# RSFNet: A White-Box Image Retouching Approach using Region-Specific Color Filters

Wenqi Ouyang[1], Yi Dong[2], Xiaoyang Kang[1], Peiran Ren[1], Xin Xu[1], Xuansong Xie[1]

[1]DAMO Academy, Alibaba Group, [2]Nanyang Technological University

wenqi.oywq@alibaba-inc.com, ydong004@ntu.edu.sg,

{peiran.rpr, kangxiaoyang.kxy, chris.xx, xingtong.xxs}@alibaba-inc.com

## Abstract

*Retouching images is an essential aspect of enhancing the visual appeal of photos. Although users often share common aesthetic preferences, their retouching methods may vary based on their individual preferences. Therefore, there is a need for white-box approaches that produce satisfying results and enable users to conveniently edit their images simultaneously. Recent white-box retouching methods rely on cascaded global filters that provide image-level filter arguments but cannot perform fine-grained retouching. In contrast, colorists typically employ a divide-and-conquer approach, performing a series of region-specific fine-grained enhancements when using traditional tools like Davinci Resolve. We draw on this insight to develop a white-box framework for photo retouching using parallel region-specific filters, called RSFNet. Our model generates filter arguments (e.g., saturation, contrast, hue) and attention maps of regions for each filter simultaneously. Instead of cascading filters, RSFNet employs linear summations of filters, allowing for a more diverse range of filter classes that can be trained more easily. Our experiments demonstrate that RSFNet achieves state-of-the-art results, offering satisfying aesthetic appeal and increased user convenience for editable white-box retouching. Code is available at https://github.com/Vicky0522/RSFNet.*

## 1. Introduction

Photos and videos recorded by the camera usually lack aesthetic quality due to poor shooting condition and inexperienced photographer. Artists often use professional-grade softwares (*e.g*., PhotoShop for image, Davinci Resolve for video) to enhance image and video quality. However, it requires professional retouching skills to conduct a series of sophisticated manual adjustments. The use of fool-proof applications that present various style templates simplifies the retouching procedure, but it is unable to achieve opti-
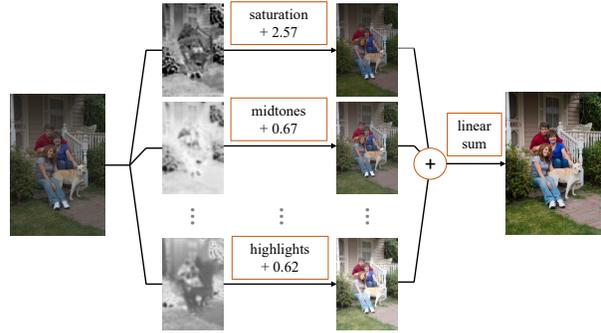


Figure 1: Architecture of our white-box retouching framework. Our model generates filter arguments (e.g., saturation, midtones, highlights) and attention maps of regions for corrsponding filter simultaneously. Final result is obtained by conducting linear summations on filtered results.

mal results due to the lack of enhancement capability.

Recent learning-based methods have demonstrated the strong capability of deep neural nets for automatic photo retouching. Automatic systems are established to generate the optimal result end-to-end. However, one of the most significant considerations is that retouching is not a problem with an exclusive solution. People have different retouching preferences. Even the same artist may retouch the same image in different styles to meet various demands. In order to provide convenience for manual edits, automatic retouching systems must provide not only the suggested result, but also the retouching strategy in a way understandable by humans.

Taking these considerations into account, we propose a white-box framework that uses the divide-and-conquer strategy employed by artists in traditional retouching tools. Our model generates attention maps of regions as well as filter arguments for traditional edits for these regions. This allows users to alter the suggested results to their preferences.

The main **contributions** of this work are as follows:

• We redefine the retouching problem according to the

divide-and-conquer strategy, focusing on finding attention maps for regions and human-understandable filter adjustments per map to achieve the best results.

- We propose RSFNet, which generates pixel-level attention maps of regions and filter arguments simultaneously. By conducting linear summations on the filtered results, our model demonstrates superior performance to global white-box retouching methods, with a wider range of filter classes for fine-grained enhancement and a simpler training procedure.
- We propose a scheme that allows users to edit RSFNet suggested results, demonstrating its effectiveness and convenience in image applications.

## 2. Related Works

**Deep Image Enhancement.** Many attempts have been made in image enhancement using learning-based methods. These methods can broadly be classified into two categories: deep-convolutional-based models and physical-inspired models. The former [3, 4, 6, 27, 32] view enhancement as a synthesis problem and use fully convolutional generators [15] to achieve dense image-to-image translation. While these methods exhibit strong capability for generating and enhancing details, they suffer from heavy structure and low inference speed. The latter category of works defines enhancement as the parameter prediction of a physical model and uses deep learning strategies to fit the model. These models include 3D-LUT [23, 31, 30], parametric filters [17], conditional sequential global modulation [8], affine color transformations [7, 22], 1D mapping curves [20, 12, 18]. Their enhancement capabilities depend on the range of transformation functions the physical model covers. Among those, models based on global 3D-LUT [30, 31] are well-designed for retouching and could achieve high performance with fast inference speed. However, they suffer from uneven transitions in smooth areas due to lack of fine-grained local adjustment. Spatial-Aware 3D-LUT [23] alleviate the problem of global 3D-LUT by computing pixel-level weight maps. All these methods have a black-box structure or unintuitive parameters that are hard for humans to understand. DeepLPF [17] has been a significant milestone in the field of region-specific color enhancement via local parametric filters but suffers from low inference speed and limited filter shapes. The three types of filters with three implementations each results in complex maps that are hard to understand.

**White-box Image Editing.** Recent white-box methods [10, 11, 29, 33] have decoupled image editing into a series of human understandable operations and deep learning strategies to predict them. Of these, our work is most related to [11, 10]. However, they train networks with cascaded filter argument prediction modules to perform global retouching step-by-step. Our model utilizes pixel-wise region-specific filters to capture local features and employs linear summations to combine filters, thus having stronger capability to cover a wider range of color transformation functions.

