discriminator and generator. In this paper, the generated image and ground truth images are augmented by Flipping, Color, Translation, Cutout as in DiffAug. The effectiveness is shown in Table 2 and Figure 2.

Channel Attention Block. Local attention, nonetheless, sacrifices the ability to model long-range dependencies. To activate more effective pixels for restoration tasks, channel atten-

tion utilizing more global information is used to further improve the performance of the Transformer. The effectiveness is shown in Table 2 and Figure 2.

3.3. Comparison with State-of-the-art Methods

Synthetic Dataset. We compare BFRFormer with several state-of-the-art face blind restoration methods, and the quantitative results are shown in Table 3. Our BFRFormer achieves better PNSR compared to other methods. Our results also obtain lower FID and LPIPS. BFRFormer can generate realistic and high-fidelity face images with details such as hair, eyes, mouth, etc.

Real World low-quality Datasets. The final target of all methods is to restore authentic world LQ face images. To evaluate the generalization ability of different methods, we also compare the performance of BFRFormer on the LFW-Test, CelebChild-Test, WebPhoto-Test, and Casia-Test for evaluating the generalization of the proposed method. As shown in Table 1, compared to previous face restoration methods on CelebAHQ-Test, BFRFormer gets better performance in FID and NIQE.

4. CONCLUSION

In this paper, we propose a Transformer-based training method embedded in the GAN prior-based framework for blind face restoration. An aggregate attention module, composed of channel attention and double attention, is proposed to further improve generation quality. Extensive experiments on synthetic data and real-world low-quality images have demonstrated that BFRFormer can restore high-quality facial details while retaining the image background properly.

References

- [1] C. Chen, X. Li, L. Yang, X. Lin, L. Zhang, and K. K. Wong, "Progressive semantic-aware style transformation for blind face restoration," in *CVPR*, 2021.
- [2] Y. Chen, Y. Tai, X. Liu, C. Shen, and J. Yang, "Fsr-net: End-to-end learning face super-resolution with facial priors," in *CVPR*, 2017.
- [3] D. Kim, M. Kim, G. Kwon, and D. S. Kim, "Progressive face super-resolution via attention to facial landmark," in *arXiv preprint arXiv:1908.08239*, 2019.
- [4] Z. Shen, W. S. Lai, T. Xu, J. Kautz, and M. H. Yang, "Deep semantic face deblurring," in *CVPR*, 2018.
- [5] X. Yu, B. Fernando, B. Ghanem, F. Porikli, and R. Hartley, "Face super-resolution guided by facial component heatmaps," in *ECCV*, 2018.
- [6] Y. Xin, B. Fernando, R. Hartley, and F. Porikli, "Super-resolving very low-resolution face images with supplementary attributes," in *CVPR*, 2018.
- [7] S. Zhu, S. Liu, C. C. Loy, and X. Tang, "Deep cascaded bi-network for face hallucination," in *ECCV*, 2016.
- [8] B. Dogan, S. Gu, and R. Timofte, "Exemplar guided face image super-resolution without facial landmarks," in *CVPRW*, 2019.
- [9] X. Li, C. Chen, S. Zhou, X. Lin, W. Zuo, and L. Zhang, "Blind face restoration via deep multi-scale component dictionaries," in *ECCV*, 2020.
- [10] Z. Wang, Z. Zhang, X. Zhang, H. Zheng, M. Zhou, Y. Zhang, and Y. Wang, "Dr2: Diffusion-based robust degradation remover for blind face restoration," in *CVPR*, 2023.
- [11] X. Li, W. Li, D. Ren, H. Zhang, M. Wang, and W. Zuo, "Enhanced blind face restoration with multi-exemplar images and adaptive spatial feature fusion," in *CVPR*, 2020.
- [12] X. Li, M. Liu, Y. Ye, W. Zuo, L. Lin, and R. Yang, "Learning warped guidance for blind face restoration," in *ECCV*, 2018.
- [13] Z. Wang, J. Zhang, R. Chen, W. Wang, and P. Luo, "Restoreformer: High-quality blind face restoration from undegraded key-value pairs," in *CVPR*, 2022.
- [14] Y. Gu, X. Wang, L. Xie, C. Dong, G. Li, Y. Shan, and M. Cheng, "Vqfr: Blind face restoration with vector-quantized dictionary and parallel decoder," in *ECCV*, 2022.
- [15] X. Li, S. Zhang, S. Zhou, L. Zhang, and W. Zuo, "Learning dual memory dictionaries for blind face restoration," in *PAMI*, 2022.
- [16] T. Yang, P. Ren, X. Xie, and L. Zhang, "Gan prior embedded network for blind face restoration in the wild," in *CVPR*, 2021.
- [17] X. Wang, Y. Li, H. Zhang, and Y. Shan, "Towards real-world blind face restoration with generative facial prior," in *CVPR*, 2021.
- [18] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *ECCVW*, 2018.
- [19] K. Chan, X. Wang, X. Xu, J. Gu, and C.C.Loy, "Glean: Generative latent bank for large-factor image super-resolution," in *CVPR*, 2021.
- [20] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *CVPR*, 2019.
- [21] B. Zhang, S. Gu, B. Zhang, Bao J, D. Chen, F. Wen, Y. Wang, and B. Guo, "Styleswin: Transformer-based gan for high-resolution image generation," in *CVPR*, 2021.
- [22] S. Qi, W. Shi, G. Ge, and L. Chang, "Degradation conditioned gan for degradation generalization of face restoration models," in *ICIP*, 2023.
- [23] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017.
- [24] R. Gal, D. C. Hochberg, A. Bermano, and D. Cohen-Or, "Swagan: A style-based wavelet-driven generative model," *ACM Trans. Graph.*, 2021.
- [25] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," in *arXiv preprint arXiv:1710.10196*, 2017.
- [26] Y. Dong, L. Zhen, S. Liao, and S. Z. Li, "Learning face representation from scratch," *Computer Science*, 2014.
- [27] L. Yang, S. Wang, S. Ma, W. Gao, P. Liu, C. and Wang, and P. Ren, "Hifacegan: Face renovation via collaborative suppression and replenishment," in *ACM MM*, 2020.
- [28] S. Zhou, K. Chan, C. Li, and C. C. Loy, "Towards robust blind face restoration with codebook lookup transformer," in *NeurIPS*, 2022.