

Modernizing Old Photos Using Multiple References via Photorealistic Style Transfer

Agus Gunawan¹ Soo Ye Kim^{1,2*} Hyeonjun Sim^{1†} Jae-Ho Lee³ Munchurl Kim^{1‡}

¹KAIST ²Adobe Research ³ETRI

{agusgun, flhy5836, mkimee}@kaist.ac.kr sooyek@adobe.com jhlee3@etri.re.kr



Figure 1. The results of old photo modernization produced by our method. Our method is able to generate more modern-looking images that resemble the style of input reference images **without the use of old photos during training**.

Abstract

This paper firstly presents old photo modernization using multiple references by performing stylization and enhancement in a unified manner. In order to modernize old photos, we propose a novel multi-reference-based old photo modernization (MROPM) framework consisting of a network MROPM-Net and a novel synthetic data generation scheme. MROPM-Net stylizes old photos using multiple references via photorealistic style transfer (PST) and further enhances the results to produce modern-looking images. Meanwhile, the synthetic data generation scheme trains the network to effectively utilize multiple references to perform modernization. To evaluate the performance, we propose a new old photos benchmark dataset (CHD) consisting of diverse natural indoor and outdoor scenes. Extensive experiments show that the proposed method outperforms other baselines in performing modernization on real old photos, even though no old photos were used during training. Moreover, our method can appropriately select styles from multiple references for each semantic region in the old photo to further improve the modernization performance.

*Soo Ye Kim is currently affiliated with Adobe Research.

†Hyeonjun Sim is currently affiliated with Qualcomm.

‡Corresponding author.

1. Introduction

Old photos taken a long time ago may contain important information that carry cultural and heritage values, e.g., photos of Queen Elizabeth II’s coronation. Such old images may contain multiple degradations, e.g., scratches, and old photo artifacts, e.g., color fading, often preventing people from understanding the scene. To restore these images, a skilled expert needs to perform laborious manual processes such as degradation restoration and modernization, i.e., colorization or enhancement, to make them look modern [44]. Consequently, early studies [8, 39] try to restore damaged old photos automatically by using traditional inpainting techniques. However, solely re-synthesizing damaged regions in the image is inadequate to ensure old photos look modern, as the overall style remains similar.

Recent work [28] formulates the task as time-travel rephotography which aims to translate old photos into a modern photos space. The authors considered a multi-task problem consisting of two main tasks: (i) restoration of old photos with both unstructured (noise, blur) and structured (scratch, crack) degradations; (ii) modernization which aims to change old photos’ characteristics to look like modern images, e.g., better color saturation and contrast by using colorization [28, 48] or enhancement [42]. However, simply using an enhancement method [42] fails to modernize old photos, as shown in Fig. 1, since the overall look

still remains similar to old photos, e.g., with a sepia color.

In this paper, we propose to modernize old color photos of natural scenes by changing their styles and enhancing them to look modern. For this, a novel unified framework is proposed which leverages multiple modern photo references in solving the modernization task of old photos by utilizing photorealistic style transfer (PST). Although one prior work [48] is also reference-based, it only relies on a single reference to colorize greyscale portrait photos. However, in natural scene cases, it is challenging to find a single modern photo as a reference that can well match the whole semantics of an old photo. Moreover, changing only the color is not sufficient to alter the overall look of an image [12]. Thus, our framework uses multiple references to modernize old photos by changing the *style* instead of only the color. Since there is no public old photos benchmark dataset of natural scenes, we propose a new Cultural Heritage Dataset (CHD) with 644 indoor and outdoor old color photos collected from various national museums in Korea.

Our multiple-reference-based old photo modernization framework (MROPM) consists of two main parts: (i) *MROPM-Net* and (ii) *a novel training strategy* that enables the network to utilize multiple references. The MROPM-Net consists of two different subnets: The first is a single stylization subnet that transfers both global and local styles without any semantic segmentation from a modern photo into an old photo; Specifically, we propose an improved version of WCT2 [51], inspired by its universal generalization, as the backbone of the single stylization subnet, and present a new architecture that can perform both global PST and local PST without requiring any semantic segmentation; The second is a merging-refinement subnet that merges multiple stylization results from multiple references based on semantic similarities and further refines the merged result to produce a modernized version of the old photo. To effectively train the MROPM-Net, we propose a synthetic data generation scheme that uses the style-variant (i.e., color jittering and unstructured degradation) and -invariant (i.e., rotation, flipping, and translation) transformations. Our MROPM can modernize old photos better than the state-of-the-art (SOTA) old photo restoration method [42], even without using any old photos during training, thanks to the generalization of PST. Our contributions are summarized as follows:

- We propose the *first* old photo modernization framework (MROPM) that allows the usage of *multiple references* to guide the modernization process.
- Our *photorealistic multi-stylization network* and training strategy enable the MROPM-Net to utilize multiple style references in modernizing old photos.
- Our training strategy based on synthetic data allows the MROPM-Net to modernize *real* old photos even without using any old photos during training.

- We propose a new old photo dataset of natural scenes, called Cultural Heritage Dataset (CHD), with 644 outdoor and indoor cultural heritage images.

2. Related Work

Reference-based color-related tasks. One way to change the overall look of an image is by changing color, which is one of the style components [12]. To change the color of old photos, one can employ two methods: *exemplar-based colorization* [10, 26, 40, 46, 49, 50] and *color transfer* or *recolorization* [1, 11, 20, 24]. However, exemplar-based colorization methods cannot utilize the color information in the input images for matching, although color is an important feature representing object semantics [38], limiting the methods for the modernization of old color photos. Color transfer aims to transfer the reference image’s color statistics into the input image. Early deep learning works [11, 24] use deep feature matching from features extracted with pre-trained VGG19 [37] to perform the color transfer, which can also be extended to multi-reference cases [11]. Due to the long execution time of the optimization process, recent works develop end-to-end networks, where Lee *et al.* [20] utilize color histogram analogy, and Afifi *et al.* [1] utilize a color-controlled generative adversarial network (GAN). However, recent works can only use a single reference, where finding a single reference image containing similar semantics as the input old photo can be challenging. Thus, from the perspective of color transfer, our work is the first end-to-end network that can utilize multiple references to handle content mismatch without any slow optimization technique.