## 3. Method

### 3.1. Design Motivation

**Traditional Retouching Strategy.** In the traditional retouching process, local adjustments for different regions are conducted separately and then aggregated together to accomplish fine-grained enhancement. All the adjustments could be accomplished in one *Layer*, where filters are all conducted on the same image instead of results from previous filters. We adopt this divide-and-conquer strategy to build our framework. We select 10 commonly used retouching filters from traditional tools(*e.g.*, Davinci Resolve) to represent adjustment manipulations, including *contrast, saturation, hue, temperature, shadows, midtones, highlights* and *shift*. For more details about filters, please refer to Appendix A . Therefore, retouching is defined as finding pixel-level attention maps and corresponding adjustments for each map to achieve the most desirable result. We adopt this convention and represent retouching result $Y$ as adding linear summations of increments resulted by filters to the original image as:

$$Y = X + \sum_m \sum_n (F_{m,n}(\theta_{m,n}, X) - X) \odot M_m \quad (1)$$

Where $\theta_{m,n}$ and $F_{m,n} \in \{contrast, ..., shift\}$ are the argument value and filter function of the $n$th filter for the $m$th attention map of input image $X$ respectively, $M_m$ is the $m$th attention map for image $X$. It should be noticed that this linear representation makes our optimization process much simpler than other white-box methods.

### 3.2. Architecture

Our model consists of two modules, the parameter predictor and the renderer. The architecture of the parameter predictor shares a similar overall structure with previous segmentation models [24, 25]. As illustrated in Figure 2, the network $g$ is composed of a backbone $g_{backbone}$, a FPN neck $g_{neck}$, a map generator $h_{map}$ for attention map prediction and an argument regressor $h_{arg}$ for arguments regression. We set the last layer of map generator $h_{map}$ to sigmoid, followed by an upsampling layer to match the size of the original input image. We use $f(X)$ to represent the output features of FPN neck. Maps generated by $h_{map}$ after upsampling are:

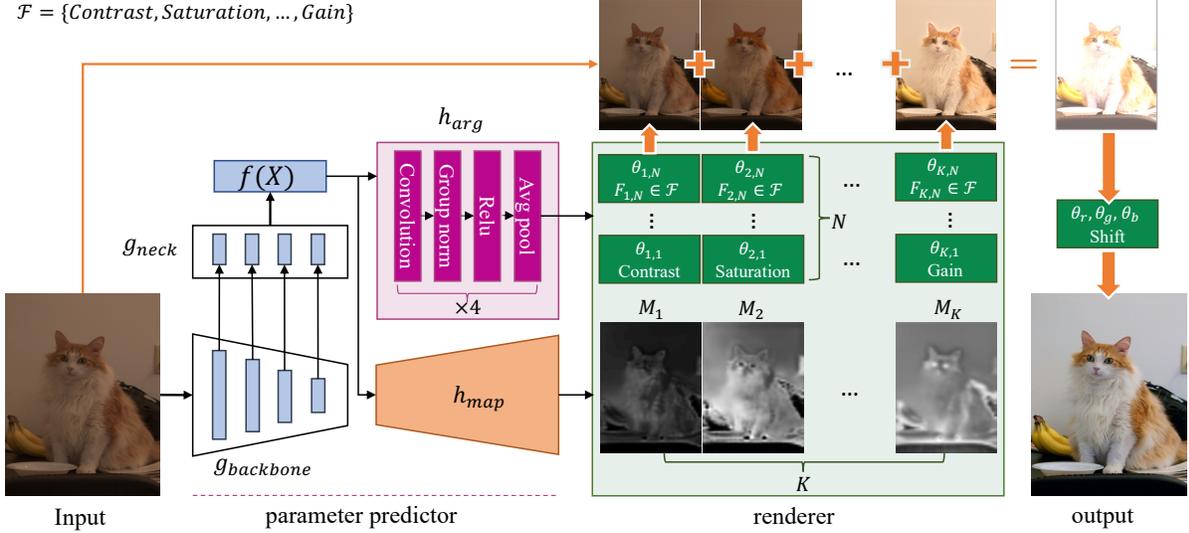$$M_m = h_{map}(f(X))_m, m \in \{1, ...., K\} \quad (2)$$

Figure 2: Our network comprises two modules: the parameter predictor and renderer. The input image is transformed into feature embeddings $f(X)$ through a backbone $g_{backbone}$ and FPN neck $g_{neck}$, which are then utilized by the map generator $h_{map}$ and argument regressor $h_{arg}$ to generate region maps and filter arguments (e.g. contrast, hue, saturation), respectively. The renderer applies the filter arguments to each region map and aggregates them to produce the final retouching result.

$K$ is the number of output channels of the last layer. The argument regressor $h_{arg}$ consists of 4 units of convolutional modules. Each module consists of a convolution layer followed by group normalization, $relu$ activation and average pooling. Output feature $f(X)$ is finally reduced by $h_{arg}$ to a $K \times N$ vector representing arguments of all filters for every map:

$$\theta_{m,n} = h_{arg}(f(X))_{m,n}, m \in \{1, ...., K\}, n \in \{1, ..., N\} \tag{3}$$

$K$ is the number of color filters per map. Given the maps and arguments from the parameter predictor, the renderer applies filter functions $F_{m,n} \in \{contrast, ..., shift\}$ to the input image and merges filtered results via linear summations of increments. The final output image $Y$ in Equation 1 is:

$$Y = X + \sum_{m=1}^{K} \sum_{n=1}^{N} (F_{m,n}(\theta_{m,n}, X) - X) \odot M_m \tag{4}$$

In our implementation, we set the number of filters per map as $N = 1$, indicating that there is one specific filter assigned to each map. This is illustrated in the second row of filter boxes depicted in Figure 2. We also consider the filter *Shift* as a global function that is applied to the entire image without an attention map. Through our experimentation, we observe that this configuration is sufficient for generating satisfactory results while also providing convenience for users to edit.

**Training Loss.** The overall framework can be trained end-to-end. Our training loss function is defined as follows:
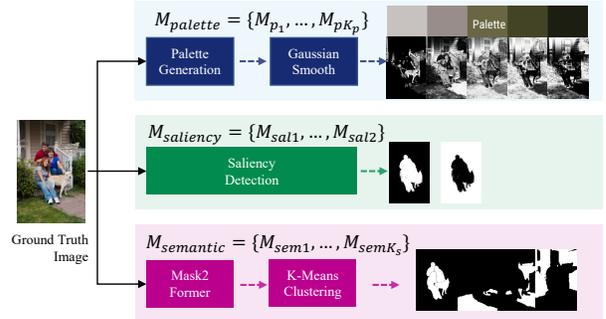


Figure 3: A variation of RSFNet could be used to generate masks with controlled shapes. We generate ground truth masks for training via three types of region masks derived from off-the-shelf models: palette-based segmentation, semantic segmentation, and saliency detection.

$$L = L_{recon}, \tag{5}$$

where $L_{recon}$ is the $l1$ loss for reconstruction.

