# Improving the Harmony of the Composite Image by Spatial-Separated Attention Module

Xiaodong Cun and Chi-Man Pun, *Senior Member, IEEE*

*Abstract*—Image composition is one of the most important applications in image processing. However, the inharmonious appearance between the spliced region and background degrade the quality of the image. Thus, we address the problem of Image Harmonization: Given a spliced image and the mask of the spliced region, we try to harmonize the "style" of the pasted region with the background (non-spliced region). Previous approaches have been focusing on learning directly by the neural network. In this work, we start from an empirical observation: the differences can only be found in the spliced region between the spliced image and the harmonized result while they share the same semantic information and the appearance in the non-spliced region. Thus, in order to learn the feature map in the masked region and the others individually, we propose a novel attention module named Spatial-Separated Attention Module (S²AM). Furthermore, we design a novel image harmonization framework by inserting the S²AM in the coarser low-level features of the Unet structure by two different ways. Besides image harmonization, we make a big step for harmonizing the composite image without the specific mask under previous observation. The experiments show that the proposed S²AM performs better than other state-of-the-art attention modules in our task. Moreover, we demonstrate the advantages of our model against other state-of-the-art image harmonization methods via criteria from multiple points of view.

*Index Terms*—Image Harmonization, Image Synthesis, Attention Mechanism

## I. INTRODUCTION

**P**Hotos play a more important role in daily life nowadays. The popularization of the smartphone and the advance of image processing applications also make image editing easily. Thus, any ordinary person, who has little knowledge of image editing, can synthesize two images together by photo manipulation software (Such as Photoshop) automatically under several easy steps [1]: The composite image is synthesized by cutting region of the donor image and then paste it to the host image with several post-processing techniques, such as Gaussian Smoothing. However, human eyes can easily identify the images are composited or not according to the color, texture or luminance differences between the splicing region and background. Meanwhile, these artifacts in the images degrade the visual quality of the whole image and make the image unrealistic. For these reasons, automatically image harmonization is an essential task in the digital image processing community. To address these problems, some previous works

Xiaodong Cun and Chi-Man Pun are with the Department of Computer and Information Science, University of Macau, Macau, 999078, China e-mail: {yb87432,cmpun}@umac.mo.

Corresponding author: Chi-Man Pun

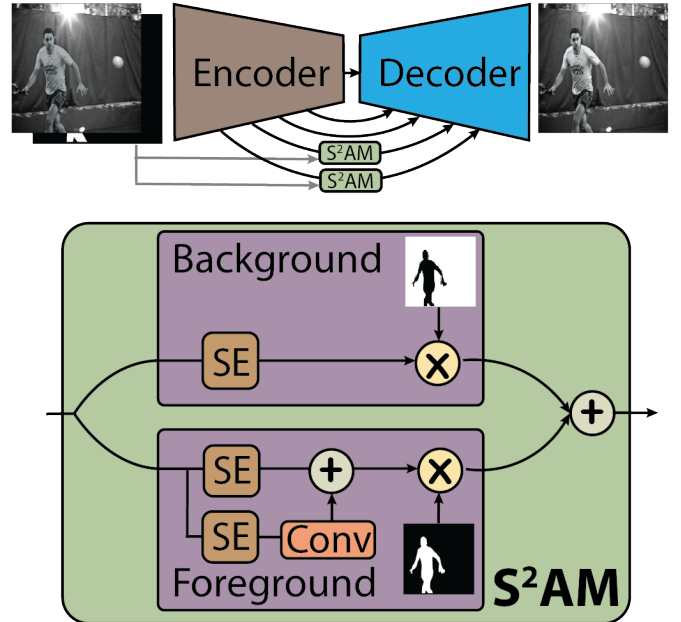Code is available at: https://github.com/vinthony/s2am

Fig. 1: The basic idea of our method. In the task of image harmonization, we argue that the spliced image (left) and the harmonized image (right) share the same high-level features as well as the appearances in the non-spliced region. The difference mainly comes from the appearance in the spliced region. So we design the S²AM module to learn these differences and make sure the consistency in the non-spliced region in low-level features. In detail, with the restrict of the hard-coded mask, we use several SE-Blocks [6] to re-weight the encoded features and learn the spatial location differently. This module is inserted in the coarser levels features of the encoder-decoder structure since the main network ensures consistency in high-level features. More detail of each component in S²AM can be found in the Section. III.

increase the quality of the synthesized image by hand-crafted features, such as color [2] and texture [3]. Nonetheless, these methods only focus on one particular kind of feature, causing unreliable results when the appearances are vastly different [4]. Recently, driven by deep learning based techniques, Zhu *et al.*[5] train a neural network to predict the realism of the composite image and optimize the color of the input image by freezing the parameters of the neural network, [4] harmonize the images by an end-to-end neural network with semantic guidance. However, their image harmonization method solves the non-convex model by multiple running [5] or rely on the power of learning directly [4].

Different from the previous methods, our key observation

is: *For image harmonization task, the model needs to learn the dissimilarity in the appearance of the spliced region and ensure the consistency in the non-spliced regions and high-level features.* As shown in the Fig.1, the spliced image and harmonized result only have differences in the appearance of the spliced region. This is because, from the viewpoint of image forgery detection, the artifacts in the composite images mainly come from the differences between various low-level features, such as camera internal features [7], different compression levels in JPEG compression [8] and noise levels [9]. Consequently, we need to design an algorithm that pays closer attention to learning the appearance differences in a certain region. Meanwhile, the spliced image and the harmonized result describe the same story by analyzing the semantic information and the sense category. So, in our proposed method, the input spliced image and the harmonized result need to share the same high-level information.

To fit previous analysis, we proposed a novel framework for image harmonization task. Firstly, following the previous work on this task [4], we regard the image harmonization as a supervised learning problem by synthesizing the composite images from some region of the natural image rather than the real copy-and-paste image. This alternative approach allows us to train the neural network in a supervised manner by considering the composite image and mask as input and natural image as the target. Then, we propose a novel attention module named Spatial-Separated Attention Module ($S^2AM$) to learn the feature changes in the masked region and the others separately. Our $S^2AM$ module uses several channel attention gates to weight features for different purposes with the help of the hard-coded mask and learns the changes in the specific masked region by a learning-able block. Moreover, we design the new framework of the image harmonization task. Our network base on the encoder-decoder [10] by inserting the $S^2AM$ module to low-level features in two different ways. Thus, the proposed $S^2AM$ module learns the regional appearance changes in low-level features, and the encoder-decoder structure will ensure the consistency in high-level features as well. Finally, we make a big step for image harmonization without giving mask as well. As far as we know, our method is the first approach for harmonizing the composite image without mask by deep learning. This goal is achieved by replacing the hard-coded mask in $S^2AM$ with the spatial attention block and we also use the additional supervision by measuring the loss between the spatial attention map and the ground truth mask.

We summarize the contributions as follows:

* We propose a novel attentive module named Spatial-Separated Attention Module ($S^2AM$) to learn the features in the foreground and background area by hard-coded masks individuality.
* We design two variants of the proposed $S^2AM$ in Unet structure for image harmonization task. Our network fits the intention of the difference in the spliced area and the similarities between high-level features and the non-spliced region. Besides, our network is fully-convolutional which is suitable for real-world image harmonization tasks.
* We generalize our method to the image harmonization

task without ground truth mask by the spatial attention module, attention loss and generative adversarial network.
* Experiments show that the proposed $S^2AM$ module gains better results than other attention mechanisms and our image harmonization method shows better performance than the state-of-the-art methods.

## II. RELATED WORKS

We do a brief review of the recent development of image harmonization and image-to-image transformation tasks. Importantly, we review the recently proposed attention mechanisms for the comparison of the proposed $S^2AM$.

**Image Harmonization** Recently, because of the efficient and impressive results, the neural network based methods have drawn much attention in image editing. However, in image harmonization field, there are still limited research works to solve this problem. Zhu *et al.* [5] propose a fine-tuned model for predicting the realism of the composite images. They also fix the parameters of the model and adjust the color of the masked region of the input image by optimizing the predicted realism score. It can also be regarded as an image harmonization task. [11] proposes a method for exchange the sky of images while this method only works for the sky. An end-to-end learning based approach of image harmonization has been proposed by [4], they trained an end-to-end neural network with the ground truth image and segmentation mask. However, their method just learns everything from the power of the neural network. Another related work is GP-GAN [12]. It utilizes the generator in GAN to synthesize the low-resolution images, and then, the high-resolution output is synthesized by the input image and the low-resolution outputs using traditional Gaussian-Poisson Equation. Their method focuses more on blending the edge of the two different images rather than the style transform in the images. Thus, this method also changes the appearance of the background while the image harmonization only focus on the spliced region. More recently, video harmonization methods have been proposed by [13] and [14]. Their methods harmonize the video and insert the video to another video using spatial and temporal information respectively.

