
Deep Reflection Prior

Yingda Yin*
Peking University
yingda.yin@gmail.com

Qingnan Fan*
Stanford University
fqchina@gmail.com

Dongdong Chen
Microsoft Cloud AI
cddlyf@gmail.com

Yujie Wang
Shandong University
yujiew.cn@gmail.com

Angelica Aviles-Rivero
University of Cambridge
ai323@cam.ac.uk

Ruoteng Li
National University of Singapore
liruoteng@gmail.com

Carola-Bibiane Schönlieb
University of Cambridge
cbs31@cam.ac.uk

Dani Lischinski
The Hebrew University of Jerusalem
danix3d@gmail.com

Baoquan Chen
Peking University
baoquan.chen@gmail.com

Abstract

Reflections are very common phenomena in our daily photography, which distract people’s attention from the scene behind the glass. The problem of removing reflection artifacts is important but challenging due to its ill-posed nature. Recent learning-based approaches have demonstrated a significant improvement in removing reflections. However, these methods are limited as they require a large number of synthetic reflection/clean image pairs for supervision, at the risk of overfitting in the synthetic image domain. In this paper, we propose a learning-based approach that captures the reflection statistical prior for single image reflection removal. Our algorithm is driven by optimizing the target with joint constraints enhanced between multiple input images during the training stage, but is able to eliminate reflections only from a single input for evaluation. Our framework allows to predict both background and reflection via a one-branch deep neural network, which is implemented by the controllable latent code that indicates either the background or reflection output. We demonstrate superior performance over the state-of-the-art methods on a large range of real-world images.

1 Introduction

Reflection is commonly observed when taking photos through a piece of glass due to the interference between reflected light and the light coming from the background scene. These reflection artifacts significantly degrade the visibility of target scene, and distract users from focusing on it. Therefore, removing reflections is a problem of great interest for many computer vision and graphics applications.

To deal with this problem, the traditional vision-based approaches either leverage the relation from multiple input images [38, 1], or rely on strong priors applied on the background or reflection layer from a single input [17, 41]. Due to the use of hand-crafted features, the traditional methods tend to fail in some challenging cases where prior assumptions are violated. Inspired by the tremendous success

*Equal Contribution

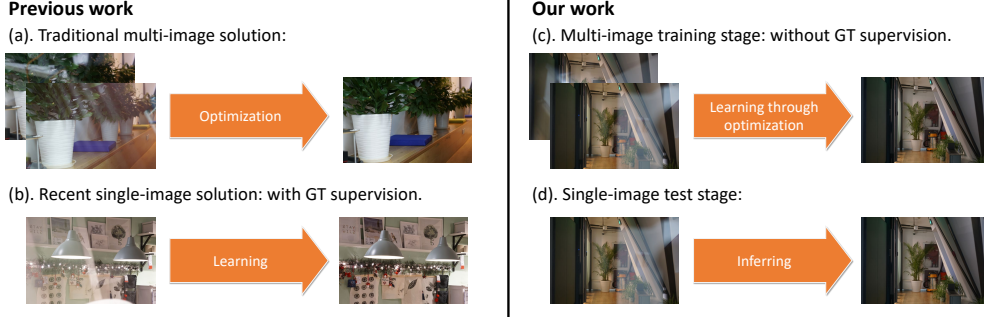


Figure 1: Advantages of our algorithms. (a) The traditional approaches reason the background more theoretically by optimizing an objective function that is limited by the requirement of multi-image inputs. (b) The recent learning methods work by memorizing the reflection data distribution from a single image with direct ground truth supervision, but enlarge the domain gap between synthetic training and real testing images. Our algorithm overcomes their difficulties, and combines their advantages. (c) In the training stage, it learns a deep network to reason the background through the optimization from multiple images without any ground truth supervision, and (d) is able to infer the background only from a single input with a better learned deep reflection prior during evaluation.

of deep learning in many image processing tasks, a variety of recent works [6, 42, 33, 40, 35, 36, 3] have focused on learning-based solutions for reflection removal by taking use of the synthetic reflections as training signals. These techniques have been proven to boost the overall performance significantly and achieve state-of-the-art results. Despite their success, by directly supervising the synthetic labels from a single input, the network may learn limited transferable knowledge and generalizes weakly in the real test data as observed in many previous work [5, 27, 15, 29].

To alleviate such a difficulty, in this paper, we propose a deep learning approach for reflection removal that requires neither ground truth data for supervision, nor hand-crafted priors designed for the non-learning approaches. Instead, it learns a better deep learning prior from a pair of constrained reflection images during training, and is able to reason the reflection removal process on a single input image for evaluation. Our framework builds on an implicit connection between multi-image and single-image methods, as shown in Figure 1.

Our framework is implemented via a single-branch fully convolutional neural network (FCN) with only one input and one output, which is however able to predict the background and reflection layer simultaneously. To this aim, a latent code within the FCN is learned dynamically to control the background or reflection prediction. As reflection removal is under-constrained, we build a joint combination of multiple label-free objective functions to reduce the solution space and help the convergence of the deep network. Optimized with our deep learning framework, we demonstrate superior quantitative and qualitative performance over the state-of-the-art reflection removal approaches on various real reflection cases.

