# Learning Pixel-Adaptive Weights for Portrait Photo Retouching

Binglu Wang    Chengzhe Lu*    Dawei Yan*    Yongqiang Zhao
* Equal Contribution

## Abstract

*Portrait photo retouching is a photo retouching task that emphasizes human-region priority and group-level consistency. The lookup table-based method achieves promising retouching performance by learning image-adaptive weights to combine 3-dimensional lookup tables (3D LUTs) and conducting pixel-to-pixel color transformation. However, this paradigm ignores the local context cues and applies the same transformation to portrait pixels and background pixels when they exhibit the same raw RGB values. In contrast, an expert usually conducts different operations to adjust the color temperatures and tones of portrait regions and background regions. This inspires us to model local context cues to improve the retouching quality explicitly. Firstly, we consider an image patch and predict pixel-adaptive lookup table weights to precisely retouch the center pixel. Secondly, as neighboring pixels exhibit different affinities to the center pixel, we estimate a local attention mask to modulate the influence of neighboring pixels. Thirdly, the quality of the local attention mask can be further improved by applying supervision, which is based on the affinity map calculated by the groundtruth portrait mask. As for group-level consistency, we propose to directly constrain the variance of mean color components in the Lab space. Extensive experiments on PPR10K dataset verify the effectiveness of our method,* e.g. *on high-resolution photos, the PSNR metric receives over 0.5 gains while the group-level consistency metric obtains at least 2.1 decreases.*

## 1. Introduction

With promising development, the image enhancement community has observed significant advances and received numerous well-performed methods [10, 23, 28, 30, 35, 38, 39]. Recently, Liang *et al.* focused on the visual quality of portrait photos and proposed the portrait photo retouching task. When recording meaningful moments by a photograph, the human often acts as an essential character in the photo, and people usually take multiple photos in similar scenes. However, the quality of raw photos taken directly by a cell phone or camera is easily affected by weather, lighting and shooting equipment, *etc*. Meanwhile, the tone of multi-
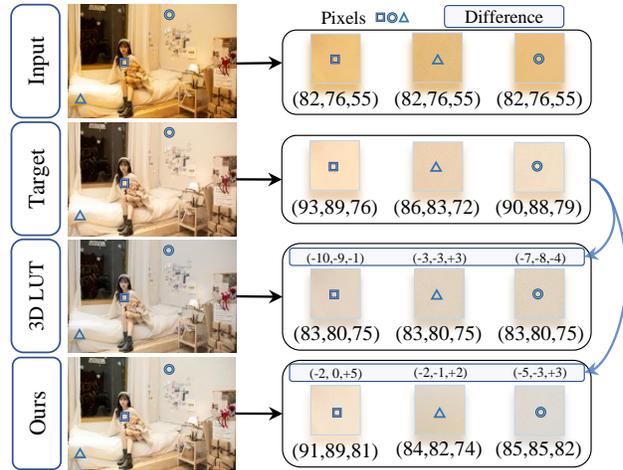


Figure 1. Comparison of retouching quality. The outputs of 3D LUT are identical for the same RGB inputs, but the outputs of our method vary according to context pixels. The RGB values of each pixel are shown in parentheses.

ple photos taken on the same subject at the same scene may also become different. Liang *et al.* [19] indicated two critical factors for portrait photo retouching: one is the human-region priority, *i.e.* the human region in the photo should receive more refined retouching compared to the background region. The other is group-level consistency, *i.e.* a group of portrait photos should be retouched to a consistent tone, even though the shooting conditions are different.

The lookup table-based method [35] exhibits significant performance for portrait photo retouching. In particular, Zeng *et al.* [35] proposed an image-adaptive 3D LUT method that combines multiple lookup tables with a convolutional neural network. Recently, based on 3D LUT, Liang *et al.* [19] weighted the MSE loss to emphasize the human region and build a baseline method for portrait photo retouching. Actually, in the photo retouching process, experts usually perform different operations on portrait and background regions. As shown in Fig. 1, even raw pixels exhibit the same RGB values, and their retouching targets could be different because each pixel has its specific contexts. However, the current method [19] tackles the whole image and only generates the image-adaptive weights to combine

the looking-up results. In contrast, when contextual information is considered, our proposed method can adaptively adjust the retouching results and achieve high-quality enhancement.

An intuitive way to consider the contextual information of a pixel is to select an image patch containing $k \times k$ pixels and predict pixel-adaptive weights to combine the looking-up results. Besides, in an image patch, the contextual pixels may exhibit different affinities with respect to the center pixel. Thus, we propose to enlarge the influence of homogeneous pixels (*e.g.* two portrait pixels) and suppress the influence of heterogeneous pixels (*e.g.* a portrait pixel and a background pixel) via learning a pixel-adaptive local attention mask. Moreover, considering each photo is accompanied by a portrait mask, we can construct the groundtruth affinity maps and apply pixel-wise supervision, which can further improve the quality of local attention masks. The above processes form the *Local-context Aware Module* (LAM). It explores contextual cues, predicts pixel-adaptive weights, and brings complementary information to image-adaptive LUT weights. Consequently, our proposed method can effectively integrate lookup table results and achieve high-quality retouching results.

In addition to the retouching of individual photos, ensuring the consistent color temperatures and tones of a group of photos is also important for the portrait photo retouching task. Liang *et al.* [19] simulate the group-level variations by cropping two patches from the same image and applying different transformations. We propose a group-style aware module, which is simple to implement but can effectively guide a group of photos to exhibit consistent color temperatures and tones. In particular, given a group of photos, we first convert the enhanced photo and all other target photos to the Lab space and then constrain the variance of mean color components to achieve the group-style aware retouching.

The contributions of this paper are summarized as follows:

- We reveal the influence of contextual cues to the portrait photo retouching task and design the local-context aware module to explicitly consider context pixels and predict pixel-adaptive weight.

- We propose the group-style aware module, which performs in a simple manner but can effectively ensure consistent color temperatures and tones among a group of enhanced photos.

- Extensive experiments on the large-scale PPR dataset demonstrate the efficacy of our proposed method, *e.g.* given high-resolution portrait photos, the PSNR metric receives over 0.5 gains, and the group-level consistency metric obtains at least 2.1 decreases.