**Attention Mechanisms** Attention mechanisms draw the attention of researchers because these modules allocate the most informative components according to analysis the features themselves. *Squeeze-and-Excitation Network* [15] insert a channel attention module in the network to make the network itself more focus on the essential features by weighting the channels with global feature analysis and *Convolutional Block Attention Module* [6] extend this method by adding a spatial attention module for weighting the features in spatial space. For choosing the feature kernel with different size, *Selective Kernel Networks* [16] is proposed. As for learning the specific region in spatial space, [17] propose *Partial Convolution* for focusing on the masked region and [18] design *Contextual Attention* to borrow the background information to hole filling. [19] use *Gated Convolution*, which can also be regarded as the spatial attention mechanism for each convolution block, to update the parameters of the network and mask (or attention map ) softly. However, these methods are designed for image
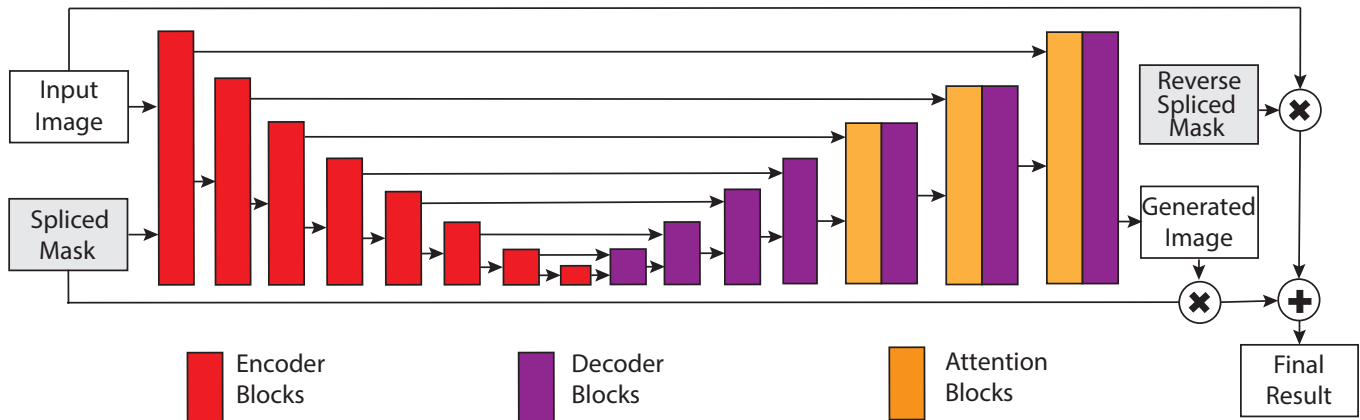
Fig. 2: The main framework of our network.

in-painting task specifically. Another noticeable direction of attention mechanism is self-attention. Self-Attention aims to capture the long-range dependencies in the image context. *Non-local Network* [20] extend the self-attention [21] as a form of non-local means [22] to capture long-range dependencies for video classification and [23] extends this method in GAN-based structure for generating divergence samples. Most recently, [24] model the *Non-Local Network* and *Squeeze-and-Excitation Network* in a new form called *Global Content Network* for various computer vision tasks. However, these attention mechanisms are focus on learning the representation on the high level computer vision task other than digital image processing with low-level features change.

**Image-to-Image Transformation** Image harmonization can also be regarded as an image-to-image transformation task which focuses on the changes the appearance between the input and target. Chen *et al.* [25] propose a neural network for fast image processing with the Dilated Convolution [26] to capture a larger reception field without reducing the feature map. By the excellent performance of guided filter [27] in edge preserving, [28] and [29] model this filter by the neural network. Taking the benefit of [25] and [27], [30] propose the deep guided filter for fast image processing. When the content of the image changes significantly, pix2pix [10] train the network for alignment datasets by GAN [31] while Cycle-GAN [32] is proposed for unsupervised learning. By combining the image-to-image transformation task with the attention network, Another noticeable direction of unsupervised image-to-image transformation is attention-based methods. [33], [34] extend CycleGAN with the attention network as an additional block and learning the foreground specifically. These methods generate a global attention mask for image-to-image transformation for specific objects only.

## III. METHODS

Following the previous work [4] on image harmonization, we create synthesized image harmonization datasets by modifying the appearance of the specific region of the original natural image. Thus, image harmonization task can be modeled as a supervised learning task by feeding the synthesized composite image and the corresponding mask to the network
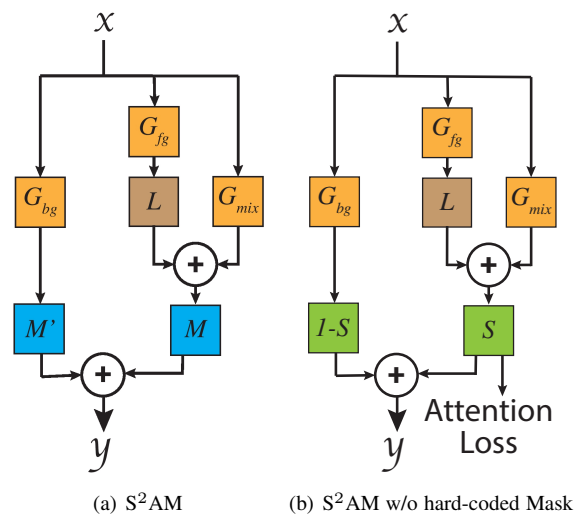


(a) S$^2$AM                    (b) S$^2$AM w/o hard-coded Mask

Fig. 3: The difference between the S$^2$AM and S$^2$AM without the hard-coded mask. We replace the original hard-coded mask by spatial attention block and attention loss from the ground truth mask.

and learning to harmonize the natural image. Fig. 2 shows an overview of our method for image harmonization. In detail, we use an encoder-decoder with skip-connection (Unet in [10]) as the backbone structure for our task. Then, with the interpolating of the proposed attention block, our network can learn the spliced region and background region separately, causing better results than previous works. In this section, we introduce the proposed Spatial-Separated Attention Module (S$^2$AM) in Section. III-A firstly. Then, for our method can learn the specific region automatically, we generalize our S$^2$AM module to image harmonization task without mask in Section. III-B. Next, we discuss the backbone network and the post-processing technique for image harmonization in Section. III-C. Finally, we discuss the loss functions in both tasks and synthesized datasets used in our framework in Section III-D and Section III-E, respectively.

## A. Spatial-Separated Attention Module ($S^2AM$)

The basic idea behind our $S^2$AM module is straightforward: the traditional convolution filter in convolution neural network share the kernel parameters over the whole image, which benefits higher-level tasks to look for a given feature everywhere in the image, rather than in just a certain area. However, for the image-to-image transformation in the certain region, such as image harmonization and inpainting, it shows a relatively poorer performance [17], [35]. Thus, with the help of the hard-coded mask, we design the attention module attempting to learn the features individually in the spatial space.

As illustrated in Fig. 3(a), given an intermediate feature map $x \in \mathbb{R}^{C \times H \times W}$ and the binary mask $M \in \mathbb{R}^{1 \times H \times W}$ of the spliced region as the inputs ( $C$ is the feature channel, $H$ and $W$ is the height and width of the feature map respectively ), our attention module utilizes three 1D channel attention gates $G_{bg}, G_{mix}, G_{fg}$ to weight the input features for different purposes with the help of hard-coded spliced region mask $M$ and non-spliced region mask $M'$. In detail, $G_{bg}$ aims to select the features which are necessary for image rebuild in the non-spliced region. $G_{fg}$ intends to weight the input features which have differences between the input and target. Moreover, we design the third channel attention module $G_{mix}$ to weight features which is unnecessary to change in the spliced region. Additionally, for the output features of $G_{fg}$, we add a learning-able block $L$ to learn the style transfer. So our RAB module between input features $x$ and output features $y$ can be formalized as:

$$y = M \times [L(G_{fg}(x)) + G_{mix}(x)] + (1 - M) \times G_{bg}(x) \quad (1)$$

Our $S^2$AM module is a general form on separated-learning the appearance changes in the specific region and others. So it is possible to design different structures of Channel Attention Modules $G_{bg}, G_{mix}, G_{fg}$ and learning-able block $L$ for the particular computer vision tasks. Here, inspired by the state-of-the-art attention modules, we design one instance of $S^2$AM for image harmonization as follows.

**Channel Attention Modules** ($G_{bg}, G_{mix}, G_{fg}$). With the help of hard-coded mask, the goal of the *Channel Attention Module* is to select the related features channels from the input by applying 1D channel weighting for input features. Here, we modify the state-of-the-art channel attention module in *Squeeze-and-Excitation Network* [6] by adding a MaxPooling layer. This channel gate is also a part of the structure in [15]. As shown in Fig. 4(a), we start from squeezing the input features $x \in \mathbb{R}^{C \times H \times W}$ to $x_{maxpooling} \in \mathbb{R}^{C \times 1 \times 1}$ and $x_{avgpooling} \in \mathbb{R}^{C \times 1 \times 1}$ by the *Global Max Pooling* and *Global Average Pooling* operators, separately. Then, we concatenate $x_{avgpooling}$ and $x_{maxpooling}$ in channel axis and learn the channel weighting by two fully-connected layers and a Sigmoid layer. Following previous work [6], [15], we reduce the size of the input feature to $\frac{C}{16}$ in the first Linear layer and increase the length to $C$ in the second and a Sigmoid layer is added at the end of Linear layer to ensure the range of scale factor under $[0, 1]$. The attentive modules show great success in image classification [6] and object detection tasks [15] with the structure of ResNet [36] or Inception [37]. As far as we
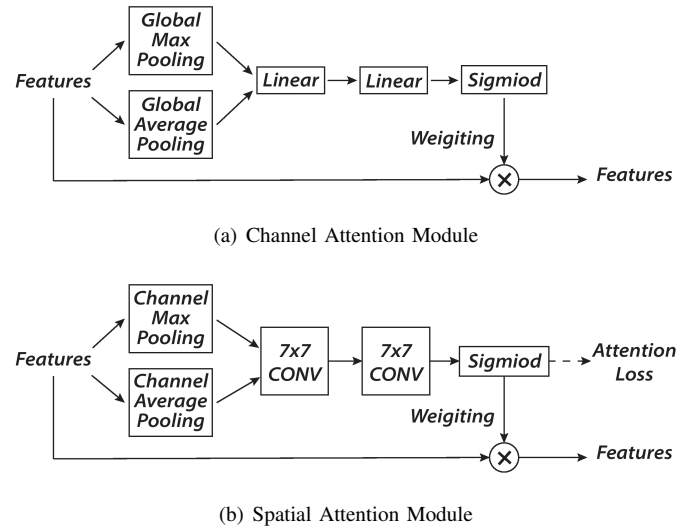


(a) Channel Attention Module



(b) Spatial Attention Module

Fig. 4: The detail of the *Channel Attention Module* and *Spatial Attention Module* in the proposed $S^2$AM.

know, we are the first to insert the attentive structure for the module with the hard-coded mask.

