REVISITING VECTOR-QUANTIZATION FOR BLIND IMAGE RESTORATION

Anonymous authors

003 004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026 027 028

029

Paper under double-blind review

ABSTRACT

Vector-Quantization (VQ) generative models are widely used to learn a highquality (HQ) codebook and a decoder as powerful generative priors for blind image restoration (BIR). In this paper, we revisit the key VQ process in VQ-based BIR methods, and provide three close observations on the side effects of VO for code index prediction: 1) confining the representational capability of HQ codebook, 2) being error-prone on code index prediction, and 3) under-valuing the low-quality (LQ) feature for BIR. These observations motivate us to replace discrete VQ selection by continuous feature transformation from input LQ image to output HQ image with the HQ codebook. To this end, in this paper, we propose a new Self-in-Cross-Attention (SinCA) module to augment the HO codebook with the LQ feature of input LQ image and perform cross-attention between LQ feature and input-augmented codebook. In this way, our SinCA extends the representational capability of the HQ codebook and effectively leverages the self-expressiveness property of input LQ image. Experiments on four typical VQ-based BIR methods demonstrate that, by replacing the VQ process with transformers using our SinCA, they achieve better quantitative and qualitative performance on blind image superresolution and blind face restoration. The code will be publicly released.

1 INTRODUCTION

Blind image restoration (BIR) aims to recover high-quality (HQ) images from the corresponding low-quality (LQ) images affected by complex degradation (Wang et al., 2021c; 2023b; Zhang et al., 2021). This ill-posed problem has been addressed by many generative models (Wang et al., 2021b; Chen et al., 2022; Lin et al., 2023; Wang et al., 2023a), under the architectures of VAEs (Kingma & Welling, 2013; Rezende et al., 2014), GANs (Goodfellow et al., 2014; Karras et al., 2019), or diffusion models (Ho et al., 2020; Song et al., 2020). Recently, with successes in applications like DALL·E (Ramesh et al., 2021), Vector-Quantization (VQ) based discrete generative models like VQVAE (Van Den Oord et al., 2017) or VQGAN (Esser et al., 2021) have been increasingly adopted in many BIR methods (Chen et al., 2022; Zhou et al., 2022; Liu et al., 2023a; TSAI et al., 2024) as robust backbones against diverse image degradations.

040 Current VQ-based BIR methods typically follow a multi-stage training scheme. First, an encoder-041 decoder model and a discrete codebook are learned to reconstruct HQ images using VQGAN (Esser 042 et al., 2021). The encoder is then fine-tuned to restore LQ images, during which the VQ process 043 replaces each pixel-wise LQ feature vector with a selected code item from the HQ codebook. Though 044 allowing for useful information recovery from the HQ codebook, the VQ process also brings three notable side effects to VQ-based blind restoration methods. Firstly, the representational range of VQ process is confined to the finite set of HQ codebook items. Secondly, two main VQ strategies, 046 *i.e.*, nearest-neighbor feature matching (Fig. 1 (a)) and transformer-based prediction (Fig. 1 (c)), are 047 error-prone on selecting code items for BIR. Thirdly, the VQ process undervalues the essential role 048 of LQ features for blind image restoration, since it simply replaces LQ features by HQ code items. 049

Considering the side effects of VQ process mentioned above, a natural question raises: is it feasible
 to replace the vulnerable discrete VQ process by continuous transformation from LQ features to HQ
 ones with the HQ codebook? In this paper, we provide positive feedback to the above question by
 implementing continuous feature transformation via cross-attention between LQ feature and HQ
 codebook. Specifically, the cross-attention computes the attention map by using the LQ feature as the



widely adopted by BIR methods (Pan et al., 2020; 2021; Tao Yang & Zhang, 2021; Wang et al.,

2021b; 2022a). For example, GFPGAN (Wang et al., 2021b) and GPEN (Tao Yang & Zhang, 2021)

employed a pre-trained StyleGAN2 (Karras et al., 2020) in a U-shaped decoder network for face

106

108 image restoration. However, GANs are prone to generate unrealistic textures due to the inherent 109 difficulty on distinguishing similar patterns (Chen et al., 2022). Recently, diffusion generative 110 priors (Ho et al., 2020; Song et al., 2020) have also been exploited by many BIR methods (Wang 111 et al., 2023a; Lin et al., 2023; Yang et al., 2023). Despite the impressive progress, these methods are 112 usually not robust to severe image degradations (Zhou et al., 2022).

113 Benefited by the power of vector-quantization (VQ) models (Van Den Oord et al., 2017; Esser et al., 114 2021) in image generation, VQ-based BIR methods (Zhao et al., 2022; Chen et al., 2022; Zhou et al., 115 2022; Liu et al., 2023a; Wang et al., 2023b) have been developed to utilize discrete HQ codebook 116 priors. In this work, we examine the side effects of discrete VQ process on feature matching and 117 replace it with our Self-in-Cross-Attention (SinCA) for continuous feature learning.

118 VQ-based Generative Models learn discrete codebook priors of images in latent space. This idea 119 is first introduced in VQVAE (Van Den Oord et al., 2017) and further enhanced by VQGAN (Esser 120 et al., 2021) with better perceptual quality induced by learning codebook priors and autoregressive 121 transformer (Vaswani et al., 2017). Built upon VQVAE and VQGAN, recent VQ-based image 122 generation methods (Cao et al., 2023; Chang et al., 2022; Lee et al., 2022; Yu et al., 2022; Zhang 123 et al., 2023; Zheng et al., 2022) primarily focus on improving the quantization and token generation 124 processes. For example, MaskGIT (Chang et al., 2022) utilized a bidirectional transformer (Kenton & 125 Toutanova, 2019) to simultaneously predict all the image tokens.

126 VO-based generative priors have also been adopted by many methods for face restoration (Zhou et al., 127 2022; Gu et al., 2022; TSAI et al., 2024; Wang et al., 2023b; Zhao et al., 2022) and image super-128 resolution (Chen et al., 2022; Liu et al., 2023a; Wu et al., 2023; Liu et al., 2023b). In particular, the 129 methods of FeMaSR (Chen et al., 2022), AdaCode (Liu et al., 2023a), and RestoreFormer++ (Wang 130 et al., 2023b) performed codebook selection via nearest-neighbor (NN) feature matching. Code-131 Former (Zhou et al., 2022) predicted the indices of code items using transformers. DAEFR (TSAI et al., 2024) used an extra HQ encoder as the prior to bridge the domain gap between LQ and HQ 132 images. AdaCode learned five categories of HQ codebooks with a weight predictor to effectively 133 restore the LQ images. In this paper, we propose to replace the discrete VQ process by continuous 134 transformation from LQ feature to HQ ones with the HQ codebook for VQ-based BIR. 135

136 **Vision Transformer** (Dosovitskiy et al., 2021) has inspired great progress in computer vision tasks (Wang et al., 2021a; Carion et al., 2020). It extends the idea of self-attention (Vaswani et al., 137 2017) by taking a sequence of image patches as input tokens. SwinIR (Liang et al., 2021) performed 138 self-attention on shifted local windows (Liu et al., 2021) and transmitted information between them. 139 Restormer (Zamir et al., 2022) exploited self-attention across feature channels for efficient image 140 restoration. Cross-attention is also developed to mix the information from two different inputs (Chen 141 et al., 2021). It is applied in RestoreFormer (Wang et al., 2022b; 2023b) to fuse the LQ and HQ 142 features for blind face restoration. In this paper, we propose a Self-in-Cross-Attention module to 143 collaboratively perform self-attention of LQ feature and cross-attention between LQ feature and HQ 144 codebook for feature learning in VQ-based BIR methods. 145

