



Progressive Recurrent Network for Shadow Removal

Yonghui Wang^a, Wengang Zhou^{a,**}, Hao Feng^a, Li Li^a, Houqiang Li^{a,**}

^aUniversity of Science and Technology of China, Hefei, China

ABSTRACT

Single-image shadow removal is a significant task that is still unresolved. Most existing deep learning-based approaches attempt to remove the shadow directly, which can not deal with the shadow well. To handle this issue, we consider removing the shadow in a coarse-to-fine fashion and propose a simple but effective Progressive Recurrent Network (PRNet). The network aims to remove the shadow progressively, enabling us to flexibly adjust the number of iterations to strike a balance between performance and time. Our network comprises two parts: shadow feature extraction and progressive shadow removal. Specifically, the first part is a shallow ResNet which constructs the representations of the input shadow image on its original size, preventing the loss of high-frequency details caused by the downsampling operation. The second part has two critical components: the re-integration module and the update module. The proposed re-integration module can fully use the outputs of the previous iteration, providing input for the update module for further shadow removal. In this way, the proposed PRNet makes the whole process more concise and only uses 29% network parameters than the best published method. Extensive experiments on the three benchmarks, ISTD, ISTD+, and SRD, demonstrate that our method can effectively remove shadows and achieve superior performance.

© 2023 Elsevier Ltd. All rights reserved.

1. Introduction

Shadow is a widespread natural phenomenon that appears when the light source is blocked. Generally, the shape, position, and intensity of shadows can aid us in understanding natural scenes (Karsch et al., 2011; Lalonde et al., 2012; Okabe et al., 2009; Panagopoulos et al., 2009). However, the existence of shadows may also degrade human perception experience as well as the performance of various computer vision tasks, such as object detection (Mikic et al., 2000; Cucchiara et al., 2003; Nadimi and Bhanu, 2004), object tracking (Khan et al., 2015; Sanin et al., 2010), and others (Levine and Bhattacharyya, 2005; Jung, 2009; Zhang et al., 2018; Sekhavat, 2016). To address this problem, shadow removal has become an essential topic in the computer vision community and been investigated for many years (Finlayson et al., 2009; Hu et al., 2019a; Chen et al., 2021; Liu et al., 2021).

Traditional methods on shadow removal mainly rely on physical models, *e.g.*, entropy minimization (Finlayson et al., 2009,

2005; Guo et al., 2011) and intrinsic priors (Choi et al., 2010; Gryka et al., 2015; Guo et al., 2012; Vicente et al., 2017; Xiao et al., 2013; Wang et al., 2019). However, due to the complexity and uncertainty of the real world, these physical models are not well applicable to the natural shadow scenes. Recently, several studies have suggested that deep learning-based methods are effective to address shadow removal (Qu et al., 2017; Xu et al., 2017; Wang et al., 2018; Hu et al., 2019b,a; Le and Samaras, 2019; H. Le and D. Samaras, 2020; Cun et al., 2020; Fu et al., 2021; Jin et al., 2021; Zhu et al., 2022b,a; Guo et al., 2023; Liu et al., 2023). The mainstream methods commit to solving this problem by designing various specialized architectures and restoring the shadow region directly (Qu et al., 2017; Le and Samaras, 2019; Hu et al., 2019a; Cun et al., 2020; Liu et al., 2021; Fu et al., 2021; Jin et al., 2021; Guo et al., 2023). Albeit the de-shadowing performance is improving, it still has blurry shadow regions with incorrect color. We argue that the illumination information lost in the shadow region can be recovered progressively. Imagine it is late at night, as the morning sun rises, the dim environment gradually becomes brighter. Intuitively, the removing of shadows can also be a progressive process. Therefore, we leverage a coarse-to-fine fashion to remove

**Corresponding author.

e-mail: zhwg@ustc.edu.cn (Wengang Zhou), lihq@ustc.edu.cn (Houqiang Li)

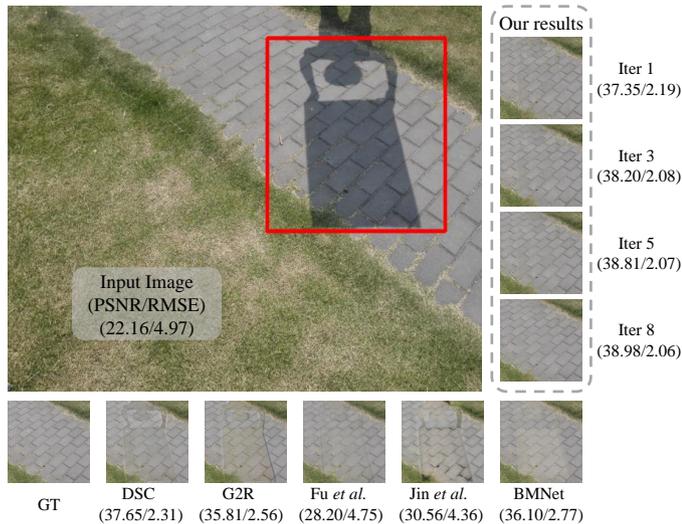


Fig. 1: Removal results of local shadow region in red bounding box by DSC (Hu et al., 2019a), G2R (Liu et al., 2021), Fu *et al.* (Fu et al., 2021), Jin *et al.* (Jin et al., 2021), BMNet (Zhu et al., 2022a) and our method (inference at iter=1, 3, 5, and 8), respectively. It can be seen that our method achieves the best performance at and after the third iteration.

the shadow gradually, which is capable of handling the shadow in different natural scenes with more or fewer iterations.

Typically, ARGAN (Ding et al., 2019) is the first method that removes shadow in a recurrent manner and achieves great success, proving the feasibility of restoring shadow region step-by-step. The core design of ARGAN (Ding et al., 2019) is that it formulates an attentive recurrent generative adversarial network to jointly detect and remove shadows. Moreover, this method is trained with an adversarial training strategy. Given that during adversarial learning the discriminator becomes harder and harder to distinguish whether the generated image is real or fake, it utilizes a semi-supervised strategy to use sufficient unsupervised shadow images available online to strengthen the training and boost the de-shadowing performance. Nevertheless, as a crucial component of ARGAN, shadow attention decoder generates attention maps that directly impact the performance of shadow removal. Additionally, the instability of adversarial training presents challenges in training the network.

Formally, we propose a new simple but effective Progressive Recurrent Network (PRNet) to iteratively recover the content of shadow regions. Our approach follows a coarse-to-fine fashion and allows for flexible adjustment of inference iterations based on demand, thereby achieving a balance between performance and time. The PRNet comprises two parts: shadow feature extraction and progressive shadow removal. The shadow feature extraction network is a shallow ResNet with six residual blocks (He et al., 2016), which extracts features from the original image size for subsequent processing. The progressive shadow removal network is parameter-shared and exhibits two main designs, *i.e.*, the re-integration module and the update module. By repeatedly feeding the refined features into a GRU-based (Cho et al., 2014a) update module, we can obtain more optimized features to achieve progressive shadow removal. Different from the GRU proposed initially (Cho et al., 2014a), we employ the ConvGRU as the recurrent unit as many

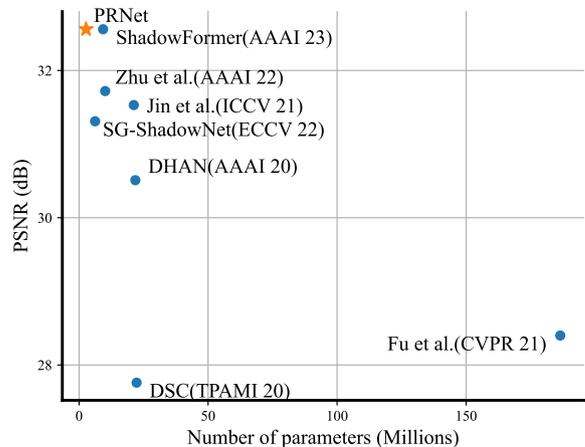


Fig. 2: The PSNR performance v.s., the number of model parameters of shadow removal models on SRD dataset (Qu et al., 2017). The metric is conducted on images with 256×256 resolution.

other tasks (Tokmakov et al., 2017; Teed and Deng, 2020). Furthermore, to better leverage the outputs of the previous iteration, we propose a re-integration module. This module can use the previous output information and guide the update module to obtain better results than the previous one. As shown in Figure 1, a recurrent structure of 8 iterations achieves better results than others (Hu et al., 2019a; Liu et al., 2021; Fu et al., 2021; Jin et al., 2021; Zhu et al., 2022a).

The network is simple and has 2.7M parameters. As shown in Figure 2, by only using 29% parameters of the state-of-the-art methods, ShadowFormer(9.3M) (Guo et al., 2023), we achieve the comparative results in terms of PSNR on SRD dataset (Qu et al., 2017). Such a progressive method makes the pipeline more concise and eases the difficulty of training CNNs directly to recover shadow-free images from shadow images.

We summarize our contributions as follows:

- We propose a new Progressive Recurrent Network to address the problem of shadow removal iteratively.
- The proposed re-integration module can efficiently integrate the last output and hidden state, and provide the refined features to the update module for better optimizing.
- Comprehensive experimental results on the three public datasets, ISTD, ISTD+, and SRD, demonstrate that the proposed method can address the shadow cases well and achieve superior performance.

2. Related Work

Traditional shadow removal. Traditional methods mainly rely on physical models with intrinsic shadow properties, *e.g.*, illumination (Shor and Lischinski, 2008; Xiao et al., 2013; Zhang et al., 2015) and regions (Yang et al., 2012; Guo et al., 2012; Vicente et al., 2017) for shadow removal.

For illumination-based methods, Shor *et al.* (Shor and Lischinski, 2008) use an illumination-invariant distance measure to identify shadow and lit areas, and then these areas are used to estimate the affine shadow information model. This method can produce shadow-free images and avoid loss of texture contrast and introduction of noise. After detecting shadows using a gaussian mixture model, Xiao *et al.* (Xiao et al., 2013) apply an adaptive illumination transfer approach to remove the shadows and leverage a multi-scale illumination transfer technique to improve the contrast and noise level. Further, the method can also extend to video dataset and achieve temporally consistent de-shadowing results. Zhang *et al.* (Zhang et al., 2015) propose a simple shadow removal framework for single natural images as well as color aerial images using an illumination recovering optimization method. The key idea of this method is constructing an optimized illumination recovering operator, which can effectively remove the shadows and recover the texture details.

For region-based methods, Yang *et al.* (Yang et al., 2012) propose a fully automatic method which does not require shadow detection. Based on the chromaticity, they extract a 2-D intrinsic image from a single RGB camera image. Then, using the bilateral filtering technique and the 2-D intrinsic image, a 3-D intrinsic image is recovered. By decomposing and combing these patch regions, they can get the correct luminance pixel values and obtain shadow-free images. Guo *et al.* (Guo et al., 2012) consider the relative illumination between the segmented regions and perform pairwise classification based on these information. Then, they apply a lighting model to relight the shadow pixels. Vicente *et al.* (Vicente et al., 2017) propose another region-based method for shadow removal. Given a pair of shadow and non-shadow regions, they use a relighting transformation method to relight the shadow pixel based on histogram matching of non-shadow pixels.