Photorealistic style transfer (PST). The PST aims at achieving photorealistic rendering of an image with the style of another image. Since the development of post-processing and regularization techniques [27, 31], PST has gained much popularity. Recent works can be categorized into architecture [2, 6, 7, 23, 35, 47, 51] and feature transformation [15, 21, 22] improvements to effectively and efficiently produce photorealistic results. Specifically, WCT2 [51] utilizes wavelet-based skip connection and progressive stylization to achieve better PST where the method can work universally without re-training to pre-defined styles. Due to these benefits, we base our network architecture on WCT2. However, WCT2 produces unnatural style transfer results when performing global and local stylization with unreliable semantic segmentation (shown in Supplementary Material), which hinders the application to old photos. Thus, our MROPM-Net is designed to enable local stylization without any semantic segmentation, which in consequence, can perform multi-style PST in one unified framework without specifying any masks. To the best of our knowledge, this is the first work in multi-style PST, although there is one work in multi-style artistic style transfer

(AST) [14]. Note that the AST is different from the PST in that it utilizes learning-dependent feature transformation, which can cause severe visual artifacts in PST.

Old photo restoration. Early works in old photo restoration focus on detecting and restoring structured degradation (scratch and crack) of images using traditional inpainting techniques [8, 39]. Besides the structured degradation, [25, 42, 48] incorporate additional spatially-uniform unstructured degradation, e.g., blur and noise, using synthetic degradation and formulate the problem as mixed degradation restoration. However, restoring mixed degradation is not enough to ensure that old photos look modern. Consequently, Luo *et al.* [28] formally introduce the time-travel rephotography problem, which aims to translate old photos to look like ones taken in the modern era. This problem adds modernization, synthesizing the missing colors and enhancing the details, on top of degradation restoration. To solve the modernization problem, Luo *et al.* [28] use a StyleGAN2-generated [18] sibling image to serve as a reference for old portrait photos. However, generating complex natural scene photos via GAN to be used as references is challenging [3], making the method unable to be applied to natural scene old photos. Another work [48] proposes to use a single reference image to colorize an old greyscale photo. However, using a single reference is not enough to cover the whole semantics of old photos (shown in Fig. 1). Thus, different from previous methods, we propose to modernize old photos by stylizing and enhancing old photos in a unified manner using *multiple references* to better cover the entire semantics of old photos.

3. Proposed Cultural Heritage Dataset (CHD)

Although some public datasets such as Historical Wiki Face Dataset (HWFD) [28] and RealOld [48] have been released recently, these datasets only contain portrait or face photos which are much simpler compared to natural scenes. In addition, these datasets only contain greyscale photos and disregard color photos produced during the 20th century using reversal films [29], which have specific degradations such as color dye fading and have not been analyzed before. Therefore, we propose a Cultural Heritage Dataset (CHD) consisting of 644 old color photos produced in the 20th century. Specifically, we collect these old photos in the form of reversal films or papers from three national museums in Korea, which are then scanned in resolutions varying from 4K to 8K. The photos have been well preserved and stored carefully due to their value, containing little structured degradation, e.g., scratches, but varying degrees of unstructured degradation, e.g., noises. These photos contain indoor and outdoor scenes of cultural heritage, such as special exhibitions and excavation ruins. After collection, all photos are divided into train and test sets by randomly splitting with a proportion of 8 (514 photos):2 (130 photos). The train set

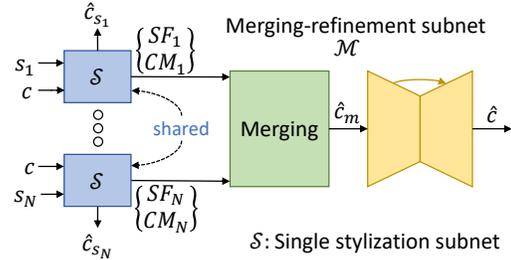


Figure 2. The overall framework of our multiple-reference-based old photo modernization network (MROPM-Net).

is only used for other baselines that need to be trained using real old photos. Since our task is reference-based old photo modernization, we further collect modern photos as references by crawling images with similar contexts from the internet. In total, we obtain 130 old photos in the test set, each of which has one or two references selected manually. Further details can be found in *Supplementary Material*.

4. Proposed Method

4.1. Overall Framework

Fig. 2 shows our proposed multi-reference-based old photo modernization network (MROPM-Net) with a shared single stylization subnet \mathcal{S} and a merging-refinement subnet \mathcal{M} . We denote an old photo input as content $c \in \mathbb{R}^{H \times W \times 3}$ and N number of modern photos as styles $s = \{s_i\}_{i=1}^N \in \mathbb{R}^{N \times H \times W \times 3}$. Our goal is to modernize c using s . In the first step, we utilize \mathcal{S} , which is built based on a photorealistic style transfer (PST) backbone, to stylize c using each s_i , yielding N stylized features and correlation matrices $\{SF_i, CM_i\}_{i=1}^N$. After having multiple stylization results, we merge the features $\{SF_i\}_{i=1}^N$ based on the semantic similarity $\{CM_i\}_{i=1}^N$ between c and s and further refine the merging result via \mathcal{M} . Specifically, \mathcal{M} selects the appropriate styles for each semantic region based on multiple stylization results $\{SF_i\}_{i=1}^N$ to produce an intermediate merging image output \hat{c}_m , e.g., selecting the most appropriate feature for a sky region from SF_1 that contains a sky style, not from $SF_{i \neq 1}$, which do not contain sky styles. Then, \hat{c}_m is further refined to get the final result \hat{c} . Given relevant references, \hat{c} becomes a modern version with a modern style and enhanced details for old photo input c .