3 PRELIMINARY

146

147

155

161

148 **Vector Quantization** (VQ) is a classical quantization technology originally developed for signal 149 compression (Linde et al., 1980). With VQ, VQVAE (Van Den Oord et al., 2017) learns an encoder 150 E, a decoder D, and a discrete visual codebook $\mathbf{C} = [\mathbf{c}_1, ..., \mathbf{c}_B]^\top \in \mathbb{R}^{B \times d}$ of images in latent 151 space with a deep neural network. Given an input image \mathbf{x} , the encoder \mathbf{E} extracts its latent 152 feature as $\mathbf{z} = \mathbf{E}(\mathbf{x}) \in \mathbb{R}^{h \times w \times d}$, which is then quantized by replacing each of its feature vector \mathbf{z}_i (i = 1, ..., hw) with the corresponding nearest code item in codebook \mathbf{C} , as follows: 153 154 $\hat{\mathbf{z}}_i =$

$$\mathbf{c}_{k_i}$$
, where $k_i = \underset{j \in \{1,...,B\}}{\arg \min} \|\mathbf{z}_i - \mathbf{c}_j\|_2$. (1)

The quantized latent feature \hat{z} including all replaced code items $\{\hat{z}_i\}_{i=1}^{hw}$ is fed into the decoder D 156 to output the reconstructed image. For model training, VQVAE utilizes three loss functions, *i.e.*, 157 a reconstruction loss \mathcal{L}_{rec} to minimize the distance between the output and the target image y, a 158 codebook loss \mathcal{L}_{code} and a commitment loss \mathcal{L}_{com} with a weighting factor β . Denoting sg(·) as the 159 stop-gradient operator (Van Den Oord et al., 2017), the overall objective function \mathcal{L}_{total} is as follows: 160

$$\mathcal{L}_{\text{total}} = \underbrace{\|\mathbf{y} - \mathbf{D}(\hat{\mathbf{z}})\|_2^2}_{\mathcal{L}_{\text{rec}}} + \underbrace{\|\text{sg}(\mathbf{E}(\mathbf{x})) - \hat{\mathbf{z}}\|_2^2}_{\mathcal{L}_{\text{code}}} + \beta \underbrace{\|\text{sg}(\hat{\mathbf{z}}) - \mathbf{E}(\mathbf{x})\|_2^2}_{\mathcal{L}_{\text{com}}}.$$
(2)



Figure 2: (a) Usage rates of codebook in FeMaSR, AdaCode, CodeFormer and DAEFR. (b) Prediction
accuracy of code index in FeMaSR and AdaCode for ×2 task on different test sets. (c) Prediction
accuracy of code index in CodeFormer and DAEFR on synthetic CelebA-Test set. (d) Inaccurate
index prediction brings performance drop. AdaCode and CodeFormer achieve better results when
using ground-truth (GT) indices (Sec.§4.2).

Since quantization is non-differentiable, VQVAE adopts straight-through gradient estimator (Huh et al., 2023) to back-propagate the gradients of the reconstruction loss \mathcal{L}_{rec} from decoder to encoder. VQGAN (Esser et al., 2021) further improves VQVAE by extraly utilizing an adversarial loss (Goodfellow et al., 2014) and a perceptual loss (Johnson et al., 2016) for better reconstruction quality. These VQ-based generative models have inspired many VQ-based blind image restoration methods (Chen et al., 2022; Zhou et al., 2022; Liu et al., 2023a; Wang et al., 2023b; TSAI et al., 2024).

179 VQ-based Blind Image Restoration (BIR) methods mainly leverage the learned codebook and decoder as a high-quality (HQ) prior robust to diverse degradation. The training pipeline of these 181 methods can be generally divided into two stages with different goals. The first Prior Learning stage aims to reconstruct the HQ image \mathbf{x}^h by learning an HQ encoder \mathbf{E}_h , an HQ decoder \mathbf{D}_h , and an HQ 182 codebook C. The second **Restoration** stage is to restore the low-quality (LQ) images along with the 183 learned HQ prior, *i.e.*, the HQ decoder D_h and HQ codebook C. To this end, these methods learn an LQ encoder \mathbf{E}_l (initialized from the HQ encoder \mathbf{E}_h) to extract from the LQ image \mathbf{x}^l its latent 185 feature $\mathbf{z}^l = \mathbf{E}_l(\mathbf{x}^l)$. Each vector \mathbf{z}_i^l in \mathbf{z}^l is replaced by a predicted code item in HQ codebook C via a VQ process, usually implemented by nearest-neighbor feature matching (Chen et al., 2022; Liu 187 et al., 2023a) or code index prediction (Zhou et al., 2022; TSAI et al., 2024). The quantized HQ 188 feature \hat{z} is fed into the HQ decoder D_h to recover the HQ image x^h . Besides the two stages, many 189 VQ-based BIR methods (Zhou et al., 2022; Wang et al., 2023b) further fuse the LQ feature from 190 encoder and the HQ feature from decoder to trade-off the restoration fidelity and quality. 191

- 192
- 193 193

195 196

197

4 OBSERVATIONS ON VQ-BASED BLIND IMAGE RESTORATION METHODS

Despite promising performance, current VQ-based methods (Chen et al., 2022; Zhou et al., 2022; Wang et al., 2022b; 2023b; Liu et al., 2023a; TSAI et al., 2024) rarely discuss the potential side effects of the essential VQ process for blind image restoration (BIR). In this section, we provide three close observations on the VQ process in the second **Restoration** stage of VQ-based BIR methods.

199 200 201

202

203

4.1 OBSERVATION 1: VQ CONFINES THE CODEBOOK'S REPRESENTATIONAL CAPABILITY

204 The high-quality (HQ) codebook serves as an expressive generative prior for VQ-based BIR (Chen 205 et al., 2022; Zhou et al., 2022; Liu et al., 2023a; TSAI et al., 2024). VQ performs one-hot code 206 selection to replace each low-quality (LQ) feature vector by a single HQ code item from the HQ 207 codebook. This, however, confines the representation range of HQ codebook to a finite set of code items. This limitation would be further amplified by low codebook usage rates of VQ-based BIR 208 methods. As illustrated in Fig. 2, though the codebook usage rate of CodeFormer (Zhou et al., 209 2022) and DAEFR (TSAI et al., 2024) are 98.73% and 100%, respectively, for blind face restoration 210 on 3,000 face images from CelebA-Test (Karras et al., 2018)). FeMaSR (Chen et al., 2022) and 211 AdaCode (Liu et al., 2023a) only used 3.32% and 20.76% of the HQ codebook vectors, respectively, 212 for $\times 2$ blind image super-resolution on the DIV2K validation set. 213

Since the representational capability of VQ process is confined by the discrete code selection of HQ
 codebook in VQ-based BIR methods, it is necessary to develop alternative solutions that can well
 utilize HQ codebook and expand the representional range of HQ codebook on diverse LQ images.



(a) Feature matching in FeMaSR. (b) Code index prediction in CodeFormer. (c) Prediction accuracy v.s. SSIM 224 Figure 3: T-SNE visualization of VQ process in FeMaSR (a) and CodeFormer (b). Different code 225 items in HQ codebook are marked by "⁺⁺⁺ in different colors. The color of LQ feature vector marked 226 by " \circ " or the HQ feature vector marked by " \triangle " is the same with the codebook item they select in 227 VQ process. Gray dashed lines "--" connects the LQ feature vector o and its selected codebook item 228 "☆" ("VQ Process"). Red dash lines "--" connects the LQ feature vector ∘ and the codebook item 229 "☆" selected by the corresponding HQ feature vector using nearest-neighbor (NN) feature matching 230 in the Prior Learning Stage ("GT Selection"). (a) NN feature matching on LQ feature vectors are 231 inconsistent with "GT Selection". For a given LQ feature vector, the codebook item selected by NN feature matching (gray dash line, "- -") is quite different from the corresponding "GT Selection" (232 Red dashed line, "--"). (b) Transformer for code index prediction is not robust to image degradation. 233 In a degraded LQ image, an "LQ Eye" patch looks like skin area and selects the code item represented 234 by many "HQ Skin" patches. The regions marked by purple dashed box and orange dashed box are 235 enlarged for better view ("Enlarge"). (c) Prediction accuracy of code indices by transformer and 236 SSIM results achieved by CodeFormer using "LQ Indices" predicted on LQ images, "HQ Indices" 237 predicted on HQ images, or "GT Indices" defined in §4.2.