### 3.3. A Variation of RSFNet with Controlled Region Shape

Our proposed framework has the capability to generate region masks with controlled shapes, with minimal modifications to the training data and loss function. In our observation, artists tend to select regions based on their main color, saliency, or other high-level semantic informa-

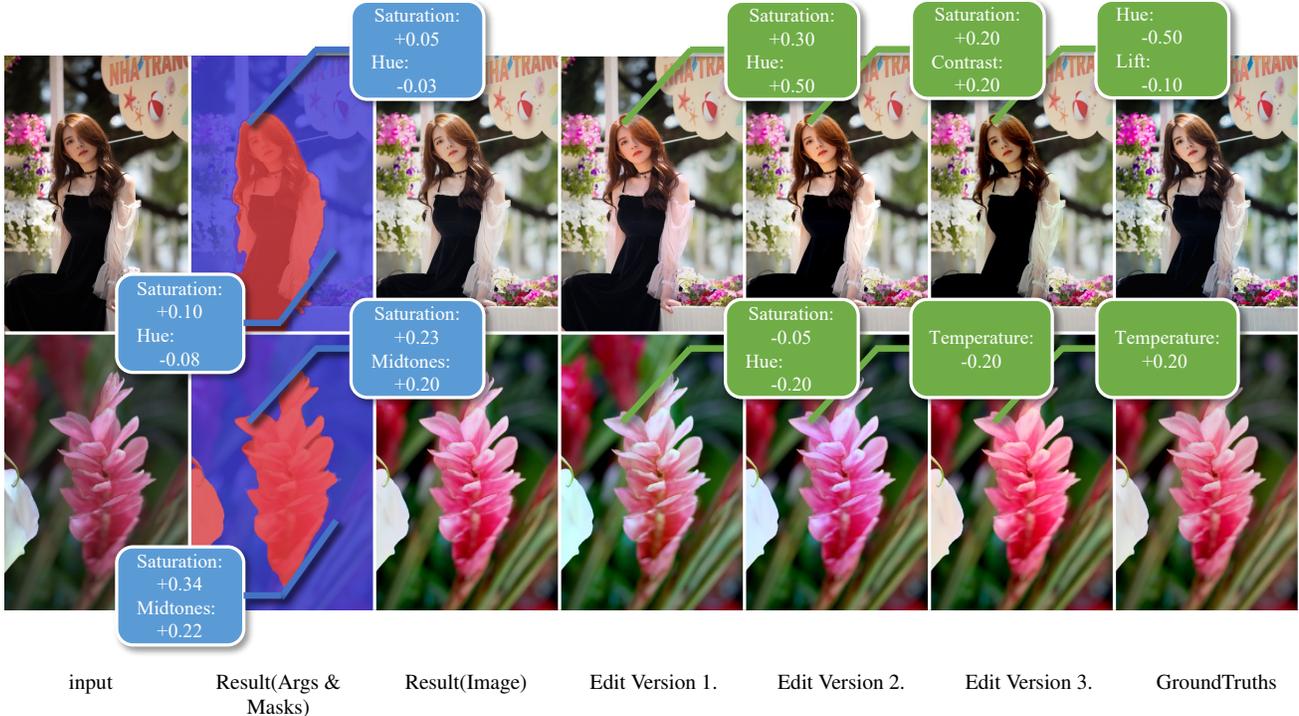| input | Result(Args & Masks) | Result(Image) | Edit Version 1. | Edit Version 2. | Edit Version 3. | GroundTruths |

Figure 4: Editable white-box retouching. The second column illustrates the arguments and masks generated by RSFNet-saliency that is trained with saliency masks. The third column displays the retouched result. On the right-hand side, the figure showcases three versions of adjustments conducted on the retouched results, with numbers denoted in green boxes indicating the relative incremental value of filters.

tion. To leverage this insight, we utilize palette-based methods [2], saliency detection [14], and panoptic segmentation networks [5] to generate palette-based, saliency, and semantic masks of the input images used for training, as depicted in Figure 3. Further implementation details can be found in Appendix B. Our training loss function in Equation 5 is modified as follows:

$$L = L_{recon} + \lambda L_{mask}, \qquad (6)$$

where $L_{recon}$ is the $l1$ loss for reconstruction. When training with palette-based masks and saliency masks, $L_{mask}$ is the Dice Loss for mask prediction following loss functions in [24, 25]. We change it to the Cross Entropy Loss when training with semantic masks, at the time map generator $h_{map}$ ends with a softmax layer.

The model generates region masks and corresponding filter arguments for the input image simultaneously. Users can select region maps and adjust filter arguments to edit results. We show an example of *RSFNet-saliency* trained with saliency masks in Figure 4.

## 4. Experiments

### 4.1. Datasets and Application Settings

We conduct experiments on two publicly available datasets: MIT-Adobe FiveK [1] and PPR10K [13]. The MIT-Adobe FiveK consists of 5,000 RAW images with their retouched versions. We follow prior works [8, 12, 31, 30] to adopt images retouched by expert C as ground truths and split the dataset into 4,500 pairs for training and 500 pairs for validation. Images are resized to 480p during training stage, whereas both of 480p resolution and original resolution are used during validation. We follow the official split [13] to split PPR10K dataset into 8,875 pairs for training and 2,286 pairs for testing. Images are resized to 360p resolution for training and validation.

For the MIT-Adobe FiveK dataset, we conduct experiments on two input settings: *input zeroed with expertC white balance* and *input zeroed as shot*. The second set of inputs is more challenging as it also deteriorates in white balance. Therefore, to restore its visual appeal, our framework's ability to perform color temperature correction and hue adjustment is required.

| mask as input | mask type | PSNR↑ | SSIM↑ |
|---|---|---|---|
| w/ | saliency(1) | 23.83 | 0.906 |
| w/ | palette(5) | 23.45 | 0.903 |
| w/ | semantic(10) | 21.65 | 0.866 |
| w/ | semantic(5) | 21.77 | 0.867 |
| w/o | saliency(1) | 24.01 | **0.909** |
| w/o | palette(5) | 24.14 | 0.906 |
| w/o | semantic(10) | 24.09 | 0.908 |
| w/o | semantic(5) | **24.17** | 0.909 |
| w/o | semantic(133) | 23.75 | 0.905 |

(a) Analysis on using mask as input and different mask types. Bracket after "mask type" indicates number of masks.

| mask first | masktype | PSNR↑ | SSIM↑ |
|---|---|---|---|
| w/ | saliency(1) | 23.26 | 0.890 |
| w/o | saliency(1) | **24.01** | **0.909** |

(b) Analysis of training order. "Mask first" refers to building model on pre-trained Pool-Net [14].

| backbone | downscaling factor | PSNR↑ | Runtime↓ |
|---|---|---|---|
| resnet18 | x8 | 24.65 | 9.98 |
| resnet18 | x4 | **24.83** | 12.38 |
| resnet10 | x8 | 24.50 | **7.98** |
| resnet10 | x4 | 24.38 | 8.71 |

(c) Trade off between speed and accuracy.