4.2. Network Architecture

Single stylization subnet \mathcal{S} . Fig. 3 shows a detailed structure of \mathcal{S} . For given multiple references, our single stylization subnet is shared for all input pairs and takes a single pair of an old photo c and a reference s_i at a time. Given a pair of (c, s_i) , \mathcal{S} stylizes c based on the style code of s_i locally and globally, resulting in a stylized feature SF_i , a stylized old photo \hat{c}_{s_i} , and a correlation matrix CM_i . This subnet \mathcal{S} consists of two main parts: (i) *an improved PST*

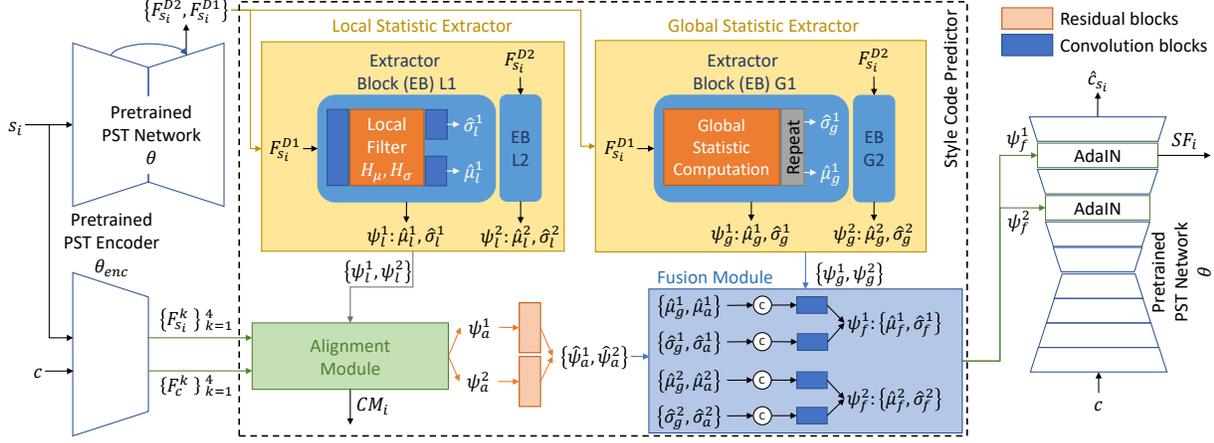


Figure 3. The architecture of the single stylization subnet \mathcal{S} .

network and (ii) a style code predictor.

For the PST network, we improve some drawbacks of the concatenated version of WCT2 [51]. We observed that the stylization only affects the last decoder block due to the design of its skip connection, where this issue is called a “short circuit” in [2]. Thus, instead of transferring three different high-frequency components as in the WCT2, we propose to simplify it by transferring a single high-frequency component in level-0 of the Laplacian pyramid representation [4]. Second, we only apply feature transformation in the network’s decoder part, especially the last two decoder blocks, which achieves the best trade-off between the stylization effect and the photorealism. Third, we use the differentiable adaptive instance normalization (AdaIN) [13] instead of the non-differentiable WCT [22] to learn and predict the local style rather than compute it.

The second part of \mathcal{S} is a style code predictor. This part aims to predict style codes $\psi = \{\mu, \sigma\}$ consisting of mean and standard deviation (std), which are statistics used to perform stylization in AdaIN [13]. We propose to predict ψ instead of computing it as in AdaIN to perform local style transfer without requiring any semantic segmentation. The first step (yellow) of the style code predictor is to extract local style codes $\psi_l^j = \{\hat{\mu}_l^j, \hat{\sigma}_l^j\}$ and global style codes $\psi_g^j = \{\hat{\mu}_g^j, \hat{\sigma}_g^j\}$ from the j -th level feature $\{F_{s_i}^{Dj}\}$ extracted by the last two decoder blocks ($j = 1, 2$) of the pre-trained PST network as shown in Fig. 3. In this regard, ψ_l^j is extracted using a local statistic extractor which consists of a local mean filter H_μ and local std filter H_σ with a kernel size of 3 and convolution blocks to refine both filtered outputs. Meanwhile, the global statistic extractor extracts ψ_g^j by computing channel-wise mean and std values, which are then spatially repeated to the same spatial size of ψ_l^j . After style code extraction, the second step (green) is to align ψ_l^j to c by using non-local attention [43]. Specifically, we extract multi-level feature maps $\{F_c^k\}_{k=1}^4$ and $\{F_{s_i}^k\}_{k=1}^4$ for both c and s_i , respectively, map them into the same feature

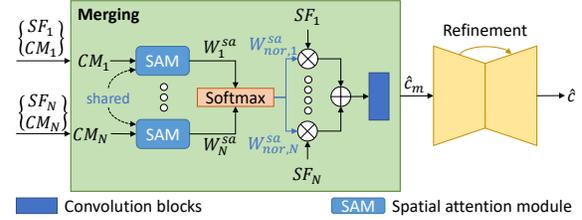


Figure 4. The architecture of the merging-refinement subnet \mathcal{M} .

space using shared convolution blocks, and perform matrix multiplication between mapped features to obtain correlation matrix CM_i . Then, we align ψ_l^j to c by using CM_i via matrix multiplication. The aligned style code ψ_a^j is further refined to prevent interpolation artifacts by using residual blocks [9], resulting in a refined version $\hat{\psi}_a^j = \{\hat{\mu}_a^j, \hat{\sigma}_a^j\}$. After obtaining $\hat{\psi}_a^j$, we fuse it with ψ_g^j via the fusion module to obtain a fused style code $\psi_f^j = \{\hat{\mu}_f^j, \hat{\sigma}_f^j\}$. The fusion module performs channel-wise concatenation for $\hat{\psi}_a^j$ and ψ_g^j , which is then fed into the following convolution blocks as shown in the blue part of Fig. 3.

Finally, after performing all the operations from the local and global statistic extractors to the fusion module, we obtain ψ_f^1 and ψ_f^2 . These fused style codes are then used for stylizing c . We use our PST network with AdaIN to perform the stylization as shown in the right part of Fig. 3.