4.2 OBSERVATION 2: THE VQ PROCESS IS ERROR-PRONE ON LQ FEATURES

240 Current VQ-based BIR methods mainly adopt two VQ strategies for code index prediction: 1) 241 nearest-neighbor feature matching independently selects a nearest code item from HQ codebook 242 for each LQ feature vector (Chen et al., 2022; Liu et al., 2023a) and 2) learning a transformer to 243 exploit global correlations of the input LQ image for code index prediction (Zhou et al., 2022; TSAI 244 et al., 2024). However, both strategies suffer from inaccurate code selection. To illustrate this point, 245 we evaluate mainstream VQ-based BIR methods on the prediction accuracy of code index, which 246 refers to the percentage that, the number of indices predicted from LQ feature vectors equal their 247 ground-truth (GT) indices. The GT indices are obtained through nearest-neighbor feature matching in Eqn. (1) using the corresponding HQ feature vectors $\mathbf{z}^{h} = \mathbf{E}_{h}(\mathbf{x}^{h})$ (Zhou et al., 2022). In Figs. 2 (b) 248 and (c), we visualize the prediction accuracies of four typical VQ-based BIR methods. The accuracies 249 of FeMaSR (Chen et al., 2022) and AdaCode (Liu et al., 2023a) are at most 30.95% on different test 250 sets, while those of CodeFormer (Zhou et al., 2022) and DAEFR (TSAI et al., 2024) are 5.63% and 251 3.42%, respectively. As shown in Fig. 2 (d), AdaCode (Liu et al., 2023a) and CodeFormer (Zhou et al., 2022) achieve higher PSNR results when using GT code indices. This demonstrates that low 253 accuracy of code index prediction degrades the performance of VQ-based BIR methods. 254

The low prediction accuracy is mainly attributed to the quality degradation of LQ images, as shown in 255 Figs. 3 (a) and (b). Figs. 3 (c) also shows that, using HQ images for index prediction in CodeFormer 256 ("HQ Indices") increases the accuracy from 5.63% to 24.02% with clear improvements on SSIM. 257 Furthermore, learning a transformer to predict code indices (Zhou et al., 2022) casts this problem as 258 a classification task. However, this is error-prone since CodeFormer has $1024^{256} \approx 10^{768}$ possible 259 prediction choices even on a 16×16 LQ feature with an HQ codebook of 1024 items. As shown in 260 Fig. 2 (c), both the image quality degradation and classification challenge conspire to the clear drops 261 in prediction accuracy of code index and SSIM results of CodeFormer. Besides, the VQ process 262 in (Chen et al., 2022; Liu et al., 2023a) also brings gradient estimation errors when back-propagate 263 the gradients from decoder to encoder (Huh et al., 2023), as described in Sec. 3. Thus, it is potentially 264 meaningful to replace the discrete VQ process by alternative solutions that are feasible to perform HQ feature learning while avoiding error-prone index prediction. 265

266 267

- 4.3 OBSERVATION 3: LQ FEATURE IS IMPORTANT FOR BIR, BUT UNDERVALUED IN VQ
- In VQ-based BIR methods (Zhou et al., 2022), the VQ process directly replace the LQ features by selected HQ code items to retrieve high-quality image information. However, this fails to establish a



Figure 4: **Importance of LQ feature for BIR**. (a) Quantitative results of FeMaSR and AdaCode w or w/o feature fusion for $\times 2$ blind super-resolution on DIV2K validation set. (b) Quantitative results of CodeFormer w or w/o feature fusion on synthetic CelebA-Test set.



Figure 5: Comparison of retrained FeMaSR with (w) or without (w/o) feature fusion module. (a) Quantitative results on DIV2K-validation set. Given an LR image (b), the retrained FeMaSR w/o feature fusion (c) loses many texture details when compared to the FeMaSR w feature fusion (d).

288 direct connection between the LQ features and the final restoration performance. To alleviate this 289 issue, many VQ-based BIR methods (Chen et al., 2022; Zhou et al., 2022) further fuse the LQ feature from encoder and the HQ features from decoder to enhance the restoration performance. To study 290 the role of LQ features, we perform experiments on the released models of FeMaSR, AdaCode, and 291 CodeFormer. We remove the corresponding feature fusion module from these models and denote the 292 corresponding variants as "FeMaSR w/o f", "AdaCode w/o f", and "CodeFormer w/o f", respectively. 293 As shown in Figs. 4 (a) and (b), it is not surprising to observe a huge drop of these variants on BIR. We also retrained the restoration stage of FeMaSR and "FeMaSR w/o f". The results in Fig. 5 show 295 that the retrained variant "FeMaSR w/o f" still suffers from clear performance drop. These results 296 validate that the feature of input LQ image is essential to final BIR performance. 297

298 Despite their efforts to preserve LQ information, these VQ-based BIR methods are still constrained 299 by the VQ bottleneck. As the LQ feature is only used for the prediction of code index, regardless of 300 how informative the LQ feature is, the information about the LQ feature transmitted to the decoder is 301 encoded in $\log_2 B$ bits (*B* is the number of items in the HQ codebook). To this end, we argue that 302 the VQ process still underestimates the importance of the LQ feature for BIR. Directly alleviating 303 this problem in the VQ process can potentially improve the performance of VQ-based BIR methods.

303 304 305

277

278

285

287

5 METHODOLOGY

306 307

5.1 Replacing Discrete VQ Selection by Continuous Feature Transformation

308 Based on the three observations analyzed in §4, VQ is a two-sided coin with clear rewards and 309 punishments for VQ-based BIR methods. To avoid the side effects of discrete VQ selection, motivated 310 by the self-attention for code index prediction (Fig. 6 (a)), it is natural to employ cross-attention 311 for continuous feature transformation from the LO feature of input LO image to HO one with the 312 HQ codebook. Specifically, as shown in Fig. 6 (b), to replace discrete VQ process, cross-attention 313 takes the input feature as the query and the HQ codebook as both the key and value. The attention 314 map correlates each LQ feature vector with HQ codebook items. Each output feature vector is an 315 adaptively weighted combination of HQ codebook. In this way, the input LQ feature is transformed 316 into HQ ones. However, the vanilla cross-attention ignores the self-expressiveness of LQ feature and 317 would fail to preserve the fidelity of diverse LQ images (Figs. 1 (b) and (d)).

318

319 5.2 PROPOSED SELF-IN-CROSS-ATTENTION (SINCA)

320

To exploit useful self-expressiveness (Elhamifar & Vidal, 2013) of LQ images for better BIR per formance, in this paper, we propose a new Self-in-Cross-Attention (SinCA) module to augment
 the HQ codebook with specific feature of input LQ image and performs cross-attention between
 input LQ feature and augmented codebook. As shown in Fig. 6 (c), given an LQ feature tensor

333 334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

355

356 357 358

367

368



Figure 6: (a) Self-attention for code index prediction. (b) Cross-attention for feature transformation. (c) Our proposed Self-in-Cross Attention (SinCA) for effective feature transformation.



Figure 7: Visualization of the attention map of the first SinCA in 'FeMaSR+SinCA' for $\times 2$ blind image super-resolution. We take inputs of size 64×64 which will be encoded into $16 \times 16 = 256$ feature vectors and visualize the attention weights by selecting the first 100 indices of self-part (0-99) and the first 100 indices of cross-part (256-355).