2 Related Work

Non-learning approaches. Since reflection removal is a highly underdetermined problem, approaches leverage on different type of priors to deal with this problem. A set of approaches are oriented to use multiple input images [7, 28, 21, 8, 16, 25, 10, 38, 39, 20, 41]. These techniques require that the scene being captured from different viewpoints. They exploit the motion cue, which assumes the background and foreground motion is usually different due to visual parallax [10, 8, 28, 38]. Some other approaches take a sequence of images using special conditions, such as shooting with polarizer [11, 23, 21], flash/non-flash image pair [1], focus/defocus image pair [22], *etc.* Due to the additional information obtained from multiple images, this problem becomes less ill-posed. However, the special data requirement limits these methods from more practical application scenarios.

Another set of techniques rely on the use of single input image [13, 14, 31, 2, 41]. These approaches leverage the gradient sparsity prior for layer decomposition [26, 12]. As the reflection and background layers are usually at different depth, the reflection plane is very likely to remain outside the depth of

field (DOF) and hence becomes blurry. The approaches of [17, 34] incorporate priors explicitly into the objective function to be optimized. A more realistic physics model takes the thickness of glass into consideration [24].

Learning approaches. [6] proposed the first deep learning approach for reflection removal. It relied on the common assumption of blurry reflection to develop a novel reflection image synthesis method. [42, 36] followed their data generation approach, and incorporates the adversarial loss to mitigate the image degradation issue. [33] acquired a reflection image dataset (RID) which demonstrates mostly sharp and weak reflections, and proposed to unify gradient inference and image reconstruction concurrently into the network design. [35, 19] collected the misaligned or unpaired reflection data to enhance the performance on real reflections. Instead of using single training image, recently [37] proposed to learn from multiple polarized images, along with a novel synthetic data generation approach, which accurately simulates reflections in polarized images. [18] incorporates a dense motion estimation module and online optimization to remove obstructions from multiple frames.

Unlike the previous non-learning or learning approaches, our algorithm takes advantage of methods of both types, that is: we leverage multiple input images in the training stage to ease the ill-posedness, while learning a general parametric solution in the absence of ground truth as direct supervision signal, which requires less hand-crafted features and only a single input during evaluation.

3 Overall Framework

The overall framework is shown in Figure 2. In the training phase, our algorithm takes two reflection images as input (I_1, I_2), which are forwarded into the same reflection removal network independently to predict their corresponding background and reflection layers ($\{\hat{B}_1, \hat{R}_1\}, \{\hat{B}_2, \hat{R}_2\}$). The reflection layers (\hat{R}_1, \hat{R}_2) are composed with every other background layer (\hat{B}_1, \hat{B}_2) to generate a set of reflection images denoted as ($\hat{I}_{11}, \hat{I}_{12}, \hat{I}_{21}, \hat{I}_{22}$).

Although this framework supports separation of background and reflection, without ground truth supervision or specific constraints on the output, this in itself cannot determine the appearance of background and reflection layer as there are arguably countless feasible solutions to the reflection image formation model as shown below.

Reflection Image Synthesis. Given a background (B) and a reflection (R) layer², the corresponding reflection image (I) is a linear combination as,

$$I_i = B_i + R_i \quad (1)$$

where i indicates the image index. Therefore to prepare for multi-image reflection removal, we build an inner connection between the different input reflection images, which reads:

$$B_1 = B_2 = \dots = B_n \quad (2)$$

where n refers to number of reflections images employed in the loss constraints. It means that all the input reflection images share the same background, and differentiate only in the reflection layer. In order to force the network to perform the desired separations, we incorporate an objective combined with four loss functions that take use of the special data configuration to jointly train our deep network.

Naive Reconstruction Loss. Given the decomposed background and reflection layer, the most straightforward supervision can be applied by minimizing the per-pixel difference between the recomposed reflection image and the original input samples as,

$$\mathcal{L}_{naive} = \|\hat{I}_{11} - I_1\|_2^2 + \|\hat{I}_{22} - I_2\|_2^2 \quad (3)$$

where \hat{I}_{ij} is synthesized via $\hat{I}_{ij} = \hat{R}_i + \hat{B}_j$. Despite this objective already reduces the solution space greatly, it is still unclear which one of the two outputs should be background, or likewise reflection.

Cross Reconstruction Loss. Following the definition of our reflection synthesis pipeline, given the predicted \hat{B}_2 and \hat{R}_1 , we should ideally be able to reconstruct the reflection input I_1 . Therefore to

²Their detailed data source is introduced in the following section.