Merging-refinement subnet \mathcal{M} . After stylizing an old photo c with N different modern photos $s = \{s_i\}_{i=1}^N$ using \mathcal{S} , we obtain multiple stylized features and correlation matrices $\{SF_i, CM_i\}_{i=1}^N$. The next step is to select the most appropriate styles from $\{SF_i\}_{i=1}^N$ for each semantic region via the merging part of \mathcal{M} , as shown in Fig. 4. For this, a spatial attention module (SAM) [45] is employed, which strengthens and dampens semantically related and unrelated spatial features, respectively, in the merging process of the stylized features. The SAM computes spatial attention weights W_i^{sa} by using CM_i for the corresponding SF_i . Then, we normalize all the spatial attention weights by

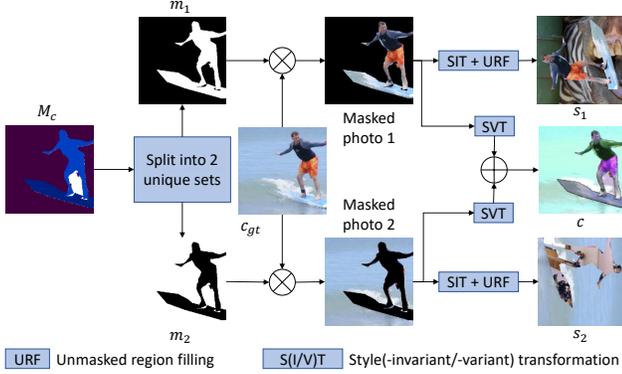


Figure 5. Our synthetic data generation pipeline.

using Softmax, thus having $\mathbf{W}_{nor}^{sa} = \{W_{nor,i}^{sa}\}_{i=1}^N$. All normalized attention weights $W_{nor,i}^{sa}$ are multiplied with their corresponding SF_i , whose results are summed and then fed into the final convolution blocks to obtain a merging result \hat{c}_m as an intermediate multi-style PST image. We further refine \hat{c}_m via the U-Net [36] based refinement subnet to produce a final modern version \hat{c} for old photo input c .

4.3. Training Strategy

Synthetic data generation. Since there is no ground truth for multiple-reference-based old photo modernization tasks, we generate synthetic data for training the network in a self-supervised manner. For this, the COCO-stuff dataset [5] is utilized, which has a semantic segmentation mask for each image. Fig. 5 shows the pipeline of generating the synthetic data where each sample consists of a *synthetic* old photo c , its corresponding N different references $s = \{s_i\}_{i=1}^N$, and its ground truth c_{gt} . We use two style references with $N = 2$ for each sample throughout our experiments. First, we randomly select a photo from the COCO-stuff dataset as ground truth c_{gt} and its corresponding semantic mask M_c at each iteration during the training. Then, all the semantic regions in M_c are randomly separated into N non-overlapping parts $\{m_i\}_{i=1}^N$. For the example shown in Fig. 5, M_c consists of two semantic regions, *surfer* and *sea*, and thus separated to: one mask m_1 with the surfer, and the other mask m_2 with the sea. The next step is to generate N different masked photos using $\{m_i\}_{i=1}^N$ by element-wise multiplication between m_i and c_{gt} . Then, we use these masked photos to generate multiple references $\{s_i\}_{i=1}^N$ and c via style-invariant (SIT) and -variant transformations (SVT) respectively.

The properties of style-variant and -invariant transformations are determined by whether the transformations alter the mean and std of any semantic region. Hence, we use random translation (only for the regions that can be translated), rotation, and flipping as our SIT. Meanwhile, random color jittering and unstructured degradation, i.e. blur, noise, resizing, and compression artifacts, are used for SVT. Other

types of degradation, e.g., scratches can be included in the SVT to make the method able to generalize to these types of degradation. To generate c , we apply different SVTs for each masked photo and sum up the results. Meanwhile, to generate $\{s_i\}_{i=1}^N$, randomly selected SITs are applied for each masked photo, and then the unmasked region is filled (URF) with another photo randomly selected from the same COCO-Stuff dataset. Our MROP-Net can work reasonably well for real old photos after training with this synthetic data. This is because the synthetic data make our MROP-Net able to (i) robustly find local semantic correspondences between degraded synthetic old photo c and semantically confusing synthetic modern photo s_i , (ii) accurately transfer the styles of each s_i to c locally, and (iii) merge and refine multiple stylization results from multiple styles $\{s_i\}_{i=1}^N$ to produce an output similar to c_{gt} . Thus, our synthetic data creation pipeline can be effectively used for multi-reference-based old photo modernization.

Our MROP-Net is trained in multiple stages: (i) **Stage 1:** Our PST network is trained using a similar training strategy to [51]; (ii) **Stage 2:** Our single stylization subnet \mathcal{S} is trained, while the pre-trained PST network is frozen; (iii) **Stage 3:** We train our merging-refinement subnet \mathcal{M} , with both the pre-trained PST network and \mathcal{S} frozen.

Loss function. In Stage 2 of training, our goal is to obtain a faithful stylization result from each of the style reference images $\{s_i\}_{i=1}^N$. Specifically, we use a weighted sum of the following different losses:

$$\mathcal{L}_{s_i}^{Stage2} = \lambda_{ML} \cdot \mathcal{L}_{ML}(\hat{c}_{s_i}, c_{gt}) + \lambda_p \cdot \mathcal{L}_p(\hat{c}_{s_i}, c_{gt}) + \lambda_{CX} \cdot \mathcal{L}_{CX}(\hat{c}_{s_i}, c_{gt}) \quad (1)$$

where \mathcal{L}_{ML} , \mathcal{L}_p and \mathcal{L}_{CX} represent masked reconstruction, perceptual [17] and contextual [32] losses respectively, and the λ 's control relative weights for their respective losses. We use the features extracted from VGG-19 [37] at layer *relu4_1* for \mathcal{L}_p , and *relu3_1* and *relu4_1* for \mathcal{L}_{CX} . Different from [32], GT image c_{gt} is used as the reference instead of style s_i for \mathcal{L}_{CX} to compare with our output \hat{c}_{s_i} because using s_i can cause severe structure distortion. \mathcal{L}_{ML} in Eq. 1 can be expressed as:

$$\mathcal{L}_{ML}(\hat{c}_{s_i}, c_{gt}) = \|(\hat{c}_{s_i} - c_{gt}) \odot m_i\|_1 \quad (2)$$

where m_i is a mask used to generate s_i in our data generation scheme as shown in Fig. 5. Correspondingly, these three losses are used to encourage \mathcal{S} (i) to faithfully stylize c at the pixel level for semantic regions that also appear in s_i and disregard other unrelated semantic regions, (ii) to faithfully stylize c at the semantic level, and (iii) to perform better semantic style transfer. In Stage 2 of training, we only use a single style reference for each c to reduce the computation complexity and stabilize the training of the single stylization subnet \mathcal{S} .

In Stage 3, we train our merging-refinement subnet \mathcal{M}

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
ExColTran [50] + OPR-R	19.5796	0.7885	0.2563
ReHistoGAN [1] + OPR-R	<u>20.0458</u>	0.7987	<u>0.2109</u>
MAST [15] + OPR-R	19.0148	0.7853	0.2270
PCAPST [6] + OPR-R	19.1731	0.7908	0.2197
Ours	21.2212	<u>0.7919</u>	0.2027

Table 1. Quantitative results of modernization on synthetic dataset.

Method	NIQE \downarrow	BRISQUE \downarrow
OPR [42]	4.8705	21.4588
ExColTran [50] + OPR	4.9415	18.8971
ReHistoGAN [1] + OPR	4.8051	26.2557
MAST [15] + OPR	4.8111	18.9555
PCAPST [6] + OPR	4.7094	18.9860
Ours - Single	<u>3.4737</u>	<u>15.5152</u>
Ours - Multiple	3.4487	15.4180

Table 2. Quantitative results of modernization on real old photos.

by using weighted sum of four different losses:

$$\mathcal{L}^{Stage3} = \lambda_{L1} \cdot \mathcal{L}_{L1}(\hat{c}, c_{gt}) + \lambda_p \cdot \mathcal{L}_p(\hat{c}, c_{gt}) + \lambda_{sm} \cdot \mathcal{L}_{sm}(\hat{c}) + \lambda_{adv} \cdot \mathcal{L}_{adv}(\hat{c}, c_{gt}) \quad (3)$$

where \mathcal{L}_{L1} , \mathcal{L}_p , \mathcal{L}_{sm} and \mathcal{L}_{adv} are reconstruction, perceptual [17], local smoothness [52] and least square adversarial [30] losses respectively, and λ 's control relative weights for corresponding losses. \mathcal{L}_p in Eq. 3 and Eq. 1 refer to the same loss function. We use these four losses accordingly to encourage the merging-refinement subnet \mathcal{M} to produce: (i) accurate merging and better refinement, (ii) perceptually plausible output, (iii) spatially smooth output, and (iv) realistic output.

5. Experiments

5.1. Experimental Settings

Training details. We use our proposed synthetic data generation scheme with the aforementioned multi-stage training strategy to train the network: (i) We train our PST network for five epochs; (ii) Then, we train our single stylization subnet \mathcal{S} based on \mathcal{L}^{Stage2} in Eq. 1 for two epochs, not to be overfitted for synthetic data, while freezing our PST network, where we set $\lambda_{ML} = 1$, $\lambda_p = 1$, and $\lambda_{CX} = 1$; (iii) Finally, we train our merging-refinement subnet \mathcal{M} for three epochs which is sufficient while freezing both \mathcal{S} and our PST network. The loss function to train \mathcal{M} is \mathcal{L}^{Stage3} in Eq. 3, where we set $\lambda_{L1} = 2$, $\lambda_p = 1$, $\lambda_{sm} = 3$, and $\lambda_{adv} = 0.2$. For all of the training, we use an ADAM optimizer [19] with a learning rate of $1e - 4$ and batch size of 1 to optimize our network and discriminator (PatchGAN discriminator [16]). In addition, we apply a linear learning decay in the last epoch of the \mathcal{M} training.

Baselines. Our work can be seen as handling a joint task of stylization and enhancement by using multiple references for old photo modernization. Since there are no baselines in reference-based old photo modernization, we compare to

Method	Top 1	Top 2	Top 3	Top 4	Top 5
OPR [42]	<u>17.44</u>	<u>39.83</u>	57.05	70.90	87.22
ExColTran [50] + OPR	1.62	5.13	10.77	24.87	47.27
ReHistoGAN [1] + OPR	7.91	32.27	<u>61.84</u>	<u>83.80</u>	<u>96.92</u>
MAST [15] + OPR	5.68	21.62	41.92	66.33	86.20
PCAPST [6] + OPR	10.98	28.50	44.87	61.97	84.66
Ours	56.37	72.69	83.55	92.14	97.74

Table 3. User study results. The percentage of user selection is shown.

sequential models consisting of stylization then enhancement, which can perform the same task. Reversing the order (enhancement and then stylization), results in worse outcomes. For stylization, we employ four state-of-the-art (SOTA) methods as baselines: (i) exemplar-based colorization: transformer-based method (ExColTran [50]); (ii) recolorization: recolorization using color-controlled GAN (ReHistoGAN [1]); photorealistic style transfer (PST): semantic PST (MAST [15]) and PCA-based knowledge distillation PST (PCAPST [6]). Meanwhile, for enhancement, we employ SOTA no-reference-based old photo restoration method (OPR [42]), which can perform similar enhancement as ours. For the stylization baselines, we use their pre-trained models in the evaluation since they cannot be re-trained with our synthetic data due to their different training strategies. However, note that these pre-trained models have already been trained to achieve the same goal of changing the overall look of input images based on the given reference image. Meanwhile, for enhancement, we use the pre-trained OPR model for real old photo evaluation, denoted as OPR, since it achieves better performance on real old photos, and retrain the OPR using our synthetic data and CHD training set for synthetic data evaluation, denoted as OPR-R. Since the four baselines can only utilize a single reference, we average the results of using different sets of references in quantitative evaluation, and randomly select one of two references for each input to the baseline networks in qualitative evaluation and user study.