 $\mathbf{z}^{l} \in \mathbb{R}^{h \times w \times d}$ of an LQ image extracted from the LQ encoder, our SinCA first reshapes it into an LQ feature matrix $\mathbf{X} \in \mathbb{R}^{hw \times d}$ and then multiplies it with a linear projection matrix $\mathbf{W}_{\mathbf{Q}}$ to obtain the query matrix \mathbf{Q} . To jointly explore the expressiveness of HQ codebook $\mathbf{C} \in \mathbb{R}^{B \times d}$ and excavate the self-expressiveness of the LQ feature itself, we concatenate the LQ feature \mathbf{X} and the HQ codebook \mathbf{C} to obtain the key matrix $\mathbf{K} \in \mathbb{R}^{(hw+B) \times d}$ and the value matrix $\mathbf{V} \in \mathbb{R}^{(hw+B) \times d}$ with the corresponding linear projection matrices $\mathbf{W}_{\mathbf{K}}$ and $\mathbf{W}_{\mathbf{V}}$, respectively, as follows:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_{\mathbf{Q}}, \mathbf{K} = \begin{bmatrix} \mathbf{X}\mathbf{W}_{\mathbf{K}} \\ \mathbf{C}\mathbf{W}_{\mathbf{K}} \end{bmatrix}, \mathbf{V} = \begin{bmatrix} \mathbf{X}\mathbf{W}_{\mathbf{V}} \\ \mathbf{C}\mathbf{W}_{\mathbf{V}} \end{bmatrix}.$$
(3)

Denoting $\mathbf{K}_{\mathbf{X}} = \mathbf{X}\mathbf{W}_{\mathbf{K}}, \mathbf{K}_{\mathbf{Code}} = \mathbf{C}\mathbf{W}_{\mathbf{K}}, \mathbf{V}_{\mathbf{X}} = \mathbf{X}\mathbf{W}_{\mathbf{V}}, \text{ and } \mathbf{V}_{\mathbf{Code}} = \mathbf{C}\mathbf{W}_{\mathbf{V}}, \text{ the attention}$ matrix $\mathbf{A} \in \mathbb{R}^{hw \times (hw+B)}$ and the output feature matrix **O** of our SinCA are computed as follows:

$$\mathbf{A} = \operatorname{Softmax}\left(\frac{1}{\sqrt{d}}\mathbf{Q}\mathbf{K}^{\top}\right) = \operatorname{Softmax}\left(\frac{1}{\sqrt{d}}\left[\mathbf{Q}\mathbf{K}_{\mathbf{X}}^{\top} \ \mathbf{Q}\mathbf{K}_{\mathbf{Code}}^{\top}\right]\right), \mathbf{O} = \mathbf{A}\begin{bmatrix}\mathbf{V}_{\mathbf{X}}\\\mathbf{V}_{\mathbf{Code}}\end{bmatrix}.$$
 (4)

As revealed by Eqn. (4), the attention map A correlates both the LQ feature and the input-augmented codebook. This is then used to adaptively weight the value feature matrices V_X and V_{Code} . In this way, our SinCA simultaneously leverages the expressive HQ codebook prior and the selfexpressiveness of input LQ image to obtain the output feature matrix O.

We replace the VQ process in VQ-based BIR methods (Chen et al., 2022; Liu et al., 2023a; Zhou et al., 2022; TSAI et al., 2024) by a transformer (Dosovitskiy et al., 2021) using our SinCA, which aims to transform from the LQ feature of input image to HQ one with the HQ codebook.

5.3 A CLOSER LOOK AT OUR SINCA

To study the working mechanism of our SinCA, in Fig. 7 we visualize the attention map of "Fe-369 MaSR+SinCA" (variant of FeMaSR (Chen et al., 2022)) for $\times 2$ blind image super-resolution. Here, 370 the VQ in FeMaSR is replaced by a transformer using our SinCA. The "Example 1" in Fig. 7 (a) 371 show highly structured with repeated patterns, which could be better represented by itself. Thus, the 372 self-part in the attention map of our SinCA exhibits dense and high attention weights. This indicates 373 that our SinCA effectively utilizes the self-expressiveness of the LQ feature itself for BIR. In contrast, 374 the "Example 2" in Fig. 7 (b) presents an animal that can be well represented by the HQ codebook 375 (cross-part). In this case, our SinCA utilizes more the HQ codebook to recover the HQ feature. 376

377 In sum, our SinCA utilizes augmented codebook to exploit the self-expressiveness of the LQ feature and the correlation between LQ feature and HQ codebook. This enables each pixel of input LQ image to be adaptively recovered by a weighted combination of augmented codebook. Compared to discrete
 VQ process relying on HQ code index selection, our SinCA extends the representational range of the
 HQ codebook and further exploits the self-expressiveness of LQ image for VQ-based BIR methods.

- 382 383 6 EXPERIMENTS
- 384 385

386

6.1 EXPERIMENTAL SETUP

Baselines. We evaluate our Self-in-Cross-Attention (SinCA) on four typical VQ-based BIR methods: 387 FeMaSR (Chen et al., 2022) and AdaCode (Liu et al., 2023a) for blind image super-resolution (BSR), 388 CodeFormer (Zhou et al., 2022) and DAEFR (TSAI et al., 2024) for blind face restoration (BFR). For 389 each baseline method, we directly use the HQ encoder, HQ codebook, and HQ decoder it learned in 390 the first **Prior Learning** stage (§3) to reconstruct the HQ images. Then we replace the VQ process 391 by a transformer with our SinCA, and fine-tuned the encoder to restore the LQ images along with 392 fixed codebook and decoder learned by each baseline. For all models in our experiments, we set the 393 number of attention heads as eight in the transformers using our SinCA and the channel dimension of both the input and the output of the transformer equal the channel dimension d of codebook. More 394 details will be provided in the Appendix. 395

Training Dataset. For BSR task, we train models on DIV2K (Agustsson & Timofte, 2017) training set, including 800 HQ images of 2K resolution. The HR images are cropped into 256×256 patches for training. Following the settings in FeMaSR and Adacode, we generate pairs of high-resolution (HR) and low-resolution (LR) training images with the degradation pipeline in BSRGAN (Zhang et al., 2021). For the BFR task, we employ the FFHQ (Karras et al., 2019) dataset with 70,000 HQ images of size 1024×1024 . All HQ face images are resized into 512×512 for training. The LQ images are synthesized by following the degradation pipeline in CodeFormer or DAEFR, respectively.

Test Set. The BSR methods are evaluated on the DIV2K validation set, Urban100, BSDS100, and
Manga109. The LR images are generated with the mixed degradation pipelines of (Zhang et al., 2021; Wang et al., 2021c) used in FeMaSR. The BFR methods are evaluated on the 3000 images used
in (Zhou et al., 2022; TSAI et al., 2024) from CelebA-Test set (Karras et al., 2018) and the real-world
dataset CelebChild-Test (Wang et al., 2021b). The LQ images from CelebA-Test set are synthesized
with the same settings used in CodeFormer and DAEFR, respectively.

Metrics. For BSR, we report PSNR and SSIM results computed on the y-channel, and LPIPS (Zhang et al., 2018b) on RGB images. For BFR, we compute PSNR, SSIM, and LPIPS on CelebA-Test, while FID (Heusel et al., 2017) and NIQE (Mittal et al., 2012) on real-world CelebChild-Test.

Training Details. For each baseline using our SinCA, we follow its original setting of codebook size, optimizer, and learning rates. For FeMaSR and AdaCode, we train the second stage with a batchsize of 16 for 100K and 160K iterations, respectively. For CodeFormer, we train the stage-2 with a batchsize of 8 and stage-3 with a batchsize of 4. For DAEFR, we train the last stage with a batchsize of 8 for 200K iterations. The models are implemented by PyTorch. The BSR methods are trained with 2 RTX 4090 GPUs, while BFR methods are trained with 1 Tesla H100 GPU.