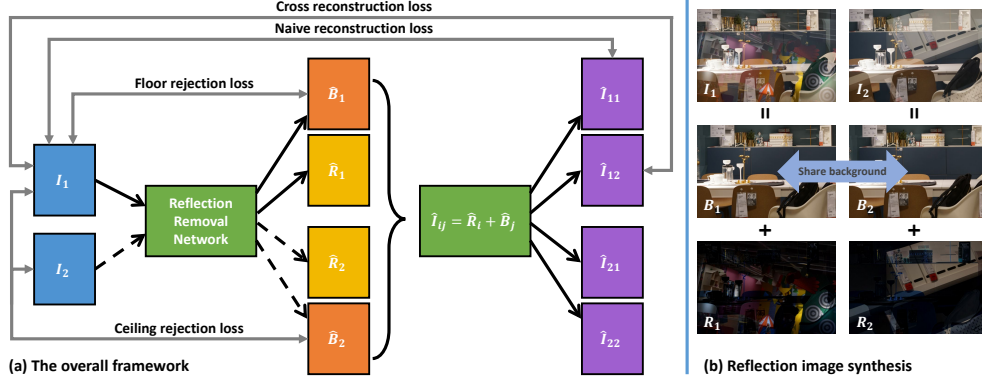


Figure 2: Our overall framework for reflection removal. **(a)** Illustration of the proposed algorithm. The deep network is fed with a pair of reflection images, and decompose them into corresponding background and reflection layers, which are further composed into new reflection images. We append a joint objective of four loss functions on all the predictions for optimization in absence of ground truth. Since the input images are forwarded into the deep network individually, during evaluation, our framework is able to take a single input to generate the clean image. Note the framework is symmetric and for clarity we omit similar lines. **(b)** The corresponding reflection image synthesis pipeline. The two inputs are composed of the same background image.

enhance such a constraint, we further reconstruct the cross compositions of background and reflection layers as

$$\mathcal{L}_{cross} = \|\hat{I}_{12} - I_1\|_2^2 + \|\hat{I}_{21} - I_2\|_2^2 \quad (4)$$

Then the network is forced to disentangle the background from reflection through the reconstruction process by specifying each output component. We demonstrate a toy example in Figure 3. The input contains a rectangle that belongs to the background, and a circle that belongs to the reflection. By simply applying the reconstruction loss, a naive solution of zero energy can be easily explored by pushing all the information from the reflection image (I) into the reflection layer (R), while having the network output a black background layer of all zero values, denoted as $\hat{B}(\mathcal{L}_r)$ in Figure 3.

The naive reconstruction loss and cross reconstruction loss can be combined as

$$\mathcal{L}_{recons} = \mathcal{L}_{naive} + \mathcal{L}_{cross} \quad (5)$$

Floor Rejection Loss. To avoid the background to degenerate, we force the output background to approach the input image. Therefore, we define:

$$\mathcal{L}_{floor} = \|\hat{B}_1 - I_1\|_1 + \|\hat{B}_2 - I_2\|_1 \quad (6)$$

Although this sounds counter-intuitive as the background may dominate the image and hurt the strong reflections, a further constraint of $\hat{B}_1 = \hat{B}_2$ limits the predicted background to only occupy the common part of I_1 and I_2 . As shown in the toy example, we demonstrate three possible solutions for both output backgrounds by simply combining the above two terms, denoted as $\hat{B}(\mathcal{L}_r + \mathcal{L}_f)$. These alternative solutions all reach the condition of zero reconstruction energy, and the same floor rejection energy, which means the strong reflections, the white circle in this case, can also be maintained with the floor rejection loss. We further demonstrate dozens of results of removing both strong and weak reflections in the main paper and supplementary material.

Ceiling Rejection Loss. The combination of the above loss functions yields a better disentanglement, which is however still insufficient and ambiguous to generate a purely clean background as shown in the toy example. As the two input reflection images are highly correlated, we take better advantage of their relationship by appending a ceiling rejection loss which prevents the background intensity from exceeding each of the two input images, reads as,

$$f(\hat{B}_1, I_1, m) = \begin{cases} \|\hat{B}_{1,m} - I_{1,m}\|_1 & \hat{B}_{1,m} > I_{1,m}, \\ 0 & \text{otherwise;} \end{cases} \quad (7)$$

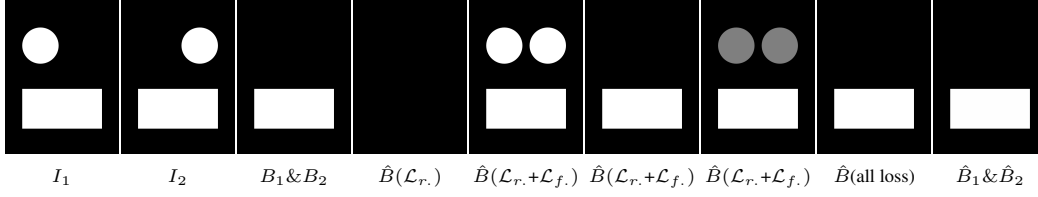


Figure 3: Analysis of the three loss terms for a toy training example. \mathcal{L}_r and \mathcal{L}_f are short for \mathcal{L}_{recons} and \mathcal{L}_{floor} separately. Please refer to the text for the explanation.

$$\mathcal{L}_{ceiling} = \sum_m (f(\hat{B}_1, I_1, m) + f(\hat{B}_1, I_2, m) + f(\hat{B}_2, I_1, m) + f(\hat{B}_2, I_2, m)) \quad (8)$$

where m indicates each image pixel. The ceiling rejection loss is a summation of f among the whole image region. It sets an upper bound for each background prediction, and punishes any background pixel whose intensity value is larger than any input. The ceiling rejection loss helps solve the ambiguity of multiple background solutions for the toy example, and leaves only the true answer as shown in Figure 3 (all loss). To justify our analysis, we further jointly optimize the three objectives on this toy example, and obtains the same result ($\hat{B}_1 \& \hat{B}_2$) as the analytical one.