Additionally, there are other physics based methods. For instance, Finlayson *et al.* (Finlayson et al., 2005, 2009) formulated a physics-based method, entropy minimization, to capture the invariant features of shadow and non-shadow regions belonging to the same surfaces in the log-chromaticity space. Such method is more insensitive to quantization and is quite reliable. However, the above methods rely heavily on the intrinsic properties of shadows. Due to the prior limitations, traditional methods are not effective to handle shadow regions in complex natural scenes.

Deep learning-based shadow removal. Recently, deep learning-based methods have shown remarkable success in the shadow removal field based on the published large-scale datasets (Qu et al., 2017; Wang et al., 2018; Hu et al., 2019b). Specifically, DshadowNet (Qu et al., 2017) designs a multi-context architecture to predict shadow matte for shadow removal. Inspired by physical models of shadow formation, Le *et al.* (Le and Samaras, 2019; H. Le and D. Samaras, 2020) formulate a linear illumination model and apply the network to predict the corresponding shadow parameters for shadow removal. Hu *et al.* (Hu et al., 2019a) propose a direction-aware method to analyze the spatial image context, and use these information for shadow removal. Zhang *et al.* (Zhang et al., 2020) propose

RIS-GAN, which conducts shadow removal in a coarse-to-fine fashion. The network predicts negative residual images and inverse illumination maps to optimize the coarse shadow-removal image, and generates the fine shadow-free image. In the same year, Cun *et al.* (Cun et al., 2020) propose a context aggregation network and hierarchically aggregate these features to produce high-quality border-free images. Fu *et al.* (Fu et al., 2021) estimate multiple over-exposure images and then compensate each pixel individually to tackle position specified color and illumination degradation problem. Niu *et al.* (Niu et al., 2022) propose a boundary-aware network to perform shadow removal and shadow boundary optimization simultaneously. Zhu *et al.* (Zhu et al., 2022b) introduce a new shadow illumination model and reformulate the shadow removal task as a variational optimization problem. The new model is effective and efficient. BMNet (Zhu et al., 2022a) leverages auxiliary shadow invariant color information for bidirectional shadow generation and removal, which can benefit from each other. Wan *et al.* (Wan et al., 2022) design a style-guided shadow removal network to restore the style consistency between shadow and non-shadow regions. Most recently, Guo *et al.* (Guo et al., 2023) propose ShadowFormer, the first transformer-based method for shadow removal. The proposed method exploits the global contextual correlation between shadow and non-shadow regions, which can effectively model the context correlation between these two regions. ST-CGAN (Wang et al., 2018) and ARGAN (Ding et al., 2019) design a novel framework to jointly perform shadow detection and removal, and use the predicted mask of shadow detection to assist shadow removal. Different from the above approaches, some unsupervised deep learning-based methods (Hu et al., 2019b; H. Le and D. Samaras, 2020; Liu et al., 2021; Jin et al., 2021) are proposed, making it possible to perform shadow removal on unpaired datasets with promising results. Furthermore, some tasks also treat shadow removal as a subtask. Zhang *et al.* (Zhang et al., 2021) propose a novel unsupervised framework called UIDNet for intrinsic image decomposition of natural images. Comprising a reflectance prediction network (RPN) and a shading prediction network (SPN), this framework can decompose images into reflectance and shading by promoting the internal self-similarity of the reflectance component. The method can be trained using individual images solely and has demonstrated superior performance. Jin *et al.* (Jin et al., 2023) propose a two-stage learning method for single-image reflectance prediction. In the first stage, an initial reflectance layer is obtained from shadow-free and specular-free priors. In the second stage, a S-Aware network is introduced to distinguish the reflectance image from the input image, further enhancing the performance of the network.

Among the methods mentioned above, ARGAN (Ding et al., 2019) is the most relevant to us. This method employs the update module to gradually optimize hidden features. However, the new features fed into the update module suffer from sub-optimal optimization, leading to poor performance. To perform semi-supervised learning, ARGAN (Ding et al., 2019) resorts to utilizing additional unlabelled data. In contrast, we propose a re-integration module to optimize the features input to the update module, which achieves SOTA performance with only

2.1% parameters of it without the need for additional data.

Progressive learning. Progressive learning mechanism has been explored in a range of computer vision tasks, such as image generation (Gregor et al., 2015; Ahn et al., 2018; Ren et al., 2019), object detection (Cai and Vasconcelos, 2018; Gidaris and Komodakis, 2015; Najibi et al., 2016), and others (Carreira et al., 2016; Liu et al., 2018). More specifically, in low-level vision tasks, Ren *et al.* (Ren et al., 2019) propose a progressive ResNet (He et al., 2016) to take advantage of recursive computation for rain streaks removal. Ahn *et al.* (Ahn et al., 2018) design a CARN module to maintain the stability of training process, and the model can increase the resolution of the output image in a recurrent manner. Moreover, Zamir *et al.* (Zamir et al., 2021) present a novel synergistic multi-stage network to progressively restore the degraded images. Jin *et al.* (Jin et al., 2022) propose a novel progressive method for removing self and soft shadows using a diffusion model (Ho et al., 2020). This method is based on self-tuned ViT feature similarity and color convergence. Additionally, a color convergence loss is introduced to mitigate color deviations, thus facilitating the proficient elimination of hard, soft, and self shadows.

In this paper, we propose a new simple progressive shadow removal method, PRNet. While many of the individual ingredients used in the progressive networks can be found in the literature, *e.g.*, ConvGRU (Cho et al., 2014b). How to make subtle modifications and combinations of them, and apply them to solve the task of shadow removal is novel. Specifically, we conduct extensive experiments on widely used datasets (Wang et al., 2018; Le and Samaras, 2019; Qu et al., 2017) and achieves comparable results with the SOTA methods.

3. METHODOLOGY

In this section, we present our PRNet for shadow removal. As shown in Figure 3, we first apply a shallow ResNet with six residual blocks (He et al., 2016) to extract shadow features from image without downsampling operation. Subsequently, the extracted features are fed into the GRU-based update module as the initial hidden state. By repeatedly feeding the refined features into the update module, we can obtain the output hidden state with more shadow-free signals rather than shadow signals which we refer to as shadow-attenuated features next. For these features, we apply a predict tail for iterative prediction. In the following, we elaborate on the components of our framework separately.

3.1. Shadow Feature Extraction

As shown in Figure 3, given the input shadow image $I_{in} \in \mathbb{R}^{H \times W \times 3}$ and corresponding shadow mask $M_{in} \in \mathbb{R}^{H \times W \times 1}$, we first use the feature extraction module E_θ to extract shadow features. Since shadow removal is a low-level vision restoration task and the downsampling operation would sacrifice the high-frequency details, the feature extraction is performed on the original input scale. Specifically, our shadow feature extraction is a residual module with six residual blocks (He et al., 2016). As shown in Table 1, we first extract features from the shadow image and its mask using a convolutional layer with a kernel

Table 1: The detailed structure of the feature extraction network E_θ in PRNet.

Layer	Output size	Operation
conv1	$64 \times 256 \times 256$	$7 \times 7, 64, s1, p3$
layer1	$64 \times 256 \times 256$	$3 \times 3, 64, s1, p1$
		$3 \times 3, 64, s1, p1$
		$3 \times 3, 64, s1, p1$
		$3 \times 3, 64, s1, p1$
layer2	$96 \times 256 \times 256$	$3 \times 3, 96, s1, p1$
		$3 \times 3, 96, s1, p1$
		$3 \times 3, 96, s1, p1$
		$3 \times 3, 96, s1, p1$
layer3	$128 \times 256 \times 256$	$3 \times 3, 128, s1, p1$
		$3 \times 3, 128, s1, p1$
		$3 \times 3, 128, s1, p1$
		$3 \times 3, 128, s1, p1$
conv2	$128 \times 256 \times 256$	$1 \times 1, 128, s1, p0$

size of 7, where the input channel is 4, similar to previous methods (Hu et al., 2019a; Cun et al., 2020; Zhu et al., 2022a,b; Guo et al., 2023). Next, the network is divided into three layers, and each layer contains two residual blocks. After each layer, we increase the number of channels by 24, and the output channel in layer three is 128. To dramatically accelerate the training speed as well as boost the network performance, instance normalization (Ulyanov et al., 2016) and ReLU activation function (Nair and Hinton, 2010) are added after every convolution operation. Finally, we use a convolution with kernel size 1 to further enhance the nonlinear ability of the network. The shadow feature extraction network E_θ produces the feature $\mathbf{h}_0 \in \mathbb{R}^{H \times W \times C}$, where $C = 128$. \mathbf{h}_0 will be integrated with the shadow image as the input of the update module, whereas it will also be used as the initial hidden state of the update module. Note that our feature extraction module E_θ will only extract features once. Then all subsequent iterative processes will be carried out in the progressive shadow removal, which is discussed in the following subsection.

3.2. Progressive Shadow Removal

Progressive shadow removal is the critical component of our method, which consists of two parts: the re-integration module and the update module. The re-integration module is applied to fuse the outputs of the last iteration and provide the input of the update module, while the update module is utilized to obtain the predicted results of each iteration.

Re-integration module. Figure 4 left shows the proposed re-integration module, which integrates the outputs of the last iteration to provide input for the next iteration. Taking the k^{th} iteration as an example, there are two outputs in the last iteration: one is the shadow-attenuated image $I_{k-1} \in \mathbb{R}^{H \times W \times 3}$, and the other is the hidden state $\mathbf{h}_{k-1} \in \mathbb{R}^{H \times W \times C}$. For shadow-attenuated image I_{k-1} , we first concatenate it with the corresponding shadow mask M_{in} to provide the shadow region information, and then extract features through the convolution operation to obtain the feature $\mathbf{F}_s \in \mathbb{R}^{H \times W \times C_1}$. For the hidden state \mathbf{h}_{k-1} , we also perform feature refinement through the convolutional layer, and obtain the feature $\mathbf{F}_l \in \mathbb{R}^{H \times W \times C_2}$, where we

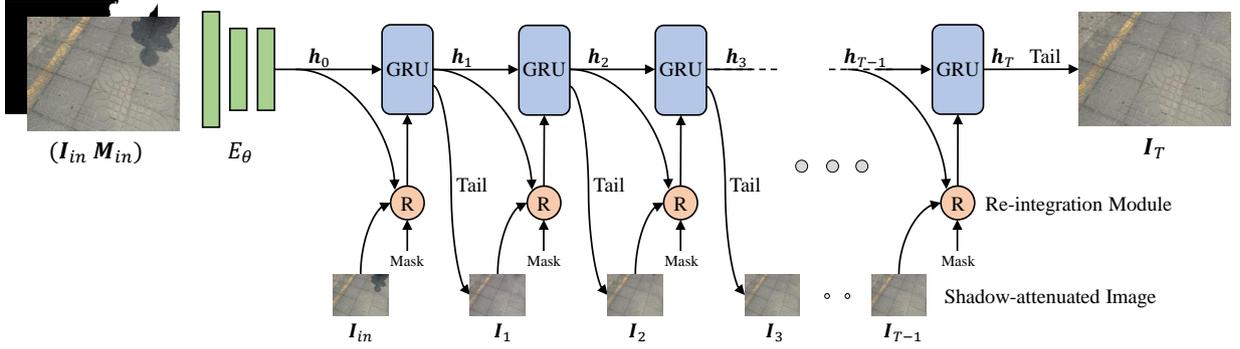


Fig. 3: An overview of the proposed PRNet. PRNet is divided into two parts: shadow feature extraction and progressive shadow removal. Shadow feature extraction is a shallow ResNet E_θ with six residual blocks. Progressive shadow removal consists of two components: the re-integration module and the GRU-based update module. The re-integration module fuses the outputs of last iteration and produces the integrated feature as the input of the next iteration. Then the update module is applied to generate shadow-attenuated features and feed to the prediction tail for prediction. We iteratively conduct the update operation to progressively improve the shadow removal result.