**Learning-able Block** ($L$) As shown in Fig. 3(a), we aim to learn the appearance changes in the spliced regions by Learning-able Block additionally. We use $3 \times 3$ convolutional layers with two CONV-BN-ELU blocks. Similar to channel attention module, we double the channel size in the first CONV-BN-ELU block and reduce the channel size in the second to expend the learning capacity. Notice that, we do not reduce the feature size in the learning-able blocks for preserving the feature details.

**Mask** ($M, M'$) Spliced region mask $M$ and non-spliced region mask $M'$ are the necessary parts for our $S^2$AM module because the channel attention modules are guided by these hard-coded masks. Furthermore, they provide a strong prior knowledge to the neural network. We also multiply a uniform $7 \times 7$ Gaussian Kernel to the mask $M$ and the reverse mask $M'$ respectively. This simple operation allows the neural network to generate the boundary by the closest object for the hard-coded masks are not always accurate. We discuss the effect of Gaussian Kernel in our attention module in the Section. IV-E2

## B. Generalizing $S^2AM$ to Image Harmonization without Mask

We generalize the $S^2$AM module for image harmonization without the hard-coded mask as input. Thus, we first replace the original hard-coded mask with the spatial attention block. We design the spatial attention module with the similar structure of the previous introduced *Channel Attention Modules*. For the *Spatial Attention Module* aims to learn the spatial position separately, we use the *Max Channel Pooling* and *Average Channel Pooling* to extract the maximum responses and medium responses from all feature channels for each pixel. Then, unlike the channel attention module which used fully-connected layer for channel weighting. It is memory-consuming *to* capture global features. Thus, we use two Convolutional layers with the kernel size $7 \times 7$ as alternative.

The first Convolutional layer broadcast the features to 16 layers while the second squeeze the channel map to fit the output of Sigmoid. Then, a Sigmoid layer is added at the end of the module to restrict the response of features weighting in $[0, 1]$. Notice that, there is only one *Spatial Attention Module* in our S$^2$AM module because we attempt to force the spatial attention map in two folds without overlap. The reverse mask $M'$ is calculated by $M' = 1 - M$. We show the structure of the spatial attention module in Fig. 4(b).

Moreover, we add a loss between the output of the sigmoid layer and the ground truth mask. This attention loss guide the focus of spatial attention module for the features from $7 \times 7$ Convolutional block is not always accurate. We review the choice of attention loss and the influence of attention loss in Section.III-D and Section. IV-D.

### C. Network Structure

As the previous analysis, we model the image harmonization task as an image-to-image transformation task. Differently, our task needs to: (1) ensuring the consistence in high level features. (2) taking more attention to the low-level appearance changes. Thus, we use the encoder-decoder structure as the backbone network to ensure the consistency in high-level features and S$^2$AM for learning appearance changes carefully. We introduce the network structure in our approach below.

**Network Backbone** Different from the encoder-decoder structure in *Deep Image Harmonization (DIH)* [4], the state-of-the-art work on image harmonization, we use Unet [10], [38] as our backbone network. Unet owns more high-level filters and treats the global features in the fully convolutional manner while DIH uses Linear Layer which will hugely increase the parameters of networks and limit the input size. Another difference between these two backbones is that DIH use ELU for non-linear activation while Unet chooses RELU/Leaky RELU. As for the operation between the skipped encoder features and decoder features, we concatenate the features from encoder rather than add or mix them. The experiment shows that the new baseline network improves the results of harmonization task to a large extent.

**The Role of S$^2$AM** The first choice of interpolating S$^2$AM to Unet is based on a general observation: The original skip-connections preserve the encoded low-level features in the original image and these features are concatenated to the low-level features of the decoder. The concatenation of the low-level features from different appearances will mislead the learning process to make a good choice. Thus, we replace the original skip-connection in the encoder-decoder by the proposed attention module. Our attention module will filter the features in the encoder and learn the appearance changes in the spliced region and non-spliced region individually. We call this type of connection *Spatial-Separated Attentive Skip-Connection (S$^2$ASC)*. Moreover, we also interpolate the proposed module after the concatenation of the shadow features from the encoder and the decoded features from high-level network (as the network structure in Fig. 2). This kind of choice will use more features (features from encoder and decoded features from high-level) for learning a better channel

attention decision to learn separately. Thus, the appearance changes will also contain the influence of the global features. We name this type of structure *Spatial-Separated Attentive Decoder (S$^2$AD)*. We show the differences between S$^2$ASC and S$^2$AD in Fig. 5.

Notice that, the proposed S$^2$ASC and S$^2$AD module only replace the original three levels of the structure in the unet for ensuring a similar high-level feature (as in Fig. 2). We also verify this assumption in Section. IV-E2.
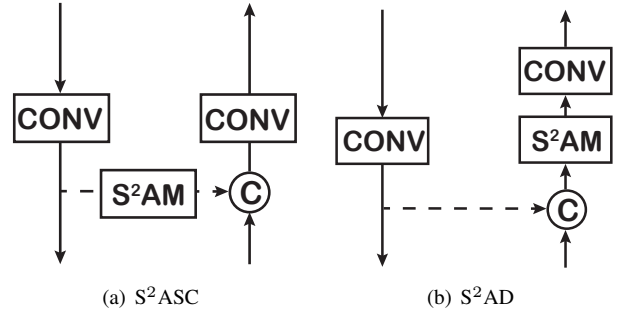


(a) S$^2$ASC      (b) S$^2$AD

Fig. 5: The positions of the proposed module.

**Post-Processing (pp)** As for image harmonization task, the background region should stay unchanged when we generate the final output. Thus, in the inference stage, we mask the harmonized region by the original ground truth mask and replace the others with the original image. For image harmonization without mask, we train the post-processing altogether. The final output $I_{final}$ is calculated by the coarser attention mask $M_{2s}$ and input image $I$ of our network as:

$$I_{final} = I_g \otimes Up(M_{2s}) + I \otimes (1 - Up(M_{2s})) \quad (2)$$

where $\otimes$ represent the pixel-wise multiplication and $I_g$ is the output of the network. $Up(\cdot)$ means the Upsample layer for the attention mask 2s is the $1/2$ size of the original image. Notice that we learn the whole network in image harmonization without mask task while we use this technique as post-processing in image harmonization task. This is because the predicted mask is not always accurate and the additional loss in image harmonization without mask will learn to make the overall image look realistic rather than the attention mask region only.

### D. Loss Function

We design different loss functions for image harmonization and image harmonization without mask since the latter task have weaker prior information.

*1) Loss Function for Image Harmonization:* Following our baseline method [4], we use the $L_2$-norm criterion as the main criterion in the pixel domain. Besides, *Perception Loss* [39] and *Adversarial Loss* [31] ( pix2pix method in Section. IV ) are added individuality as the additional loss to constrict the accuracy of results in the feature domain and distribution. These two criteria have shown great success in various image generation tasks. However, we do not observe a significant improvement in our task which focuses on low-level feature processing. Similar experiment results have been observed in image processing [25].

*2) Loss Function for Image Harmonization w/o Mask:*
Different from previous image harmonization task, we use multiple criteria for Image harmonization without the mask. Firstly, as the same loss function in image harmonization, we set the pixel-level criterion by $L_2$ norm for the harmonized image should have the same pixel values with the target in the corresponding spatial position. We mark the output image as $I_{predict}$ and the target image as $I$. The pixel level loss function can be written as:

$$L_{pixel} = ||I_{predict} - I||_2 \qquad (3)$$

Although our spatial attention module can learn the foreground features and background features individually, it is still a challenging task when there is no ground truth mask as input. However, for the ground truth spliced mask is easy to get by the synthesized dataset, we add the attention loss in the spatial attention module of the S$^2$AM module for guiding the image synthesis. In specific, we resize the original ground truth mask to constrict the output of the sigmoid in every level of spatial attention mask. It is obvious that the gated attention can be considered as a supervised binary classification task, we use the Binary Cross Entropy as the criterion firstly. Then, every pixel in the image has spatial connections. Thus, $L_2$ Loss is also a satisfying option. We evaluate all the choices of the loss function in Section IV-D. Finally, we choice $L_2$ loss as the attention loss in each spatial attention block:

$$L_{att} = \sum_{i=1}^{3} ||M_i^{'} - M_i||_2 \qquad (4)$$

$M_i^{'}$ is the $i$-th level attention mask of S$^2$AM while $M_i$ is the corresponding ground truth spliced mask.