The overall objective function is a weighted combination of the above four loss terms, defined as

$$\mathcal{L}_{unsuper} = \lambda_1 \mathcal{L}_{recons} + \lambda_2 \mathcal{L}_{floor} + \mathcal{L}_{ceiling} \quad (9)$$

where the specific value for λ_1 and λ_2 are addressed in the experiment section. Note the network training procedure is essentially a continuous coordination and competition process of all the loss terms. The reconstruction loss sets the lower bound of the background, and the floor and ceiling rejection loss set the upper bound of the background. For more ablation study about the designed loss terms, please refer to the supplementary material.

Multi-input or single-input? Please note our algorithm only takes advantage of the supervision on multiple outputs, and it still takes single image as input for each forward pass. Hence it is fairly compared with previous single-image approaches via training on the same synthetic data as shown in the experimental session.

4 Reflection Removal Network

In this section, we introduce the details of our reflection removal network that unifies background and reflection prediction within a one-branch pipeline. Our network structure is shown in Figure 4.

For our specialized framework, one network of multiple output branches or even two independent neural networks are usually required for both background and reflection prediction. In this paper, we develop a more elegant way to couple the background and reflection within a single network, and provide a new understanding of the learning based reflection removal approach.

Following the basic network structure in [6], our neural network is a fully convolutional neural network (FCN), which contains 32 convolution layers, where the middle 26 ones are organized into 13 residual blocks to accelerate convergence. All the convolution layers use 3×3 kernels with 64 output channels, and are followed by instance normalization [30] and ReLU layer except for the last one. The input and output of the network are both a three channel color image.

However, training such a one-branch fully convolutional neural network is challenging, as it does not support two concurrent outputs for a given single input image. To transform this network in a more suitable way, we introduce the latent code that explicitly learns the background or reflection information.

This latent code is implemented as the learnable parameters (scale and shift) in the first instance normalization layer, which is decoupled from the network as an independent parameterized vector. During the training process, two latent codes that represent the background and reflection are independently learned. Given a single input image, the output of the reflection removal network is determined by switching the two controllable latent codes. As the instance normalization layer takes 64 feature maps as input, the latent code is only defined in a tiny vector of size 128. Benefited

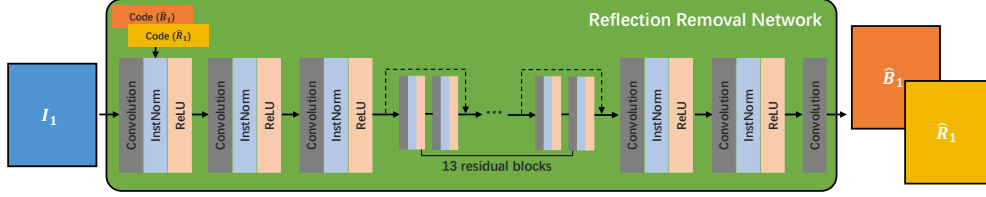


Figure 4: Our detailed reflection removal network structure. This is a one-branch fully-convolutional neural network. To enable both background and reflection prediction within such a network, we introduce the latent code defined by the learnable parameters in the first instance normalization layer to encode the output image information. To be specific, two latent codes that represent background and reflection are learned separately, and fed to the network independently to generate the corresponding output.

from such a design, the specific background or reflection information are all encoded into these 128 learnable parameters, while all the other network weights are shared for them.

Experimentally, we observe the best performance by selecting the first instance normalization layer, since in this manner, most subsequent layers can leverage the latent code to better differentiate the background and reflection.

Note there is only one reflection removal network that exists in our full pipeline, and all the different inputs are fed through the same network individually. Such a design support further extension to three or more inputs without modifying the architecture. In this paper, we take only two reflection images as input for simplicity. Benefited from such a design, our framework can also be naturally switched for single image reflection removal during evaluation. We demonstrate more inspiring experiments regarding the learned latent code in the supplementary material.

5 Experiments

Due to the lack of large-scale high-quality reflection image dataset, the learning based approaches mostly synthesize reflection images for training. Currently, there are two mainstream data synthesis methods, which mainly differentiate in the logic of computing the reflection layer: (1) One of the methods is developed by [6] and followed by [42, 35]. Their reflection layer is generated via subtraction and clipping operations. It reflects some physical properties observed in natural scenes. The generated images tend to contain strong and blurry reflections. (2) The other method is the common linear additive image mixing model used by [33, 40]. Their reflection layer is only scaled by a small factor for direct composition. The resultant reflections are weaker yet sharper.

In order to conduct a fair comparison with corresponding previous works, and as these approaches have their specialties in image formation process and bias towards target images, we split our experiments by different training data generated from the above two synthesis methods for our algorithm³.