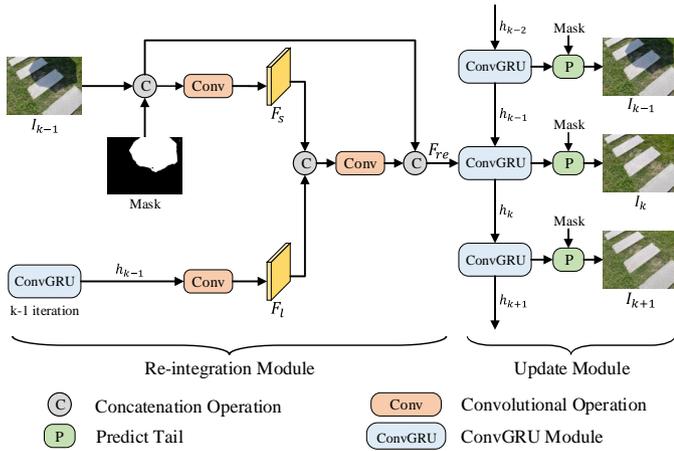


Fig. 4: Illustration of the process of k^{th} iteration. Given the hidden state h_{k-1} and the prediction results I_{k-1} of the last iteration, first they are fed into the re-integration module to produce the integrated features F_{re} . Then both h_{k-1} and F_{re} are put into the ConvGRU (Cho et al., 2014b) operator and output the k^{th} hidden state h_k . Subsequently, h_k is sent to the predict tail for k^{th} prediction.

set $C_1 = 192$ and $C_2 = 64$, respectively. Next, we concatenate both of them followed by another convolution operation to generate the final integrated feature $F_r \in \mathbb{R}^{H \times W \times C}$. This re-integrated feature combined with the prediction of last iteration is set as the input of the update module. The whole process can be viewed as using the prediction results of the last iteration to enhance the current iteration. This way, the update module can flexibly update and reset the upcoming hidden features.

Update module. Figure 5 shows our update module. The core component of it is a ConvGRU (Cho et al., 2014b) block, which has been used in many other computer vision tasks (Tokmakov et al., 2017; Teed and Deng, 2020). ConvGRU is a variant of the original GRU (Cho et al., 2014a), in which the fully connected layers are replaced by the convolutional layers. The whole pro-

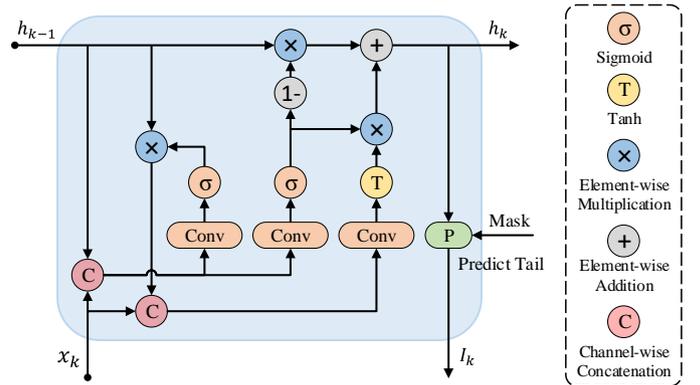


Fig. 5: Illustration of the proposed update module, which consists of a ConvGRU (Cho et al., 2014b) operator and a prediction tail.

cess can be formulated as follows,

$$\begin{aligned}
 z_k &= \sigma(\text{Conv}([h_{k-1}, x_k], W_z)), \\
 r_k &= \sigma(\text{Conv}([h_{k-1}, x_k], W_r)), \\
 \tilde{h}_k &= \tanh(\text{Conv}([r_k \odot h_{k-1}, x_k], W_h)), \\
 h_k &= (1 - z_k) \odot h_{k-1} + z_k \odot \tilde{h}_k,
 \end{aligned} \tag{1}$$

where x_k is the output of the re-integration module, which is the fusion of the shadow-attenuated image I_{k-1} and the hidden state h_{k-1} of the last iteration. W is the learnable parameter of a convolutional layer. The other part of update module is the shadow prediction tail which has two convolutional operation. The first convolution is followed by a ReLU activation function (Nair and Hinton, 2010) and the output channel is set as 256 and the second convolution predicts the output directly. After the hidden state is renewed by the ConvGRU (Cho et al., 2014b), it is passed through the prediction tail to produce the shadow-attenuated image I_k . Subsequently, this image is passed to the re-integration module for next iteration.

3.3. Loss Function

During the training phase, we supervise our progressive recurrent network with the L_1 distance loss between the predicted shadow-attenuated image and the ground truth shadow-

free image over all T iterations. With exponentially increasing weights, the loss is formulated as follows

$$\mathcal{L} = \sum_{i=1}^T \gamma^{T-i} \|I_{gt} - I_i\|_1, \quad (2)$$

where I_{gt} and I_i denote ground-truth shadow-free image and i^{th} iteration shadow-attenuated image, respectively. Empirically, we set $\gamma = 0.8$ in our experiment.

4. Experiments

4.1. Experiment Setup

Benchmark datasets. We train and evaluate the proposed method on three public datasets: ISTD (Wang et al., 2018), adjusted ISTD (ISTD+) (Le and Samaras, 2019), and SRD (Qu et al., 2017). ISTD dataset (Wang et al., 2018) consists of 1870 image triples (shadow images, shadow-free images, and shadow masks), which are divided into 1330 training triplets and 540 testing triplets. Adjusted ISTD (ISTD+), presented in (Le and Samaras, 2019), has the same number of triples as ISTD. By applying the proposed color adjustment algorithm, the color inconsistency between the shadow and shadow-free images is decreased. SRD dataset (Qu et al., 2017) consists of 2680 training pairs and 408 testing pairs of shadow and shadow-free images without ground-truth shadow masks. Since SRD (Qu et al., 2017) does not provide the ground-truth shadow masks, we utilize the public SRD shadow masks provided by DHAN (Cun et al., 2020) during training and testing, following the previous methods (Cun et al., 2020; Fu et al., 2021; Zhu et al., 2022a,b; Wan et al., 2022; Guo et al., 2023).

Evaluation metrics. We employ the root mean square error (RMSE) between the predicted shadow removal image and the ground-truth shadow-free image in LAB color space. For the RMSE metric, the lower the better. We also adopt the Peak Signal-to-Noise Ratio (PSNR) and the structural similarity (SSIM) (Wang et al., 2004) to measure the de-shadowing performance in the RGB color space. The higher the PSNR and SSIM, the better the performance. Following the previous works (Fu et al., 2021; Jin et al., 2021; Liu et al., 2021; Zhu et al., 2022a), we conduct the evaluation with a resolution of 256×256 and compare our method with several state-of-the-art methods on the ISTD, ISTD+, and SRD datasets (Wang et al., 2018; Le and Samaras, 2019; Qu et al., 2017) in both quantitative and qualitative ways.

Implementation details. Our proposed method is implemented by PyTorch 1.8 on the linux platform with NVIDIA RTX 2080Ti GPUs. During training, we randomly crop the images into 256×256 patches. For the three benchmarks, the total training epochs and mini-batch size are set to 300 and 4, respectively. An Adam (Kingma and Ba, 2014) optimizer with an initial learning rate of 2×10^{-4} is applied to optimize the network, and the learning rate will be linearly decayed to 0 in the last 250 epochs. For ISTD+ dataset (Le and Samaras, 2019), we set the training iteration $T = 7$, while for SRD dataset (Qu et al., 2017)

which contains more training samples, we set $T = 8$. During inference, we take the same number of iterations as during training. In practical application, the number of inference iterations can be flexibly adjusted based on specific requirements to strike a balance between the performance and time.

4.2. Comparison with State-of-the-art Methods

Shadow removal evaluation on ISTD dataset. We first report the quantitative shadow removal results of our method on ISTD dataset (Wang et al., 2018). As shown in Table 2 and Table 3, to validate the scalability of the method, we conduct evaluation on both 256×256 resolution and the original size. We compare the proposed method with the state-of-the-art algorithms: Guo *et al.* (Guo et al., 2012), ShadowGAN (Hu et al., 2019b), ST-CGAN (Wang et al., 2018), ARGAN (Ding et al., 2019), DSC (Hu et al., 2019a), DHAN (Cun et al., 2020), G2R (Liu et al., 2021), Fu *et al.* (Fu et al., 2021), Jin *et al.* (Jin et al., 2021), Zhu *et al.* (Zhu et al., 2022b), BMNet (Zhu et al., 2022a), SG-ShadowNet (Wan et al., 2022), and ShadowFormer (Guo et al., 2023). Note that different from other deep learning-based methods, Guo *et al.* (Guo et al., 2012) is the traditional shadow removal method. The results of the state-of-the-art methods are directly provided by the authors or obtained from the original paper. However, the code of ARGAN (Ding et al., 2019) is not publicly available, so we carefully calculate it based on the details provided in the original paper. Our method performs better than ARGAN (Ding et al., 2019) in terms of PSNR, SSIM, and RMSE value, indicating the effectiveness of the progressive method. Additionally, we only use 2.1% of its parameters. Compared to G2R (Liu et al., 2021), Fu *et al.* (Fu et al., 2021), and Jin *et al.* (Jin et al., 2021), our methods also achieves the best performance among all metrics. Compared to the two papers by Zhu *et al.* (Zhu et al., 2022b,a), most results of our method are superior to them. ShadowFormer (Guo et al., 2023), which is the first transformer-based shadow removal method, achieves the state-of-the-art performance on this task. Our method also obtains the competitive performance with it.

Figure 6 illustrates the visualization comparison results of the shadow removal from other state-of-the-art methods and our method on ISTD dataset (Wang et al., 2018). As mentioned in the original paper (Wang et al., 2018), there are a slightly inconsistent colors between shadow and shadow-free images in this dataset, which is caused by the different capturing times of the day. We can see that for the traditional method, Guo *et al.* (Guo et al., 2012) can not remove the shadow effectively due to the limited modeling capacity in the relatively complex scenes. In the third and fourth columns, DHAN (Cun et al., 2020) and Fu *et al.* (Fu et al., 2021) tend to generate blurry images, and they also contain random artifacts and incorrect colors. For the results of Zhu *et al.* (Zhu et al., 2022b), due to the inability of physical model to adapt to various environments, it contains artifacts around the shadow region. ShadowFormer (Guo et al., 2023) performs best among the above methods, and is able to restore more realistic color in the shadow area. Compared to them, by removing the shadows progressively, our method can maintain the color consistency between the shadow and non-shadow regions.