Evaluation metrics. We evaluate all the methods on synthetically degraded and real old photos (CHD testing set). In synthetic degraded photos evaluation, we employ: 1) peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) to measure the pixel-level discrepancy between output and ground truth, 2) learned perceptual image patch similarity (LPIPS) [53] to measure the perceptual quality of the output. For evaluation on real old photos, we employ no-reference image quality assessment metrics such as NIQE [34] and BRISQUE [33], similar to [28, 41, 42], since the modernization ground truth photos do not exist.

5.2. Experimental Results

Quantitative comparison. We evaluate our method and baselines on a synthetic dataset and real-world old photos. In synthetic dataset evaluation, we evaluate all methods,

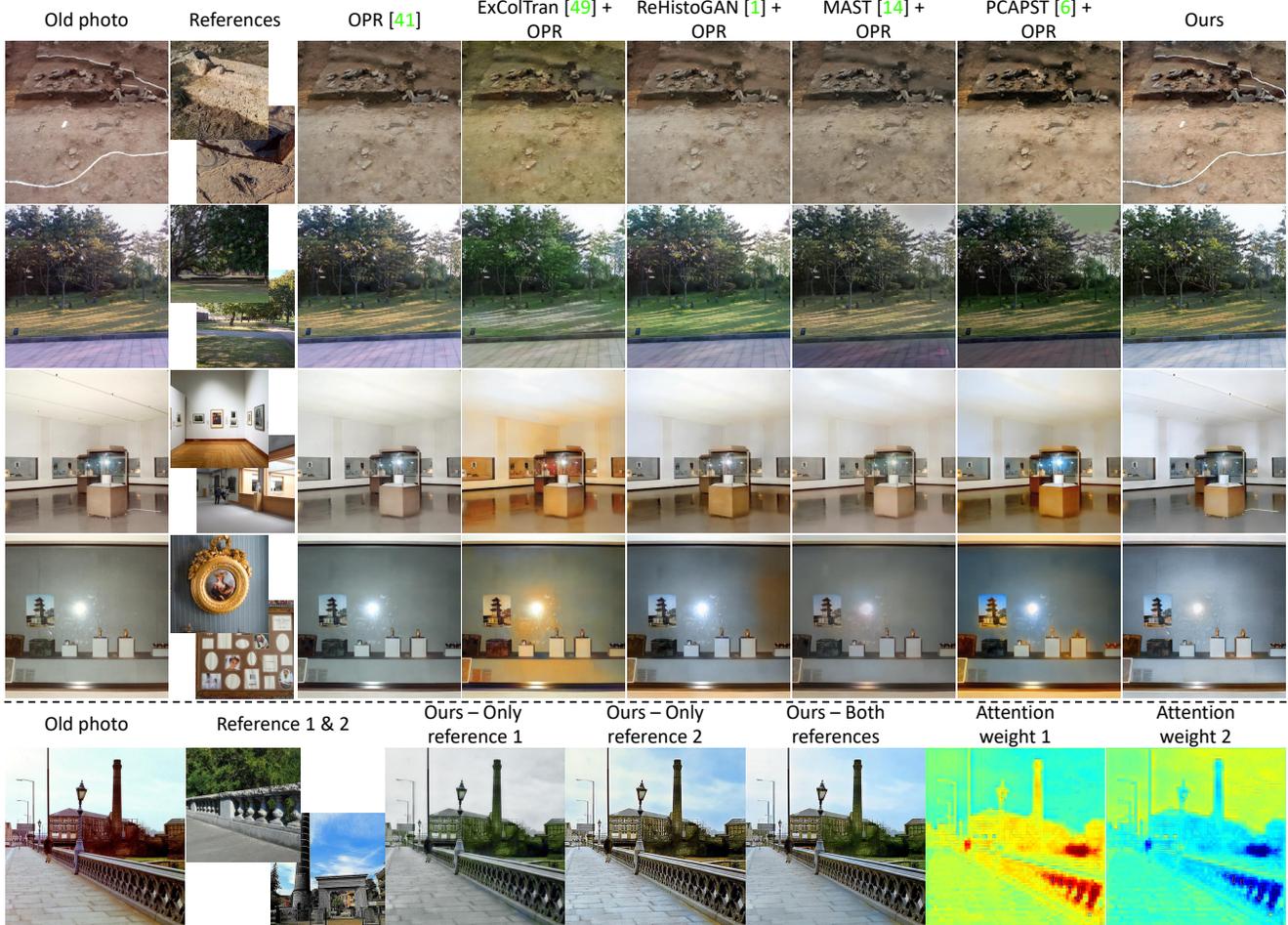


Figure 6. **Top**: qualitative results of modernization on real outdoor and indoor old photos. The baselines use top-left reference as their reference. **Bottom**: Attention weight visualization, blue (lowest)-red (highest) color coded. Our method can select appropriate styles from multiple references depending on the availability of similar objects to achieve better modernization. (Zoom in for a better view)

including ours, in a single-reference-based scenario since the baselines cannot utilize more than one reference by using ADE20K validation set [54] that includes semantic segmentation masks. Specifically, we generate 1,000 evaluation pairs, each consisting of a synthetic degraded photo and a reference image, by randomly degrading half set of the semantic regions using our synthetic data generation scheme, e.g., only the surfer in Fig. 5. Table. 1 shows our method achieves the best PSNR and LPIPS score, which means our method can effectively utilize the references to jointly stylize and enhance the synthetically degraded images, thus generating an output similar to the ground truth both in the pixel and semantic levels. In terms of SSIM, we achieve the second-best compared to recolorization (ReHistoGAN [1] + OPR-R) since our method can change the other aspects besides color, such as texture and luminance, which may result in a lower SSIM score. Interestingly, compared to other PST baselines, especially MAST [15], which is designed to perform semantic style transfer, our method

achieves the most accurate stylization (PSNR and LPIPS) while still preserving the structure (SSIM), which are two important aspects in PST. A similar observation can be seen for real old photos evaluation shown in Table. 2, where our method outperforms other baselines significantly using a single reference and further improves the performance by using multiple references.