Moreover, we use the adversarial loss [10] for removing the artifacts in the predicted image. Adversarial loss has shown much success in many image-to-image transformation tasks and image generation tasks. Different from original works in Pix2pix [10], we use the Least-Square GAN as the adversarial loss following the work of CycleGAN [32]. For image harmonization task, we regard the Unet with S$^2$AD is the generator of GAN and we use the same Patch-based discriminator structure of pix2pix [10].

Then, the total loss for image harmonization without mask can be written as:

$$L = \alpha L_{attention} + \beta L_{pixel} + L_{gan} \qquad (5)$$

We experimentally set the $\alpha = 90$ and $\beta = 100$ in all experiments. We also evaluate the different settings of hyperparameters $\alpha, \beta$ in Section. IV-D.

### E. Data Acquisition

Since there is no publicly available dataset for image harmonization, we build the below two datasets to compare for different reasons.

**Synthesized COCO (S-COCO)** To test the robustness of approaches over the tampered images with ground truth instance mask, we build the S-COCO dataset. COCO [40] contains various natural images with ground truth instance

mask annotated by the human. Different from DIH [4], we tamper the images by the state-of-the-art style transfer methods [30]. As shown in Fig. 6, firstly, we temper the original natural image by the pre-trained *auto-ps* model. Then, for further simulating the copy-and-paste image in the real world, we randomly enhance the appearance of the transformed region with brightness, color and contrast. Finally, we crop the tampered image by the ground truth instance mask and paste the cropped region to the corresponding position of the natural image. We use the Train2017 and Val2017 split in COCO for training and testing, respectively. We only select the tampered image with the spliced region larger than 15% of the whole image since smaller regions may harm the consistency of the prediction. Overall, there are about 43k and 1.7k images for training and testing respectively.

**Synthesized Adobe5K (S-Adobe5K)** We evaluate the robustness of the image harmonization methods on inaccurate ground truth mask and manually copy-and-paste under Adobe5K dataset [41]. Adobe5K contains 5k raw photos, and each image is dedicated to photo adjustment by five experts using Adobe Lightroom software. So each image has five different natural styles. To get the ground truth image spliced binary mask, we segment image from one specific style by the pre-trained state-of-the-art semantic segmentation model, DeepLab-V3 [42]. We generate the binary mask by each class of segmentation results. Finally, the samples are created by pasting the mask region from other different styles to the original image. We split the original dataset for learning randomly and filter out the smaller region as S-COCO. There are 36k images for training and 2k images for testing in S-Adobe5K.

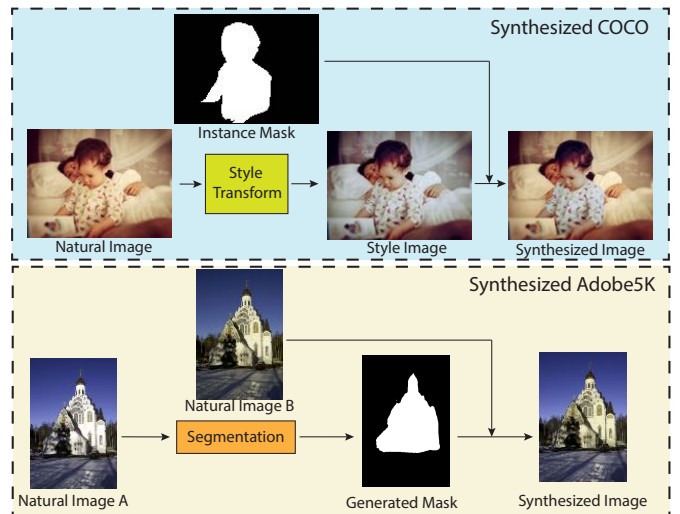The overall progress is shown in Fig 6. We will public these two synthesized datasets for further research.



Fig. 6: The workflow of building synthesized datasets.

## IV. EXPERIMENTS

In this section, we conduct experiments on our methods with state-of-the-art methods on image harmonization. Firstly, we illustrate all the training/testing details of our methods

and the baselines for comparisons in Section. IV-A and Section. IV-B respectively. Then, for the comparison, we evaluate our method for image harmonization and image harmonization without the mask on several datasets in Section. IV-C and Section. IV-D. Finally, a detailed analysis of the proposed $S^2$AM is reported in Section. IV-E.

### A. Implementation detail

All the models are trained with PyTorch [43] v1.0 and CUDA v9.0. We train the image with the resolution of $256 \times 256$ and run the model 80 and 60 epochs for converging on S-COCO and S-Adobe5K, respectively. All the optimizers are Adam [44] with the learning rate of 0.001, and there is no learning rate adjustment schedule for a fair comparison. For testing, our model runs at 0.012 seconds pre-image on a single NVIDIA 1080 GPU with a resolution of $256 \times 256$. Since the synthesized datasets have the ground truth target, we evaluate our approach with other methods on multiple popular numerical criteria, such as Mean Square Error (MSE), Structural Similarity (SSIM) and Peak Signal-to-Noise Ratio (PSNR). All the results in the synthesized datasets are the mean of the 1.7k and 2.4k test images in S-COCO and S-Adobe5K, separately. Moreover, following the metrics in previous works [4], [12], a pre-trained realism prediction model [5] is used to evaluate all the results from the viewpoint of the neural network. Finally, We evaluate the state-of-the-art methods on the real composite images which are collected from the real world with a user study.

### B. Baselines

We list all the image harmonization baselines for comparison on the synthesized datasets as follows:

**Deep Image Harmonization (DIH) [4]** DIH harmonizes images in a supervised manner with an encoder-decoder structure. Additionally, they use the ground truth segmentation mask as a redundant decoder since the harmonized images should also share the same segmentation results with the target. We train the DIH on the same synthesized datasets for fair comparison. As our $S^2$AM can be plugged into DIH easily by replacing the original skip-connections in color and semantic branch, we modify the original model for comparison on S-COCO and S-Adobe5K. Notice that, because S-Adobe5K has no ground truth segmentation mask, we only evaluate the color branch.

**RealismCNN (R-CNN) [5]** RealismCNN uses the pre-trained VGG16 network as the basic feature extractor and fine-tune the model in the fully-connection layer for classifying the realism of the input image. Additionally, they freeze the parameters in the realism model and optimize the color of the input image by a novel loss function. However, this method tries to solve a non-convex function by selecting the minimal cost from multiple running. We compare the results on our datasets with the official implementation[1].

**pix2pix [10]** Pix2pix use GAN in the image-to-image transformation task firstly. GAN has significant benefits synthesizing the meaningful photo-realistic images because the

| Method | MSE↓ | SSIM↑ | PSNR↑ | LPIPS↓ |
|---|---|---|---|---|
| S-COCO dataset | | | | |
| Original Copy-and-Paste | 25.94 | 0.9469 | 26.62 | 0.0507 |
| R-CNN [5] | 31.21 | 0.9509 | 27.17 | 0.0458 |
| DIH [4] | 20.88 | 0.9662 | 31.73 | 0.0308 |
| **DIH + S²ASC** | 18.09 | 0.9734 | 33.55 | 0.0174 |
| pix2pix [10] | 21.24 | 0.9633 | 32.13 | - |
| **pix2pix [10] + S²ASC** | 17.67 | 0.9707 | 33.32 | - |
| Unet [10] | 17.28 | 0.9757 | 33.55 | 0.0162 |
| **Unet + S²ASC** | 15.59 | 0.9790 | 34.74 | 0.0144 |
| **Unet + S²AD** | 14.88 | 0.9790 | 35.23 | **0.0128** |
| **Unet + S²AD + pp** | **14.26** | **0.9811** | **35.38** | **0.0128** |
| S-Adobe5K dataset | | | | |
| Original Copy-and-Paste | 37.43 | 0.9401 | 26.60 | 0.0501 |
| R-CNN [5] | 41.29 | 0.9338 | 26.11 | 0.0506 |
| DIH[4] | 34.71 | 0.9243 | 27.94 | 0.0506 |
| **DIH + S²ASC** | 26.96 | 0.9586 | 31.45 | 0.0282 |
| pix2pix [10] | 32.26 | 0.9560 | 30.01 | - |
| **pix2pix [10] + S²ASC** | 22.67 | 0.9702 | 32.59 | - |
| Unet [10] | 23.18 | 0.9693 | 32.46 | 0.0232 |
| **Unet + S²ASC** | 21.83 | 0.9711 | 33.31 | 0.0210 |
| **Unet + S²AD** | 21.52 | 0.9743 | 33.67 | **0.0186** |
| **Unet + S²AD + pp** | **20.47** | **0.9757** | **33.94** | 0.0195 |

TABLE I: The numerical comparison on synthesized datasets. Here, $S^2$ASC and $S^2$AD are the different interpolate method of our $S^2$AM. **pp** represent the post-processing in Section. III-C .

discriminator in GAN distinguishes the real or fake image by the neural network itself. However, it fails to get the accurate numerical results. Since we use unet as our backbone network, we compare the unet, unet with $S^2$AM as Skip-connection (Unet+$S^2$ASC) and unet with $S^2$AM as Decoder (Unet+$S^2$AD). Furthermore, we compare pix2pix by considering Unet and Unet+$S^2$ASC as the generators of GAN respectively. Notice that, we train the GAN framework under the default configuration by the author's implementation[2].

### C. Comparison

*1) Comparison on Synthesized Datasets.:* For fair compare the effect of the backbone network and our attention module, we mainly compare the state-of-the-art methods on Unet+$S^2$ASC. Additionally, we compare the Unet+$S^2$ASC with Unet+$S^2$AD individually.