Table 2: Quantitative comparison of our method with the state-of-the-art methods on ISTD dataset (Wang et al., 2018). The best and the second results are highlighted in bold and underlined, respectively. “↑” indicates the higher the better and “↓” indicates the lower the better. S, NS, and ALL indicate the shadow region, non-shadow region, and all the image, respectively. T represents the number of iterations. All metrics are conducted on images with 256×256 resolution.

Method	Params	Flops	Shadow Region (S)			Non-Shadow Region (NS)			All Image (ALL)		
			PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓
Input Image	-	-	22.40	0.936	32.10	27.32	0.976	7.09	20.56	0.893	10.88
Guo <i>et al.</i> (Guo et al., 2012)	-	-	27.76	0.964	18.65	26.44	0.975	7.76	23.08	0.919	9.26
ShadowGAN (Hu et al., 2019b)	11.4M	56.8G	-	-	12.67	-	-	6.68	-	-	7.41
ST-CGAN (Wang et al., 2018)	29.2M	<u>17.9G</u>	33.74	0.981	9.99	29.51	0.958	6.05	27.44	0.929	6.65
ARGAN (Ding et al., 2019)	125.8M	-	-	-	6.65	-	-	5.41	-	-	5.89
DSC (Hu et al., 2019a)	22.3M	123.5G	34.64	0.984	8.72	31.26	0.969	5.04	29.00	0.944	5.59
DHAN (Cun et al., 2020)	21.8M	262.9G	35.53	0.988	7.73	31.05	0.971	5.29	29.11	0.954	5.66
G2R (Liu et al., 2021)	22.8M	113.9G	32.66	0.984	10.47	26.27	0.968	7.57	25.07	0.946	7.88
Fu <i>et al.</i> (Fu et al., 2021)	143.0M	160.3G	34.71	0.975	7.91	28.61	0.880	5.51	27.19	0.945	5.88
Jin <i>et al.</i> (Jin et al., 2021)	21.2M	105.0G	31.69	0.976	11.43	29.00	0.958	5.81	26.38	0.922	6.57
Zhu <i>et al.</i> (Zhu et al., 2022b)	10.1M	56.1G	36.95	0.987	8.29	31.54	0.978	4.55	29.85	0.960	5.09
BMNet (Zhu et al., 2022a)	0.4M	11.0G	35.61	0.988	7.60	<u>32.80</u>	<u>0.976</u>	4.59	30.28	0.959	5.02
SG-ShadowNet (Wan et al., 2022)	6.2M	39.7G	36.03	0.988	7.30	<u>32.56</u>	<u>0.978</u>	4.38	30.23	0.961	4.80
ShadowFormer (Guo et al., 2023)	9.3M	100.9G	38.19	0.991	5.96	34.32	0.981	3.72	32.21	0.968	4.09
Ours	<u>2.7M</u>	<u>73.7+88.5T</u>	36.47	<u>0.990</u>	6.43	<u>32.80</u>	<u>0.978</u>	4.26	<u>30.57</u>	<u>0.964</u>	4.57

Table 3: Quantitative comparison of our method with the state-of-the-art methods on ISTD dataset (Wang et al., 2018). The best and the second results are highlighted in bold and underlined, respectively. “↑” indicates the higher the better and “↓” indicates the lower the better. S, NS, and ALL indicate the shadow region, non-shadow region, and all the image, respectively. All metrics are conducted on images with the original size.

Method	Shadow Region (S)			Non-Shadow Region (NS)			All Image (ALL)		
	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓
Input Image	22.34	0.935	33.23	26.45	0.947	7.25	20.33	0.874	11.35
ARGAN (Ding et al., 2019)	-	-	9.21	-	-	6.27	-	-	6.63
DSC (Hu et al., 2019a)	33.45	0.967	9.76	28.18	0.885	6.14	26.62	0.845	6.67
DHAN (Cun et al., 2020)	34.79	0.983	8.13	29.54	0.941	5.94	27.88	0.921	6.29
G2R (Liu et al., 2021)	32.31	<u>0.978</u>	11.18	25.51	0.941	8.10	24.40	0.915	8.42
Fu <i>et al.</i> (Fu et al., 2021)	33.59	0.958	8.73	27.01	0.794	6.24	25.71	0.745	6.62
Jin <i>et al.</i> (Jin et al., 2021)	30.59	0.949	12.43	25.88	0.785	7.11	24.16	0.724	7.79
Zhu <i>et al.</i> (Zhu et al., 2022b)	33.78	0.956	9.44	27.39	0.786	6.23	26.06	0.734	6.68
BMNet (Zhu et al., 2022a)	34.84	<u>0.983</u>	8.31	31.14	0.949	5.16	29.02	0.929	5.59
SG-ShadowNet (Wan et al., 2022)	35.17	<u>0.982</u>	8.21	30.86	0.950	5.04	28.95	0.928	5.48
ShadowFormer (Guo et al., 2023)	37.03	0.985	6.76	32.20	0.953	4.44	30.47	0.935	4.79
Ours	<u>35.65</u>	0.985	<u>7.12</u>	<u>31.17</u>	<u>0.951</u>	<u>4.85</u>	<u>29.29</u>	<u>0.933</u>	<u>5.17</u>

Table 4: Quantitative comparison of our method with the state-of-the-art methods on SRD dataset (Qu et al., 2017). The best and the second results are highlighted in bold and underlined, respectively. “↑” indicates the higher the better and “↓” indicates the lower the better. S, NS, and ALL indicate the shadow region, non-shadow region, and all the image, respectively. All metrics are conducted on images with 256×256 resolution.

Method	Shadow Region (S)			Non-Shadow Region (NS)			All Image (ALL)		
	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓
Input Image	18.96	0.871	36.69	31.47	0.975	4.83	18.19	0.830	14.05
Guo <i>et al.</i> (Guo et al., 2012)	-	-	29.89	-	-	6.47	-	-	12.60
DeShadowNet (Qu et al., 2017)	-	-	11.78	-	-	4.84	-	-	6.64
DSC (Hu et al., 2019a)	30.65	0.960	8.62	31.94	0.965	4.41	27.76	0.903	5.71
ARGAN (Ding et al., 2019)	-	-	6.35	-	-	4.46	-	-	5.31
DHAN (Cun et al., 2020)	33.67	0.978	8.94	34.79	0.979	4.80	30.51	0.949	5.67
Fu <i>et al.</i> (Fu et al., 2021)	32.26	0.966	8.55	31.87	0.945	5.74	28.40	0.893	6.50
Jin <i>et al.</i> (Jin et al., 2021)	34.00	0.975	7.70	35.53	0.981	3.65	31.53	0.955	4.65
Zhu <i>et al.</i> (Zhu et al., 2022b)	34.94	0.980	7.44	35.85	0.982	3.74	31.72	0.952	4.79
BMNet (Zhu et al., 2022a)	35.05	0.981	6.61	36.02	0.982	3.61	31.69	0.956	4.46
SG-ShadowNet (Wan et al., 2022)	36.55	0.981	7.56	34.23	0.961	3.06	31.31	0.927	4.30
ShadowFormer (Guo et al., 2023)	36.91	0.989	<u>5.90</u>	<u>36.22</u>	0.989	3.44	32.90	<u>0.958</u>	<u>4.04</u>
Ours	36.30	<u>0.984</u>	5.66	36.56	<u>0.983</u>	<u>3.34</u>	<u>32.56</u>	0.960	3.99

Shadow removal evaluation on ISTD+ dataset. We report the shadow removal performance of our method on the adjusted ISTD (ISTD+) dataset (Le and Samaras, 2019). As shown in Table 6 and Table 7, we compare the proposed method with several state-of-the-art algorithms: Guo *et al.* (Guo et al., 2012),

ST-CGAN (Wang et al., 2018), DeShadowNet (Qu et al., 2017), Mask-ShadowGAN (Hu et al., 2019b), Param+M+D-Net (H. Le and D. Samaras, 2020), G2R (Liu et al., 2021), SP+M-Net (Le and Samaras, 2019), Fu *et al.* (Fu et al., 2021), Jin *et al.* (Jin et al., 2021), SG-ShadowNet (Wan et al., 2022), BM-

Table 5: Quantitative comparison of our method with the state-of-the-art methods on SRD dataset (Qu et al., 2017). The best and the second results are highlighted in bold and underlined, respectively. “↑” indicates the higher the better and “↓” indicates the lower the better. S, NS, and ALL indicate the shadow region, non-shadow region, and all the image, respectively. All metrics are conducted on images with the original size.

Method	Shadow Region (S)			Non-Shadow Region (NS)			All Image (ALL)		
	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓	PSNR↑	SSIM↑	RMSE↓
Input Image	19.00	0.871	39.23	28.41	0.949	5.86	17.87	0.804	14.62
DSC (Hu et al., 2019a)	25.95	0.912	20.40	22.46	0.748	16.89	20.15	0.642	17.75
DHAN (Cun et al., 2020)	32.21	0.969	8.39	30.58	0.943	5.02	27.70	0.898	5.88
Fu <i>et al.</i> (Fu et al., 2021)	31.19	0.955	9.65	28.10	0.894	6.63	25.83	0.825	7.32
Jin <i>et al.</i> (Jin et al., 2021)	31.21	0.955	9.23	28.62	0.896	5.91	26.18	0.827	6.72
Zhu <i>et al.</i> (Zhu et al., 2022b)	28.25	0.930	11.57	23.83	0.803	8.48	21.95	0.700	9.19
BMNet (Zhu et al., 2022a)	<u>33.28</u>	<u>0.973</u>	<u>7.84</u>	<u>32.71</u>	0.963	<u>4.38</u>	<u>29.21</u>	<u>0.923</u>	<u>5.24</u>
Ours	34.07	0.975	6.93	32.98	<u>0.962</u>	4.16	29.70	0.925	4.85

Table 6: Quantitative comparison of our method with the state-of-the-art methods on ISTD+ datasets (Le and Samaras, 2019). The best and the second results are highlighted in bold and underlined, respectively. “↓” indicates the lower the better. All metrics are conducted on images with 256×256 resolution.

Method	RMSE↓		
	Shadow	Non-Shadow	All Image
Input Images	39.0	2.6	8.4
Guo <i>et al.</i> (Guo et al., 2012)	22.0	3.1	6.1
ST-CGAN (Wang et al., 2018)	13.4	7.7	8.7
DeshadowNet (Qu et al., 2017)	15.9	6.0	7.6
Mask-ShadowGAN (Hu et al., 2019b)	12.4	4.0	5.3
Param+M+D-Net (H. Le and D. Samaras, 2020)	9.7	3.0	4.0
G2R (Liu et al., 2021)	7.3	2.9	3.6
SP+M-Net (Le and Samaras, 2019)	7.9	3.1	3.9
Fu <i>et al.</i> (Fu et al., 2021)	6.7	3.8	4.2
Jin <i>et al.</i> (Jin et al., 2021)	10.4	3.6	4.7
SG-ShadowNet (Wan et al., 2022)	5.9	2.9	3.4
BMNet (Zhu et al., 2022a)	5.6	2.5	3.0
ShadowFormer (Guo et al., 2023)	5.2	2.3	2.8
Ours	<u>5.5</u>	2.3	<u>2.9</u>

Table 7: Quantitative comparison of our method with the state-of-the-art methods on ISTD+ datasets (Le and Samaras, 2019). The best and the second results are highlighted in bold and underlined, respectively. “↓” indicates the lower the better. All metrics are conducted on images with the original size.

