# Mobile Image Restoration via Prior Quantization

Shiqi Chen, Jinwen Zhou, Menghao Li, Yueting Chen, Tingting Jiang

*Abstract*—In digital images, the performance of optical aberration is a multivariate degradation, where the spectral of the scene, the lens imperfections, and the field of view together contribute to the results. Besides eliminating it at the hardware level, the post-processing system, which utilizes various prior information, is significant for correction. However, due to the content differences among priors, the pipeline that aligns these factors shows limited efficiency and unoptimized restoration. Here, we propose a prior quantization model to correct the optical aberrations in image processing systems. To integrate these messages, we encode various priors into a latent space and quantify them by the learnable codebooks. After quantization, the prior codes are fused with the image restoration branch to realize targeted optical aberration correction. Comprehensive experiments demonstrate the flexibility of the proposed method and validate its potential to accomplish targeted restoration for a specific camera. Furthermore, our model promises to analyze the correlation between the various priors and the optical aberration of devices, which is helpful for joint soft-hardware design.

*Index Terms*—image processing, neural networks, optical aberration correction, priority

## I. INTRODUCTION

**A**NY digital imaging system suffers from optical aberration, and thus correcting aberration is necessary for accurate measurements [1]. Unfortunately, the optical aberration in image is affected by many factors, which is formulated as:

$$J_e = \int \mathcal{C}_e(\lambda) \cdot [I_e(h, w, \lambda) * L_e(h, w, \lambda)] d\lambda + N_e(h, w), \quad (1)$$

here $(h, w)$ indicates the coordinates of the pixel, $\lambda$ is the wavelength, $\mathcal{C}_e$ and $I_e$ denote the spectral response and the point spread function (PSF), $L_e$, $J_e$, and $N_e$ are the latent sharp image, the observed image, and the measurement noise, respectively. Note that the subscript $e$ indicates the measurement that represents the energy received by the sensor. Therefore, energy dispersion and field-of-view (FoV) clue of the optical system, spectral properties, and sensor noise together contribute to the optical aberration expression on digital images. In other words, these factors serve as the priors to facilitate the correction.

Although there are many algorithms for optical aberration correction, it still faces a few challenges for widespread application. One issue is that the existing methods generally have a limited application scope, *e.g.*, the deconvolution is

Shiqi Chen, Jinwen Zhou, and Yueting Chen are with the College of Optical Science and Engineering, Zhejiang University, Hangzhou 310000, China (e-mail: chenshiqi@zju.edu.cn).

Tingting Jiang is with the Research Center for Intelligent Sensing Systems, Zhejiang Laboratory, Hangzhou 311100, China (e-mail: eager-jtt@zhejianglab.com).

inefficient in handling spatial-varying kernels, and the deep learning method is trained for a specific device according to the data [2]. Another challenge is the quality of restoration, which is generally unsatisfactory due to the lack of sufficient information. An inherent defect of these unflexible methods is that they are incapable of mining the interaction between the digital pixel and other optical priors [3].

As mentioned above, the expression of optical aberration correlates with multiple factors. Designing a general method to utilize these priors for targeted correction is an issue worth discussing. Meanwhile, this baseline can help analyze the correlation between different priors and optical aberration correction, aiming to guide the co-design of the hardware configuration and the post-processing pipeline in the high-end imaging system.

In this letter, we propose a prior quantization model to correct the optical degradation influenced by multiple factors, where different priors corresponding to each factor are fed to the model. However, due to the content differences among the multimodal priors, they cannot directly integrate into the model for efficient post-processing. To this end, we encode the various priors into a high-dimensional latent space and characterize it by a learnable codebook. The learned code reduces the dimension of multimodal priors representation and models the global interrelations of auxiliary information. Then the model integrates the quantized prior representation into the image restoration branch by fusing it with features of different scales. Finally, we supervise the restoration in multiple scales to ensure that the cross-scale information used for fusion is accurate and valid. Instead of a black box, the model can adjust the input pixel-level prior according to the demands of the user, aiming to achieve a targeted development. Moreover, the learned codebook bridges the gap between different priors and restoration. Therefore, the proposed model has the potential to analyze the utilization of priors, which is meaningful for imaging system design.