Table 1: Ablation analysis of our framework. Models are trained on the FiveK dataset [1] with inputs zeroed as shot.
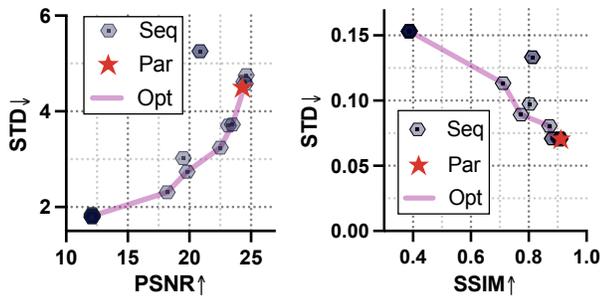


Figure 5: Comparison between sequential and parallel methods under the same filter functions and network structures.

## 4.2. Implementation Details

To implement our proposed framework as described in Section 3.2, we employ ResNet18 [9] as the backbone and utilize the output features of layer $(0, 1, 2, 3)$ as concatenated input with $256 \times 4$ channels for the FPN neck. The FPN neck reduces the features to $64 \times 4$ channels and feeds them into the map generator and argument regressor. We set the output numbers $K$ of region maps to 10 and 16 for the Adobe-MIT FiveK and PPR10K datasets, respectively. The number of filters per map $N$ is set to $N = 1$, with 10 filters used for MIT-Adobe FiveK, 1 each for *shadow, midtones, highlights*, where arguments for RGB channels are set equal. For PPR10K, we use 16 filters, with 3 each for *shadow, midtones, highlights*, where arguments for RGB channels are set differently.

To implement the variation of RSFNet with controlled region shapes, as described in Section 3.3, we generate ground truth masks for the MIT-Adobe FiveK dataset. Specifically, we generate five palette-based masks, one saliency mask, and five semantic masks for each image. For the PPR10K dataset, we directly use the human-region masks provided in the dataset as ground truth masks for

training. The map generator is designed to output 5 channels for the palette-based version, 2 channels for the saliency version, and 5 channels for the semantic version. We use 10 filters for the MIT-Adobe FiveK and 16 filters for the PPR10K, as mentioned above.

We use the standard Adam optimizer [19] to minimize the loss function in Equation 5 and 6. The mini-batch size is 16. We train models for 1,000,000 iterations with initial learning rate $l_r = 2 \times 10^{-4}$. We use a cosine strategy to gradually decay the learning rate every 250,000 iterations. All experiments are conducted on NVIDIA Tesla V100 GPU.

In the following sections, we denote our main model in Section 3.2 as *RSFNet-map*. Variations depicted in Section 3.3 are denoted as *RSFNet-palette, RSFNet-saliency* and *RSFNet-semantic*.

## 4.3. Ablation Studies

In this section, we evaluate the ability of our frameworks in different settings.

**Comparative study of parallel and sequential approaches.** We conduct experiments to compare parallel and sequential approaches under the same set of filter functions and backbone structures. For sequential implementations, we modify the network structure to mimic artistic workflows and resemble Harmonizer [11]. We train the network using randomly shuffled 20 unique sequences of filter functions. The results are shown in Figure 5. The parallel approach shows Pareto-optimality with near-maximum PSNR and a lower PSNR standard deviation (std). It also outperforms all sequential approaches in terms of SSIM and SSIM std. We find that altering the sequence of applying filters can lead to considerable changes in performance. Nearly half of these sequences fail to produce satisfactory results. Moreover, during training, the parallel approach converges faster than its sequential counterparts. In addition, the parallel method can be faster than the sequential one.

| Method | 480p zeroed with C's | | | 480p zeroed as shot | | | | Full Resolution zeroed as shot | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | $\triangle E_{ab}$ ↓ | PSNR↑ | SSIM↑ | $\triangle E_{ab}$ ↓ | Runtime↓ | PSNR↑ | SSIM↑ | Runtime↓ |
| UPE [22] | 21.88* | 0.853* | 10.80* | / | / | / | / | / | / | / |
| DPE [4] | 23.75* | 0.908* | 9.34* | / | / | / | / | / | / | / |
| HDRNet [7] | 24.66* | 0.915* | 8.06* | / | / | / | / | / | / | / |
| DeepLPF [17] | 24.73* | 0.916* | 7.99* | 23.38 | 0.880 | 10.03 | 44.25 | 23.40 | 0.863 | 1133.90 |
| CSRNet [8] | 25.17* | 0.924* | 7.75* | 24.24 | 0.910 | 9.70 | 3.49 | 23.04 | 0.874 | 80.6 |
| SA-3DLUT [23] | 25.50* | / | / | / | / | / | / | / | / | / |
| 3D-LUT [31] | 25.29* | 0.923* | 7.55* | / | / | / | / | / | / | / |
| 3D-LUT+AdaInt [30] | 25.49* | 0.926* | 7.47* | 24.50 | 0.912 | 9.22 | 1.59 | 24.24 | 0.857 | 1.80 |
| Harmonizer [11] | 24.11 | 0.904 | 8.23 | 23.23 | 0.893 | 10.14 | 16.74 | 22.57 | 0.870 | 27.98 |
| RSFNet-map | 25.49 | 0.924 | 7.23 | 24.64 | 0.915 | 9.16 | 9.98 | 24.39 | 0.894 | 12.35 |
| RSFNet-palette | 25.01 | 0.914 | 7.62 | 24.22 | 0.911 | 9.52 | 19.76 | 23.88 | 0.888 | 68.57 |
| RSFNet-saliency | 24.78 | 0.916 | 7.86 | 24.20 | 0.912 | 9.42 | 14.09 | 24.00 | 0.890 | 59.42 |
| RSFNet-semantic | 24.76 | 0.915 | 7.77 | 24.19 | 0.912 | 9.45 | 24.07 | 23.89 | 0.891 | 71.18 |
| RSFNet-global | 24.31 | 0.911 | 8.21 | 23.43 | 0.904 | 10.16 | 9.42 | 23.26 | 0.885 | 10.99 |

Table 2: Quantitative comparisons for retouching tasks on FiveK dataset [1]. Runtime is measured in miliseconds. "*" means the result is adopted from the paper [30]. "/" means the result is not available. Runtime shown in previous works are not adopted due to different hardware settings. Results with the first and second performance are colored as red and blue respectively. White-box methods are colored as violet. Full resolution images results are evaluated using Python, while others are evaluated on Matlab.

| Method | 360p expert a | | 360p expert b | | 360p expert c | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| DeepLPF [17] | 23.47 | 0.892 | 22.77 | 0.875 | 23.73 | 0.896 |
| CSRNet [8] | 24.01 | 0.936 | 23.91 | 0.938 | 24.31 | 0.931 |
| 3D-LUT+AdaInt [30] | 25.98 | 0.947 | 24.91 | 0.936 | 25.48 | 0.919 |
| Harmonizer [11] | 24.76 | 0.929 | 22.79 | 0.885 | 24.56 | 0.902 |
| RSFNet-map | 25.58 | 0.949 | 24.81 | 0.945 | 25.52 | 0.939 |
| RSFNet-saliency | 25.53 | 0.946 | 24.72 | 0.944 | 25.11 | 0.939 |

Table 3: Quantitative comparisons for retouching tasks on PPR10K dataset [13]. All the models are trained on data without augmentations. White-box methods are colored as violet.