## 2. Related work

**Learning-based image restoration.** Image restoration [2, 5, 8, 20, 31, 34, 42] aims to remove image distortions caused by various situations. Since manual retouching requires a mass of expert knowledge and subjective judgment, many learning-based image restoration methods come out recently. For image exposure and low-light issues, Wang *et al.* [31] propose a neural network for photo enhancement using deep lighting estimation to fix underexposed photos. Recently Mahmoud *et al.* [1] propose a model that can perform multi-scale exposure correction for images. Guo *et al.* [7] propose a model based on a luminance enhancement curve, which is iterated continuously to achieve gradual contrast and luminance enhancement of low-light images. As for the noise reduction task, Huang *et al.* [11] present a model which generates two independent noise-bearing images of similar scenes to train the noise reduction network.

**Learning-based image retouching.** Image retouching aims to improve the quality of raw images and has received a large number of significant advances [3, 6, 9, 12, 14–16, 22, 24, 27, 41]. Focusing on real-time processing of high-resolution images (*e.g.* over 24 megapixels), Gharbi *et al.* proposed the HDRNet [6] that put most of the computation on downsampled images. He *et al.* [9] proposed the CSRNet that extracted global features with normal CNN and modulated the features after $1 \times 1$ convolution. Some researchers, such as Sean *et al.* [22] and Han-Ul *et al.* [15], adjusted the research paradigm and introduced residuals into the model in order to pursue an enhancement over the original image instead of training the image directly. In addition to expert retouching, Han-Ul *et al.* [16] developed PieNet to solve the problem of personalization in image enhancement. These excellent works demonstrate that image enhancement has splendid results based on paired datasets. However, the collection of paired datasets is expensive. Therefore, image enhancement methods based on unpaired datasets started to receive attention. There are some GAN-based methods [3, 12], while do not require pairs of input images and target images. Meanwhile, Jongchan *et al.* [24] proposed the Distort-and-Recover method focusing on image color enhancement, which also only required high-quality reference images for training. Although existing image enhancement models have good generalizability, our work focuses on portrait photo retouching, which emphasizes both the human-region priority and group-level consistency, leaving it an under-explored task.

In addition, HDRNet [6] can be regarded as a masterpiece along the lines of bilateral filter. For image processing, the features extracted by the network are complete. The applicability of this method is very broad, considering that the human photographer's retouching process also considers only these features. CSRNet [9] used an end-to-end approach that is easy to train. By utilizing an expanded con-
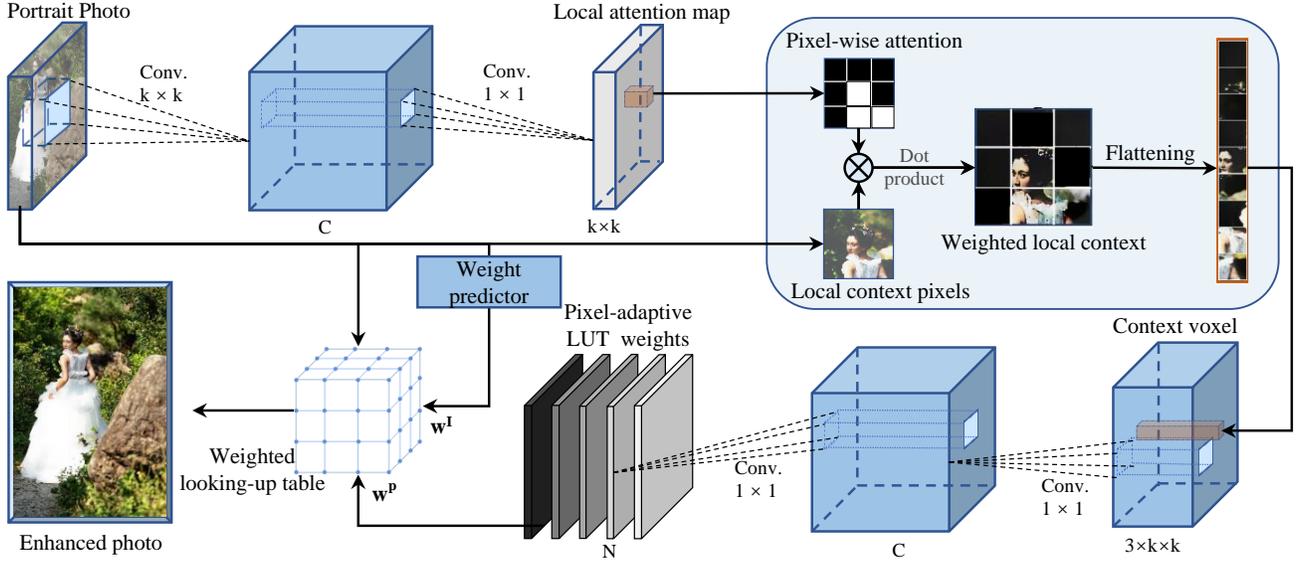
Figure 2. Framework of the proposed local-context aware module. Given a portrait photo, we first elaborately estimate pixel-wise attention and then adaptively modulate local context pixels. Afterward, neighboring pixels are flattened to form the context voxel, which are used to predict pixel-adaptive LUT weights $\mathbf{w}^p$. Meanwhile, a weight predictor estimates the image-adaptive weights $\mathbf{w}^I$. Finally, the weighted looking-up table operation generates the enhanced photo.

volutional layer, CSRNet can extend the field of perception without the loss of the receptive field.

**Deep learning methods with 3D LUT.** In previous studies, Zeng *et al*. [35] combined the deep learning-based weight predictor and 3D LUT for the first time, which achieved outstanding results. In particular, 3D LUT learned a weight predictor based on the convolutional network, which tackled low-resolution images and estimated image-adaptive weights. Then, multiple lookup tables were dynamically combined to capture image characteristics and perform the enhancement. The 3D LUT was a representative work that combined the deep learning paradigm with the traditional image enhancement paradigm, which presented the community with a new perspective on how to implement image enhancement. Subsequently, Jo *et al*. [13] implemented a similar lookup process from input to output values by training a deep super-resolution network and transferring the output values of the deep learned model into the LUT. A lot of works show that the traditional LUTs can be well integrated with image enhancement tasks. In addition, due to the hardware-friendly nature, LUTs have been used in camera imaging pipelines to retrieve precomputed output from memory.