## II. THE PROPOSED METHOD

As mentioned above, multiple factors contribute to the expression of optical aberration. Thus the model must possess the ability to perceive a complex combination and transform it into a simpler form. Directly engaging the image signal with individual priors is inefficient, where the model will spend substantial computing overhead on the auxiliary information. Therefore, we encode these priors into a high-dimension space and represent them with a learnable codebook. Since the aberration is relatively constant in the neighborhood, the optical prior of an image $x \in \mathcal{R}^{H \times W \times 3}$ can be represented by a spatial collection of codebook $z_q \in \mathcal{R}^{h \times w \times d_z}$, where $d_z$ is the dimension of the code. To effectively learn such
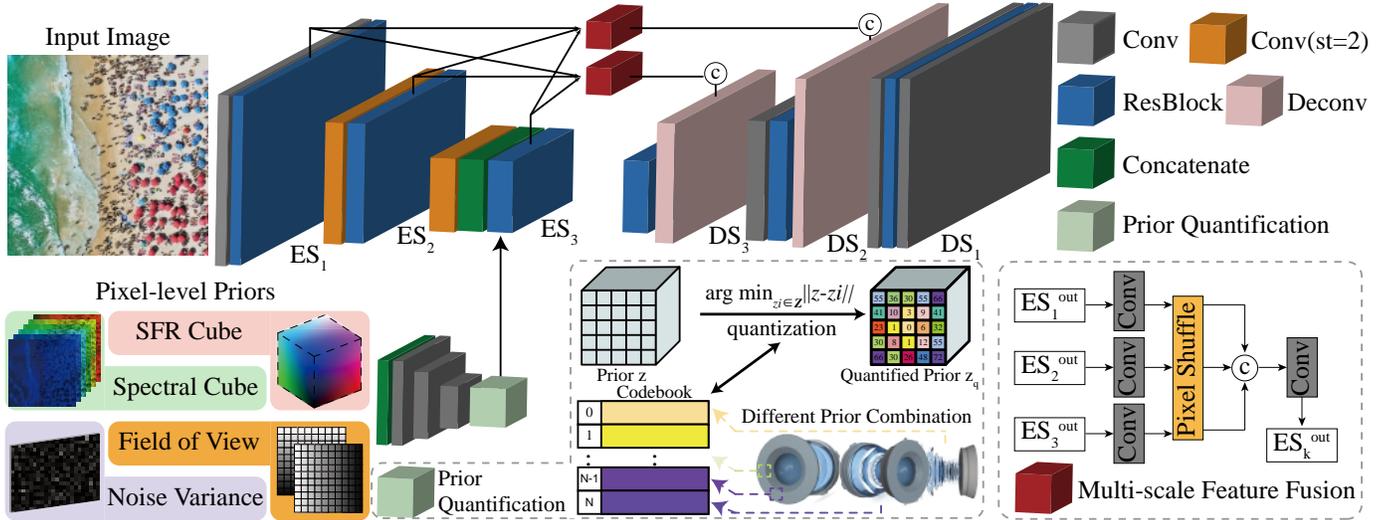
Fig. 1. The detail of the proposed prior quantization model. The layer configurations are illustrated with different colored blocks.

a codebook, we propose to exploit the spatial extraction capability of convolution and incorporate with the neural discrete representation learning [4]. First, we use CNNs to encode the priors into a latent space with the same spatial resolution ($h \times w$). Then, we embed the encoded features into the spatial code $\hat{z}_{ij} \in \mathcal{R}^{d_z}$ and quantify each code onto its closest entry in the discrete codebook $\mathcal{Z} = \{z_k\}_{k=1}^K \subset \mathcal{R}^{n_z}$:

$$z_q = \left( \arg\min_{z_k \in \mathcal{Z}} ||\hat{z}_{ij} - z_k|| \right) \in \mathcal{R}^{h \times w \times d_z}, \quad (2)$$

Third, we concatenate the quantized prior representation with the features in the restoration branch. And the subsequent ResBlock actively fused the image feature, where the prior after quantization provides clues of the spatial/channel importance to the reconstruction branch.