**Training map generation and arguments regression separately.** To train models with controlled region shapes, we experiment with using pre-trained PoolNet [14] and adding an FPN neck followed by an argument regressor to the backbone instead of training from scratch. During training, we fix the weights of the backbone to keep the saliency mask generated by the network unchanged. The quantitative results are shown in Table 1b, which indicates that there is no superiority over training from scratch. The reason for this might be that a network trained only on semantic data lacks sufficient information for retouching. Therefore, it is better to train the map generation and arguments regression tasks simultaneously.

**Using Masks as Input.** For training models with controlled region shapes, we can use masks concatenated with the image as input for our model. We train three models using three sets of masks as input, respectively. The results are shown in Table 1a. The retouching results evaluated on PSNR and SSIM [26] are lower than in other settings. This may be because the ground truth masks are generated by off-the-shelf models, which may differ from the underlying real masks of the dataset. Masks are slightly modified by the map generator in the inference stage to achieve better retouching results. Therefore, simple concatenation of masks with the input image can lead to worse performance compared to our design of the variant of RSFNet, which generates masks and filter arguments simultaneously.
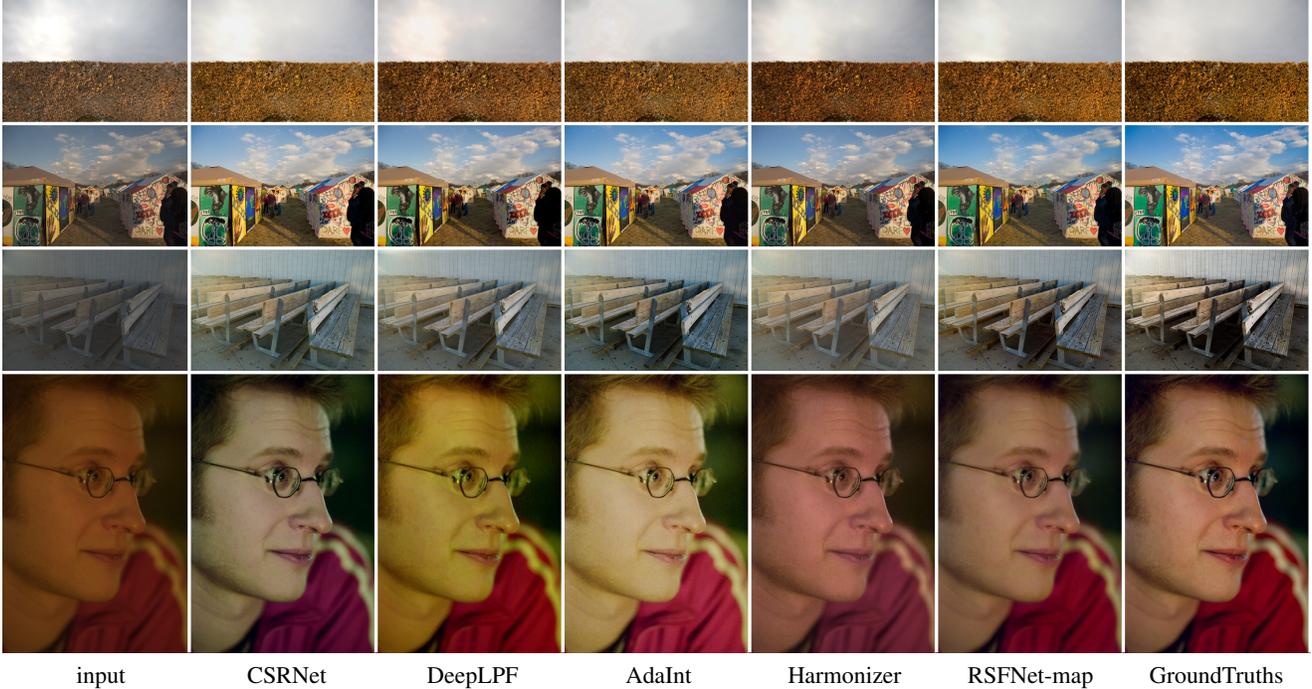
| input | CSRNet | DeepLPF | AdaInt | Harmonizer | RSFNet-map | GroundTruths |

Figure 6: Qualitative comparison of selected methods. "RSFNet-map" means our main model explained in Section 3.2
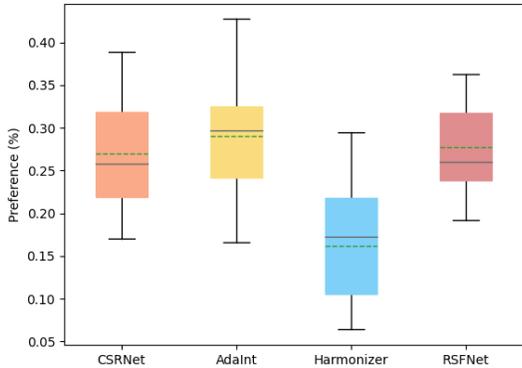


Figure 7: Boxplot of user study. The dashed green line and the solid black line inside the box are the mean and the median preference percentage respectively.

**Trade Off Between Speed and Accuracy.** As we decrease the number of downsample operations in the backbone, the feature size grows larger and time costs increase. However, quantitative results of retouching also become better, as shown in 1c. Speed is measured in miliseconds.

## 4.4. Comparison with State-of-the-Art

We compare our methods with the state-of-the-art for black-box and white-box image retouching tasks. The selected methods are compared on PSNR, SSIM [26], the $L_2$-distance in CIE LAB color space ($\triangle E_{ab}$) and the inference speed. We follow the practice in [30] to measure the GPU inference time on 100 images and report the average. Quantitative results are shown in Table 2 and 3. *RSFNet-map* refers to our main model in Section 3.2, *RSFNet-saliency*, *RSFNet-palette* and *RSFNet-semantic* refers to models trained with three sets of masks respectively in Section 3.3. In addition, we evaluate a global retouching model without the map generator, denoted as *RSFNet-global*. Other models are trained using their official public codes and default configurations. All experiments are executed on an NVIDIA Tesla V100 GPU.