- 418 419
- 6.2 COMPARISON RESULTS
- 420

421 Blind Image Super-Resolution. In Table 1, we summarize the quantitative results of comparison 422 methods on four benchmarks. One can see that the FeMaSR and AdaCode using our SinCA (denoted as "FeMaSR+SinCA" and "AdaCode+SinCA", respectively) outperform their corresponding baselines 423 (retrained under the same settings) by $0.70 \sim 1.10$ dB in PSNR, $0.01 \sim 0.03$ in SSIM, and generally 424 better results in LPIPS. For reference, in Table 1 we also report the results of the released models, 425 denoted as "FeMaSR (release)" and "AdaCode (release)". Note that albeit being trained with much 426 larger datasets, the released models still achieve worse results than the corresponding methods using 427 our SinCA in terms of PSNR and SSIM. 428

Blind Face Restoration. In Table 2, we provide the quantitative results of blind face restoration on synthetic and real-world datasets. One can see that, on CelebA-Test set, compared with retrained ones, CodeFormer using our SinCA obtains a gain of 0.23 dB and 0.0129 on PSNR and SSIM, respectively, while DAEFR using our SinCA achieves an improvement of 1.92 dB, 0.0636 and 0.0069 on PSNR,



"+SinCA": Figure 8: Comparison on blind image super-resolution and blind face restoration. employing a transformer using our SinCA for continuous feature transformation.

Scale	Method	١	Urban10	0]	BSDS100			Manga10	19	DIV2K valid		
Scale		PSNR↑	SSIM↑	LPIPS↓	PSNR ↑	SSIM↑	$\text{LPIPS}{\downarrow}$	PSNR↑	SSIM↑	$\text{LPIPS}{\downarrow}$	PSNR↑	SSIM↑	LPIPS↓
	FeMaSR (release)	20.11	0.5769	0.3847	21.90	0.5189	0.4225	22.14	0.7075	0.3358	23.44	0.6443	0.3878
	FeMaSR (re-train)	19.61	0.5607	0.4103	21.25	0.5052	0.4315	21.85	0.7092	0.3505	22.76	0.6311	0.4129
~ 2	FeMaSR + SinCA	20.59	0.5853	0.3958	22.43	0.5335	0.4387	22.61	0.7222	0.3428	23.84	0.6498	0.4026
~2	AdaCode (release)	20.46	0.5886	0.3808	22.03	0.5173	0.4199	22.35	0.7097	0.3226	23.38	0.6467	0.3759
	AdaCode (re-train)	19.47	0.5565	0.4124	21.28	0.5014	0.4422	21.81	0.7038	0.3522	22.55	0.6231	0.4102
	AdaCode + SinCA	20.46	0.5924	0.3940	22.44	0.5440	0.4293	22.74	0.7306	0.3326	23.75	0.6621	0.3917
	FeMaSR (release)	18.52	0.4891	0.4358	20.49	0.4528	0.4647	18.85	0.6107	0.3945	21.72	0.5626	0.4418
	FeMaSR (re-train)	18.41	0.4729	0.4759	20.82	0.4585	0.4950	18.86	0.5999	0.4269	21.72	0.5634	0.4715
$\times 4$	FeMaSR + SinCA	19.11	0.4887	0.4707	20.80	0.4477	0.4928	19.47	0.6168	0.4267	22.30	0.5703	0.4673
×4	AdaCode (release)	18.71	0.4875	0.4444	20.71	0.4495	0.4752	19.00	0.6067	0.3955	21.80	0.5638	0.4432
	AdaCode (re-train)	17.94	0.4644	0.4796	19.75	0.4237	0.5025	18.62	0.5936	0.5607	20.73	0.5404	0.4760
	AdaCode + SinCA	18.72	0.4780	0.4660	21.15	0.4617	0.4882	19.51	0.6093	0.4229	22.09	0.5658	0.4587

Table 1: Results of blind image super-resolution methods on four benchmark test sets.

Table 2: Results of blind face restoration methods on two synthetic and real-world test sets.

477	Method	CelebA-Test			CelebChild-Test		Mathod	CelebA-Test			CelebChild-Test	
478		PSNR ↑	SSIM↑	LPIPS↓	FID↓	NIQE↓	Method	PSNR ↑	SSIM↑	LPIPS↓	FID↓	NIQE↓
479	CodeFormer (release)	22.19	0.5957	0.3152	116.23	4.983	DAEFR (release)	19.92	0.5534	0.3880	105.70	4.143
/190	CodeFormer (re-train)	22.66	0.6248	0.3100	116.53	4.883	DAEFR (re-train)	19.65	0.5456	0.3675	105.23	4.220
400	CodeFormer + SinCA	22.89	0.6377	0.3108	121.50	5.112	DAEFR + SinCA	21.57	0.6092	0.3606	104.56	4.097
481					-						-	

SSIM, LPIPS, respectively. On CelebChild-Test dataset, we cannot evaluate their performance on restoration fidelity since there is no ground-truth image. The CodeFormer using our SinCA obtain similar results of FID and NIQE to the released model, while DAEFR using our SinCA achieves an improvement in FID and NIQE. These results demonstrate that our SinCA serves as a promising replacement for the discrete VQ process for blind face restoration.

>	Blind Face Restoration									
Method	DIV2K valid			Urban100			Method	CelebA-Test		
Wiethou	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	Wieulou	PSNR ↑	SSIM↑	LPIPS↓
FeMaSR (index)	22.85	0.6399	0.4124	19.78	0.5666	0.4101	CodeFormer (index)	22.63	0.6243	0.3061
FeMaSR (feature)	23.84	0.6498	0.4026	20.59	0.5853	0.3958	CodeFormer (feature)	22.89	0.6377	0.3108
AdaCode (index)	22.17	0.5542	0.4992	20.04	0.5635	0.4045	DAEFR (index)	19.85	0.5551	0.3635
AdaCode (feature)	23.75	0.6621	0.3917	20.46	0.5924	0.3940	DAEFR (feature)	21.57	0.6092	0.3606

Table 3: Comparison of using our SinCA for code index prediction or feature fusion.

Table 4: Comparison of Self-Attention (SA), Cross-Attention (CA), and our Self-in-Cross-Attention (SinCA) used by transformers for feature fusion in VQ-based BIR methods.

95	×	2 Blind	Image S	Blind Face Restoration							
0	Method	D	IV2K va	lid	i I	Urban10	0	Method	CelebA-Test		
)7	Wiethou	PSNR ↑	SSIM↑	LPIPS↓	PSNR ↑	SSIM↑	LPIPS↓	Wiethou	PSNR ↑	SSIM↑	LPIPS↓
3	FeMaSR + SA	23.01	0.6352	0.4035	19.87	0.5689	0.3990	CodeFormer + SA	22.85	0.6341	0.3165
	FeMaSR + CA	23.22	0.6424	0.3970	20.08	0.5740	0.3957	CodeFormer + CA	22.83	0.6310	0.3154
	FeMaSR + SinCA	23.84	0.6498	0.4026	20.59	0.5853	0.3958	CodeFormer + SinCA	22.89	0.6377	0.3108
	AdaCode + SA	21.88	0.6053	0.4185	19.06	0.5379	0.4228	DAEFR + SA	21.53	0.5987	0.4118
	AdaCode + CA	22.67	0.6349	0.3654	19.53	0.5583	0.4247	DAEFR + CA	21.47	0.6001	0.3770
	AdaCode + SinCA	23.75	0.6621	0.3917	20.46	0.5924	0.3940	DAEFR + SinCA	21.57	0.6092	0.3606

505 Visual Comparisons in Fig. 8 show that, the four VQ-based BIR methods using our SinCA consis-506 tently preserves the colors and geometric shapes. For example, "FeMaSR + SinCA" and "AdaCode + SinCA" restore the steel ropes on the bridge and colors of bricks in the 1-st and 2-nd rows, respectively, while "CodeFormer + SinCA" and "DAEFR + SinCA" properly recover the skin tones and 508 509 grin in the 3-rd and 4-th rows, respectively.