In most reconstruction models, the skip connections are only processed in one scale when the coarse-to-fine architecture is applied. Inspired by the dense connection between multi-scale features, we implement a multi-scale feature fusion (MFF) module to fuse the quantified representation with the features from other scales [5]. As shown in the right-bottom of Fig. 1, this module receives the outputs of different encoding scales ($ES_i$) as inputs. After adjusting the channels of each $ES_i^{out}$ by convolution, we use pixel shuffle to transform the encoded information into the same spatial resolution and then perform the fusion [6]. The output of the MFF is delivered to its corresponding decoding scales. In this way, each scale can perceive the encoded information of other scales, especially the lowest-scale features filtered by the quantified priors, resulting in improved restoration quality.

Here we discuss the loss function of our model. Due to the quantization operation in Eq. 2 is non-differentiable in backpropagation, the CNN encoder of priors cannot receive a gradient to optimize its parameters if the entire network is directly trained with end-to-end supervision. Fortunately, the strategy of the codebook alignment allows the gradient propagated from decoders to update these CNN encoders. Specifically, we keep the encoded priors approaching the

vectors of the learnable codebook. In this way, the quantization progress like a gradient estimator, allowing the CNN encoder to estimate the backpropagation and update the parameters even when the gradient is truncated in training. Thus, we apply the codebook alignment loss to realize the optimization of encoders [7]:

$$\mathcal{L}_{align} = ||sg[\hat{z}_{ij}] - z_q||_2^2 + ||\hat{z}_{ij} - sg[z_q]||_2^2, \quad (3)$$

here $\hat{z}_{ij} = \{E_1(p_1), E_2(p_2), \ldots, E_n(p_n)\}$, where $p_i$ and $E_i$ are the $i^{th}$ prior and encoder, respectively. $\{\cdot\}$ is the operation to concatenate the encoded features. $sg[\cdot]$ denotes the stop-gradient operation to ensure the loss function only guide the encoded priors and the codebook vectors approaching.

Because our network relies on the refinement in a coarse-to-fine manner, we supervise the image reconstruction on different scales. For the supervision of image content, we find that the L1 loss performs better on quantitative metrics for optical aberration correction. However, since the content loss only measures the pixel-level difference, the model does a good job restoring the low-frequency information (such as brightness and color, *etc.*) while performing poorly in restoring the textures of the scene. To prevent the limitation of pixel-level supervision, we adopt the supervision in the Fourier domain as an auxiliary loss [5]. Compared with the gradient constraints (*e.g.*, total variation) or the feature similarity on the pre-trained model (*e.g.*, perceptual loss), this item provides more information on different spatial frequencies. The loss $\mathcal{L}_{content}$ is formulated as follows:

$$\mathcal{L}_{content} = \frac{1}{t_k} \sum_{k=1}^K ||\hat{I}_k - I_k||_1 + \lambda ||\mathcal{F}(\hat{I}_k) - \mathcal{F}(I_k)||_1, \quad (4)$$

here the output of the k-th scale is $I_k$ and its corresponding ground-truth is $\hat{I}_k$, where $K$ is the number of scales. The content loss in each scale is the average on the total elements, whose number is denoted as $t_k$. $\mathcal{F}$ is the fast Fourier transform (FFT) operation that performs on different scales and the $\lambda$ is experimentally set to 0.1. The overall loss function is the sum of $\mathcal{L}_{content}$ and $\mathcal{L}_{align}$.