Our proposed method, *RSFNet-map*, demonstrates its effectiveness in temperature correction and color enhancement, especially for the second set of inputs under more deteriorating shooting conditions. Although *RSFNet-map* achieves state-of-the-art results with negligible increase, while other RSFNet variations lag behind the state of the art, it is important to note that our white-box framework provides human-understandable ways for retouching, which makes it more convenient for users to edit and assess retouching results compared to other black-box methods. When compared with state-of-the-art white-box methods

Figure 8: Editable image retouching. "RSFNet-saliency", "RSFNet-semantic" means RSFNet trained with saliency masks and semantic masks respectively.

such as Harmonizer [11], all of our models exhibit better performance with faster speed. Furthermore, using 3D-LUTs to encode filter functions could potentially accelerate our models, a technique commonly implemented in traditional retouching tools. For PPR10K [13], we also evaluate our methods using a random split setup. The results are provided in Appendix C, consistently reinforcing the conclusions drawn from previous experiments.

The qualitative results presented in Figure 6 demonstrate that our proposed method, *RSFNet-map*, exhibits superior performance in handling color transitions across regions, particularly in highlight areas, compared to 3D-LUT based methods such as AdaInt, as shown in the first row of the figure. Furthermore, we conduct a comparison with Harmonizer [11] on editable retouching, as shown in Figure 8. Our white-box framework offers more degrees of freedom than global-retouching manipulations to achieve region-specific retouching. For instance, when the temperature is adjusted globally using Harmonizer, the sky in the first image turns yellow. On the other hand, our *RSFNet-saliency* model can modify the temperature of the foreground girl while leaving other regions unaffected. In the second image, global temperature adjustment by Harmonizer turns the lake and mountains blue. In contrast, *RSFNet-semantic* only modifies the sky temperature, leaving the mountains and lake unchanged.

To further validate the effectiveness of our proposed framework, we conduct a user study to evaluate human preferences for RSFNet and other state-of-the-art methods, including CSRNet [8], AdaInt [30], and Hormonizer [11]. To form the test set, we randomly select 58 images from the validation set of Adobe-MIT FiveK using the *input zeroed with expertC white balance* and *input zeroed as shot* settings, as well as an additional 20 images downloaded from the internet. During the experiment, we display retouched images produced by all methods, including the original input, to 22 participants. Participants are asked to select the best result from a randomly shuffled set of retouched images generated by all methods. We calculate the preferences of participants for each method and plot the results using a box plot, as shown in Figure 7.

Among the compared methods, AdaInt [30] achieves

higher mean and median percentages. RSFNet exhibits more stable performance with a higher bottom percentage and the smallest standard variance percentage than the other methods. In comparison to the white-box method Harmonizer [11], RSFNet demonstrates superior user preference percentages. The boxplot illustrates the varying preferences of different users, highlighting the significance of editable capability. As RSFNet utilizes traditional color filters for retouching, it can be readily integrated into traditional retouching tools. Users can further enhance the visual appeal of the retouched results according to their own preferences through an interface.

**Video Retouching.** Our proposed RSFNet model also demonstrates applicability to video retouching, as shown in Figure 8. The filter arguments remain constant within a single video clip for all frames, ensuring consistency throughout the retouching process. However, for RSFNet variations with controlled region shapes, such as *RSFNet-saliency*, the retouching consistency across frames relies on the consistency of region masks across frames. Although filter arguments can alleviate the inconsistency caused by masks, severe inconsistency problems in masks may still affect retouching results. With the assistance of a more robust tracking algorithm, these models can achieve better performance in editable video retouching, highlighting the potential of RSFNet in various retouching applications.

## 5. Limitation and Conclusion

We develop our framework under the assumption that all adjustments are done in one *layer*, which is then combined via linear summations of increments of all filters to obtain the final result. Therefore, although our model provides a white-box retouching framework that aligns with the intuition of human artists, it is unable to cover all adjustments that artists could conduct using professional software. However, we believe that if the cache data in the retouching process of artists is available, such as the region masks used for region-specific retouching, our model equipped with a more complicated structure, such as a cascaded structure in [10, 11], could behave more similarly to a real human artist by learning from this data.

This paper introduces RSFNet, a white-box framework for image retouching that utilizes a divide-and-conquer strategy to generate region maps and human-understandable filter arguments. The resulting filtered images are combined using linear summations, allowing for a wider range of filter classes and achieving fine-grained enhancement with superior performance. A variation of RSFNet with controlled shape region masks is also proposed for user convenience. Extensive experiments are presented to demonstrate the effectiveness of RSFNet.

## References

[1] Vladimir Bychkovsky, Sylvain Paris, Eric Chan, and Frédo Durand. Learning photographic global tonal adjustment with a database of input / output image pairs. In *CVPR 2011*, pages 97–104, 2011.

[2] Huiwen Chang, Ohad Fried, Yiming Liu, Stephen DiVerdi, and Adam Finkelstein. Palette-based photo recoloring. *SIGGRAPH*, 34(4), July 2015.

[3] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, pages 3291–3300, 2018.

[4] Yu Sheng Chen, Yu Ching Wang, Man Hsin Kao, and YungYu Chuang. Deep photo enhancer: Unpaired learning for image enhancement from photographs with gans. In *CVPR*, pages 6306–6314, 2018.

[5] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, pages 1290–1299, 2021.

[6] Yubin Deng, Chen Change Loy, and Xiaoou Tang. Aestheticdriven image enhancement by adversarial learning. In Susanne Boll, Kyoung Mu Lee, Jiebo Luo, Wenwu Zhu, Hyeran Byun, Chang Wen Chen, Rainer Lienhart, and Tao Mei, editors, *ACM MM*, pages 870–878. ACM, 2018.

[7] Michaël Gharbi, Jiawen Chen, Jonathan T. Barron, Samuel W. Hasinoff, and Frédo Durand. Deep bilateral learning for real-time image enhancement. *ACM TOG*, 36(4):118:1–118:12, 2017.

[8] Jingwen He, Yihao Liu, Yu Qiao, and Chao Dong. Conditional sequential modulation for efficient global image retouching. In *ECCV*, pages 679–695, 2020.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2022.

[10] Yuanming Hu, Hao He, Chenxi Xu, Baoyuan Wang, and Stephen Lin. Exposure: A white-box photo post-processing framework. *ACM TOG*, 37(2):26, 2018.