510 511

512

493

494

504

507

6.3 ABLATION STUDY

513 **Necessity of Replacing Index Prediction with Feature Transformation**. To study this aspect, we 514 additionally design a variant of transformer using our SinCA for discrete index prediction in VO-based 515 BIR methods. As shown in Table 3, the four methods using our SinCA for feature transformation 516 ("feature") obviously outperform those for code index prediction ("index"), especially on PSNR and 517 SSIM. This validates the effectiveness of replacing the VQ process for code index prediction by 518 feature transformation using our SinCA in VQ-based BIR methods.

519 Effectiveness of our SinCA. Here, we compare VQ-based BIR methods with transformers using 520 self-attention (SA) (Dosovitskiy et al., 2021), cross-attention (CA) (Chen et al., 2021), or our SinCA 521 for continuous feature learning. For CA, we generate the query/value matrix from HQ codebook and 522 the key matrix from LQ feature for consistent input and output dimensions. As shown in Table 4, 523 compared with those using SA and CA, the baselines using our SinCA achieve superior results on 524 objective metrics in most cases. This validates the effectiveness of SinCA in VQ-based BIR methods.

525 526

7 CONCLUSION

527 528 529

In this paper, we revisited the key VQ process in VQ-based blind image restoration (BIR) methods 530 and provided three close observations on the side-effects of VQ on code index prediction. We 531 revealed that discrete VQ limits the representational capability of HQ codebook prior, is error-prone 532 on code index prediction, and under-values the important LQ feature for BIR. Based on these 533 observations, we proposed to replace the discrete VQ selection by continuous feature transformation 534 from LQ feature to HQ ones via cross-attention of LQ feature and HQ codebook. We further proposed a Self-in-Cross-Attention (SinCA) module to augment HQ codebook with LQ feature and 536 perform cross-attention between LQ feature and input-augmented codebook. Our SinCA extends 537 the representational capability of HQ codebook and well leverages the self-expressiveness of input LQ image. Experiments demonstrated that, the four VQ-based BIR methods replacing the discrete 538 VQ process with a transformer using our SinCA achieve better performance on blind image superresolution and blind face restoration. Ablation studies also validated the effectiveness of our SinCA.

540 8 ETHICS STATEMENT

542

543

544

546

547 548

549

550

551

552 553 554

555

556

559

560

592

Our work adheres to ethical standards. We ensure that our work does not perpetuate bias or harm. All datasets used in our experiments are publicly available and ethically sourced, with proper consideration for privacy and consent.

9 REPRODUCIBILITY STATEMENT

We are committed to reproducibility. We provide the working mechanism of our proposed module in Sec. §5, our experimental setup in Sec. §6, and implementation details in Appendix A and Appendix B. All datasets used in our experiments are publicly available, and we will also release our code to facilitate further research and reproducibility.

REFERENCES

- Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pp. 126–135, 2017.
- Shiyue Cao, Yueqin Yin, Lianghua Huang, Yu Liu, Xin Zhao, Deli Zhao, and Kaigi Huang. Efficientvqgan: Towards high-resolution image generation with efficient vision transformers. In *Int. Conf. Comput. Vis.*, pp. 7368–7377, 2023.
- 561
 562 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
 563 Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis.*, pp. 213–229.
 564 Springer, 2020.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. Maskgit: Masked generative image transformer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2022.
- Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo.
 Real-world blind super-resolution via feature matching with implicit high-resolution priors. In
 ACM Int. Conf. Multimedia, pp. 1329–1338, 2022.
- Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Int. Conf. Comput. Vis.*, pp. 357–366, 2021.
- Yu Chen, Ying Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. Fsrnet: End-to-end learning face
 super-resolution with facial priors. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 2492–2501, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In Int. Conf. Learn. Represent., 2021. URL https://openreview.net/forum? id=YicbFdNTTy.
- Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications.
 IEEE Transactions on Pattern Analysis and Machine Intelligence, 35(11):2765–2781, 2013. doi: 10.1109/TPAMI.2013.57.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 12873–12883, 2021.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Adv. Neural Inform. Process. Syst.*, 27, 2014.
- Jinjin Gu, Hannan Lu, Wangmeng Zuo, and Chao Dong. Blind super-resolution with iterative kernel correction. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 1604–1613, 2019.

594 595 596	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 <i>Eur. Conf. Comput. Vis.</i> , 2022.
597 598 599 600	Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In <i>Adv. Neural Inform. Process. Syst.</i> , 2017.
601 602	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Adv. Neural Inform. Process. Syst., 33:6840–6851, 2020.
604 605	Yan Huang, Shang Li, Liang Wang, Tieniu Tan, et al. Unfolding the alternating optimization for blind super resolution. <i>Adv. Neural Inform. Process. Syst.</i> , 33:5632–5643, 2020.
606 607 608	Minyoung Huh, Brian Cheung, Pulkit Agrawal, and Phillip Isola. Straightening out the straight- through estimator: Overcoming optimization challenges in vector quantized networks. In <i>Interna-</i> <i>tional Conference on Machine Learning</i> , pp. 14096–14113. PMLR, 2023.
609 610 611	Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In <i>Eur. Conf. Comput. Vis.</i> , pp. 694–711. Springer, 2016.
612 613 614	Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In <i>Int. Conf. Learn. Represent.</i> , 2018. URL https://openreview.net/forum?id=Hk99zCeAb.
615 616 617	Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 4401–4410, 2019.
618 619 620	Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 8110–8119, 2020.
621 622 623	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of NAACL-HLT</i> , pp. 4171–4186, 2019.
624 625 626	Diederik P Kingma and Max Welling. Auto-encoding variational bayes. <i>arXiv preprint arXiv:1312.6114</i> , 2013.
627 628 629	Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image generation using residual quantization. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 11523– 11532, 2022.
630 631 632 633	Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In <i>International Conference on Computer Vision Workshops (ICCVW)</i> , pp. 1833–1844, 2021.
634 635 636	Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Ben Fei, Bo Dai, Wanli Ouyang, Yu Qiao, and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. <i>arXiv</i> preprint arXiv:2308.15070, 2023.
637 638 639	Y. Linde, A. Buzo, and R. Gray. An algorithm for vector quantizer design. <i>IEEE Transactions on Communications</i> , 28(1):84–95, 1980. doi: 10.1109/TCOM.1980.1094577.
640 641	Kechun Liu, Yitong Jiang, Inchang Choi, and Jinwei Gu. Learning image-adaptive codebooks for class-agnostic image restoration. In <i>Int. Conf. Comput. Vis.</i> , pp. 5373–5383, 2023a.
642 643 644 645	Yunlong Liu, Tao Huang, Weisheng Dong, Fangfang Wu, Xin Li, and Guangming Shi. Low-light image enhancement with multi-stage residue quantization and brightness-aware attention. In <i>Int.</i> <i>Conf. Comput. Vis.</i> , pp. 12140–12149, 2023b.
646 647	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In <i>Int. Conf. Comput.</i> <i>Vis.</i> , pp. 10012–10022, 2021.