Method	RMSE↓		
	Shadow	Non-Shadow	All Image
Input Images	38.5	3.3	9.2
Fu <i>et al.</i> (Fu et al., 2021)	10.4	8.0	8.4
Jin <i>et al.</i> (Jin et al., 2021)	11.9	5.3	6.3
SG-ShadowNet (Wan et al., 2022)	7.3	3.8	4.3
BMNet (Zhu et al., 2022a)	6.6	3.2	3.7
ShadowFormer (Guo et al., 2023)	6.2	<u>3.2</u>	3.6
Ours	<u>6.3</u>	3.1	3.6

Net (Zhu et al., 2022a), and ShadowFormer (Guo et al., 2023). Unlike other methods, Mask-ShadowGAN (Hu et al., 2019b) adopts unpaired shadow and shadow-free images for training. For ARGAN (Ding et al., 2019), due to the using of extra on-line data for semi-supervised learning and the hyperparameters of the training detail are unknown, we can not reproduce it. The results show that our method outperforms most previous methods. For instance, compared to the second best method BMNet (Zhu et al., 2022a), our method outperforms it by reducing the RMSE from 3.0 to 2.9 for the whole image, indicating the effectiveness of our method. Compared to the latest method ShadowFormer (Guo et al., 2023), our method obtains the same shadow removal results with the lowest RMSE in the non-shadow region. While for the shadow and the whole image, our method also has competitive results. Specifically, our PR-

Net is 0.1 worse than ShadowFormer (Guo et al., 2023) in the whole image in terms of RMSE on the images with 256×256 resolution. We argue that the performance can be further improved by increasing the number of iterations.

Figure 7 illustrates the visualization comparison results of the shadow removal from other state-of-the-art methods and our method on ISTD+ dataset (Le and Samaras, 2019). G2R (Liu et al., 2021) and Jin *et al.* (Jin et al., 2021) tend to generate blurry images, and they also contain random artifacts and incorrect colors. For the results of Param+M+D-Net (H. Le and D. Samaras, 2020), due to the simplified linear shadow model, it contains artifacts around the shadow region. Although the methods (Wan et al., 2022; Zhu et al., 2022a) can remove most of the shadows, they still suffer from the inconsistent color and shadow boundaries between the restored shadow region and the non-shadow region. In contrast, ShadowFormer (Guo et al., 2023) and our method perform well in these cases.

Shadow removal evaluation on SRD dataset. As shown in Table 4 and Table 5, we report the comparison results with other state-of-the-art methods on SRD dataset (Qu et al., 2017), including Guo *et al.* (Guo et al., 2012), DeshadowNet (Qu et al., 2017), DSC (Hu et al., 2019a), ARGAN (Ding et al., 2019), DHAN (Cun et al., 2020), Fu *et al.* (Fu et al., 2021), Jin *et al.* (Jin et al., 2021), Zhu *et al.* (Zhu et al., 2022b), BMNet (Zhu et al., 2022a), SG-ShadowNet (Wan et al., 2022), and ShadowFormer (Guo et al., 2023). In terms of RMSE value, the proposed method obtains the best shadow removal performance in the all image. Specifically, our method outperforms the ARGAN (Ding et al., 2019) in the shadow, non-shadow, and the whole image. Compared to Fu *et al.* (Fu et al., 2021), Jin *et al.* (Jin et al., 2021) and Zhu *et al.* (Zhu et al., 2022b), our method performs best among all the metrics, including PSNR, SSIM and RMSE values. In addition, our method outperforms the method BMNet (Zhu et al., 2022a) by 14.4%, 7.5%, and 10.5% RMSE decreasing in the shadow region, non-shadow region, and the whole image, respectively. Compared to SG-ShadowNet (Wan et al., 2022), our method reduces the RMSE from 4.30 to 3.99, achieving 7.21% decreasing in the whole image. While for the transformer-based method ShadowFormer (Guo et al., 2023), we still have competitive results. In addition, we also provide the visual comparison results in Figure 8. For the first row, the PRNet can well restore the orig-

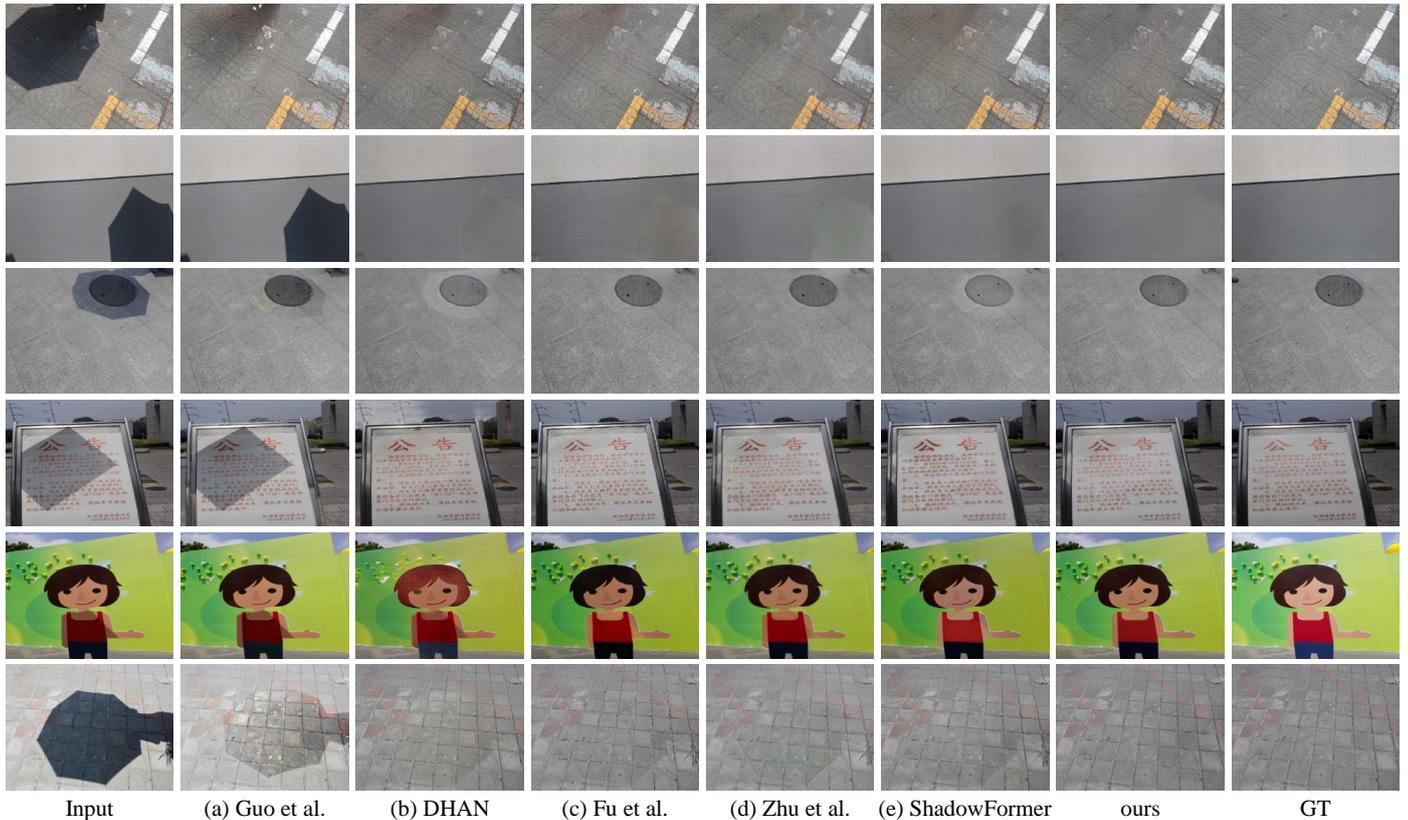


Fig. 6: Visual comparison results of shadow removal on the ISTD dataset (Wang et al., 2018). (a) to (f) are the predicted results from SOTA methods: Guo *et al.* (Guo et al., 2012), DHAN (Cun et al., 2020), Fu *et al.* (Fu et al., 2021), Zhu *et al.* (Zhu et al., 2022b), and ShadowFormer (Guo et al., 2023), respectively.

Table 8: Ablation study of the component on SRD dataset (Qu et al., 2017). The best result is highlighted in bold.

	re-integration	update	RMSE↓
Basic	×	×	6.32
Basic+re	✓	×	4.61
Basic+up	×	✓	4.50
Ours	✓	✓	3.99

inal color information of shadow region, avoiding color-bias effect. Other visual comparison results show that our method can well remove the shadows, and have good visual perception effect. Additionally, as shown in Figure 10, we present more visual results of our method and show various types of shadows, including small shadows, soft shadows, and dark shadows on black objects.

4.3. Ablation Studies of Network Component

We conduct ablation studies on the re-integration module and update module to verify the effectiveness of our network design, and all experiments are conducted on the SRD dataset (Qu et al., 2017). Here we consider three baseline networks. The first baseline network (denoted as "Basic") only has feature extraction network. The second (denoted as "Basic+re") and the third (denoted as "Basic+up") consider the re-integration module and update module, respectively. Table 8 shows the quantitative comparison results. Both the re-integration module and

Table 9: Ablation study of the number of training iteration T on SRD dataset (Qu et al., 2017). Empirically, we set $T = 8$ in our paper.

Iteration	Metrics		
	PSNR↑	SSIM↑	RMSE↓
1	31.13	0.952	4.57
2	31.48	0.954	4.38
3	31.98	0.956	4.25
4	32.19	0.958	4.19
5	32.33	0.958	4.11
6	32.35	0.958	4.06
7	32.44	0.959	4.01
8	32.56	0.960	3.99
9	32.58	0.960	3.98
10	32.60	0.960	3.97

the update module can boost the shadow removal performance. More specifically, with the re-integration module and the update module, the RMSE value is improved from 6.32 to 4.61 and 6.32 to 4.50, respectively. By using both the two modules, the RMSE value reaches 3.99, demonstrating the importance of each component for shadow removal.