Qualitative comparison. As shown in Fig. 6, no-reference OPR [42] can restore both structured (SD) and unstructured degradations (UD). However, SD restoration cannot generalize well to real old photos since it significantly degrades the important regions of the original photos. In addition, it fails to modernize some old photos because the overall styles still remain similar to the original old photos. ExColTran [50], which can only use luminance information for semantic matching, fails to locally change the color of old photos, thus producing unnatural results. Meanwhile, ReHistoGAN [1] can better recolorize old photos, producing more modern-looking images than only OPR.



Figure 7. Ablation study on the single stylization subnet.



Figure 8. Ablation study on the merging-refinement subnet.

Compared to other PST methods combined with OPR, and other baselines, our method achieves better local and global PST and yields enhancement, consequently achieving better modernization. Moreover, our method can utilize multiple references better in all examples, e.g., the second row of Fig. 6, where styles of tree leaves and road come from the first and second references, respectively. Meanwhile, the fourth row shows the generalization of our method, which can handle unrelated references well. The bottom part under the dotted line in Fig. 6 shows the visualization of spatial attention weights, where our method can select appropriate styles for each semantic object in an old photo from multiple references to achieve better modernization, e.g., the bridge and sky style in the first and second reference respectively. All in all, our method can modernize old photos better than the baselines by leveraging multiple modern photo references, even though it has not been trained with any old photos. These results also show that restoring the degradation of old photos cannot guarantee the outputs to look modernized, but changing their styles can contribute to the *modern* look more than restoring the degradation.

User study. We conduct a user study to compare the modernization results of our method with those of the baselines. Specifically, we select 130 photos from the CHD testing set and ask 18 users to rank the modernization results. As shown in Table. 3, our method outperforms other baselines with 56.37% chance selected as the best method.

5.3. Analysis

Ablation study on the single stylization subnet. We analyze the contribution of each module in the single stylization subnet \mathcal{S} . As shown in Fig. 7, the subnet fails to accurately transfer the local styles of objects, e.g., the styles of the blue building and grass, when the alignment module is removed. Even though the subnet can perform better local style transfer of building and grass regions with the alignment module, the stylization results are not smooth, which may produce unnatural results. Thus, adding a fusion module that merges global and local styles can produce smoother stylization locally and globally.

Ablation study on the merging-refinement subnet. To



Figure 9. Limitation of our method.

evaluate the contribution of the merging-refinement subnet \mathcal{M} , we change this subnet to a simple concatenation between multiple stylization features and feed the concatenated features into several convolution blocks (denoted as *w/o* \mathcal{M}). As shown in Fig. 8, without \mathcal{M} , the network cannot select appropriate styles from different references and fail to enhance the results. In addition, retraining the network is required to use a different number of references between inference and training. With our \mathcal{M} , we can adaptively choose the number of references without retraining. Results of other ablation studies and modernization using more than two references are in *Supplementary Material*.

Limitation. The limitation of our work mainly comes from the selection of references. As shown in Fig. 9, our method may produce unsatisfying modernization when a related object in the references has a style that does not enhance the old photo. However, finding references that contain a similar local object with a modern style in an automated way is highly challenging using existing image retrieval methods. Moreover, using VGG feature space matching similar to [10] fails to produce semantically similar references due to the domain gap between old and modern photos.

6. Conclusion

In this paper, we first proposed old photo modernization by using multiple references. In order to perform modernization, we proposed MROPM, which performs old photo stylization using multiple references via photorealistic style transfer and enhancement in one unified framework. Thanks to the generalization of PST and our synthetic data generation scheme, our work outperforms baselines for real-world old photos, even without using any old photos during the training. Furthermore, we analyze that our method can select appropriate styles from multiple references, further improving the modernization performance. Also, we propose an old color photos dataset CHD consisting of natural indoor and outdoor scenes to spur future research in the domain.

Acknowledgment. This research was supported by Culture, Sports and Tourism R&D Program through the Korea Creative Content Agency grant funded by the Ministry of Culture, Sports and Tourism in 2020 (Project Name: CHIC, Project Number: R2020040045, Contribution Rate: 100%). We would like to thank Gimhae, Jeju, and National Museum of Korea for the old photos.