648 649	Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a "completely blind" image quality analyzer. <i>IEEE Signal processing letters</i> , 20(3):209–212, 2012.
651 652	Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. In <i>Eur. Conf. Comput. Vis.</i> , 2020.
653 654 655	Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , pp. 1–1, 2021. doi: 10.1109/TPAMI.2021.3115428.
656 657 658	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. pp. 8821–8831. Pmlr, 2021.
659 660	Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. pp. 1278–1286. PMLR, 2014.
662 663	Ziyi Shen, Wei-Sheng Lai, Tingfa Xu, Jan Kautz, and Ming-Hsuan Yang. Deep semantic face deblurring. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 8260–8269, 2018.
664 665	Assaf Shocher, Nadav Cohen, and Michal Irani. "zero-shot" super-resolution using deep internal learning. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 3118–3126, 2018.
667 668	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In Interna- tional Conference on Learning Representations, 2020.
669 670	Xuansong Xie Tao Yang, Peiran Ren and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , 2021.
671 672 673 674	YU-JU TSAI, Yu-Lun Liu, Lu Qi, Kelvin C.K. Chan, and Ming-Hsuan Yang. Dual associated encoder for face restoration. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=gwDuW70k5f.
675 676	Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Adv. Neural Inform. Process. Syst., 30, 2017.
678 679 680	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Adv. Neural Inform. Process. Syst.</i> , 30, 2017.
681 682 683	Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to- end panoptic segmentation with mask transformers. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 5463–5474, 2021a.
684 685 686	Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. <i>arXiv preprint arXiv:2305.07015</i> , 2023a.
687 688	Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 9168–9178, 2021b.
689 690 691 692	Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super- resolution with pure synthetic data. In <i>International Conference on Computer Vision Workshops</i> (<i>ICCVW</i>), 2021c.
693 694 695	Yinhuai Wang, Yujie Hu, and Jian Zhang. Panini-net: Gan prior based degradation-aware feature interpolation for face restoration. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pp. 2576–2584, 2022a.
696 697 698 699	Zhouxia Wang, Jiawei Zhang, Runjian Chen, Wenping Wang, and Ping Luo. Restoreformer: High- quality blind face restoration from undegraded key-value pairs. In <i>IEEE Conf. Comput. Vis. Pattern</i> <i>Recog.</i> , pp. 17512–17521, 2022b.
700 701	Zhouxia Wang, Jiawei Zhang, Tianshui Chen, Wenping Wang, and Ping Luo. Restoreformer++: Towards real-world blind face restoration from undegraded key-value pairs. <i>IEEE Transactions on</i> <i>Pattern Analysis and Machine Intelligence</i> , 2023b.

702 703 704	Rui-Qi Wu, Zheng-Peng Duan, Chun-Le Guo, Zhi Chai, and Chongyi Li. Ridcp: Revitalizing real image dehazing via high-quality codebook priors. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 22282–22291, 2023.
706 707	Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic image super-resolution and personalized stylization. <i>arXiv preprint arXiv:2308.14469</i> , 2023.
708 709 710 711	Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved VQGAN. In <i>International Conference on Learning Representations</i> , 2022. URL https://openreview.net/forum?id=pfNyExj7z2.
712 713 714	Xin Yu, Basura Fernando, Bernard Ghanem, Fatih Porikli, and Richard Hartley. Face super-resolution guided by facial component heatmaps. In <i>Eur. Conf. Comput. Vis.</i> , pp. 217–233, 2018.
715 716 717	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 <i>IEEE Conf.</i> <i>Comput. Vis. Pattern Recog.</i> , pp. 5728–5739, 2022.
718 719 720 721	Jiahui Zhang, Fangneng Zhan, Christian Theobalt, and Shijian Lu. Regularized vector quantization for tokenized image synthesis. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 18467–18476, 2023.
722 723	Kai Zhang, Wangmeng Zuo, and Lei Zhang. Learning a single convolutional super-resolution network for multiple degradations. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 3262–3271, 2018a.
724 725 726	Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In <i>Int. Conf. Comput. Vis.</i> , pp. 4791–4800, 2021.
727 728 729	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 586–595, 2018b.
730 731 732 732	Yang Zhao, Yu-Chuan Su, Chun-Te Chu, Yandong Li, Marius Renn, Yukun Zhu, Changyou Chen, and Xuhui Jia. Rethinking deep face restoration. In <i>IEEE Conf. Comput. Vis. Pattern Recog.</i> , pp. 7652–7661, 2022.
734 735 736	Chuanxia Zheng, Tung-Long Vuong, Jianfei Cai, and Dinh Phung. Movq: Modulating quantized vectors for high-fidelity image generation. <i>Adv. Neural Inform. Process. Syst.</i> , 35:23412–23425, 2022.
737 738 739	Shangchen Zhou, Kelvin Chan, Chongyi Li, and Chen Change Loy. Towards robust blind face restoration with codebook lookup transformer. <i>Adv. Neural Inform. Process. Syst.</i> , 35:30599–30611, 2022.
740 741 742 743 744	Feida Zhu, Junwei Zhu, Wenqing Chu, Xinyi Zhang, Xiaozhong Ji, Chengjie Wang, and Ying Tai. Blind face restoration via integrating face shape and generative priors. In <i>IEEE Conf. Comput. Vis.</i> <i>Pattern Recog.</i> , pp. 7662–7671, 2022.
745 746 747 748	
749 750 751 752	
753 754	

756 MORE ARCHITECTURE DETAILS А

758 We develop transformers (Dosovitskiy et al., 2021) using our SinCA for feature transformation to 759 replace the discrete VQ process for code index prediction in four VQ-based BIR methods. To verify 760 the effectiveness of our SinCA, we apply the transformers using our SinCA on FeMaSR (Chen et al., 761 2022) and AdaCode (Liu et al., 2023a) for blind super-resolution as well as CodeFormer (Zhou et al., 762 2022) and DAEFR (TSAI et al., 2024) for blind face restoration. The transformers contain multiple layers. Each transformer layer incorporates one SinCA to excavate the useful information from input 764 LQ feature itself and HQ codebook prior. Following our SinCA, each layer employs a Gated-Dconv Feed-Forward Network from (Zamir et al., 2022). We also use a skip connection to concatenate 765 the input LQ feature of the transformer with the output HQ feature, followed by a linear layer for 766 feature fusion. In all experiments, the input and output of the transformers maintain the same channel 767 dimension as that of the HQ codebook in different baseline methods. 768

769 **FeMaSR** (Chen et al., 2022). The transformer using our SinCA for FeMaSR contains nine layers. 770 Since the HQ codebook in FeMaSR is of size 1024×512 , where the number of code items in codebook is B = 1024. The input and output channel dimension in our SinCA should be d = 512. 771

772 AdaCode (Liu et al., 2023a) has five categories of basis codebooks in its network backbone: archi-773 tecture, indoor objects, natural scenes, street views and portraits, with codebooks of size 512×256 , 774 256×256 , 512×256 , 256×256 , and 256×256 , respectively. We replace the VQ process 775 for each category by a transformer using our SinCA with a channel dimension of d = 256, and B = 512, 256, 512, 256, 256. We set the number of transformer layers for these five categories to be 776 5, 4, 3, 4, and 4, respectively. The number of transformer layers is determined by the average value 777 of weight maps from the pre-trained weight predictor. 778

779 CodeFormer (Zhou et al., 2022) uses a transformer (Dosovitskiy et al., 2021) consists of nine 780 transformer layers with a channel dimension of d = 512, followed by a linear projection for code 781 index prediction. To equip CodeFormer with our SinCA, we directly replace its transformer for code index prediction by our alternative transformer containing nine transformer layers with a channel 782 dimension of d = 256, and the number of code items in HQ codebook is B = 1024. 783

784 **DAEFR** (TSAI et al., 2024) also employs a nine-layer transformer for code index prediction. Similar 785 to our practice on CodeFormer (Zhou et al., 2022), we replace its transformer with the transformer 786 containing nine layers with a channel dimension of d = 256, and the number of code items in HQ 787 codebook is B = 1024.