In S-COCO, we train all the learning-based methods under the same framework except the model structure. As shown in Table. I, RASC module get better numerical results when it is plugged into three baseline structures: DIH, unet and pix2pix. Moreover, our full method (Unet+$S^2$ASC) shows the best results. These experiments confirm that our method is more suitable for image harmonization task comparing with others because their approaches only gain the results relying on larger datasets. Additionally, our method achieves better visual effects. We plot the harmonized results in Figure. 7 and calculate the absolute difference colormap between the target and result for better visualization. It is clear that our method outperforms other methods to a large extent. Figure. 7, our method shows better results in various image content and color form. It is obvious that our approach gains better results when

(a) Input/Mask    (b) R-CNN[25]    (c) DIH[4]    (d) DIH+    (e) pix2pix    (f) pix2pix+    (g) UNET[10]    (h) UNET+    (i) Target
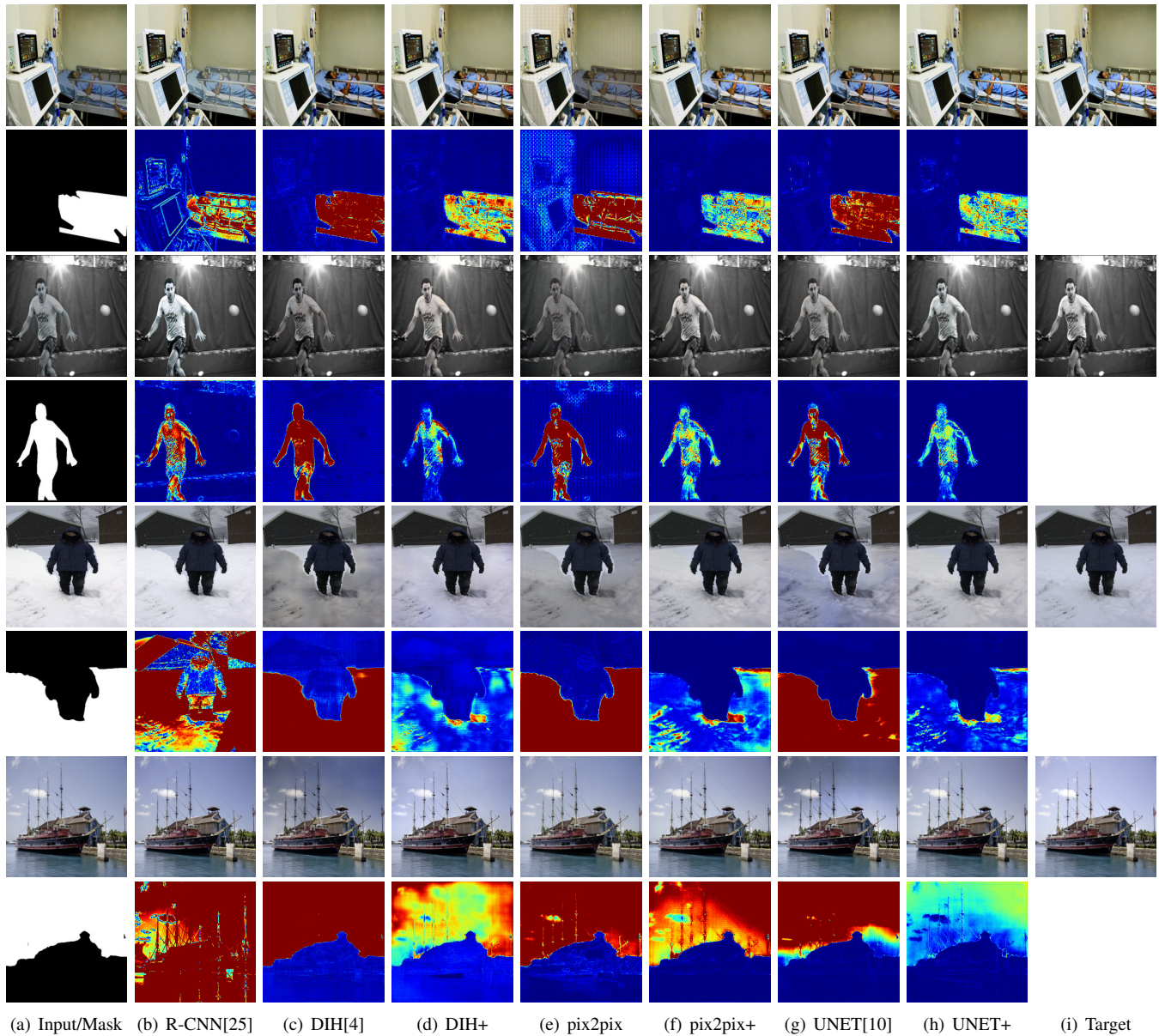
Fig. 7: Comparison on synthesized datasets. The two top samples are selected from S-COCO while the two bottom images come from S-Adobe5K. The methods with $+$ represent we replace the original coarser skip-connection with our $S^2ASC$. To visualize the differences better, we plot the absolute differences colormap between the harmonized image and target under each result. It is apparent that our method shows better results in both image quality and difference colormap. We enlarge the colormap $10\times$ for better visualization.

the bed or person is the spliced region and performs well in both gray images and color images.

Similar improvements are also observed in S-Adobe5K. This dataset is more complicated than S-COCO for the data is limited, and the mask is inaccurate. We show the visual comparisons in Figure. 7 and report the numerical results in Table. I. Both results indicate that our methods gain better performances than others. As shown in the third example in Figure. 7, the three baseline methods with $S^2AM$ perform better with inaccurate ground truth mask because our $S^2AM$ filter the regional features separately in the skip-connection.

Besides the numerical and visual comparison, we also evaluate all the methods on the pre-trained composite realism

model in R-CNN [5]. Notice that, the pre-trained model has not been trained on any images in our dataset before. The predicted score in Figure. 8 shows the predicted score of the composite images. It is clear that our $S^2AM$ improves the predicted scores from the viewpoint of the pre-trained model in most cases. Meanwhile, our Unet+$S^2AM$ method shows the best performance results comparing with others. However, it is not surprising that R-CNN performed better when the experiment is performed on its own pre-trained model. Interestingly, for R-CNN, we do not observe a similar phenomenon in S-Adobe5K. It is probably because the pre-trained model optimizes the difference in color by semantic while the mask is not always meaningful in S-Adobe5K.
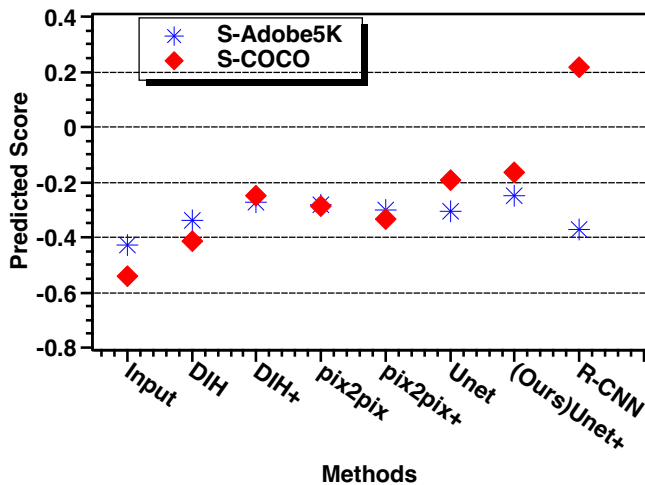
Fig. 8: The predicted score on the pre-trained model of R-CNN [5]. The methods with $+$ represent we replace the original coarser skip-connection with our $S^2$ASC. A higher score means a more realistic result.

We compare the results of $S^2$ASC and $S^2$AD on two synthesized datasets. We argue that, in Unet+$S^2$AD, our $S^2$AM learns from the high-level information and low-level information altogether. Besides the numerical analysis in Table. I, $S^2$AD show better global consistence that $S^2$ASC as shown in Fig. 9. It is clear that $S^2$AD gets the benefits from the upsampled high-level features and generates more realistic results from the global context. For example, the shoulder of the human in the third example and the hat of the child in the last sample are not synthesized well in $S^2$ASC while we get the better results from the interpolate of high-level features and low-level information from encoder.