5.1 Subtraction and Clipping Image Model

Implementation Details. Following [6], we collect around 17,000 natural images from PASCAL VOC dataset as the data source of the image synthesis approach, while 5% of them are treated as test data and the left are training data. During the network training process, we randomly sampled two images from the training split to synthesize the reflection images online. We directly adopt [6]’s method to generate our training data. Our framework is implemented in PyTorch. Our deep network is optimized by Adam with mini-batch size as 8. Some specifics about our training setting: initial learning rate (0.01), training epoch number (60), λ_1 (15), λ_2 (20).

Results. To evaluate our algorithm on the real images, we randomly collect around 30 reflection images that feature strong and blurry reflections, and conduct a user study among 40 users. They are asked to pick the best visual result among 8 approaches for each example, and we evaluate the user

³Note [36] learns the alpha mask in the linear composition model for reflection synthesis, which we also experiment with, but doesn’t work better than the image model in our problem setting.



Figure 5: Visual comparison on the strong real-world reflection images from [6]. These images are featured by subtraction and clipping image model. Ours* refers to the version with ground truth supervision.



Figure 6: Visual comparison on strong real-world reflection images collected by ourselves. These images are featured by subtraction and clipping image model. Ours* refers to the version with ground truth supervision.

selections for each method in Table 1⁴. We also experiment by replacing the proposed loss functions with the supervision losses for training on the ground truth data following [42]. The corresponding qualitative results are demonstrated in Figure 5 ([6]’s test images) and 6 (self-collected images).

For both numerical and visual results, our proposed pipeline outperforms the others. Note on the visual side, our algorithm learns to remove more reflections and generate cleaner background, even compared with [42, 6, 35] that share the similar data generation approach as ours. We further test our model on the synthetic test data, and interestingly observe degraded performance compared to our fully supervised alternative. This is due to the fact that our algorithm enables deeper understanding of the theoretical reflection removal process by optimizing the proposed objective function. Hence it overfits less on the synthetic training data, and generalizes better on the real cases.

5.2 Linear Addition Image Model

Implementation Details. Similarly to the subtraction and clipping image model, we directly take their train/test split in the PASCAL VOC dataset as our data source. For fair comparison with the previous work, we follow [40]’s approach to linearly mix two natural images with a constant scale factor to synthesize the reflection image, where the scale weight is within [0.6, 0.8] for the reflection

⁴Due to inaccessibility to [33]’s codes or trained model, their results are not present for the subtraction and clipping image model.

Data source	Collected Real	SIR ²			
Error metric	%	PSNR	SSIM	SSIM _r	SSIM _{r,s}
Input	-	26.19	0.916	0.822	0.822
Wan <i>et al.</i> 2018 [33]	-	22.16	0.828	0.821	0.829
Wen <i>et al.</i> 2019 [36]	0.28	21.23	0.854	0.779	0.798
Yang <i>et al.</i> 2018 [40]	9.07	22.25	0.853	0.832	0.834
Li <i>et al.</i> 2014 [17]	2.97	18.87	0.782	0.761	0.790
Fan <i>et al.</i> 2017 [6]	8.50	20.97	0.839	0.775	0.805
Zhang <i>et al.</i> 2018 [42]	19.83	21.20	0.862	0.807	0.817
Wei <i>et al.</i> 2019 [35]	6.52	24.14	0.879	0.819	0.826
Ours*	18.13	22.45	0.778	0.765	0.773
Ours	34.70	22.72	0.870	0.833	0.841

Table 1: SOTA comparison on the user study of strong real reflections with our subtraction model (left), and image quality evaluation of the wild scene images of SIR² dataset with our linear addition model (right). Ours* refers to the version with ground truth supervision. The numbers in red and blue are the best and second best results respectively.

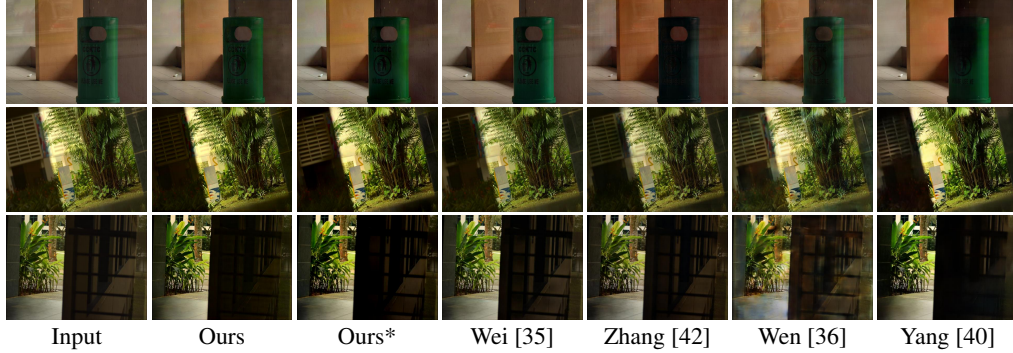


Figure 7: Visual comparison on the wild scene images of SIR² dataset. Ours* refers to the version with ground truth supervision. **Zoom in to see the details.**

layer. We adopt the similar training setting as in the subtraction and clipping image model. The difference lies in: training epoch number (30), λ_1 (80), λ_2 (50). The parameter change is caused by the significantly different data generation approach.