References

- [1] Mahmoud Afifi, Marcus A Brubaker, and Michael S Brown. Histogan: Controlling colors of gan-generated and real images via color histograms. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7941–7950, 2021. 2, 6, 7
- [2] Jie An, Haoyi Xiong, Jun Huan, and Jiebo Luo. Ultrafast photorealistic style transfer via neural architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10443–10450, 2020. 2, 4
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018. 3
- [4] Peter J Burt and Edward H Adelson. The laplacian pyramid as a compact image code. In *Readings in computer vision*, pages 671–679. Elsevier, 1987. 4
- [5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocostuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 5
- [6] Tai-Yin Chiu and Danna Gurari. Pca-based knowledge distillation towards lightweight and content-style balanced photorealistic style transfer models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7844–7853, 2022. 2, 6
- [7] Tai-Yin Chiu and Danna Gurari. Photowct2: Compact autoencoder for photorealistic style transfer resulting from blockwise training and skip connections of high-frequency residuals. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2868–2877, 2022. 2
- [8] Ioannis Giakoumis, Nikos Nikolaidis, and Ioannis Pitas. Digital image processing techniques for the detection and removal of cracks in digitized paintings. *TIP*, 15(1):178–188, 2005. 1, 3
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [10] Mingming He, Dongdong Chen, Jing Liao, Pedro V Sander, and Lu Yuan. Deep exemplar-based colorization. *ACM Transactions on Graphics (TOG)*, 37(4):1–16, 2018. 2, 8
- [11] Mingming He, Jing Liao, Dongdong Chen, Lu Yuan, and Pedro V Sander. Progressive color transfer with dense semantic correspondences. *ACM Transactions on Graphics (TOG)*, 38(2):1–18, 2019. 2
- [12] Zhiyuan Hu, Jia Jia, Bei Liu, Yaohua Bu, and Jianlong Fu. Aesthetic-aware image style transfer. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3320–3329, 2020. 2
- [13] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 4
- [14] Zixuan Huang, Jinghuai Zhang, and Jing Liao. Style mixer: Semantic-aware multi-style transfer network. In *Computer Graphics Forum*, volume 38, pages 469–480. Wiley Online Library, 2019. 3
- [15] Jing Huo, Shiyin Jin, Wenbin Li, Jing Wu, Yu-Kun Lai, Yinghuan Shi, and Yang Gao. Manifold alignment for semantically aligned style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14861–14869, 2021. 2, 6, 7
- [16] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 6
- [17] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5, 6
- [18] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [20] Junyong Lee, Hyeongseok Son, Gunhee Lee, Jonghyeop Lee, Sunghyun Cho, and Seungyong Lee. Deep color transfer using histogram analogy. *The Visual Computer*, 36(10):2129–2143, 2020. 2
- [21] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3809–3817, 2019. 2
- [22] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. *Advances in neural information processing systems*, 30, 2017. 2, 4
- [23] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 453–468, 2018. 2
- [24] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics (TOG)*, 36(4):1–15, 2017. 2
- [25] Jixin Liu, Rui Chen, Shipeng An, and Heng Zhang. Cg-gan: Class-attribute guided generative adversarial network for old photo restoration. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 5391–5399, 2021. 3
- [26] Peng Lu, Jinbei Yu, Xujun Peng, Zhaoran Zhao, and Xiaojie Wang. Gray2colormnet: Transfer more colors from reference image. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3210–3218, 2020. 2
- [27] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4990–4998, 2017. 2

- [28] Xuan Luo, Xuaner Zhang, Paul Yoo, Ricardo Martin-Brualla, Jason Lawrence, and Steven M Seitz. Time-travel rephotography. *ACM Transactions on Graphics (TOG)*, 40(6):1–12, 2021. 1, 3, 6
- [29] Octavian-Mihai Machidon and Mihai Ivanovici. Digital color restoration for the preservation of reversal film heritage. *Journal of Cultural Heritage*, 33:181–190, 2018. 3
- [30] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017. 6
- [31] Roey Mechrez, Eli Shechtman, and Lihi Zelnik-Manor. Photorealistic style transfer with screened poisson equation. *arXiv preprint arXiv:1709.09828*, 2017. 2
- [32] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The contextual loss for image transformation with non-aligned data. In *Proceedings of the European conference on computer vision (ECCV)*, pages 768–783, 2018. 5
- [33] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing*, 21(12):4695–4708, 2012. 6
- [34] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a “completely blind” image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012. 6
- [35] Ying Qu, Zhenzhou Shao, and Hairong Qi. Non-local representation based mutual affine-transfer network for photorealistic stylization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):7046–7061, 2021. 2
- [36] 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. 5
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2, 5
- [38] Aditya Singh, Alessandro Bay, and Andrea Mirabile. Assessing the importance of colours for cnns in object recognition. In *NeurIPS 2020 Workshop SVRHM*, 2020. 2
- [39] F Stanco, Giovanni Ramponi, and A De Polo. Towards the automated restoration of old photographic prints: a survey. In *The IEEE Region 8 EUROCON 2003. Computer as a Tool.*, volume 2, pages 370–374. IEEE, 2003. 1, 3
- [40] Jheng-Wei Su, Hung-Kuo Chu, and Jia-Bin Huang. Instance-aware image colorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7968–7977, 2020. 2
- [41] Ziyu Wan, Bo Zhang, Dongdong Chen, and Jing Liao. Bringing old films back to life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17694–17703, 2022. 6
- [42] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Bringing old photos back to life. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2747–2757, 2020. 1, 2, 3, 6, 7
- [43] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 4
- [44] Phillip Whitt. *Beginning photo retouching and restoration using GIMP*. Springer, 2014. 1
- [45] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 4
- [46] Yanze Wu, Xintao Wang, Yu Li, Honglun Zhang, Xun Zhao, and Ying Shan. Towards vivid and diverse image colorization with generative color prior. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14377–14386, 2021. 2
- [47] Xide Xia, Meng Zhang, Tianfan Xue, Zheng Sun, Hui Fang, Brian Kulis, and Jiawen Chen. Joint bilateral learning for real-time universal photorealistic style transfer. In *European Conference on Computer Vision*, pages 327–342. Springer, 2020. 2
- [48] Runsheng Xu, Zhengzhong Tu, Yuanqi Du, Xiaoyu Dong, Jinlong Li, Zibo Meng, Jiaqi Ma, Alan Bovik, and Hongkai Yu. Pik-fix: Restoring and colorizing old photo. *arXiv preprint arXiv:2205.01902*, 2022. 1, 2, 3
- [49] Zhongyou Xu, Tingting Wang, Faming Fang, Yun Sheng, and Guixu Zhang. Stylization-based architecture for fast deep exemplar colorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9363–9372, 2020. 2
- [50] Wang Yin, Peng Lu, Zhaoran Zhao, and Xujun Peng. “Yes,” attention is all you need”, for exemplar based colorization. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 2243–2251, 2021. 2, 6, 7
- [51] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9036–9045, 2019. 2, 4, 5
- [52] Bo Zhang, Mingming He, Jing Liao, Pedro V Sander, Lu Yuan, Amine Bermak, and Dong Chen. Deep exemplar-based video colorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8052–8061, 2019. 6
- [53] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6
- [54] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 7