4.4. Generalization ability

To verify the generalization ability of our method, we conduct experiments on the SBU-Timelapse dataset (Le and Samaras, 2021), and compare it with the state-of-the-art methods,

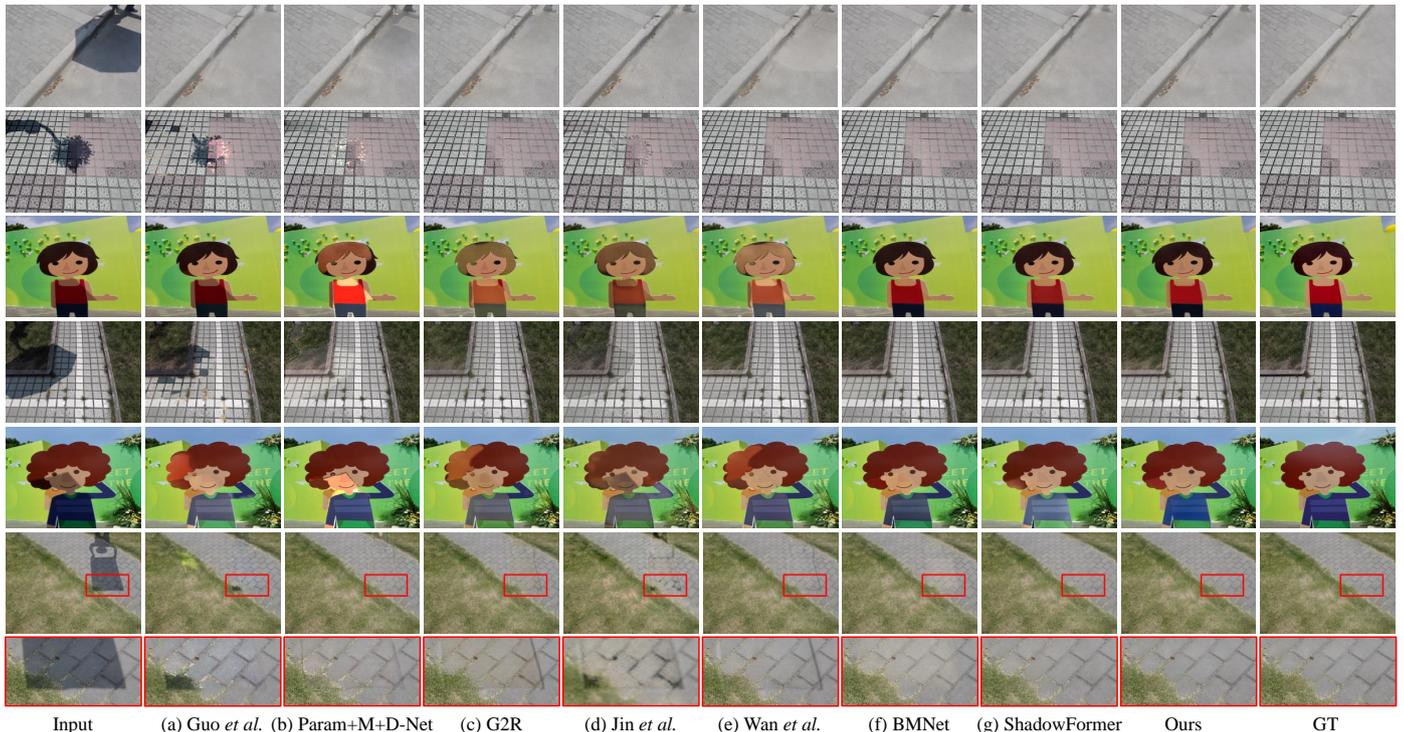


Fig. 7: Visual comparison results of shadow removal on the ISTD+ dataset (Le and Samaras, 2019). (a) to (f) are the predicted results from SOTA methods: Guo *et al.* (Guo *et al.*, 2012), Param+M+D-Net (H. Le and D. Samaras, 2020), G2R (Liu *et al.*, 2021), Jin *et al.* (Jin *et al.*, 2021), Wan *et al.* (Wan *et al.*, 2022), BMNet (Zhu *et al.*, 2022a), and ShadowFormer (Guo *et al.*, 2023), respectively.

Table 10: Quantitative comparison of our method with the state-of-the-art methods on the SBU-Timelapse dataset (Le and Samaras, 2021). The best and the second results are highlighted in bold and underlined, respectively. “ \uparrow ” indicates the higher the better and “ \downarrow ” indicates the lower the better. The evaluation is conducted in the shadow region.

Method	RMSE \downarrow	PSNR \uparrow	SSIM \uparrow
SID (Le and Samaras, 2021)	18.2	<u>20.54</u>	0.893
Fu <i>et al.</i> (Jin <i>et al.</i> , 2021)	19.0	19.63	0.893
SG-ShadowNet (Wan <i>et al.</i> , 2022)	<u>17.5</u>	20.33	<u>0.894</u>
Ours	16.0	21.66	0.903

including SID (Le and Samaras, 2021), Fu *et al.* (Jin *et al.*, 2021), and SG-ShadowNet (Wan *et al.*, 2022). As shown in Table 10, our method outperforms the other three methods in all metrics. Compared to SG-ShadowNet (Wan *et al.*, 2022), we decrease the RMSE from 17.5 to 16.0, and achieve an increase in PSNR from 20.33 to 21.66 and SSIM from 0.894 to 0.903 for shadow regions. Compared to SID (Le and Samaras, 2021) and Fu *et al.* (Jin *et al.*, 2021), our method also exhibits best performance. As shown in Figure 11, through progressive learning, our method achieves acceptable perceptual performance in complex environments, demonstrating the generalization ability of our method.

4.5. Discussions about Network Iterations

Analysis of the training iterations. Following previous methods (Ding *et al.*, 2019), we conduct experiments with training iteration $T = 1, 2, \dots, 10$ to explore how the training iteration impacts the performance. We choose the SRD dataset (Qu *et al.*,

2017) which contains more samples than ISTD (Wang *et al.*, 2018) and it can well evaluate the algorithm capability to handle various natural scenes. In Table 9, it can be concluded that when the number of iterations is increasing gradually, the performance of our method first has a significant improvement and then tends to be stable. In our experiments, we observe that the training iteration $T = 8$ is a good trade-off between computational cost and performance. In order to clearly show how does our proposed progressive shadow removal method work, we take $T = 8$ and present the visual results for different inference iterations. As shown in Figure 9, after eight iterations, our network can deal with the problem of shadow boundary and restore its original color. Specifically, PRNet mainly recovers the color of the shadow region for the first three iterations, and for the following iterations, it aims to refine the shadow boundary traces.

Analysis of the inference iterations. To provide a more specific view of the progressive shadow removal, we select the training iteration $T = 8$ on SRD dataset (Qu *et al.*, 2017) and adopt different iterations for inference. As shown in Table 11, we can see that the RMSE value is improving continuously in the top 1 ~ 4 iterations, while for the later iterations, the performance slightly increases until iteration 7. Note that when the inference iteration is 7, we can obtain the best performance. After that, the performance remains stable even if we continue to iterate. Through this, we conclude that the best results can be reached during inference by setting the same iteration as training. Therefore, we can avoid extra computational overhead caused by additional iterations.

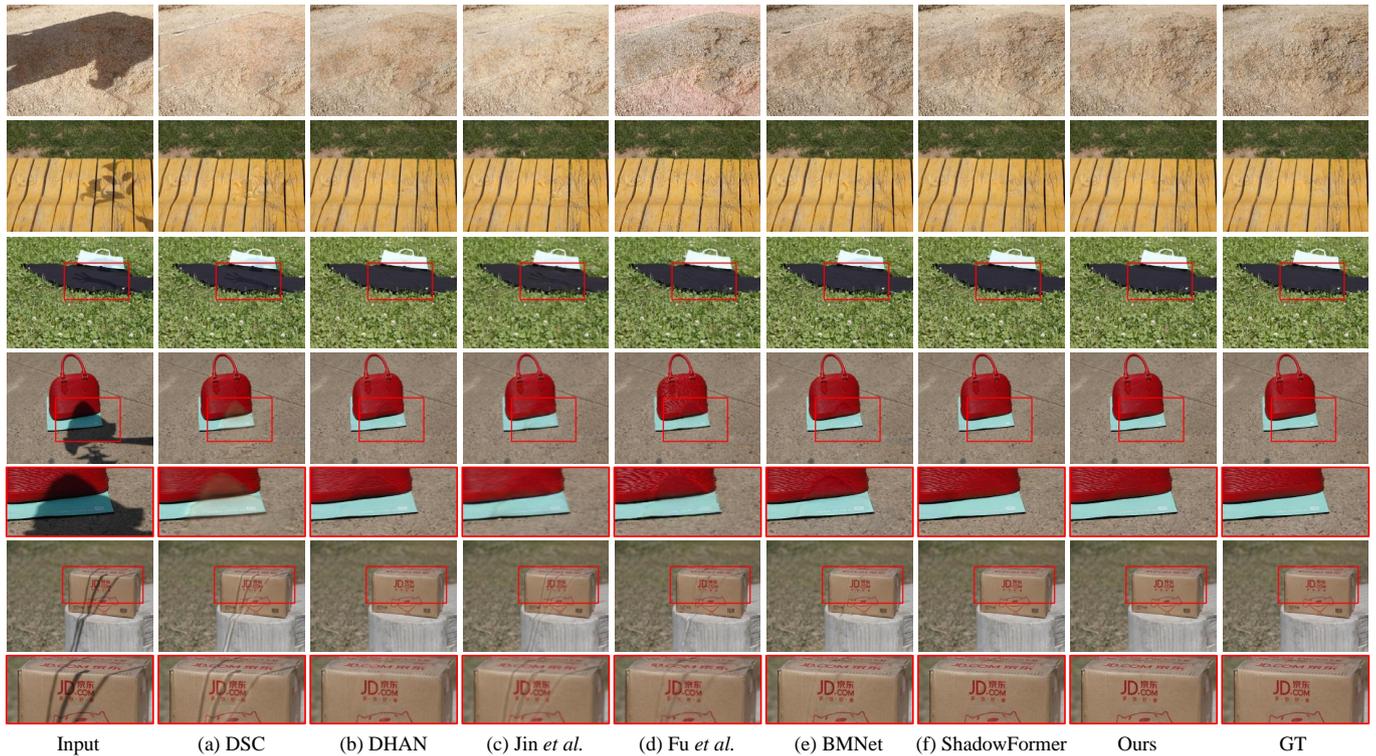


Fig. 8: Visual comparison results of shadow removal on the SRD dataset (Qu et al., 2017). (a) to (f) are the predicted results from SOTA methods: DSC (Hu et al., 2019a), DHAN (Cun et al., 2020), Jin *et al.* (Jin et al., 2021), Fu *et al.* (Fu et al., 2021), BMNet (Zhu et al., 2022a), and ShadowFormer (Guo et al., 2023), respectively.

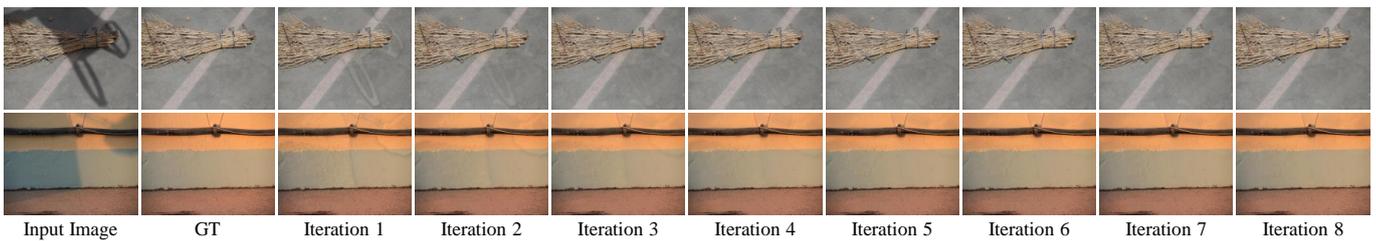


Fig. 9: Visualization of the process of our method for different inference iterations, where we set training iteration $T = 8$. From the results we can conclude that for the first three iterations, our method mainly focuses on recovering the intrinsic color of the shadow region, while for the following iterations, it aims to refine the shadow boundary traces and finally produces more realistic shadow-free images.