Results. SIR² [32] is a well-known reflection removal benchmark that contains relatively weak and sharp reflections. For a fair comparison, we follow SIR²’s subsequent work [33], led by the same researchers, to utilize the linear addition image model to prepare for the training data, and take the wild scene images in SIR² dataset.

We demonstrate the visual results in Figure 7. The reflections in SIR² dataset are weaker yet sharper, compared to the previous section. However, our approach is still able to generate cleaner background, and also preserves better image structure and color.

The numerical performances are listed in Table 1. We first test on the common PSNR and SSIM error metric. Our approach overwhelms most previous approaches except for [35]. However, interestingly we observe the input reflection images achieve the best performance, as also mentioned in [35]. This is due to the fact that most algorithms tend to touch the image as a whole, and accidentally cause errors in the reflection-free regions. Therefore, [33] manually labels the reflection regions in the SIR² dataset and proposes SSIM_r to evaluate the SSIM error only within the region of interest. Our algorithm outperforms all the existing approaches under this metric, including the input image.

We further analyze the statistics of the SIR² data, and observe that 39% of the background pixels are bigger than the input for intensity comparison. This is mainly because that the input and background images in the SIR² benchmark are both captured. Without the interception of reflections, the resultant background intensity is possible to become larger. However, such a fact is contradictory with the common assumption that background is a subset of input image in many previous work, as shown in Equation 1. To tackle this issue, we propose to minimize the channel-wise color difference between

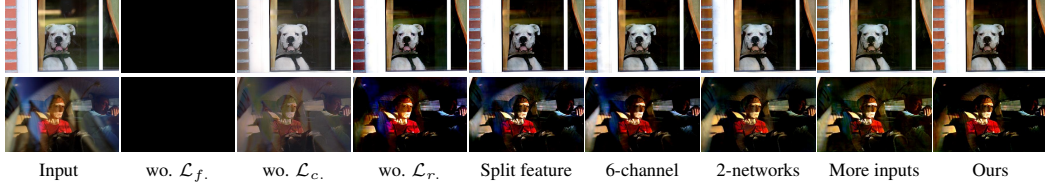


Figure 8: Visual comparison of different alternatives of our algorithm. \mathcal{L}_f , \mathcal{L}_c and \mathcal{L}_r are short for \mathcal{L}_{floor} , $\mathcal{L}_{ceiling}$ and \mathcal{L}_{recons} separately. **Zoom in to see the difference.**

	wo. \mathcal{L}_f	wo. \mathcal{L}_c	wo. \mathcal{L}_r	Split feature	6-channel	2-networks	More inputs	Ours
PSNR	6.20	16.44	16.97	16.98	17.06	17.76	17.35	17.92
SSIM	0.012	0.804	0.796	0.824	0.813	0.823	0.841	0.842

Table 2: Ablation study on the synthetic test data generated by the subtraction and clipping image model. \mathcal{L}_f , \mathcal{L}_c and \mathcal{L}_r are short for \mathcal{L}_{floor} , $\mathcal{L}_{ceiling}$ and \mathcal{L}_{recons} separately.

the output and ground truth background, similar to [9]. Specifically,

$$SSIM_{rs} = SSIM_r(B, \hat{\alpha}\hat{B}) \quad (10)$$

where $\hat{\alpha} = \argmin_{\alpha} \|B - \alpha\hat{B}\|^2$. When tested under this error metric, our algorithm is more superior than others.

We further test our model on the synthetic data, and also observe degradation compared to the fully supervised version. This justifies again that our algorithm overfits less on the synthetic training domain, and generalizes better in the real data.

6 Ablation Study and Analysis

6.1 How important is each loss?

To analyze the effectiveness of each loss function, we train our network with three alternatives, by removing (1) the floor rejection loss (wo. \mathcal{L}_f), (2) the ceiling rejection loss (wo. \mathcal{L}_c), and (3) the reconstruction loss (wo. \mathcal{L}_r).

The results are shown in Figure 8 and Table 2. (1) When training without the floor rejection loss, the network learns a naive solution by simply setting the background to be completely black, and achieves the global minimum of zero energy for the combined objective of reconstruction and ceiling rejection loss. (2) When removing the supervision of the ceiling rejection loss, the learned background tends to approach both input images, and confuses the network with the attempt to maintain reflections without punishment. Hence in the visual results, some obvious reflections remain and color degradation effect appears. (3) When removing the reconstruction loss, the predicted reflections are left without any supervision, and the two predicted backgrounds have no constraints to be the same. Hence the multi-inputs constraint is not fully leveraged in this condition. In the according visual results, some reflections are not removed correctly. On the other side, justification of the reconstruction loss also proves the importance of separate reflection and background estimations, instead of computing the reflection by simply subtracting input and predicted background. All in all, compared to the three variants of our algorithm, our full pipeline achieves the best numerical and visual results.