*2) Comparison on Real Dataset with User Study.:* Although we train the model on the synthesized dataset, our method can also harmonize the real samples of the spliced image. We compare our method with the two most relevant state-of-the-art methods on the real composite images with provided spliced mask. This dataset contains 99 images with various different scenes and it is collected by DIH [4]. As shown in Figure. 10, our method out-perfumes other states-of-the-art methods with a larger margin on this dataset. For there is no available ground truth target image for the harmonized version of these images, we conduct a user study to evaluate our method on the real dataset and synthesized datasets by Amazon Mechanical Turk. Our user study template derivates from the web page in [12]. Specifically, for each task, giving the original copy-and-paste image and corresponding mask, the users need to choose the best realistic image from all the harmonized results created by different algorithms. In user study, we use the structure of Unet+$S^2$AD structure and train our model on two synthesized datasets with 30 epochs from scratch. As for the comparison, we compare our method with the pre-trained DIH model provided by author[3]. This model pre-trained on the semantic segmentation task and fine-tune the network on three datasets synthesized by their own. As shown in the Table. II, our method gains much more votes
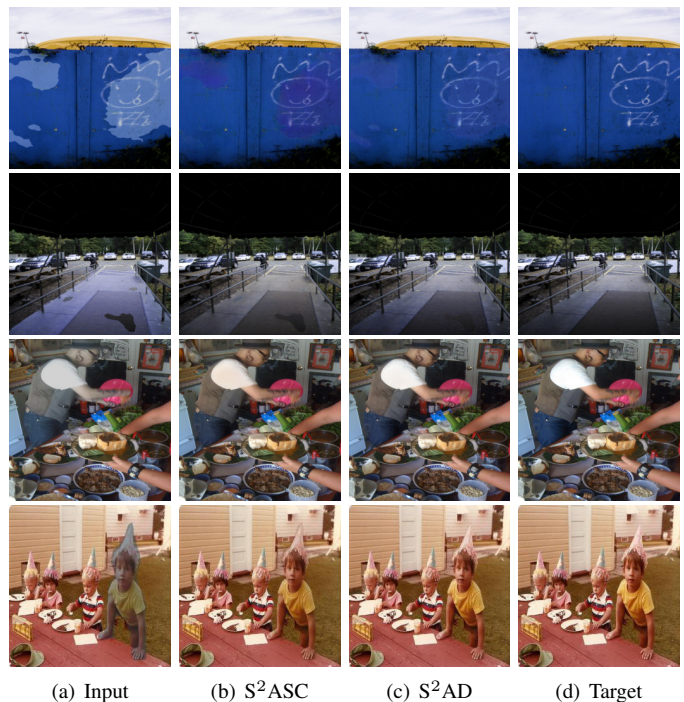
[3]https://github.com/wasidennis/DeepHarmonization



(a) Input    (b) $S^2$ASC    (c) $S^2$AD    (d) Target

Fig. 9: Comparison between the $S^2$ASC and $S^2$AD. Although $S^2$ASC harmonizes the tampered region, there are still some fake details while $S^2$AD show realistic global results.



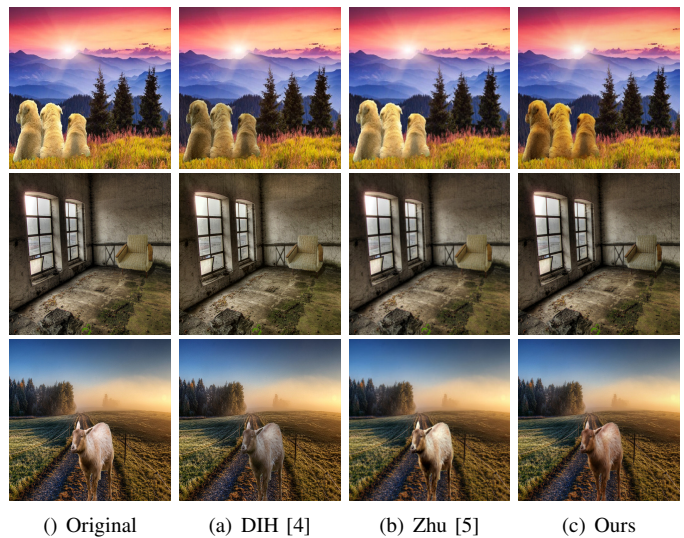() Original    (a) DIH [4]    (b) Zhu [5]    (c) Ours

Fig. 10: Comparison on the real dataset. As shown in the figure, our method gain better results than DIH [4] and Zhu [5] with the pre-trained model provided by the author.

than others with a larger margin.

*D. Image Harmonization without mask*

*1) Comparison on Synthesized Dataset:* As described in Section. III-C, our method can be adapted to the image harmonization task without the ground truth as input. So in this task, we regard the previous Unet as the baseline network structure by feeding the color image to the network only. We evaluate our method on the S-COCO dataset for comparison.

| Method | Total votes | Average votes |
|--------|-------------|---------------|
| Zhu [5] | 120 | 21% |
| DIH [11] | 212 | 35% |
| Ours | 263 | 44% |

TABLE II: The comparison between our method and other state-of-the-art methods under user study.

As shown in Fig. 12, even without the feeding of the tampered mask, our method can still predict the realistic harmonized image and achieve much more better performance than the baseline unet [10] both visual and numerical quality. We also plot the intermediate attention map (coarser level of the attention loss) created by our methods in Fig. 12. The mask predicted by our method is reasonable. We also report the numerical results under this task in Table. III, surprisedly, our full method get comparable results with the DIH methods (which need the mask as input) under the same datasets. Table. III also report the intermediate results of our methods. Each experiment shows the necessity in the network.

| Method | MSE↓ | SSIM↑ | PSNR↑ |
|--------|------|-------|-------|
| Original Copy-and-Paste | 25.94 | 0.9469 | 26.62 |
| Unet (w/o mask)  [10] | 29.86 | 0.9518 | 28.96 |
| Ours w/o Attention Loss | 26.08 | 0.9608 | 30.10 |
| Ours w/o GAN Loss | 23.93 | 0.9622 | 30.89 |
| Ours w/o post-processing | 23.86 | 0.9590 | 30.61 |
| Ours (Unet+S$^2$AD) | **20.73** | **0.9657** | **31.15** |

TABLE III: We evaluate the proposed method for image harmonization without the ground truth mask as input on the S-COCO datasets.

*2) Evaluation of Attention Loss:* The key component of our method for image harmonization without the mask is attention loss. We tuning the choice of attention loss of $L_2$, *Binary Cross Entropy* under a subset of S-COCO dataset. This subset contains 10k training images and we test it on the same test dataset as the full S-COCO dataset. We choose this subset S-COCO for evaluating the attention loss because it trains $4\times$ faster than the original dataset. As shown in the red line and blue line in Fig. 11, $L_2$ loss is a more suitable choice for attention loss than *Binary Cross Entropy* under the same ratio between the pixel loss. We also try to find the best proportion by rescaling the percentage of the attention loss to $0.1\times$, $0.01\times$, $10\times$. From the Fig. 11, it is clear that the best proportion between attention loss and pixel-level loss is $1 : 10$.

### E. Attention Mechanisms Analysis

We first evaluate the proposed S$^2$AM against other state-of-the-art attention mechanisms: *Squeeze-and-Excitation Network (SE-Block)* [6] and *Convolution Block Attention Module (CBAM)* [15] for image harmonization task in Section. IV-E1. Their methods are proposed for image classification tasks, so we replace the original skip-connection in Unet for comparison with the proposed S$^2$ASC. Besides, we increase the channels in each levels of the original Unet to match the parameters of our Unet+S$^2$AM model for a fair comparison (+ More channels in Table. IV). Finally, we compare the attention modules with the basic CONV-Block. Each
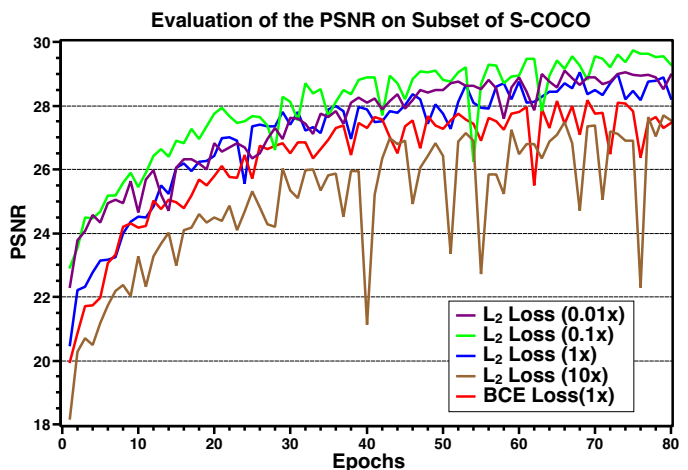


Fig. 11: The Evaluations of the attention loss on the Subset of S-COCO. We plot the PSNR on the full test set in each epoch when we train the network with different attention loss. We also tune the hyper-parameters $\alpha, \beta$ in Equation.5. $X\times$ means the proportion between the pixel level loss parameters $\alpha$ and our attention loss $\beta$. This figure is best view in color.

CONV-Block block has the structure of CONV $(3 \times 3)$-ELU-BN-CONV $(3\times3)$-ELU-BN, which is the same as the learnable block in S$^2$AM. Then, we conduct a detailed ablation study of our module in Section. IV-E2 and Section. IV-E3. All the methods are experimented on S-COCO dataset in the same framework with default configurations.