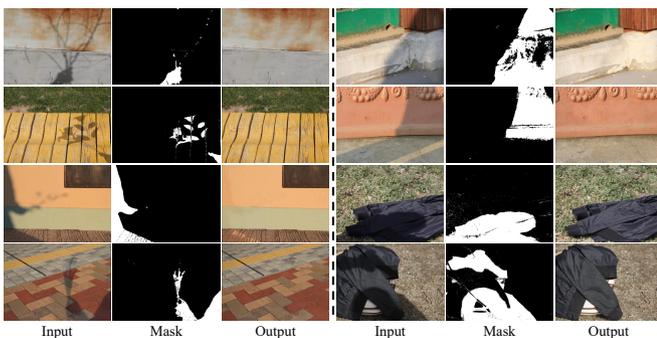


Fig. 10: More visualization results of our method on the SRD dataset (Qu et al., 2017), including small shadows, soft shadows, and dark shadows on black objects.

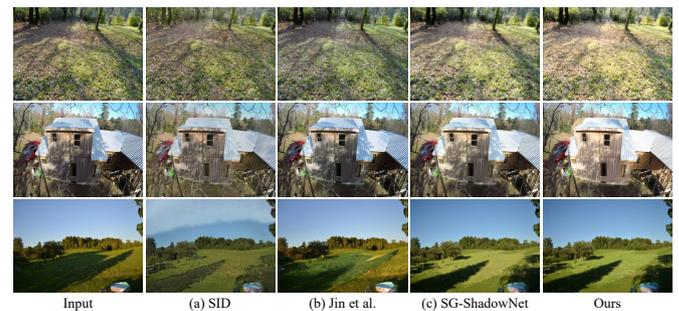


Fig. 11: Visual comparison results of shadow removal on the SBU-Timelapse dataset (Le and Samaras, 2021). (a) to (c) are the predicted results from state-of-the-art methods: SID (Le and Samaras, 2021), Fu *et al.* (Jin et al., 2021), and SG-ShadowNet (Wan et al., 2022), respectively.

Analysis of the output of each iteration. Different from the progressive optical flow prediction (Teed and Deng, 2020) and

other image restoration tasks (Wang et al., 2022; Zamir et al., 2021), which aim to learn residual signals, our update module

Table 11: Ablation study of the number of inference iteration T on SRD dataset (Qu et al., 2017). Empirically, we set $T = 8$ in our paper.

	Iterations of Inference				
RMSE↓	1	2	3	4	5
	4.37	4.16	4.09	4.05	4.02
	6	7	8	9	10
	4.00	3.99	3.99	3.99	3.99

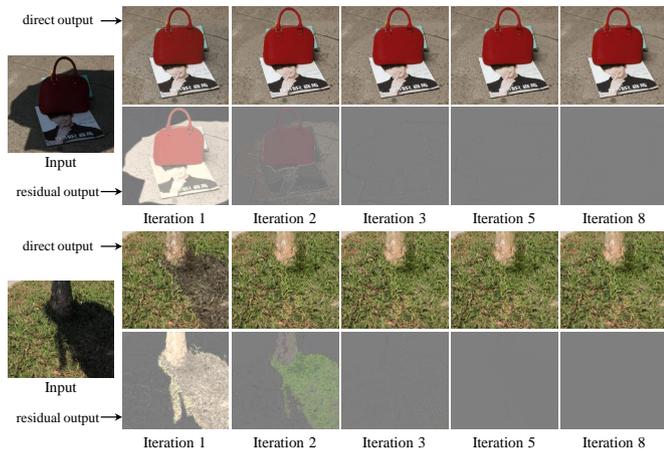


Fig. 12: Visualization of direct output and residual output at different iterations, where we set training iteration $T = 8$.

directly outputs the results in the current stage. Referring to their method, we change the output of each iteration into residual output, and then add it to the original shadow image to obtain the shadow-attenuated image. We report the RMSE value as 4.02 which is close to our result 3.99. As shown in Figure 12, we also provide the visualization results of the residual learning method. We can clearly see that for the first three iterations, the residual images are changing rapidly. In the fifth and eighth iterations, the residual images nearly turn into the same color over all regions, indicating most regions in the image has already been recovered. Through the residual image, we can also come to the same conclusion as before, *i.e.*, the color information is recovered in the first few iterations and the shadow boundary traces are refined in the following iterations.

4.6. More Discussions about Our Method

Analysis of the effectiveness of shadow masks. The shadow only occupies a part of the image. It is crucial to know the location of the shadow because it can provide the shadow information and help the network pay more attention to the shadow region. Shadow detection is another important and challenging task. The results of shadow detection can be used to provide auxiliary information for shadow removal. Here we use the results of the latest shadow detection method, FDRNet (Zhu et al., 2021), as auxiliary information to remove shadows on the ISTD+ dataset (Le and Samaras, 2019). With the detected shadow masks, the de-shadowing performance of our method is slightly decreased to 3.3, but it can still outperform most existing methods in Table 6. Further, the mask for SRD dataset (Qu et al., 2017) is from DHAN (Cun et al., 2020) and the pro-

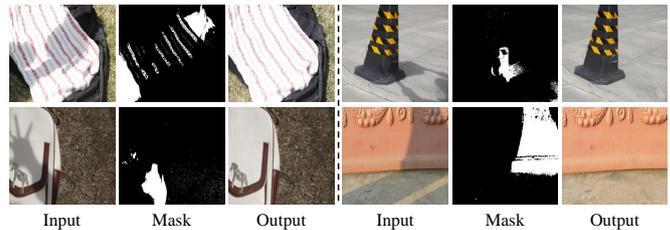


Fig. 13: Visualization results of our method with inaccurate masks, which can also remove the shadow successfully.

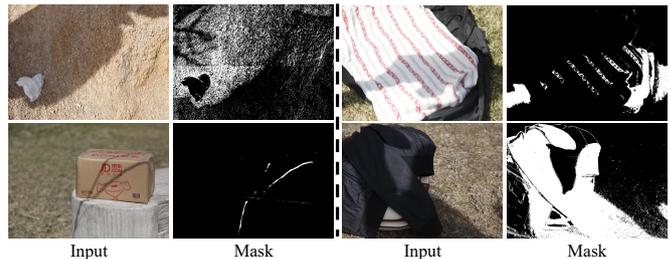


Fig. 14: Inaccurate masks on the SRD dataset (Qu et al., 2017).

vided masks are noisily-annotated. As shown in Figure 13, our method can still robustly remove the shadows even though the shadow masks are inaccurate. We analyze the reason that the training dataset includes some noisy data. As shown in Figure 14, some masks within the dataset are discontinuous or inaccurate. Consequently, during the training of the model on such noisily-annotated data, the model implicitly learns the ability to accommodate inaccurate masks, thereby enhancing its robustness.

Analysis of the parameter-shared update module. Different from the previous method (Ding et al., 2019), our update module is parameter-shared and has no extra parameter cost when we conduct more iterations in both training and testing phases. To evaluate the effectiveness of the parameter-shared update module, we conduct another experiment on SRD dataset (Qu et al., 2017) that the parameter of each update module is independent. In this way, the parameters of the network will increase linearly with the number of iterations. We report that the RMSE value of shared and not shared modules over all the images are 3.99 and 4.00, respectively. The result shows that our parameter-shared model performs similarly to the independent one, but our parameter-shared update module can reduce the number of network parameters which simplifies the structure.

Analysis of the images to be best performance. We calculate how many images can reach the best performance before the pre-defined 8 iterations on SRD dataset (Qu et al., 2017). The results are shown in Figure 15. From the statistical results we can see that most shadow images obtain the best performance in the pre-defined 8 iterations. In addition, about a quarter of images reach the optimal performance before 8 iterations.

Analysis of the loss function. In our experiments, we perform L_1 distance loss over all iterations, and the loss is exponentially increasing through the iteration. Here, we conduct another experiment that only calculates the loss at the last iteration. We report the RMSE result is 4.47 in SRD dataset (Qu et al., 2017),

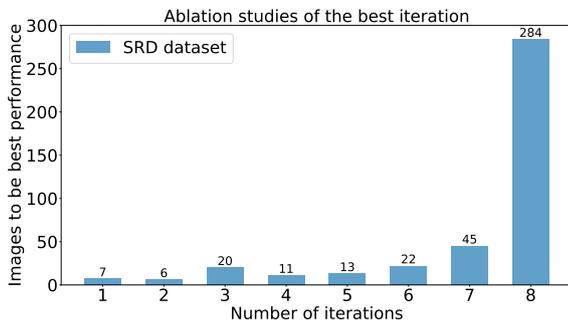


Fig. 15: The number of images to be best performance in the inference stage.

Table 12: Comparison of inference time and parameters with previous methods on the 3090Ti GPU device.

Method	Time(s)	Params(M)
DHAN	0.117	21.8
Param+M+D-Net	0.105	141.2
G2R	0.254	22.8
ShadowFormer	0.127	9.3
Ours iter2	0.090	2.7
Ours iter4	0.135	2.7
Ours iter6	0.181	2.7
Ours iter8	0.226	2.7

which is worse than ours (3.99).

4.7. Computational efficiency

In terms of computational efficiency, we compare our method with previous methods: DHAN (Cun et al., 2020), Param+M+D-Net (H. Le and D. Samaras, 2020), G2R (Liu et al., 2021), and ShadowFormer (Guo et al., 2023). We employ an NVIDIA GTX 3090Ti GPU and test on an image with the resolution of 480×640 . The overall comparison is shown in Table 12, When we set the iteration as two, the inference time is 0.090s per image. In order to achieve better shadow removal results, in this paper, we set the number of iteration as 8, which takes 0.226s to process an image. Compared to the other methods, we argue that the cost is also acceptable. In addition, users can choose the appropriate number of iterations based on their needs or computational resources.

5. Conclusion

In this work, we present a simple Progressive Recurrent Network (PRNet), which aims to address the de-shadowing problem iteratively. The key idea of our method is to apply a parameter-shared GRU-based update module and removes the shadow progressively. The results show that our method restores the color information of the shadow region in the first few iterations and refine to eliminate the shadow boundary traces in the following iterations. The results produced by our method are inconsistent in color and do not suffer from artifacts between shadow and non-shadow regions, resulting in a superior shadow removal performance. Extensive experiments on the

three datasets with both quantitative and qualitative results validate the effectiveness of our method.

In the future, we will explore the potential of our PRNet and further improve its modeling capability. Besides, we also plan to apply this progressive method to other computer vision applications, such as detection and tracking.