**Contextual information in image enhancement.** Contextual information can be taken into account in image enhancement to get better results. Its importance has been demonstrated by a number of distinguished researchers in recent work. By effectively modeling contextual information, many deep learning-based methods [4, 18, 33] were

able to remove rain marks, distortions, and shadows from images. Zhang *et al*. [40] proposed a context reasoning attention network to adaptively modulate the convolution kernel according to the global context and apply it to super-resolution tasks. Besides, there was also some work that applied contextual information to complete a missing region [36, 37]. The context information also plays an important role on other computer vision domains, *e.g*. image classification [29], saliency detection [21]. In particular, Liu *et al*. [21] proposed PiCANet to learn informative contextual features of each pixel and then embed them into a saliency detection network. In the photo retouching domain, the contextual information is still under-explored. In this work, we make an early exploration to study the influence of contextual information on the portrait photo retouching task and observe sufficient improvements.

**Spatial-aware 3D LUT.** Spatial-aware 3D LUT [32] is a most relevant work to ours, though our work is independently developed with [32]. The insight of learning pixel-adaptive weights is similar. However, we precisely model local context cues by a convolutional layer with kernel size $k \times k$, while the UNet architecture in [32] considers the entire image and is affected by both local contexts and global information. As a result, in [32], the global information in the UNet branch is somewhat redundant with the image-adaptive weights, also has the risk of bringing noises to the retouching process. Besides, [32] employed the UNet architecture and maintained multiple lookup tables. It is more complex than our local-context aware module.

# 3. Methodology

## 3.1. Image enhancement based on 3D LUT

3D LUT defines a 3D grid consisting of $M^3$ elements, where $M$ indicates the number of bins in each color channel. Given an input image $I$, the looking-up table method makes transformations and generates the enhanced output $O$. The transformation operation can be formulated as:

$$O_{(i,j,k)} = \varphi(I^R_{(i,j,k)}, I^G_{(i,j,k)}, I^B_{(i,j,k)}), \tag{1}$$

where $\varphi(\cdot)$ defines a pixel-to-pixel mapping via a looking-up table, in practice, $M{=}33$ is a common setting, which contains 108K parameters and achieves a good balance between inference speed and enhancement quality.

To tackle dramatic variations among multiple photos, Zeng *et al.* [35] capture the holistic cues within an image and predict image-adaptive weight $\mathbf{w}^I = [\mathbf{w}^I_1, \mathbf{w}^I_2 \ldots \mathbf{w}^I_n]$ to adjust the looking-up table results. Thus, the transformation operation with LUT weight can be formulated as follows:

$$
\begin{aligned}
O_{(i,j,k)} &= \Phi(I^R_{(i,j,k)}, I^G_{(i,j,k)}, I^B_{(i,j,k)}, \mathbf{w}^I_1 \ldots \mathbf{w}^I_n) \\
&= (\sum_{n=1}^{N} \mathbf{w}^I_n \times \varphi_n)(I^R_{(i,j,k)}, I^G_{(i,j,k)}, I^B_{(i,j,k)}),
\end{aligned} \tag{2}
$$

where $\Phi(\cdot)$ indicates the image enhancement process and $\varphi_n$ is the mapping operation of the $n^{th}$ looking-up table.

## 3.2. Local-context aware module

In the portrait photo retouching task, experts usually apply different adjustments about color temperatures and tones for human regions and backgrounds. Thus, the context cues should be carefully considered in the retouching process.

As shown in Fig. 2, we propose the local-context aware module to predict pixel-adaptive weights and precisely modulate the looking-up table results. Given a portrait photo $\mathbf{x}$, each pixel considers its $k \times k$ neighboring pixels to perceive local context cues. To this end, we first employ a convolutional layer with kernel size $k \times k$ to extract context features $\mathbf{f}_{\text{ctx}} = \text{ReLU}(\mathbb{C}^{k \times k}(\mathbf{x}))$. Then, we predict local attention map $\mathbf{a} \in \mathbb{R}^{(k \times k) \times H \times W}$:

$$\mathbf{a} = \text{Softmax}(\mathbb{C}^{1 \times 1}(\mathbf{f}_{\text{ctx}})), \tag{3}$$

where "Softmax" indicates softmax normalization among $k^2$ attention elements. $\mathbb{C}^{k \times k}(\cdot)$ means a $k \times k$ convolutional layer.

Considering a pixel $x_{i,j}$, the pixel-wise attention weight $\mathbf{a}_{i,j}$ can be used to modulate its neighboring context pixels via dot product. Then, similar context pixels would be strengthened while the influence of intrusive pixels would be weakened. Next, all neighboring pixels are flattened to a vector, and context vectors from all spatial locations form
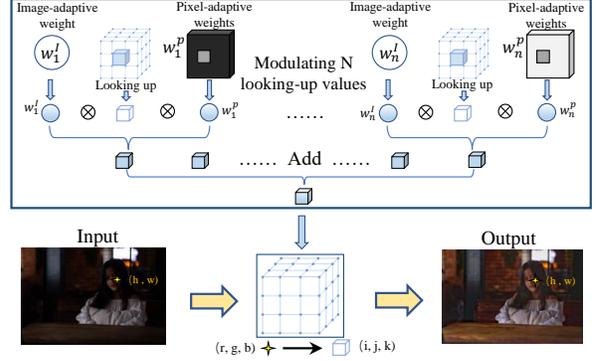


Figure 3. Portrait photo retouching process with pixel-adaptive weights and image-adaptive weights. The lookup table generates a retouching result for a pixel, which is then jointly modulated by pixel-adaptive weight $\mathbf{w}^p$ and image-adaptive weight $\mathbf{w}^I$. Finally, multiple retouching results are added to enhance the input photo.

the neighboring context voxel $\mathbf{v} \in \mathbb{R}^{(3 \times k \times k) \times H \times W}$. Afterwards, we employ a convolutional layer to tackle the context voxel and obtain feature $\mathbf{f}_{\text{vox}} = \text{ReLU}(\mathbb{C}^{1 \times 1}(\mathbf{v}))$. In the end, a convolutional layer with kernel size 1 can adaptively predict LUT weights for each pixel:

$$\mathbf{w}^p = \mathbb{C}^{1 \times 1}(\mathbf{f}_{\text{vox}}). \tag{4}$$

At each spatial location, the pixel-adaptive weights $\mathbf{w}^p$ contains $N$ elements, being responsible for zooming or shrinking the looking-up table results.