6.2 What does the latent code learn?

In our implementation, the latent code encodes the unique information for background or reflection, while the rest of network weights are shared. It means the layer separation task is conducted by the specific instance normalization layer where our learned latent code is embedded. Then we are interested in the difference between the feature maps generated by this instance normalization layer with the background and reflection latent codes.

Analysis of Active Features. Inspired by the sparsity property of convolution features [4], we study the active feature channel for background and reflection. To achieve it, we deactivate each channel

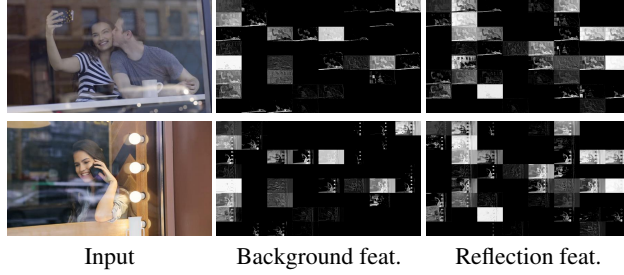


Figure 9: Visualization of background and reflection features after the first instance normalization layer where the latent code is embedded. 64 feature maps are allocated into an 8×8 block for visualization.

of the feature map by setting its value to all zero, and compute the difference (MSE) between the manipulated and untouched output image. We conduct such an experiment among 100 real reflection images. As long as the average MSE is larger than 0, this feature channel is considered “active” for output generation.

As a result, among all 64 feature maps, interestingly we find that only 35 channels are active for background, and 40 channels for reflection, while eliminating all the other channels does not influence the output generation. Moreover, the background and reflection share only 20 common active channels. However, active channel is not always effective. When filtering out the “ineffective” ones (average MSE < 1), there are only 11 overlapped channels. It means that the learned latent code differentiates background and reflection by allocating different active feature maps.

To further justify the above observation, we visualize the normalized feature maps of two randomly picked real images in Figure 9. As shown in this figure, many feature maps are totally blank while only part of them contain color intensities. And the active background and reflection feature maps are mostly different from each other. Similar phenomenon is consistently observed in many more examples.

Feature Splitting. Inspired by this observation, we learn to unify background and reflection predictions into one network by explicitly splitting features, instead of learning latent code. Specifically, we assign half of the 64 feature maps after the first convolution layer for background by setting the other half features to all zero, and vice versa for reflection. Therefore, during training, latent code is not learned, but replaced by allocating features explicitly.

We train such a network using the subtraction and clipping image model, and demonstrate its visual results in Figure 8. Compared to our result, this feature splitting network removes much less reflections. We believe this may be because the half-half splitting of features does not reflect the ideal feature distribution, while our approach is able to adaptively learn the best feature amount for both background and reflection. A corresponding numerical comparison is shown in Table 2.

6.3 Replacement for latent code?

Instead of learning latent code within one single network, one alternative solution is to directly learn two networks for background and reflection separately. As shown in Figure 8, this 2-network solution shows very similar visual results, in comparison to our proposed one-network solution. It demonstrates our efficiency by learning a single network to achieve similar performance with two networks.

Another naive replacement is to directly learn a single network with a 6-channel output, where the half are for background and the other half are for reflection. The predicted background from this trained network still contains many obvious reflections, which do not exist in our result. It justifies the effectiveness of our specific network design.

6.4 More inputs to the network?

During the implementation of our specialized framework, we take only two inputs for simplicity. But theoretically, our framework is able to take unlimited inputs. To explore its performance, we

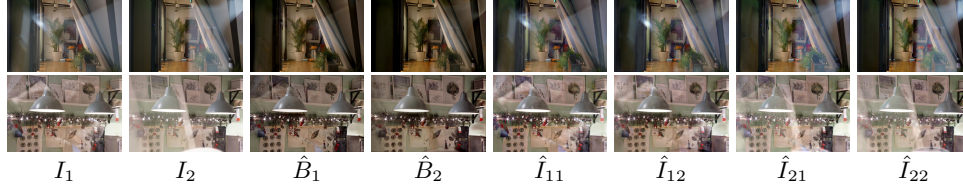


Figure 10: Visual results of our algorithm learned on the real label-free reflection data.

experiment by feeding three reflection images to our neural network in the training phase, whose results are shown in Figure 8.

It also demonstrates similar visual results to ours. It means that our two-input pipeline learns sufficiently good reflection removal effects, while feeding more inputs does not have significant impact in the output. We also fed more input images to the network, however it does not make substantial difference in the result.

7 Conclusion

In this paper, we propose a learning-based approach for reflection removal problem. Our deep network learns to optimize a combination of multiple objective functions, which take advantage of relationship between multiple input images during the training phase, and transfers the learned deep reflection prior in the evaluation stage with only a single input image.

8 Appendix

A. Robustness of hyper-parameters?

In order to test the robustness of hyper-parameters in the objective function, we test the numerical results on the SIR² dataset by randomly sampling λ_1 between [50, 110], and λ_2 between [30, 70], whose values are set as 80 and 50 by default. The resultant PSNR is between [22.523, 23.120], and SSIM is between [0.829, 0.864], which shows robustness to the change of hyper-parameters.