[11] Zhanghan Ke, Chunyi Sun, Lei Zhu, Ke Xu, and Rynson W.H. Lau. Harmonizer: Learning to perform white-box image and video harmonization. In *ECCV*, pages 690–706, 2022.

[12] Hanul Kim, Su-Min Choi, Chang-Su Kim, and Yeong Jun Koh. Representative color transform for image enhancement. In *ICCV*, pages 4459–4468, 2021.

[13] Jie Liang, Hui Zeng, Miaomiao Cui, Xuansong Xie, and Lei Zhang. Ppr10k: A large-scale portrait photo retouching dataset with human-region mask and group-level consistency. In *CVPR*, pages 653–661, 2021.

[14] Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Jiashi Feng, and Jianmin Jiang. A simple pooling-based design for real-

time salient object detection. In *CVPR*, pages 3917–3926, 2019.

[15] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015.

[16] Caron Mathilde, Touvron Hugo, Misra Ishan, Jégou Hervé, Mairal Julien, Bojanowski Piotr, and Joulin Armand. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 3917–3926, 2021.

[17] Sean Moran, Pierre Marza, Steven McDonagh, Sarah Parisot, and Gregory G. Slabaugh. Deeplpf: Deep local parametric filters for image enhancement. In *CVPR*, pages 12823–12832, 2020.

[18] Sean Moran, Steven McDonagh, and Gregory G. Slabaugh. Curl: Neural curve layers for global image enhancement. In *ICPR*, pages 9796–9803. IEEE, 2020.

[19] Kingma Diederik P. and Ba Jimmy. Adam: A method for stochastic optimization. In Bengio Yoshua and LeCun Yann, editors, *ICLR*, 2015.

[20] Yuda Song, Hui Qian, and Xin Du. Starenhancer: Learning real-time and style-aware image enhancement. In *ICCV*, pages 4106–4115, 2021.

[21] Lloyd Stuart. Least squares quantization in pcm. *IEEE TIT*, 28(2):129–137, 1982.

[22] Ruixing Wang, Qing Zhang, Chi-Wing Fu, Xiaoyong Shen, Wei-Shi Zheng, and Jiaya Jia. Underexposed photo enhancement using deep illumination estimation. In *CVPR*, pages 6849–6857, 2019.

[23] Tao Wang, Yong Li, Jingyang Peng, Yipeng Ma, Xian Wang, Fenglong Song, and Youliang Yan. Real-time image enhancer via learnable spatial-aware 3d lookup tables. In *ICCV*, pages 2451–2460, 2021.

[24] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *ECCV*, pages 649–665, 2020.

[25] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *NeurIPS*, 33:17721–17732, 2020.

[26] Zhou Wang, Bovik A.C., Sheikh H.R., and Simoncelli E.P. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004.

[27] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. In *BMVC*, page 155. BMVA Press, 2018.

[28] Van Gansbeke Wouter, Vandenhende Simon, and Van Gool Luc. Discovering object masks with transformers for unsupervised semantic segmentation. *arXiv preprint arXiv:2206.06363*, 2022.

[29] Zhicheng Yan, Hao Zhang, Baoyuan Wang, Sylvain Paris, and Yizhou Yu. Automatic photo adjustment using deep neural networks. *ACM TOG*, 35(2):11:1–11:15, 2016.

[30] Canqian Yang, Meiguang Jin, Xu Jia, Yi Xu, and Ying Chen. Adaint: Learning adaptive intervals for 3d lookup tables on real-time image enhancement. In *CVPR*, pages 17501–17510, 2022.

[31] Hui Zeng, Jianrui Cai, Lida Li, Zisheng Cao, and Lei Zhang. Learning image-adaptive 3d lookup tables for high performance photo enhancement in real-time. *IEEE TPAMI*, 2020.

[32] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In Laurent Amsaleg, Benoit Huet, Martha A. Larson, Guillaume Gravier, Hayley Hung, Chong-Wah Ngo, and Wei Tsang Ooi, editors, *ACM MM*, pages 1632–1640. ACM, 2019.

[33] Zhengxia Zou, Tianyang Shi, Shuang Qiu, Yi Yuan, and Zhenwei Shi. Stylized neural painting. In *CVPR*, pages 15689–15698, 2021.

# Appendices

## A. Filter Function

The retouched result is represented as equation:

$$Y = X + \sum_m \sum_n (F_{m,n}(\theta_{m,n}, X) - X) \odot M_m, \quad (7)$$

For different filters, $F_{m,n}(\theta_{m,n}, X)$ has different expressions:

$$
\begin{aligned}
F_{contrast}(\theta, X) - X &= \theta(X - mean(X)) \\
F_{saturation}(\theta, X) - X &= \theta(X - L(X)) \\
F_{hue}(\theta, X_c) - X_c &= \begin{cases} \alpha_h \theta X_c, X_c \in \{X_R, X_G\} \\ -\dfrac{1}{2}\alpha_h \theta X_c, X_c = X_B \end{cases} \\
F_{temperature}(\theta, X_c) - X_c &= \begin{cases} \alpha_{t,1}\theta X_c, X_c = X_R, \theta \geq 0 \\ \alpha_{t,2}\theta X_c, X_c = X_R, \theta < 0 \\ 0.0, X_c = X_G, \theta \geq 0 \\ \alpha_{t,3}\theta X_c, X_c = X_G, \theta < 0 \\ \alpha_{t,4}\theta X_c, X_c = X_B, \theta \geq 0 \\ \alpha_{t,5}\theta X_c, X_c = X_B, \theta < 0 \end{cases} \\
F_{shadows}(\theta, X) - X &= \theta(1 - X) \\
F_{midtones}(\theta, X) - X &= \theta(0.25 - (X - 0.5)^2) \\
F_{highlights}(\theta, X) - X &= \theta X \\
F_{shift}(\theta, X_c) - X_c &= \begin{cases} \alpha_{s,1}, X_c = X_R \\ \alpha_{s,2}, X_c = X_G \\ \alpha_{s,3}, X_c = X_B \end{cases}
\end{aligned}
$$

$$(8)$$

Where $mean(X)$ denotes the mean value of the entire image, while $L(X)$ represents the L channel of the image in the CIE LAB color space. Additionally, $X_c \in X_R, X_G, X_B$ denotes the RGB color channels of the image. The adjustment factor, denoted by $\alpha$, is a scalar that satisfies $\alpha > 0$. In our experiments, values of $\alpha$ are determined according to traditional color grading tools.