## 3.3. Fuse the multiple 3D LUTs.

**Image-adaptive Weight.** Given $N$ 3D LUTs, the fusion is performed according to the weight $\mathbf{w}^I$ output by the image classifier. In our work, the number of 3D LUTs is set to 5 based on ablation experiments.

**Pixel-adaptive Weight.** The fuse effect based on classifier and 3D LUTs is useful, but it can be further improved by considering the contextual information of each pixel. In this paper, we consider the contextual information for each pixel and estimate pixel-adaptive attention weights $\mathbf{w}^p$. By jointly considering image-adaptive attention weights $\mathbf{w}^I$ and pixel-adaptive attention weights $\mathbf{w}^p$, our photo retouching method can perceive both the image-level holistic information and pixel-level contextual cues and achieve high-quality image enhancement.

Fig. 3 illustrates the fuse process by tackling a certain pixel located at (h,w) within the photo. The program sends the photo to a group of 3D LUTs and obtains the converted result. The result is then integrated with the $\mathbf{w}^I$ and $\mathbf{w}^p$ to obtain the RGB value at the corresponding pixel point (h, w) of the output photo. This procedure is performed for each pixel, which generates the complete enhanced photo.
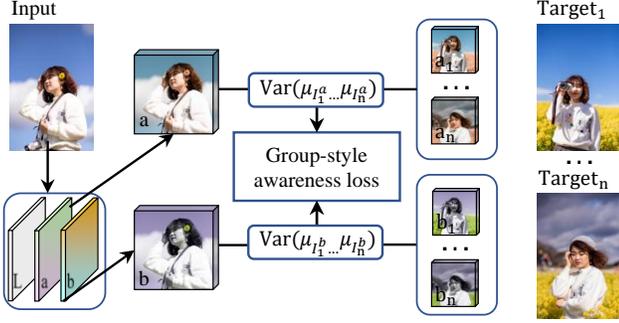
Figure 4. Calculation process of the group-style aware loss. The retouched photo and all other target photos from the same group are first converted to Lab space. Then, we calculate the variance of mean color components in a-channel and b-channel to obtain the group-style aware loss.

## 3.4. Group-style aware module

The group-level consistency characteristic emphasizes more the uniformity of style among a group of photos rather than the retouching effect of a single photo. In [19], given a group of enhanced image, Liang *et al.* [19] first convert each image from RGB to Lab color space, and then select the $a$ and $b$ channels to measure the variance of mean color components. Based on the combination of $a$ and $b$ channels, Liang *et al.* [19] defines the GLC metric.

The group-style aware module strategy is to enumerate the color of all pixels of the photo and calculate the mean value $\mu$, calculate the variance between $\mu$ and the value of the same group of target photos. The smaller the variance, the closer the values in the two channels, and the more consistent the color temperatures and tones of the image.

As shown in Fig. 4, given a group of retouched photos $[\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2, ..., \hat{\mathbf{p}}_n]$, we first calculate their mean values from the a-channel and b-channel, and then do the same for images generated by the network model retouching. Calculating the variance of these values can further obtain group-style aware loss.

Given a group of photos, we find it is inconvenient and inefficient to simultaneously forward all photos through the network and calculate the group-style aware loss $\mathcal{L}_{GAM}$, because each group contains a different number of photos. Instead, we propose to compare a retouched photo with all other target photos from the same group, where the mean value of a-channel and b-channel can be calculated in advance.

## 3.5. Loss function

As we adopt the lookup table paradigm to retouch photos, we follow [35], adopt the MSE loss $\mathcal{L}_{mse}$, the smooth regularization loss $\mathcal{R}_s$ and monotonicity regularization loss $\mathcal{R}_m$ as the basic loss terms, which can be calculated as:

$$\mathcal{L}_{LUT} = \mathcal{L}_{mse} + \lambda_s \mathcal{R}_s + \lambda_m \mathcal{R}_m, \tag{5}$$
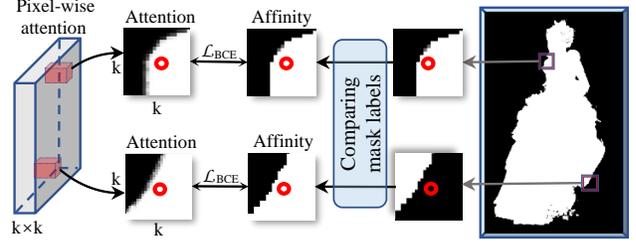


Figure 5. Calculation process of the edge supervision loss. Considering a pixel seated on the portrait border area, we compare the mask labels between the center pixel and the neighboring pixels, obtain the affinity map. Afterward, the affinity map can guide the pixel-wise attention map via calculating the BCELoss $\mathcal{L}_{BCE}$.

where $\lambda_s$ and $\lambda_m$ are trade-off coefficients, we follow [35] and set $\lambda_s = 1 \times 10^{-4}$, $\lambda_m = 10$. The MSE Loss $\mathcal{L}_{mse}$ ensures the content consistency between enhancing result and target photo. Smooth Regularization $\mathcal{R}_s$ and Monotonicity Regularization $\mathcal{R}_m$ are used to ensure the output values of the LUT are smoothed and monotonic, and the relative brightness and saturation of the input RGB values are maintained, ensuring natural enhancement results.

**Edge supervision loss.** In the proposed local-context aware module, we predict the local attention mask **a** and employ the pixel-wise attention $a_{i,j}$ to modulate neighboring context pixels. We find apply supervision can improve the quality of local attention mask. As shown in Fig. 5, given the human region mask of each photo, we can obtain the affinity of a pixel and its neighboring pixels via pairwise inclusive-or operation. This affinity map can serve as the learning target for the local attention mask, and we guide the learning process by calculating the BCELoss. Given a mask, most pixels belong to the human region or background region, where the inclusive-or operation would generate an all-one mask. This would cause dramatically imbalanced loss values and impact the quality of the local attention mask. To alleviate this issue, we only calculate the loss on edge pixels, whose contextual pixels include both foreground and background, and name this loss edge supervision loss $\mathcal{L}_{edge}$.