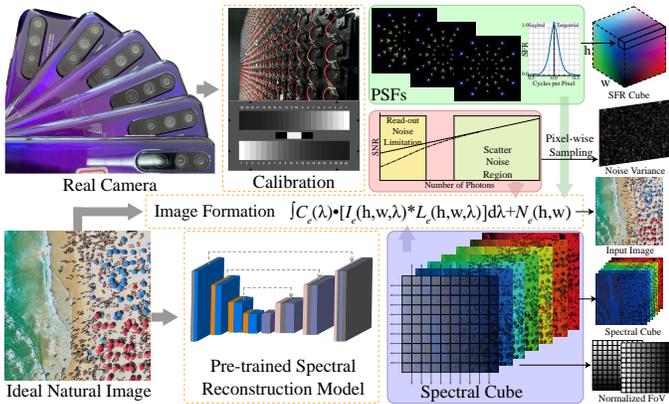| TestSet | Synthetic Evaluation | | | | | | | | Real Evaluation | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | PSNR↑ | | SSIM↑ | | VIF↑ | | LPIPS↓ | | BRISQUE↓ | | NIQE↓ | |
| SRN | 29.79 | (56.6%) | 0.9247 | (29.1%) | 0.8221 | (11.3%) | 3.428 | (48.2%) | 54.49 | (16.1%) | 5.686 | (21.3%) |
| IRCNN | 30.58 | (48.0%) | 0.9289 | (24.5%) | 0.8314 | (10.0%) | 3.151 | (43.7%) | 53.47 | (14.5%) | 5.572 | (19.7%) |
| Self-Deblur | 32.23 | (24.0%) | 0.9353 | (17.5%) | 0.8544 | (7.05%) | 2.766 | (35.8%) | 50.28 | (9.03%) | 5.323 | (16.0%) |
| GLRA | 32.47 | (19.6%) | 0.9347 | (18.2%) | 0.8662 | (5.60%) | 2.457 | (27.8%) | 49.74 | (8.04%) | 5.476 | (18.3%) |
| FDN | 33.27 | (3.39%) | 0.9402 | (12.2%) | 0.8943 | (2.28%) | 1.944 | (8.69%) | 47.42 | (3.54%) | 5.038 | (11.2%) |
| Deep Wiener | 32.02 | (27.6%) | 0.9463 | (5.7%) | 0.9007 | (1.55%) | 1.928 | (7.94%) | 46.52 | (16.8%) | 4.943 | (9.44%) |
| Ours | 33.42 | (0.0%) | 0.9517 | (0.0%) | 0.9147 | (0.0%) | 1.775 | (0.0%) | 45.74 | (0.0%) | 4.476 | (0.0%) |



Fig. 2. The synthetic flow of the training data and the priors.

## A. Synthetic Flow for training data and prior

Since optical aberration is highly related to multiple factors, we construct a comprehensive dataset for correction. The detailed synthetic flow of the training data and the priors is shown in Fig. 2. First, we calibrate the PSFs and the noise factors of many mobile cameras, where the PSFs are used to simulate the degradation and calculate the SFR prior. Second, we use the pre-trained MST++ [8] to convert RGB images to hyperspectral data for the spectral prior. Third, the FoV prior is consist of the $(h, w)$ pixel coordinates that normalized to $[-1, 1]$ [9], [10]. Finally, the multispectral data is engaged with optical aberration, where the procedure is the same as Eq. 1.

## III. EXPERIMENTAL RESULTS

### A. Datasets, Metrics, and Training Settings

As illustrated above, we use the natural images in DIV8K to synthetic the training data-pairs and priors [11]. The training dataset consists of 800 image pairs. And the resolution of image is rescaled to $3000 \times 4000$ to align with the real cameras. In the case of real-world evaluation, we capture many photographs with multiple mobile terminal of Huawei Honor 20 and iPhone 12. As for the metrics, the PSNR, SSIM, VIF [12], and LPIPS [13] evaluate the model with reference. The BRISQUE [14], NIQE [15] are used to assess the restoration of the captrued photographs. In the training, we crop the whole image to $256 \times 256$ pixels and form minibatches of 8 images. For a fair comparisons, the priors are concatnated with the images and fed into the competing models. We optimize the

model with Adam in $\beta_1 = 0.5$ and $\beta_2 = 0.999$. The initial learning rate is $10^{-4}$ and then halved every 50 epoch.