788 789

В **EXPERIMENTAL DETAILS**

790 791 792

MORE DETAILS ON MAIN EXPERIMENTS B.1

793 In this section, we elaborate more details on the main experiments. All of our implementations are 794 built upon publicly released codes of the baseline methods. For better clarity, we denote \hat{z} as the quantized HQ feature obtained by VQ or the output HQ feature by the transformer using our SinCA. 796

FeMaSR (Chen et al., 2022) adopts a two-stage training pipeline. We directly take the pre-trained 797 codebook and decoder from the first stage. Then we replace the VQ process in FeMaSR with the 798 transformer using our SinCA, and train the modified FeMaSR for the second stage. 799

800 AdaCode (Liu et al., 2023a) employs a three-stage training pipeline. The first two stages aim to 801 obtain high-quality codebook and decoder prior as well as a weight predictor, while the last stage fine-tunes the encoder and weight predictor for image restoration. We replace the VQ process in 802 AdaCode with the transformer using our SinCA, and train the modified AdaCode on the third stage. 803 We take the original loss functions used in AdaCode for training. 804

805 CodeFormer (Zhou et al., 2022) utilizes a three-stage training pipeline. The first stage learns the 806 HQ generative prior of codebook and decoder. The second stage learns a transformer to predict 807 code indices. The third stage aims for feature fusion of LQ and HQ features in the decoder. We replace the VQ process in CodeFormer with the transformer using our SinCA, and train the modified 808 CodeFormer on the second and third stages. The second training stage of vanilla CodeFormer has two loss functions: a cross-entropy loss $\mathcal{L}_{code}^{token}$ to supervise the training of transformer for index prediction and an $\ell_2 \log \mathcal{L}_{code}^{feat'}$ to align the extracted LQ feature z^l and the quantized HQ feature \hat{z} . The training objective function \mathcal{L}_{tf} of the second stage can be written as follows:

812 813 814

815 816

817

818

819

820

827

828

829

830

831

 $\mathcal{L}_{\text{code}}^{\text{token}} = \sum_{i=0}^{hw-1} -s_i \log\left(\hat{s}_i\right), \quad \mathcal{L}_{\text{code}}^{\text{feat}'} = \left\| \mathbf{z}^l - \operatorname{sg}\left(\hat{\mathbf{z}}\right) \right\|_2^2, \quad \mathcal{L}_{tf} = \lambda_{\text{token}} \, \mathcal{L}_{\text{code}}^{\text{token}} + \mathcal{L}_{\text{code}}^{\text{feat}'}, \quad (5)$

where s_i represents the *i*-th element of the GT indices (defined in §4.1 while \hat{s}_i denotes the *i*-th element of indices predicted by the vanilla transformer in CodeFormer and λ_{token} is a hyper-parameter used to balance the two loss functions. Since our transformer performs feature transformation instead of code index prediction, we replace the cross-entropy loss $\mathcal{L}_{\text{code}}^{\text{token}}$ with a feature matching loss $\mathcal{L}_{\text{matching}} = \|\mathbf{z}^t - \text{sg}(\hat{\mathbf{z}})\|_2^2$, where \mathbf{z}^t denotes the LQ feature \mathbf{z}^l after transformation. Setting equal weight to $\mathcal{L}_{\text{matching}}$ and $\mathcal{L}_{\text{cod}}^{\text{tead'}}$, our training objective function in the second stage is

$$\mathcal{L}_{tf} = \underbrace{\left\| \hat{\mathbf{z}} - \mathbf{z}^{t} \right\|_{2}^{2}}_{\mathcal{L}_{\text{matching}}} + \underbrace{\left\| \mathbf{z}^{l} - \operatorname{sg}\left(\hat{\mathbf{z}} \right) \right\|_{2}^{2}}_{\mathcal{L}_{\text{code}}^{\text{feat}'}}.$$
(6)

The training of the third stage in CodeFormer inherits the loss functions used in the second stage, with additional image-level ℓ_1 loss function and a perceptual loss (Johnson et al., 2016). Following this setting, we keep the loss functions adopted in training the second stage of the CodeFormer with the transformer using our SinCA with the ℓ_1 loss and the perceptual loss as well.

DAEFR (TSAI et al., 2024) is trained in three stages. Similar to the second stage in Code-Former (Zhou et al., 2022), the third stage of DAEFR aims to align the LQ feature with HQ codebook, by using a cross-entropy loss function and an ℓ_2 loss function of Eqn. (6) as the training objective. We replace the VQ process in DAEFR with the transformer using our SinCA for feature learning. We train the revised DAEFR in the third stage using our feature matching loss $\mathcal{L}_{matching}$ defined above and the loss function $\mathcal{L}_{code}^{feat'}$ defined in Eqn. (6), with equal weights.

837 838 839

840

B.2 MORE DETAILS ON ABLATION STUDY

Transformer Using Our SinCA for Code Index Prediction. As described in §6.3 of our main 841 paper, we also apply the transformer using our SinCA for code index prediction in VQ-based BIR 842 methods. To this end, we use the linear layer following the transformer developed by our SinCA to 843 project the output feature map of size $hw \times d$ into the size of $hw \times B$, with B is the number of code 844 items in HQ codebook. The output feature matrix of the linear layer is transformed into a probability 845 matrix **P** via a softmax operation, where p_{ij} (i = 1, ..., hw, j = 1, ..., B) represents the probability 846 that the *i*-th LQ feature vector select the *j*-th code item of HQ codebook. Then we use top-1 selection 847 for each LQ feature vector. For FeMaSR (Chen et al., 2022) and AdaCode (Liu et al., 2023a), besides their original loss functions, we introduce the cross-entropy loss $\mathcal{L}_{code}^{token}$ to supervise the learning of 848 code index prediction. For CodeFormer (Zhou et al., 2022) and DAEFR (TSAI et al., 2024), since 849 the code index prediction is incorporated in their original implementations, we directly use the loss 850 functions used in the second **Restoration** stage of the corresponding VQ-based BIR methods. 851

Self-Attention (SA) for Feature Transformation. We replace our SinCA by SA (Dosovitskiy et al., 2021) for feature transformation in the transformer of VQ-based BIR methods. The experimental setups remain the same with those methods using our SinCA, with the exception that the SA module only takes LQ feature as the input to obtain query, key, and value. The input and output channel dimensions of SA in FeMaSR (Chen et al., 2022), AdaCode (Liu et al., 2023a), CodeFormer (Zhou et al., 2022), and DAEFR (TSAI et al., 2024) are set as 512, 256, 256, and 256, respectively.

Cross-Attention (CA) for Feature Transformation. Similar to the experiments on SA-based feature transformation mentioned above, we conduct experiments on CA (Chen et al., 2021) by replacing our SinCA in transformer with CA (Chen et al., 2021) for feature transformation. CA uses LQ feature to obtain the query and HQ codebook to obtain the key and the value. For FeMaSR (Chen et al., 2022), the input LQ feature and HQ codebook have 512 channels. For the transformers used in AdaCode (Liu et al., 2023a), CodeFormer (Zhou et al., 2022), and DAEFR (TSAI et al., 2024), both the input LQ feature and HQ codebook have 256 channels.

C MORE VISUAL COMPARISON RESULTS

Here, we provide more visual comparison results of different methods on blind image super-resolution (SR) in Figs. $9 \sim 11$ for $\times 2$ SR task and in Figs. $12 \sim 13$ for $\times 4$ SR task. More visual comparison results of blind face restoration are provided in Figs. $14 \sim 16$.



Figure 9: More visual results on $\times 2$ blind image super-resolution task. "+SinCA": employing the transformer using our SinCA for feature transformation.



Figure 10: More visual results on $\times 2$ blind image super-resolution task. "+SinCA": employing the transformer using our SinCA for feature transformation.



the transformer using our SinCA for feature transformation.



Figure 12: More visual results on ×4 blind image super-resolution task. "+SinCA": employing the transformer using our SinCA for feature transformation.



Figure 13: More visual results on ×4 blind image super-resolution task. "+SinCA": employing 1057 the transformer using our SinCA for feature transformation.



1078 Figure 14: More visual results of blind face restoration task on CelebA-Test. "+SinCA": employ-1079 ing the transformer using our SinCA for feature transformation.



Figure 15: More visual results of blind face restoration task on CelebA-Test. "+SinCA": employing the transformer using our SinCA for feature transformation.



Figure 16: More visual results of blind face restoration task on CelebChild-Test. "+SinCA":
 employing the transformer using our SinCA for feature transformation.

1110