**Group-style aware loss.** The GAM-loss can be calculated as follows:

$$\mathcal{L}_{GAM} = \lambda_{GAM}(Var(\mu_{I_1^a} \ldots \mu_{I_i^a} \ldots \mu_{I_n^a})$$
$$+ Var(\mu_{I_1^b} \ldots \mu_{I_i^b} \ldots \mu_{I_n^b})), \tag{6}$$

where $\lambda_{GAM}$ is the coefficient of $\mathcal{L}_{GAM}$. In experiments, we find $\lambda_{GAM} = 0.001$ is a proper choice. $\mu_{I_n^a}$ and $\mu_{I_n^b}$ indicates the average value of the $n^{th}$ input photo in a-channel and b-channel, respectively. $Var(\cdot)$ stands for calculating the variance of all elements.

In summary, the total loss is given by

$$\mathcal{L} = \mathcal{L}_{LUT} + \mathcal{L}_{edge} + \mathcal{L}_{GAM}. \tag{7}$$

Table 1. Comparison of portrait photo retouching results on PPR10K dataset. a/b/c indicate the dataset retouched by three experts. The "LR" column and the "HR" column report the test results of each item on 360P and 4K ~ 8K photos, respectively.

| Method | Dataset | $PSNR \uparrow$ | | $\Delta E_{ab} \downarrow$ | | $PSNR^{HC} \uparrow$ | | $\Delta E_{ab}^{HC} \downarrow$ | | $M_{GLC} \downarrow$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LR | HR | LR | HR | LR | HR | LR | HR | LR | HR |
| HDRNet [6] | PPR10K-a | 23.93 | 23.06 | 8.70 | 9.13 | 27.21 | 26.58 | 5.65 | 5.84 | 14.83 | 14.37 |
| CSRNet [9] | PPR10K-a | 22.72 | 22.01 | 9.75 | 10.20 | 25.90 | 25.19 | 6.33 | 6.73 | 12.73 | 12.66 |
| 3D LUT [35] | PPR10K-a | 25.64 | 25.15 | 6.97 | 7.25 | **28.89** | 28.39 | 4.53 | 4.71 | 11.47 | 11.05 |
| Liang *et al.* [19] | PPR10K-a | **25.99** | **25.55** | **6.76** | **7.02** | 28.29 | **28.83** | **4.38** | **4.55** | 10.81 | 10.32 |
| Liang *et al.* [19]+GLC | PPR10K-a | 25.31 | 24.60 | 7.30 | 7.75 | 28.56 | 27.86 | 4.75 | 5.03 | **9.95** | **9.68** |
| Ours (LAM) | PPR10K-a | **26.23** | **26.14** | **6.62** | **6.66** | 29.53 | 29.41 | 4.29 | 4.32 | 10.77 | 10.71 |
| Ours (LAM+GAM) | PPR10K-a | 25.66 | 25.57 | 7.42 | 7.47 | 28.96 | 28.85 | 4.82 | 4.85 | **6.67** | **6.66** |
| HDRNet [6] | PPR10K-b | 23.96 | 23.51 | 8.84 | 9.13 | 27.21 | 26.55 | 5.74 | 5.92 | 13.21 | 13.04 |
| CSRNet [9] | PPR10K-b | 23.76 | 23.29 | 8.77 | 9.28 | 27.01 | 26.62 | 5.68 | 5.90 | 11.82 | 11.73 |
| 3D LUT [35] | PPR10K-b | 24.70 | 24.30 | 7.71 | 7.97 | 27.99 | 27.59 | 4.99 | 5.16 | 9.90 | 9.52 |
| Liang *et al.* [19] | PPR10K-b | **25.06** | **24.66** | **7.51** | **7.73** | **28.36** | **27.93** | **4.85** | **5.00** | 9.87 | 9.60 |
| Liang *et al.* [19]+GLC | PPR10K-b | 24.52 | 23.81 | 7.93 | 8.42 | 27.82 | 27.12 | 5.12 | 5.44 | **9.01** | **8.73** |
| Ours (LAM) | PPR10K-b | **25.35** | **25.28** | **7.31** | **7.34** | 28.63 | 28.55 | 4.73 | 4.75 | 9.29 | 9.25 |
| Ours (LAM+GAM) | PPR10K-b | 24.97 | 24.90 | 8.10 | 8.13 | 28.27 | 28.19 | 5.24 | 5.26 | **6.56** | **6.57** |
| HDRNet [6] | PPR10K-c | 24.08 | 23.66 | 8.87 | 9.05 | 27.32 | 26.93 | 5.76 | 5.99 | 14.76 | 14.28 |
| CSRNet [9] | PPR10K-c | 23.17 | 22.85 | 9.45 | 9.87 | 26.47 | 26.09 | 6.12 | 6.54 | 14.64 | 14.22 |
| 3D LUT [35] | PPR10K-c | 25.18 | 24.78 | 7.58 | 7.85 | 28.49 | 28.09 | 4.92 | 5.09 | 13.51 | 13.16 |
| Liang *et al.* [19] | PPR10K-c | **25.46** | **25.05** | **7.43** | **7.69** | **28.80** | **28.38** | **4.82** | **4.98** | 13.49 | 13.06 |
| Liang *et al.* [19]+GLC | PPR10K-c | 24.59 | 24.01 | 8.02 | 8.39 | 27.92 | 27.33 | 5.20 | 5.43 | **12.76** | **12.79** |
| Ours (LAM) | PPR10K-c | **25.65** | **25.56** | **7.39** | **7.43** | 28.95 | 28.84 | 4.80 | 4.83 | 14.62 | 14.56 |
| Ours (LAM+GAM) | PPR10K-c | 25.31 | 25.22 | 8.00 | 8.03 | 28.61 | 28.51 | 5.19 | 5.21 | **8.87** | **8.84** |

## 4. Experiments

### 4.1. Experiment setups

**Datasets.** All our experiments are performed on the PPR10K dataset [19]. It provides 4K ~ 8K high-resolution (HR) photos and corresponding 360P low-resolution (LR) photos with the retouched results by three experts. We use LR photos for training with randomly disrupted, while all of the LR and HR photos are available for testing. Following previous work [19], we divided the PPR10K dataset into a training set with 1,356 groups and 8,875 photos, and a testing set with 325 groups and 2,286 photos.

**Implementation details.** We perform our experiments on RTX 2080Ti GPUs with Pytorch framework [25]. We evaluate the performance of different settings and determine the size of the image patch is $3 \times 3$. All of the training photos are LR photos, so as to improve the training speed, and the testing photos include HR photos (4K ~ 8K) and LR photos (360P). We follow the data augment setting in [19] and train 200 epochs with batch size 16. We use Adam [17] as the network optimizer and the learning rate is fixed at $1 \times 10^{-4}$ following [35].