*1) Evaluation of Attention Modules:* As illustrated in Table. IV, Unet+CONV-Block and Unet+SE-Block show poor performance than the baseline method for similar reasons. In these methods, the convolution operators are performed on the whole image and their method cannot learn the disparities in the spliced region and preserve the features in the non-spliced region simultaneously. SE-Block shows better results than CONV-Block because SE-Block learns the channel attention additionally in the skip-connection. Moreover, it is not surprising that CBAM gets better results than the baseline since the spatial attention part in CBAM can learn the specific region automatically, and the input mask of our network makes it easy. However, our unet+S$^2$ASC method gets a better result than other attention modules with a large margin. This is because, with the help of the hard-coded mask, we can design more channel attention modules for different purposes and learning specifically while the SE-Block and CBAM can only focus on certain channels or certain regions by learning from data. Besides, the baseline network with more parameters also improve the performance, however, it still far worse than our method. Overall, the benefits of our attention module are obvious. From Fig.13, our method gets a cleaner difference colormap in the background than other attention modules for their methods only consider soft attention. Another improvement between Figure.13(e) to Figure.13(f) is the boundary of the spliced mask. From the figure, we can find our method gets better harmonized edges while the ground mask is not always true while other methods fail. It is probably because the soft spatial attentions in CBAM can not always get an accurate attention map even with the input ground truth mask

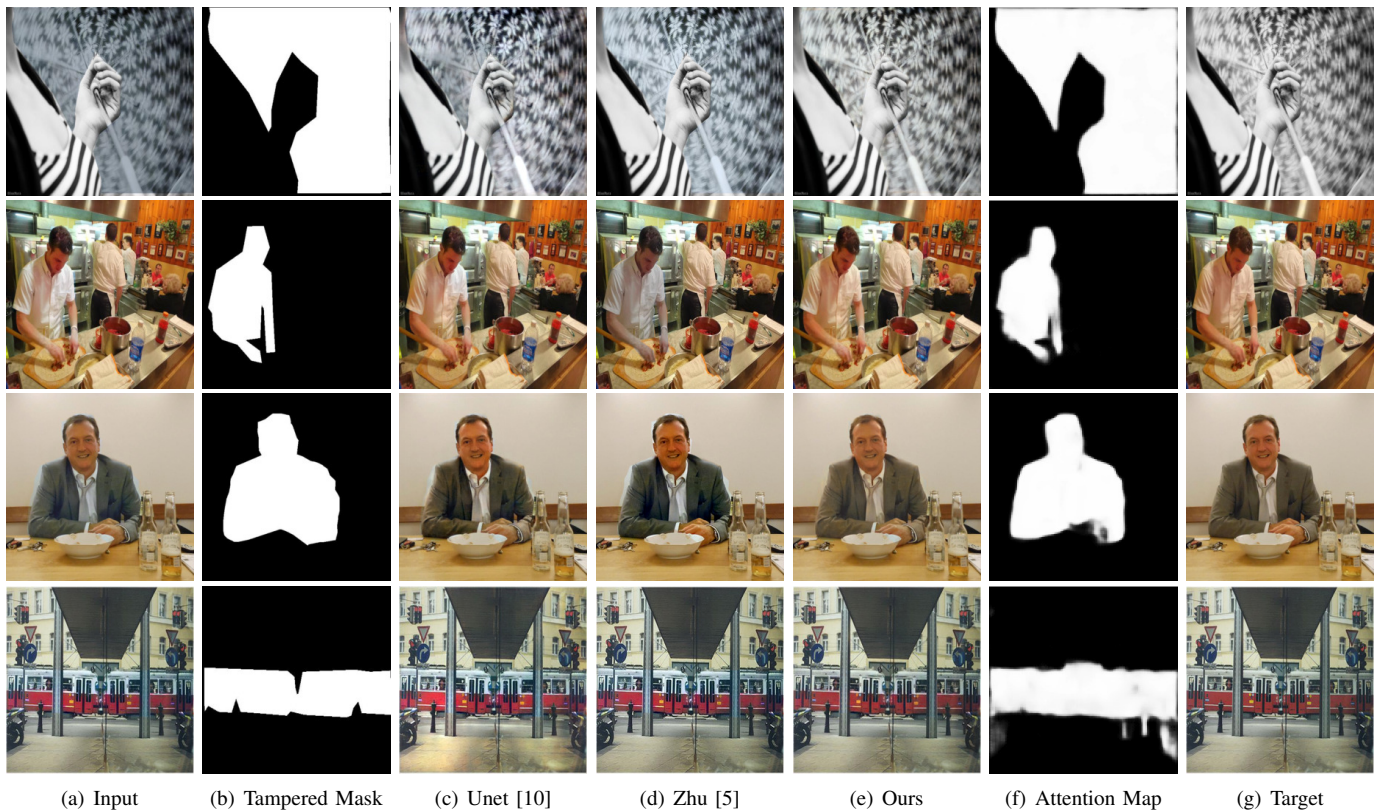| (a) Input | (b) Tampered Mask | (c) Unet [10] | (d) Zhu [5] | (e) Ours | (f) Attention Map | (g) Target |

Fig. 12: The Comparison results of our methods and other baseline methods on the S-COCO datasets without the ground truth mask as input. We also plot the attention map created by our method and the ground truth map here for comparison.

while we use the hard-coded mask and Gaussian Filter in the module.

| Method | MSE↓ | SSIM↑ | PSNR↑ |
|---|---|---|---|
| baseline unet [10] | 17.28 | 0.9757 | 33.55 |
| + More Channels | 16.93 | 0.9770 | 33.80 |
| + CONV-Block | 17.85 | 0.9757 | 33.34 |
| + SE-Block [6] | 17.45 | 0.9762 | 33.42 |
| + CBAM [15] | 16.50 | 0.9768 | 34.32 |
| unet with S$^2$ASC (Ours) | 15.59 | **0.9790** | **34.74** |
| − $G_{bg}, G_{mix}, G_{fg}$ | 17.93 | 0.9758 | 33.32 |
| − $G_{mix}$ | 15.80 | 0.9779 | 34.49 |
| − Gaussian Mask | 15.67 | 0.9788 | 34.72 |
| + 6 layers S$^2$ASC | 15.66 | 0.9787 | 34.73 |
| − Learning block($L$) | **15.49** | 0.9781 | 34.73 |

TABLE IV: The evaluation of S$^2$AM with the different attention modules and the ablation study of S$^2$AM on S-COCO dataset.

*2) Ablation Study:* We conduct ablation experiments on the inner structure of the proposed S$^2$AM model under ***S$^2$ASC***. All the experiments are performed on the S-COCO dataset with the same configuration. We illustrate all the numerical results in Table. IV.

**Without $G_{bg}, G_{mix}, G_{fg}$.** By comparison with the baseline, Unet w/o $G_{bg}, G_{mix}, G_{fg}$ gets slightly worse results because the hard-coded mask will force the learning-able block $L$ to learn the necessary and unnecessary features altogether without any identification.

**Without $G_{mix}$.** We remove the $G_{mix}$ in S$^2$ASC module for necessity. As shown in Table. IV, compared with S$^2$ASC, the method without $G_{mix}$ performs slightly worse because

$G_{mix}$ selects the original features in the spliced region which is unnecessary to transfer.

**Without Gaussian Filter.** A detailed explanation of Gaussian Filter has been introduced in the Section.III-A. Here, we report the results of experiments. As shown in Fig. 14, the Gaussian Filter shows better performance in the boundary of spliced edges. We also report the numerical results in Table. IV, it is obvious that the Gaussian Filter improves the performance of our module.

**With 6 layers S$^2$ASC.** The model gets a slightly worse performance when we replace all the skip-connection with S$^2$ASC in Unet (Unet + 6 layers S$^2$ASC in Table. IV). This experiment also fit our assumption: the style is much more related to low-level features, and the high-level features should be the same. However, in all 6 layers S$^2$ASC, our method still gains much better results than the baseline.

**Without learning-able Block.** We evaluate the importance of our learning-able block. As shown in Table.IV, the methods without learning-able block $L$ show worse results than our full method, especially in SSIM. This is because SSIM measures the image quality from the luminance, contrast, and structure while the small convolutions in the learning-able block will learn more detail.

*3) Attention Analysis:* We also analyze the detail responses of all three *Channel Attention Modules* in our S$^2$AM modules. We use ***S$^2$ASC*** in unet structure for evaluating for the features used in S$^2$ASC all comes from the encoder. As shown in Fig. 15, we visualize the weighting of $G_{fg}, G_{mix}, G_{bg}$ in the three layers of S$^2$ASC modules at the bottom of each

(a) Input      (b) Mask      (c) Target

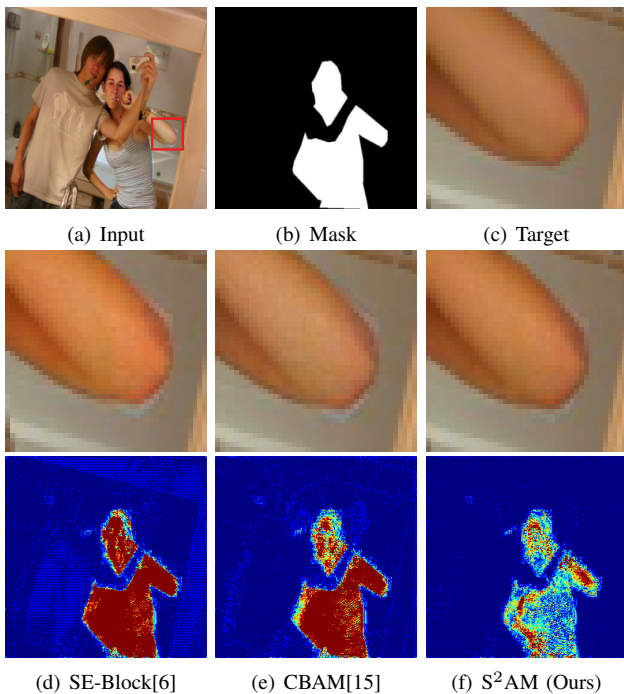(d) SE-Block[6]    (e) CBAM[15]    (f) S$^2$AM (Ours)

Fig. 13: Comparisons on different attention modules. We plot the results of Unet+SE-Block, Unet+CBAM and Unet+S$^2$AM (Ours) for comparison. For each method, we show the harmonized results of the red region and the absolute difference colormap between the result and target image. It is clear that our Unet + S$^2$AM method gains better results than others in the boundary of arms (best view with zoom in) and the absolute difference colormaps. Moreover, our method also shows a better prediction in the non-spliced region. The colormaps are enlarged $30\times$ for better visualize the differences in the non-spliced region.