## B. Variation of RSFNet with Controlled Region Shape

**Ground Truths Mask Generation.** We adopt the palette-based method proposed in [2] to generate the main colors $C = C_1, ..., C_n$ of an image, along with distance maps from pixels to $N$ color centers. We obtain region masks by applying a Gaussian smoothing function to these distance maps, resulting in $M_{palette} = M_{p1}, ...M_{pK_p}$. We also predict saliency masks using pre-trained networks from [14], yielding $M_{saliency} = M_{sal1}, M_{sal2}$.

Since previous works on panoptic segmentation split objects into more than one hundred classes, which is redundant for our task, we aim to identify the most significant pixel groupings. To accomplish this, we follow the practice in [16, 28] and train a self-attentioned network with pairwise retouching data. First, we predict semantic masks using the networks presented in [5]. We then apply a clustering algorithm (e.g., K-means [21]) to the output features of the self-attentioned network with masked images as input. Masks assigned with the same cluster index are merged, resulting in $M_{semantics} = M_{sem1}, ...M_{semK_s}$. For each of the three sets of masks, we train a separate model with the corresponding output channel numbers. The entire process is illustrated in Figure 3.

**Differentiable Adaptive Smooth Kernel.** To ensure smooth transition across mask edges, we have incorporated a differentiable adaptive smooth kernel module into our main network. We fix the Gaussian smooth kernel to a suitable size $\sigma_{max}$, such as $51 \times 51$ for the original input with a resolution of $256 \times 256$. The standard variance of the kernel is a learnable parameter, which adapts to the inputs.

## C. Additional Results

**Quantitative Results.** We valuate our methods using a random split setup for PPR10K [13]. The results are demonstrated in Table 4. For more implementation details, please refer to our codebase at https://github.com/Vicky0522/RSFNet.

**Qualitative Results.** We present more results of RSFNet-saliency, RSFNet-palette and RSFNet-map in Figure 9 and 10, including generated masks and corresponding filter arguments.

| Method | 360p expert a | | 360p expert b | | 360p expert c | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| DeepLPF [17] | 24.97 | 0.939 | 24.33 | 0.930 | 24.65 | 0.926 |
| CSRNet [8] | 24.38 | 0.938 | 24.41 | 0.940 | 24.53 | 0.931 |
| 3D-LUT+AdaInt [30] | 27.31 | 0.954 | 26.62 | 0.945 | 26.67 | 0.929 |
| Harmonizer [11] | 25.02 | 0.916 | 23.84 | 0.895 | 25.22 | 0.920 |
| RSFNet-map | 27.25 | 0.956 | 26.61 | 0.954 | 26.76 | 0.945 |
| RSFNet-saliency | 25.98 | 0.946 | 25.87 | 0.948 | 25.88 | 0.937 |

Table 4: Quantitative comparisons for retouching tasks on PPR10K dataset [13]. All the models are trained on data without augmentations.White-box methods are colored as violet.



Figure 9: Editable white-box retouching. Arguments and masks generated by RSFNet-saliency trained with saliency masks are shown in the second column. Retouched result is shown in the third column. Three versions of adjustments conducted on the retouched results are shown in the three columns on the right. Ground truths is shown in the right-most column. Numbers in green boxes indicate relative variation.
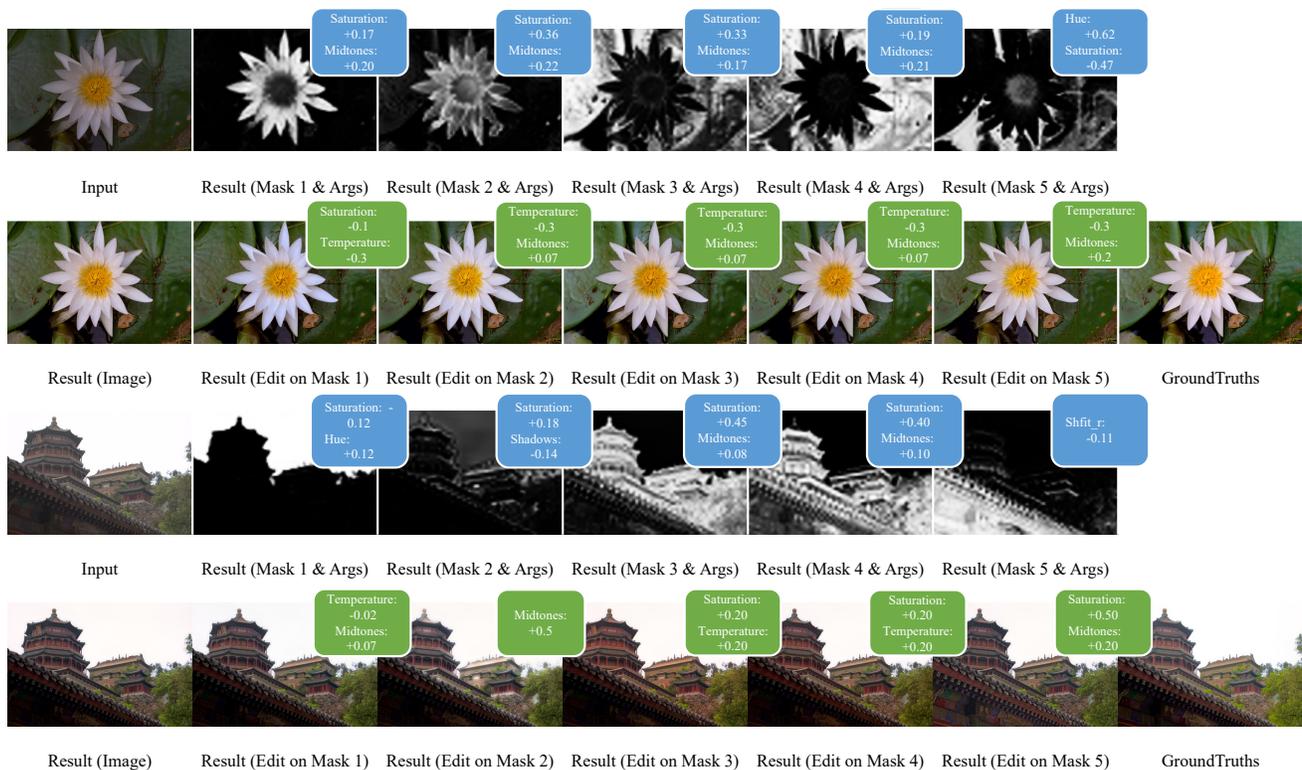
Figure 10: Editable white-box retouching. Arguments and masks generated by RSFNet-palette trained with palette-based masks are shown in the first row. Only two of the most significant arguments are presented. Retouched result is shown in the first column of the second row. Five versions of adjustments conducted on the retouched results and corresponding masks are shown in the rest columns of the second row. Ground truths is shown in the right-most column. Numbers in green boxes indicate relative variation.