**Evaluation metrics.** Five metrics are employed to evaluate the performance of different methods quantitatively. Besides basic $PSNR$, the color difference between the retouched photo $R$ and the target photo $T$ is defined as the L2-distance in CIELAB color space with

$$\Delta E_{ab} = \| R^{Lab} - T^{Lab} \|_2 . \tag{8}$$

The $PSNR^{HC}$ and $\Delta E_{ab}^{HC}$ mean human-centered PSNR and color difference. We follow the study of Liang *et al.* [19] and adopt the same evaluation metric of group-style $\mathcal{M}_{GLC}$ to evaluate the performance of our model.

### 4.2. Comparisons with state-of-the-art methods

Table 1 reports the comparison results with multiple state-of-the-art methods. The comparison results show that our proposed method has made obvious progress in both HR and LR photos. Using a specialist's retouching example, the PSNR of the HR photo is improved by 0.60, and the addition of the GAM led to a significant reduction in the $M_{GLC}$ metric, which demonstrates the effectiveness of our method.

The qualitative results are shown in Fig. 6, which not only demonstrate the photo retouching results of each network component but also show the input photo and target photo together. It can be intuitively seen that our method achieves the closest retouching results to the target photo. The results of retouching the same set of input photos with the local-context aware module are better than 3D LUT [35]. While the output photos with the group-style

Table 2. Ablation experiments. (a) studies the influence of different spatial sizes when considering the local context. (b) reports the performance when varying the coefficient of group-style aware loss. (c) and (d) verify the influence of 3D LUT's number and the dimension in each LUT unit, respectively. In addition, we recorded the average running time (in milliseconds) of these models. All methods were tested on an NVIDIA 2080Ti GPU.

| (a) | $PSNR$ | $\Delta E_{ab}$ | $PSNR^{HC}$ | $\Delta E_{ab}^{HC}$ | $M_{GLC}$ | $Time$ |
|---|---|---|---|---|---|---|
| $3 \times 3$ | **26.23** | **6.62** | **29.53** | **4.29** | 11.08 | **4.66** |
| $5 \times 5$ | 26.19 | 6.67 | 29.50 | 4.32 | 11.39 | 4.80 |
| $7 \times 7$ | 26.05 | 6.73 | 29.33 | 4.36 | **11.06** | 4.86 |

| (c) | $PSNR$ | $\Delta E_{ab}$ | $PSNR^{HC}$ | $\Delta E_{ab}^{HC}$ | $M_{GLC}$ | $Time$ |
|---|---|---|---|---|---|---|
| 3 | 25.91 | 6.84 | 29.18 | 4.44 | **10.84** | **4.56** |
| 5 | **26.23** | **6.62** | **29.53** | **4.29** | 11.08 | 4.66 |
| 7 | 26.06 | 6.70 | 29.34 | 4.34 | 11.38 | 5.82 |

| (b) | $PSNR$ | $\Delta E_{ab}$ | $PSNR^{HC}$ | $\Delta E_{ab}^{HC}$ | $M_{GLC}$ | $Time$ |
|---|---|---|---|---|---|---|
| $1\times10^{-2}$ | 24.78 | 8.47 | 28.10 | 5.48 | **5.99** | 4.66 |
| $5\times10^{-3}$ | 25.07 | 8.06 | 28.37 | 5.22 | 6.13 | 4.66 |
| $1\times10^{-3}$ | 25.66 | 7.42 | 28.96 | 4.82 | 6.67 | 4.66 |
| $1\times10^{-4}$ | **26.05** | **6.99** | **29.34** | **4.52** | 7.87 | 4.66 |

| (d) | $PSNR$ | $\Delta E_{ab}$ | $PSNR^{HC}$ | $\Delta E_{ab}^{HC}$ | $M_{GLC}$ | $Time$ |
|---|---|---|---|---|---|---|
| 17 | 25.76 | 6.85 | 29.05 | 4.44 | **10.78** | **4.59** |
| 33 | **26.23** | **6.62** | **29.53** | **4.29** | 11.08 | 4.66 |
| 49 | 26.12 | 6.72 | 29.41 | 4.36 | 11.51 | 4.70 |
| 65 | 26.21 | 6.70 | 29.50 | 4.35 | 11.75 | 4.74 |

Table 3. Ablation experiments on each component of our proposed module. 3D LUT was used as the baseline, and the effectiveness of each module is verified by step-to-step ablation.

| No. | Model | $PSNR$ | $PSNR^{HC}$ | $M_{GLC}$ |
|---|---|---|---|---|
| #1 | 3D LUT [35] | 25.99 | 28.29 | 10.81 |
| #2 | #1 + LAM | 26.13 | 29.42 | 11.58 |
| #3 | #2 + Full Sup. | 26.13 | 29.40 | 11.26 |
| #4 | #2 + Edge Sup. | **26.23** | **29.53** | 11.08 |
| #5 | #4 + GAM | 25.66 | 28.96 | **6.67** |

Table 4. Ablation experiments that replace the local-context aware module with UNet [26] to generate pixel-adaptive weights.

| Model | $PSNR$ | $\Delta E_{ab}$ | $PSNR^{HC}$ | $\Delta E_{ab}^{HC}$ | $M_{GLC}$ | $Time$ |
|---|---|---|---|---|---|---|
| LAM | **26.23** | **6.62** | **29.53** | **4.29** | **11.08** | 4.66 |
| UNet | 25.72 | 6.92 | 28.97 | 4.49 | 11.15 | 5.88 |

aware module obviously maintain a more uniform retouching style, but the PSNR value only drops a bit. Meanwhile, the output of the photos with our proposed method has better scores.

### 4.3. Ablation study

We conducted extensive ablation experiments on the PPR10K dataset to determine the validity of the components and verify the influence of each parameter.

**Efficacy of each component.** We conduct extensive experiments to analyze the key components of our model, as shown in Table 3. We regard the model in Sec. 3.1 as the baseline model and add the local-context aware module on top of it. From the ablation experiments, we can find that all the metrics have been improved. After that, we design two different supervision strategies. Both of them are groundtruth mask oriented. They are full supervision corresponding to the whole photo and edge supervision for the attention mask. Both experiments are conducted on the basis of the local-contextual aware module. The results demonstrate that full supervision does not bring a significant improvement, almost any difference, while the edge supervision can consistently improve the performance.