### B. Quantitative Assessment to SOTA Methods

We compare the proposed prior quantization model with the competing algorithms designed for optical aberration correction. All these methods are retrained with the same dataset until convergence. In the assessment of spatial-various aberration correction, we evaluate the quantitive indicators on the synthetic dataset. Tab. I reports the performance of various approaches on the synthetic dataset. The blind methods (SRN [16], IRCNN[17], SelfDeblur[18], GLRA[19]) successfully deal with the optical degradation of one camera. But there are various mobile cameras with different optical aberrations, and the blind manner fails to acquire the precise degradation clue only from the feature of the image. The non-blind approaches (FDN[20], Deep Wiener[21]), which feed the pre-trained models with PSFs to adapt to a specific camera, have a similar idea to ours. However, since the PSFs is a high-dimensional representation of degradation, only a few representative PSFs are selected to be fed into the model (only 5 PSFs in [21]), where the spatial relationship of the PSFs across the whole FoV is abandoned. Otherwise, the complexity will be extremely high. Different from this method, we use the rearranged SFR cube to represent the dispersion for each FoV, which is a lower-dimensional representation of degradation. Therefore, with relatively lower computational overhead, our model can obtain pixel-level guidance to achieve better restoration.

### C. Real Restoration Comparisons

Moreover, we test these algorithms with the degraded photographs taken from real cameras, and the results are visualized in Fig. 3. Since the priors acquired by the Deep Wiener do not correspond to the pixel neighborhood, the model needs attention mechanism to determine the specific degradation of the input. These overheads increase the processing burden, resulting in the failure to obtain efficient restoration. Compared with other algorithms and the built-in ISP, our model efficiently integrates various priors and characterizes them with the quantized codebook vectors. Therefore, our method reduces the complexity of multi-task post-processing, and achieves a comprehensive improvement in image quality. Another advantage of our non-blind model is that it can be designed as a post-processing system with better generalization.
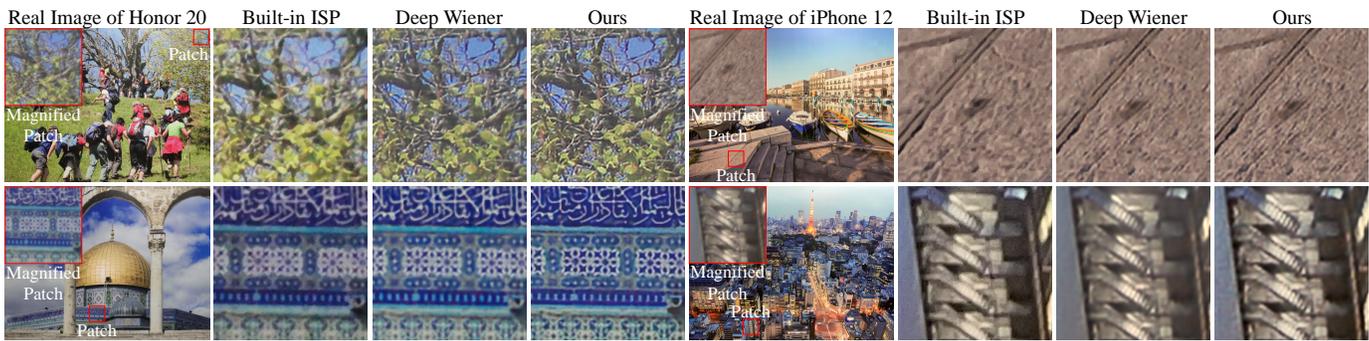
Fig. 3. Real image restoration comparison, where the position is highlighted in red. See more experimental results in Supplement 1.
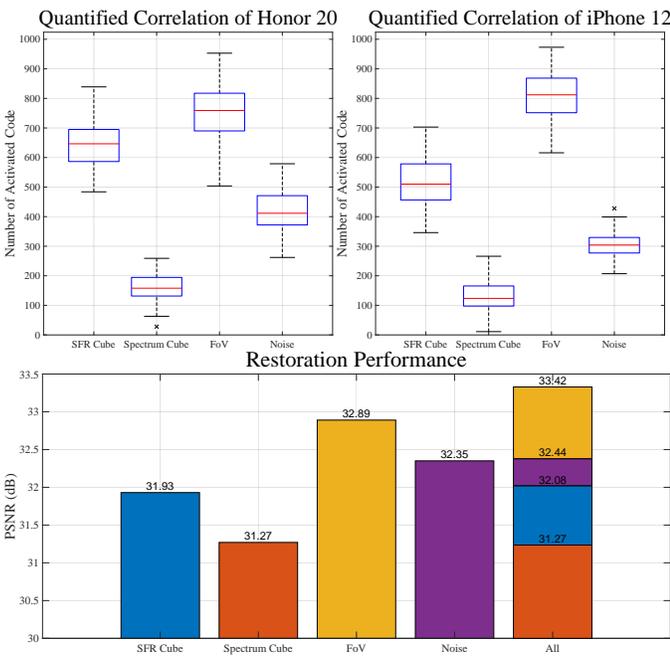