References

- Ahn, N., Kang, B., Sohn, K.A., 2018. Image super-resolution via progressive cascading residual network, in: CVPRW, pp. 791–799.
- Cai, Z., Vasconcelos, N., 2018. Cascade R-CNN: Delving into high quality object detection, in: CVPR, pp. 6154–6162.
- Carreira, J., Agrawal, P., Fragkiadaki, K., Malik, J., 2016. Human pose estimation with iterative error feedback, in: CVPR, pp. 4733–4742.
- Chen, Z., Long, C., Zhang, L., Xiao, C., 2021. CANet: A context-aware network for shadow removal, in: ICCV, pp. 4743–4752.
- Cho, K., Van Merriënboer, B., Bahdanau, D., Bengio, Y., 2014a. On the properties of neural machine translation: Encoder-decoder approaches. arXiv preprint arXiv:1409.1259 .
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014b. Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078 .
- Choi, J., Yoo, Y.J., Choi, J.Y., 2010. Adaptive shadow estimator for removing shadow of moving object. CVIU 114, 1017–1029.
- Cucchiara, R., Grana, C., Piccardi, M., Prati, A., 2003. Detecting moving objects, ghosts, and shadows in video streams. TPAMI 25, 1337–1342.
- Cun, X., Pun, C.M., Shi, C., 2020. Towards ghost-free shadow removal via dual hierarchical aggregation network and shadow matting gan, in: AAAI, pp. 10680–10687.
- Ding, B., Long, C., Zhang, L., Xiao, C., 2019. ARGAN: Attentive recurrent generative adversarial network for shadow detection and removal, in: ICCV, pp. 10213–10222.
- Finlayson, G.D., Drew, M.S., Lu, C., 2009. Entropy minimization for shadow removal. IJCV 85, 35–57.
- Finlayson, G.D., Hordley, S.D., Lu, C., Drew, M.S., 2005. On the removal of shadows from images. TPAMI 28, 59–68.
- Fu, L., Zhou, C., Guo, Q., Juefei-Xu, F., Yu, H., Feng, W., Liu, Y., Wang, S., 2021. Auto-exposure fusion for single-image shadow removal, in: CVPR, pp. 10571–10580.
- Gidaris, S., Komodakis, N., 2015. Object detection via a multi-region and semantic segmentation-aware cnn model, in: ICCV, pp. 1134–1142.
- Gregor, K., Danihelka, I., Graves, A., Rezende, D., Wierstra, D., 2015. DRAW: A recurrent neural network for image generation, in: ICML, pp. 1462–1471.
- Gryka, M., Terry, M., Brostow, G.J., 2015. Learning to remove soft shadows. TOG 34, 1–15.
- Guo, L., Huang, S., Liu, D., Cheng, H., Wen, B., 2023. ShadowFormer: Global context helps image shadow removal, in: AAAI.
- Guo, R., Dai, Q., Hoiem, D., 2011. Single-image shadow detection and removal using paired regions, in: CVPR, pp. 2033–2040.
- Guo, R., Dai, Q., Hoiem, D., 2012. Paired regions for shadow detection and removal. TPAMI 35, 2956–2967.
- H. Le and D. Samaras, 2020. From shadow segmentation to shadow removal, in: ECCV, pp. 264–281.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: CVPR, pp. 770–778.
- Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models. NIPS 33.
- Hu, X., Fu, C.W., Zhu, L., Qin, J., Heng, P.A., 2019a. Direction-aware spatial context features for shadow detection and removal. TPAMI 42, 2795–2808.
- Hu, X., Jiang, Y., Fu, C.W., Heng, P.A., 2019b. Mask-ShadowGAN: Learning to remove shadows from unpaired data, in: ICCV, pp. 2472–2481.
- Jin, Y., Li, R., Yang, W., Tan, R.T., 2023. Estimating reflectance layer from a single image: Integrating reflectance guidance and shadow/specular aware learning, in: AAAI, pp. 1069–1077.
- Jin, Y., Sharma, A., Tan, R.T., 2021. DC-ShadowNet: Single-image hard and soft shadow removal using unsupervised domain-classifier guided network, in: ICCV, pp. 5027–5036.

- Jin, Y., Yang, W., Ye, W., Yuan, Y., Tan, R.T., 2022. Des3: Attention-driven self and soft shadow removal using vit similarity and color convergence. arXiv preprint arXiv:2211.08089 .
- Jung, C.R., 2009. Efficient background subtraction and shadow removal for monochromatic video sequences. *TMM* 11, 571–577.
- Karsch, K., Hedau, V., Forsyth, D., Hoiem, D., 2011. Rendering synthetic objects into legacy photographs. *TOG* 30, 1–12.
- Khan, S.H., Bennamoun, M., Sohel, F., Togneri, R., 2015. Automatic shadow detection and removal from a single image. *TPAMI* 38, 431–446.
- Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 .
- Lalonde, J.F., Efros, A.A., Narasimhan, S.G., 2012. Estimating the natural illumination conditions from a single outdoor image. *IJCV* 98, 123–145.
- Le, H., Samaras, D., 2019. Shadow removal via shadow image decomposition, in: *ICCV*, pp. 8578–8587.
- Le, H., Samaras, D., 2021. Physics-based shadow image decomposition for shadow removal. *TPAMI* 44, 9088–9101.
- Levine, M.D., Bhattacharyya, J., 2005. Detecting and removing specularities in facial images. *CVIU* 100, 330–356.
- Liu, C., Zoph, B., Neumann, M., Shlens, J., Hua, W., Li, L.J., Fei-Fei, L., Yuille, A., Huang, J., Murphy, K., 2018. Progressive neural architecture search, in: *ECCV*, pp. 19–34.
- Liu, J., Wang, Q., Fan, H., Li, W., Qu, L., TangMember, Y., 2023. A decoupled multi-task network for shadow removal. *TMM* .
- Liu, Z., Yin, H., Wu, X., Wu, Z., Mi, Y., Wang, S., 2021. From shadow generation to shadow removal, in: *CVPR*, pp. 4927–4936.
- Mikic, I., Cosman, P.C., Kogut, G.T., Trivedi, M.M., 2000. Moving shadow and object detection in traffic scenes, in: *ICPR*, pp. 321–324.
- Nadimi, S., Bhanu, B., 2004. Physical models for moving shadow and object detection in video. *TPAMI* 26, 1079–1087.
- Nair, V., Hinton, G.E., 2010. Rectified linear units improve restricted boltzmann machines, in: *ICML*, pp. 807–814.
- Najibi, M., Rastegari, M., Davis, L.S., 2016. G-CNN: an iterative grid based object detector, in: *CVPR*, pp. 2369–2377.
- Niu, K., Liu, Y., Wu, E., Xing, G., 2022. A boundary-aware network for shadow removal. *TMM* , 1–13.
- Okabe, T., Sato, I., Sato, Y., 2009. Attached shadow coding: Estimating surface normals from shadows under unknown reflectance and lighting conditions, in: *ICCV*, pp. 1693–1700.
- Panagopoulos, A., Samaras, D., Paragios, N., 2009. Robust shadow and illumination estimation using a mixture model, in: *CVPR*, pp. 651–658.
- Qu, L., Tian, J., He, S., Tang, Y., Lau, R.W., 2017. DeshadowNet: A multi-context embedding deep network for shadow removal, in: *CVPR*, pp. 4067–4075.
- Ren, D., Zuo, W., Hu, Q., Zhu, P., Meng, D., 2019. Progressive image deraining networks: A better and simpler baseline, in: *CVPR*, pp. 3937–3946.
- Sanin, A., Sanderson, C., Lovell, B.C., 2010. Improved shadow removal for robust person tracking in surveillance scenarios, in: *ICPR*, pp. 141–144.
- Sekhvat, Y.A., 2016. Privacy preserving cloth try-on using mobile augmented reality. *TMM* 19, 1041–1049.
- Shor, Y., Lischinski, D., 2008. The shadow meets the mask: Pyramid-based shadow removal, in: *CGF*, pp. 577–586.
- Teed, Z., Deng, J., 2020. RAFT: Recurrent all-pairs field transforms for optical flow, in: *ECCV*, pp. 402–419.
- Tokmakov, P., Alahari, K., Schmid, C., 2017. Learning video object segmentation with visual memory, in: *ICCV*, pp. 4481–4490.
- Ulyanov, D., Vedaldi, A., Lempitsky, V., 2016. Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 .
- Vicente, T.F.Y., Hoai, M., Samaras, D., 2017. Leave-one-out kernel optimization for shadow detection and removal. *TPAMI* 40, 682–695.
- Wan, J., Yin, H., Wu, Z., Wu, X., Liu, Y., Wang, S., 2022. Style-guided shadow removal, in: *ECCV*, pp. 361–378.
- Wang, B., Zhao, Y., Chen, C.P., 2019. Moving cast shadows segmentation using illumination invariant feature. *TMM* 22, 2221–2233.
- Wang, J., Li, X., Yang, J., 2018. Stacked conditional generative adversarial networks for jointly learning shadow detection and shadow removal, in: *CVPR*, pp. 1788–1797.
- Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P., 2004. Image quality assessment: from error visibility to structural similarity. *TIP* 13, 600–612.
- Wang, Z., Cun, X., Bao, J., Zhou, W., Liu, J., Li, H., 2022. Uformer: A general u-shaped transformer for image restoration, in: *CVPR*, pp. 17683–17693.
- Xiao, C., She, R., Xiao, D., Ma, K.L., 2013. Fast shadow removal using adaptive multi-scale illumination transfer, in: *CGF*, pp. 207–218.
- Xu, M., Zhu, J., Lv, P., Zhou, B., Tappen, M.F., Ji, R., 2017. Learning-based shadow recognition and removal from monochromatic natural images. *TIP* 26, 5811–5824.
- Yang, Q., Tan, K.H., Ahuja, N., 2012. Shadow removal using bilateral filtering. *TIP* 21, 4361–4368.
- Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H., Shao, L., 2021. Multi-stage progressive image restoration, in: *CVPR*, pp. 14821–14831.
- Zhang, L., Long, C., Zhang, X., Xiao, C., 2020. RIS-GAN: Explore residual and illumination with generative adversarial networks for shadow removal, in: *AAAI*, pp. 12829–12836.
- Zhang, L., Zhang, Q., Xiao, C., 2015. Shadow Remover: Image shadow removal based on illumination recovering optimization. *TIP* 24, 4623–4636.
- Zhang, Q., Zhou, J., Zhu, L., Sun, W., Xiao, C., Zheng, W.S., 2021. Unsupervised intrinsic image decomposition using internal self-similarity cues. *TPAMI* 44, 9669–9686.
- Zhang, W., Zhao, X., Morvan, J.M., Chen, L., 2018. Improving shadow suppression for illumination robust face recognition. *TPAMI* 41, 611–624.
- Zhu, L., Xu, K., Ke, Z., Lau, R.W., 2021. Mitigating intensity bias in shadow detection via feature decomposition and reweighting, in: *CVPR*, pp. 4702–4711.
- Zhu, Y., Huang, J., Fu, X., Zhao, F., Sun, Q., Zha, Z.J., 2022a. Bijective mapping network for shadow removal, in: *CVPR*, pp. 5627–5636.
- Zhu, Y., Xiao, Z., Fang, Y., Fu, X., Xiong, Z., Zha, Z.J., 2022b. Efficient model-driven network for shadow removal, in: *AAAI*, pp. 3635–3643.