We think the reason is due to the existence of the large pure background or pure foreground region in the groundtruth mask, which led to dramatically imbalanced supervision. Thus the generated loss is not accurate enough

to learn. Edge supervision has a better learning efficiency because it focuses more on the demarcation of the photo. Finally, we add the group-style aware module and analyze that the previous strategy does not work well for group-level consistency, whose performance under the metric $M_{GLC}$ are greater than 10. However, after introducing our group-style aware module, the indicator can be reduced to 6.67. The substantial numerical improvement demonstrates the effectiveness of the proposed group-style aware module.

**Comparison to spatial-aware 3D LUT [32].** We follow spatial-aware 3DLUT [32] and employ the same UNet architecture to replace our local-context aware module and predict pixel-adaptive LUT weights. As shown in Table 4, under the metric PSNR, our local-context aware module exceeds the UNet architecture with 0.51 points, and the group-level consistency and running time also verify the superiority of our method. We think there are two potential reasons. Firstly, UNet analyzes the contextual information for the whole image instead of a controlled receptive field, thus introducing some noise in the portrait photo retouching process. Secondly, the UNet architecture contains more parameters than our local attention mask module, which may cause the UNet architecture not easy to train.

**The parameter of 3D LUTs.** In order to ensure that the experiment settings are appropriate, we change the parameters for different experiments while keeping other conditions constant. Results of the experiments are shown in Table 2(c) and Table 2(d). The ablation results indicate that the most suitable number of 3D LUTs is 5, while the suitable dimension is 33.

Figure 6. Qualitative effect of the LAM and LAM+GAM strategy. Using the LAM improves the PSNR, and adding the GAM improves the tonal uniformity of a group of photos.

**The size of the receptive field.** We consider that the processing of low-frequency information (brightness, color, *etc.*) in PPR tasks relies on a large receptive field, but a too-large receptive field may introduce extra noise and also lead to a larger amount of parameters in the convolutional block. To investigate the most suitable value for the patch (*e.g.* receptive field) size of the LAM, we carry out experiments and report the results in Table 2(a). The experimental results indicate that the size of $3 \times 3$ is a proper choice for our local-attention aware module.

**Determining the coefficient of GAM.** In the practical experiments, the introduction of the GAM module leads to a slight decrease in the retouching effect. Thus, it is unreasonable to purely optimize the group-level consistency and lead to large deviations in the retouching results. As shown in Table 2(b), we conduct ablation experiments and set the coefficient $\lambda_{GLC}$ to $1 \times 10^{-3}$. This makes a balance between the single-photo retouching quality (measured by PSNR) and group-level consistency (measured by $M_{GLC}$).

By integrating the above experimental results, we ensure that parameters in our network are appropriate for the model. Due to the effectiveness of the local-context aware module and the group-style aware module, our proposed method achieves promising results. We are forming a concise and efficient solution for portrait photo retouching.

## 5. Conclusion

As the existing LUT-based photo retouching paradigm [35] only conducts pixel-to-pixel transformation but ignores the context cues, this paper designs a lightweight local-context aware module to incorporate context pixels and adaptively estimate LUT weights when tackling each pixel. Besides, a local attention mask is learned to distinguish the affinity of neighboring pixels further. Moreover, we propose to constrain the group-level consistency in the Lab space and guide the color temperatures and tones for a group of photos to be consistent. As we only employ four convolutional layers to capture local context cues explicitly, our method still maintains the lightweight superiority of the LUT-based photo retouching paradigm. Extensive experiments on the large-scale PPR10K dataset demonstrate the efficacy of the proposed local-context aware module and group-style aware module. However, as we only verify our method on a benchmark dataset, it may observe performance drop when tackling practical photos. In addition, the unintended usage of our method for surveillance may violate individual privacy. In the future, it is a promising direction to introduce the local-context aware module to other works, *e.g.* super-resolution [40], image classification [29], image inpainting [36, 37].

# References

[1] Mahmoud Afifi, Konstantinos G Derpanis, Bjorn Ommer, and Michael S Brown. Learning multi-scale photo exposure correction. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9157–9167, 2021. 2

[2] Jianrui Cai, Shuhang Gu, and Lei Zhang. Learning a deep single image contrast enhancer from multi-exposure images. *IEEE Trans. Image Process.*, 27(4):2049–2062, 2018. 2

[3] Yu-Sheng Chen, Yu-Ching Wang, Man-Hsin Kao, and Yung-Yu Chuang. Deep photo enhancer: Unpaired learning for image enhancement from photographs with gans. In *Int. Conf. Comput. Vis.*, pages 6306–6314, 2018. 2

[4] Zipei Chen, Chengjiang Long, Ling Zhang, and Chunxia Xiao. Canet: A context-aware network for shadow removal. In *Int. Conf. Comput. Vis.*, pages 4743–4752, 2021. 3

[5] Lan Fu, Changqing Zhou, Qing Guo, Felix Juefei-Xu, Hongkai Yu, Wei Feng, Yang Liu, and Song Wang. Auto-exposure fusion for single-image shadow removal. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10571–10580, 2021. 2

[6] Michaël Gharbi, Jiawen Chen, Jonathan T Barron, Samuel W Hasinoff, and Frédo Durand. Deep bilateral learning for real-time image enhancement. *ACM Trans. Graph.*, 36(4):1–12, 2017. 2, 6

[7] Chunle Guo Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1780–1789, June 2020. 2

[8] Samuel W Hasinoff, Dillon Sharlet, Ryan Geiss, Andrew Adams, Jonathan T Barron, Florian Kainz, Jiawen Chen, and Marc Levoy. Burst photography for high dynamic range and low-light imaging on mobile cameras. *ACM Trans. Graph.*, 35(6):1–12, 2016. 2

[9] Jingwen He, Yihao Liu, Yu Qiao, and Chao Dong. Conditional sequential modulation for efficient global image retouching. In *Eur. Conf. Comput. Vis.*, pages 679–695. Springer, 2020. 2, 6

[10] Jie Huang, Pengfei Zhu, Mingrui Geng, Jiewen Ran, Xingguang Zhou, Chen Xing, Pengfei Wan, and Xiangyang Ji. Range scaling global u-net for perceptual image enhancement on mobile devices. In *Eur. Conf. Comput. Vis.*, pages 0–0, 2018. 1