Fig. 4. The correlation between priors and optical aberration correction in different mobile cameras.

into analyze the proceeds of introducing each prior to the restoration. When evaluating a specific prior, we zero out all other auxiliaries except this one and test the learned model with the raw photographs taken by various mobile cameras, counting the number of the activated codes in inference process. The larger number of activated codebooks indicates the more relevance between this prior and optical aberration. The result of the assessment is shown at the top of Fig. 4. We find that the SFR and FoV priors will activate more codes in inference, which means they are highly correlated with aberration and play critical roles in correction. We note that the restoration of the Honor 20 pro activates more code in the SFR and noise priors evaluation when compared with the iPhone 12, which attributes to its uneven optical aberration of the optics and lower SNR of the sensor. This experiment also demonstrates that the number of activation hints at the utilization of different priors, providing a brand new issue for imaging system assessment. Moreover, we evaluate the correlation between the restoration indicators and the priors on the synthetic testset (shown at the bottom of Fig. 4). Different from the number of activated codes, we note that the noise is more substantial on the restoration performance, which may be put down to the PSNR that calculates the absolute error in the pixel scale. On the other hand, this phenomenon indicates that we can strip the noise and the optical aberration to balance the model with restoration indicators and computational efficiency.

In our experiment, we use the model trained on the simulation dataset to post-process the real images taken by various mobile cameras. The results in the left and right of Fig. 3 demonstrate that the model pre-trained on synthetic data achieves perfect generalization ability on specific devices (Honor 20 pro and iPhone 12), where fine-tuning is not necessary. Therefore, the plug-and-play feature indicates that the learned model is promising to replace the camera-specific ISP system with a flexible model. For the detailed comparisons, we refer the readers to the supplementary for more restorations on different manufactured samples and more the quantitive results of system imaging quality.

## D. Prior Correlation Analyse

As mentioned above, the optical aberration is explicitly correlated to multiple factors when imaging at a fixed distance and this auxiliary information can benefit the restoration. However, obtaining so much assistance in real-time imaging comes at an unaffordable expense. We use the learned model

## IV. CONCLUSION

In summary, we develop a deep learning model to utilize multiple priors for optical aberration correction. Even though the content of the priors varies a lot from each other, the proposed quantization strategy efficiently fuses these factors and guides the correction of optical aberration. Comprehensive experiments show that integrating all the related information does benefit the restoration of optical aberration, and the proposed model can generalize to different mobile devices without finetuning. Moreover, the learned model quantifies the correlation and significance of different priors for the correction procedure. Thus, the aberration correction in the post-processing system can be more efficient when the critical clues are obtained. In the future, we will put efforts into engaging the proposed model with different imaging devices, aiming to realize efficient correction and better generalization.