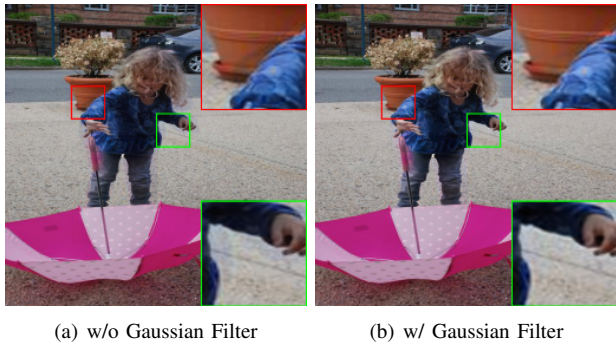(a) w/o Gaussian Filter      (b) w/ Gaussian Filter

Fig. 14: Gaussian Filter smooth the binary mask to fit object. For example, the little girl in the image is the tampered region with inaccuracy spliced mask. From two zoom-in regions in the image, the method with Gaussian Filter show better harmonization boundary.

sample. Particularly, the blacker block indicates the channel weighting is low while the whiter block means it is important for filtered features. So our first observation is that our S$^2$ASC module allocates and weights the features according to the input images by comparing the response maps between two images. Moreover, Another interesting observation is that there

are more white blocks in $G_{fg}$ than $G_{mix}$ and $G_{bg}$ in the $S^2ASC_1$ module while the opposite phenomenon is shown in $S^2ASC_3$. This fact perfectly explains our assumption: the coarser features in the spliced region need to learn specifically while the high-level features are almost the same.
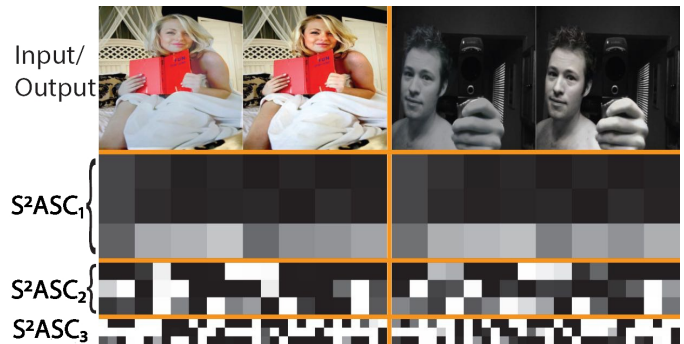


Fig. 15: From coarse-to-fine, we plot the first $\frac{1}{8}$ attention channels response on three different levels of S$^2$ASC in network structure named $S^2ASC_1$, $S^2ASC_2$ and $S^2ASC_3$, respectively. For each level of attention, we visualize the Channel Attention Module with the order $G_{bg}$, $G_{mix}$ and $G_{fg}$ from the top to the bottom.

## V. CONCLUSION

In this paper, we propose a novel method for image harmonization. Specifically, we design a new attention module called S$^2$AM to fit the image harmonization task with Unet backbone network. Additionally, we harmonize the image without mask by interpolating spatial attention module and the attention loss to S$^2$AM module. This idea comes from the original intention of the similarity in content and difference in the spliced region between the spliced image and harmonized target. The experiments show that our proposed method achieves better results than others both in quality and quantity.

Besides image harmonization, the proposed attention module can also be adapted to other computer vision tasks easily with regional differences, such as *Free-Form Image Inpainting* [18] and *Semantic Image Synthesis* [45]. We believe it is a promising direction for our future work.

## REFERENCES

[1] X. Cun and C.-M. Pun, "Image Splicing Localization via Semi-global Network and Fully Connected Conditional Random Fields." *ECCV Workshops*, 2018.
[2] J.-F. Lalonde and A. A. Efros, "Using Color Compatibility for Assessing Image Realism." *ICCV*, 2007.
[3] K. Sunkavalli, M. K. Johnson, W. Matusik, and H. Pfister, "Multi-scale image harmonization," *ACM Trans. Graph.*, vol. 29, no. 4, p. 1, Jul. 2010.

[4] Y.-H. Tsai, X. Shen, Z. L. 0001, K. Sunkavalli, X. Lu, and M.-H. Y. 0001, "Deep Image Harmonization." *CVPR*, 2017.

[5] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros, "Learning a Discriminative Model for the Perception of Realism in Composite Images." *ICCV*, 2015.

[6] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-Excitation Networks," *CVPR*, Sep. 2018.

[7] L. Bondi, S. Lameri, D. Guera, P. Bestagini, E. J. Delp, and S. Tubaro, "Tampering Detection and Localization Through Clustering of Camera-Based CNN Features," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2017, pp. 1855–1864.

[8] S. Ye, Q. Sun, and E.-C. Chang, "Detecting Digital Image Forgeries by Measuring Inconsistencies of Blocking Artifact." *ICME*, 2007.

[9] C.-M. Pun, B. Liu, and X. Yuan, "Multi-scale noise estimation for image splicing forgery detection." *J. Visual Communication and Image Representation*, vol. 38, pp. 195–206, 2016.

[10] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," *CVPR*, Nov. 2017.

[11] Y.-H. Tsai, X. Shen, Z. L. 0001, K. Sunkavalli, and M.-H. Y. 0001, "Sky is not the limit - semantic-aware sky replacement." *ACM Trans. Graph.*, vol. 35, no. 4, pp. 1–11, 2016.

[12] H. Wu, S. Zheng, J. Zhang, and K. Huang, "GP-GAN: Towards Realistic High-Resolution Image Blending," Mar. 2017.

[13] H. Huang, S. Xu, J. Cai, W. Liu, and S. Hu, "Temporally Coherent Video Harmonization Using Adversarial Networks." *CoRR*, 2018.

[14] D. Lee, T. Pfister, and M.-H. Yang, "Inserting Videos into Videos," *CVPR*, Mar. 2019.

[15] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional Block Attention Module," *ECCV*, 2018.

[16] X. Li, W. Wang, X. Hu, and J. Yang, "Selective Kernel Networks," *CVPR*, Mar. 2019.

[17] Liu, Guilin, Reda, Fitsum A, Shih, Kevin J, Wang, Ting-Chun, Tao, Andrew, and Catanzaro, Bryan, "Image Inpainting for Irregular Holes Using Partial Convolutions," *ECCV*, Apr. 2018.

[18] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-Form Image Inpainting with Gated Convolution," *arXiv.org*, Jun. 2018.

[19] J. Yu, Z. L. 0001, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative Image Inpainting With Contextual Attention." *CVPR*, pp. 5505–5514, 2018.

[20] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local Neural Networks," *CVPR*, 2018.

[21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is All you Need." *NIPS*, 2017.

[22] A. Buades, B. Coll, and J.-M. Morel, "A Non-Local Algorithm for Image Denoising." vol. 2, pp. 60–65, 2005.

[23] T. Chen, X. Zhai, M. Ritter, M. Lucic, and N. Houlsby, "Self-Supervised Generative Adversarial Networks," *arXiv.org*, Nov. 2018.

[24] Y. Cao, J. Xu, S. Lin, F. Wei, and H. Hu, "GCNet: Non-local Networks Meet Squeeze-Excitation Networks and Beyond," *arXiv.org*, Apr. 2019.

[25] Q. Chen, J. Xu, and V. Koltun, "Fast Image Processing with Fully-Convolutional Networks," *ICCV*, Sep. 2017.

[26] F. Yu and V. Koltun, "Multi-Scale Context Aggregation by Dilated Convolutions." *ICLR*, 2016.

[27] He, Kaiming, Sun, Jian, and Tang, Xiaoou, "Guided Image Filtering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 6, pp. 1397–1409, 2013.

[28] M. Gharbi, J. Chen, J. T. Barron, S. W. Hasinoff, and F. Durand, "Deep bilateral learning for real-time image enhancement." *ACM Trans. Graph.*, vol. 36, no. 4, pp. 1–12, 2017.

[29] J. Chen, A. Adams, N. Wadhwa, and S. W. Hasinoff, "Bilateral guided upsampling," *ACM Trans. Graph.*, vol. 35, no. 6, pp. 1–8, Nov. 2016.

[30] H. Wu, S. Zheng, J. Zhang, and K. Huang, "Fast End-to-End Trainable Guided Filter," *CVPR*, Mar. 2018.

[31] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Networks," *NIPS*, Jun. 2014.

[32] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks," *ICCV*, Mar. 2017.

[33] X. Chen, C. Xu, X. Yang, and D. Tao, "Attention-GAN for Object Transfiguration in Wild Images," *NeurIPS*, 2018.

[34] Y. A. Mejjati, C. Richardt, J. Tompkin, D. Cosker, and K. I. Kim, "Unsupervised Attention-guided Image-to-Image Translation." *NeurIPS*, 2018.

[35] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Deep image prior," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9446–9454.

[36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *CVPR*, Dec. 2015.

[37] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going Deeper with Convolutions," *arXiv.org*, Sep. 2014.

[38] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," May 2015.

[39] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution," *ECCV*, Mar. 2016.

[40] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft COCO: Common Objects in Context," *arXiv.org*, May 2014.

[41] V. Bychkovsky, S. Paris, E. Chan, and F. Durand, "Learning photographic global tonal adjustment with a database of input / output image pairs." 2011.

[42] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking Atrous Convolution for Semantic Image Segmentation," *arXiv.org*, Jun. 2017.

[43] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.

[44] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *ICLR*, Dec. 2014.

[45] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic Image Synthesis with Spatially-Adaptive Normalization," *CVPR*, Mar. 2019.