[11] Tao Huang, Songjiang Li, Xu Jia, Huchuan Lu, and Jianzhuang Liu. Neighbor2neighbor: Self-supervised denoising from single noisy images. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 14781–14790, June 2021. 2

[12] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. Dslr-quality photos on mobile devices with deep convolutional networks. In *Int. Conf. Comput. Vis.*, pages 3277–3285, 2017. 2

[13] Younghyun Jo and Seon Joo Kim. Practical single-image super-resolution using look-up table. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 691–700, 2021. 3

[14] Hanul Kim, Su-Min Choi, Chang-Su Kim, and Yeong Jun Koh. Representative color transform for image enhancement. In *Int. Conf. Comput. Vis.*, pages 4459–4468, 2021. 2

[15] Han-Ul Kim, Young Jun Koh, and Chang-Su Kim. Global and local enhancement networks for paired and unpaired image enhancement. In *Eur. Conf. Comput. Vis.*, pages 339–354. Springer, 2020. 2

[16] Han-Ul Kim, Young Jun Koh, and Chang-Su Kim. Pienet: Personalized image enhancement network. In *Eur. Conf. Comput. Vis.*, pages 374–390. Springer, 2020. 2

[17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *Int. Conf. Learn. Represent.*, 2015. 6

[18] Xia Li, Jianlong Wu, Zhouchen Lin, Hong Liu, and Hongbin Zha. Recurrent squeeze-and-excitation context aggregation net for single image deraining. In *Eur. Conf. Comput. Vis.*, pages 254–269, 2018. 3

[19] Jie Liang, Hui Zeng, Miaomiao Cui, Xuansong Xie, and Lei Zhang. Ppr10k: A large-scale portrait photo retouching dataset with human-region mask and group-level consistency. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 653–661, 2021. 1, 2, 5, 6

[20] Hongyu Liu, Ziyu Wan, Wei Huang, Yibing Song, Xintong Han, and Jing Liao. Pd-gan: Probabilistic diverse gan for image inpainting. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9371–9381, 2021. 2

[21] Nian Liu, Junwei Han, and Ming-Hsuan Yang. Picanet: Learning pixel-wise contextual attention for saliency detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3089–3098, 2018. 3

[22] Sean Moran, Pierre Marza, Steven McDonagh, Sarah Parisot, and Gregory Slabaugh. Deeplpf: Deep local parametric filters for image enhancement. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 12826–12835, 2020. 2

[23] Sean Moran, Steven McDonagh, and Gregory Slabaugh. Curl: Neural curve layers for global image enhancement. In *Int. Conf. Pattern Recog.*, pages 9796–9803. IEEE, 2021. 1

[24] Jongchan Park, Joon-Young Lee, Donggeun Yoo, and In So Kweon. Distort-and-recover: Color enhancement using deep reinforcement learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5928–5936, 2018. 2

[25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform. Process. Syst.*, 32:8026–8037, 2019. 6

[26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 7

[27] Yuda Song, Hui Qian, and Xin Du. Starenhancer: Learning real-time and style-aware image enhancement. In *Int. Conf. Comput. Vis.*, pages 4126–4135, 2021. 2

[28] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *Int. Conf. Comput. Vis.*, pages 4539–4547, 2017. 1

[29] Kevin Tang, Manohar Paluri, Li Fei-Fei, Rob Fergus, and Lubomir Bourdev. Improving image classification with location context. In *Int. Conf. Comput. Vis.*, pages 1008–1016, 2015. 3, 8

[30] Thang Vu, Cao Van Nguyen, Trung X Pham, Tung M Luu, and Chang D Yoo. Fast and efficient image quality enhancement via desubpixel convolutional neural networks. In *Eur. Conf. Comput. Vis.*, pages 0–0, 2018. 1

[31] Ruixing Wang, Qing Zhang, Chi-Wing Fu, Xiaoyong Shen, Wei-Shi Zheng, and Jiaya Jia. Underexposed photo enhancement using deep illumination estimation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6849–6857, 2019. 2

[32] Tao Wang, Yong Li, Jingyang Peng, Yipeng Ma, Xian Wang, Fenglong Song, and Youliang Yan. Real-time image enhancer via learnable spatial-aware 3d lookup tables. In *Int. Conf. Comput. Vis.*, pages 2471–2480, 2021. 3, 7

[33] Xiaoqing Yin, Xinchao Wang, Jun Yu, Maojun Zhang, Pascal Fua, and Dacheng Tao. Fisheyerecnet: A multi-context collaborative deep network for fisheye image rectification. In *Eur. Conf. Comput. Vis.*, pages 469–484, 2018. 3

[34] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 14821–14831, 2021. 2

[35] Hui Zeng, Jianrui Cai, Lida Li, Zisheng Cao, and Lei Zhang. Learning image-adaptive 3d lookup tables for high performance photo enhancement in real-time. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020. 1, 3, 4, 5, 6, 7, 8

[36] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for high-quality image inpainting. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1486–1494, 2019. 3, 8

[37] Yu Zeng, Zhe Lin, Huchuan Lu, and Vishal M Patel. Cr-fill: Generative image inpainting with auxiliary contextual reconstruction. In *Int. Conf. Comput. Vis.*, pages 14164–14173, 2021. 3, 8

[38] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Trans. Image Process.*, 26(7):3142–3155, 2017. 1

[39] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image restoration. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(7):2480–2495, 2020. 1

[40] Yulun Zhang, Donglai Wei, Can Qin, Huan Wang, Hanspeter Pfister, and Yun Fu. Context reasoning attention network for image super-resolution. In *Int. Conf. Comput. Vis.*, pages 4278–4287, 2021. 3, 8

[41] Zhaoyang Zhang, Yitong Jiang, Jun Jiang, Xiaogang Wang, Ping Luo, and Jinwei Gu. Star: A structure-aware lightweight transformer for real-time image enhancement. In *Int. Conf. Comput. Vis.*, pages 4106–4115, 2021. 2

[42] Zhuoran Zheng, Wenqi Ren, Xiaochun Cao, Tao Wang, and Xiuyi Jia. Ultra-high-definition image hdr reconstruction via collaborative bilateral learning. In *Int. Conf. Comput. Vis.*, pages 4449–4458, 2021. 2