## REFERENCES

[1] T. Yue, J. Suo, J. Wang, X. Cao, and Q. Dai, "Blind optical aberration correction by exploring geometric and visual priors," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1684–1692, 2015.

[2] T. Eboli, J.-M. Morel, and G. Facciolo, "Fast two-step blind optical aberration correction," *European Conference on Computer Vision*, pp. 693–708, 2022.

[3] S. Chen, T. Lin, H. Feng, Z. Xu, Q. Li, and Y. Chen, "Computational optics for mobile terminals in mass production," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–16, 2022.

[4] A. Van Den Oord, O. Vinyals *et al.*, "Neural discrete representation learning," *Advances in neural information processing systems*, vol. 30, 2017.

[5] S.-J. Cho, S.-W. Ji, J.-P. Hong, S.-W. Jung, and S.-J. Ko, "Rethinking coarse-to-fine approach in single image deblurring," *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4641–4650, 2021.

[6] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1874–1883, 2016.

[7] J. Duan, L. Chen, S. Tran, J. Yang, Y. Xu, B. Zeng, and T. Chilimbi, "Multi-modal alignment using representation codebook," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15 651–15 660, 2022.

[8] Y. Cai, J. Lin, Z. Lin, H. Wang, Y. Zhang, H. Pfister, R. Timofte, and L. Van Gool, "Mst++: Multi-stage spectral-wise transformer for efficient spectral reconstruction," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 745–755, 2022.

[9] S. Chen, H. Feng, K. Gao, Z. Xu, and Y. Chen, "Extreme-quality computational imaging via degradation framework," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2632–2641, 2021.

[10] Q. Sun, C. Wang, Q. Fu, X. Dun, and W. Heidrich, "End-to-end complex lens design with differentiate ray tracing," *ACM Trans. Graph.*, vol. 40, no. 4, jul 2021. [Online]. Available: https://doi.org/10.1145/3450626.3459674

[11] S. Gu, A. Lugmayr, M. Danelljan, M. Fritsche, J. Lamour, and R. Timofte, "Div8k: Diverse 8k resolution image dataset," in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2019, pp. 3512–3516.

[12] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on image processing*, vol. 15, no. 2, pp. 430–444, 2006.

[13] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[14] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on image processing*, vol. 21, no. 12, pp. 4695–4708, 2012.

[15] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.

[16] X. Tao, H. Gao, X. Shen, J. Wang, and J. Jia, "Scale-recurrent network for deep image deblurring," *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8174–8182, 2018.

[17] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep cnn denoiser prior for image restoration," *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3929–3938, 2017.

[18] D. Ren, K. Zhang, Q. Wang, Q. Hu, and W. Zuo, "Neural blind deconvolution using deep priors," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3341–3350, 2020.

[19] W. Ren, J. Zhang, L. Ma, J. Pan, X. Cao, W. Zuo, W. Liu, and M.-H. Yang, "Deep non-blind deconvolution via generalized low-rank approximation," *Advances in neural information processing systems*, vol. 31, 2018.

[20] J. Kruse, C. Rother, and U. Schmidt, "Learning to push the limits of efficient fft-based image deconvolution," *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4586–4594, 2017.

[21] T. Lin, S. Chen, H. Feng, Z. Xu, Q. Li, and Y. Chen, "Non-blind optical degradation correction via frequency self-adaptive and finetune tactics," *Opt. Express*, vol. 30, no. 13, pp. 23 485–23 498, Jun 2022. [Online]. Available: https://opg.optica.org/oe/abstract.cfm?URI=oe-30-13-23485