B. What if training on real images?

Our algorithm can be potentially improved with real image training pairs. However, due to the difficulty of collecting large-scale real reflection data, the experiments we conduct until now all leverage the ground truth to synthesize training reflections. In order to explore the possibility of our algorithm on the real label-free reflection training data, we capture some image sequences where the background holds mostly steady and the reflection differs. Since the real image number is limited, our network is trained on each image pair iteratively, and generates the corresponding background (\hat{B}_1, \hat{B}_2) in Figure 10. In this case, even without ground truth to synthesize perfect training pairs, our algorithm is still able to remove reflections very well. More results are given in the supplemental materials.

In Figure 10, we also further demonstrate the reconstructed reflection images ($\hat{I}_{11}, \hat{I}_{12}, \hat{I}_{21}, \hat{I}_{22}$), which are visually consistent with the corresponding input image pairs (I_1, I_2).

C. Joint training of different image mixing models

Instead of training separate networks for different image models, in this section, we experiment by incorporating both the subtraction and clipping image model, and linear addition image model into joint training of a single network. To achieve such a goal, we learn four different latent codes, which correspond to the background and reflection for either of the two image models. We mix the image models randomly into the same mini-batch to train the neural network.

We test this jointly trained network on the SIR² dataset featured by the linear addition image model, and achieve 0.814 and 0.822 for SSIM_r and SSIM_{rs} metric separately, which both lag behind the

separately trained network (0.833 and 0.841). Regarding the strong reflection images featured by the subtraction and clipping image model, we demonstrate some visual comparisons with the separately trained network in Figure 11, whose results are also consistent with the numerical comparison. This is probably caused by the confusion of different training data, and leads to degraded performance on both test data.

D. Joint training with both the fully supervised loss and the proposed loss

We further experiment by training our network with both the fully supervised loss using the ground truth label and our proposed loss not using the ground truth label, but don't observe better visual and numerical results, compared to our original pipeline.



Figure 11: Visual comparison between the results obtained by jointly or separately training on different image mixing models.

References

- [1] Amit Agrawal, Ramesh Raskar, Shree K Nayar, and Yuanzhen Li. Removing photography artifacts using gradient projection and flash-exposure sampling. *ACM Transactions on Graphics (TOG)*, 24(3):828–835, 2005.
- [2] Nikolaos Arvanitopoulos, Radhakrishna Achanta, and Sabine Susstrunk. Single image reflection suppression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4498–4506, 2017.
- [3] Yakun Chang and Cheolkon Jung. Single image reflection removal using convolutional neural networks. *IEEE Transactions on Image Processing*, 28(4):1954–1966, 2019.
- [4] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1897–1906, 2017.
- [5] Yuhua Chen, Wen Li, Xiaoran Chen, and Luc Van Gool. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1841–1850, 2019.
- [6] Qingnan Fan, Jiaolong Yang, Gang Hua, Baoquan Chen, and David Wipf. A generic deep architecture for single image reflection removal and image smoothing. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3238–3247, 2017.
- [7] Hany Farid and Edward H Adelson. Separating reflections and lighting using independent components analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 262–267, 1999.
- [8] Kun Gai, Zhenwei Shi, and Changshui Zhang. Blind separation of superimposed moving images using image statistics. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 34(1):19–32, 2012.
- [9] Roger Grosse, Micah K Johnson, Edward H Adelson, and William T Freeman. Ground truth dataset and baseline evaluations for intrinsic image algorithms. In *2009 IEEE 12th International Conference on Computer Vision*, pages 2335–2342. IEEE, 2009.
- [10] Xiaojie Guo, Xiaochun Cao, and Yi Ma. Robust separation of reflection from multiple images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2187–2194, 2014.
- [11] Naejin Kong, Yu-Wing Tai, and Joseph S Shin. A physically-based approach to reflection separation: from physical modeling to constrained optimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 36(2):209–221, 2014.
- [12] Anat Levin and Yair Weiss. User assisted separation of reflections from a single image using a sparsity prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 29(9), 2007.
- [13] Anat Levin, Assaf Zomet, and Yair Weiss. Learning to perceive transparency from the statistics of natural scenes. In *Advances in Neural Information Processing Systems*, pages 1271–1278, 2003.
- [14] Anat Levin, Assaf Zomet, and Yair Weiss. Separating reflections from a single image using local features. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 1, pages I–I. IEEE, 2004.
- [15] Siyuan Li, Iago Breno Araujo, Wenqi Ren, Zhangyang Wang, Eric K Tokuda, Roberto Hirata Junior, Roberto Cesar-Junior, Jiawan Zhang, Xiaojie Guo, and Xiaochun Cao. Single image deraining: A comprehensive benchmark analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3838–3847, 2019.
- [16] Yu Li and Michael S Brown. Exploiting reflection change for automatic reflection removal. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2432–2439, 2013.

- [17] Yu Li and Michael S Brown. Single image layer separation using relative smoothness. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2752–2759, 2014.
- [18] Yu-Lun Liu, Wei-Sheng Lai, Ming-Hsuan Yang, Yung-Yu Chuang, and Jia-Bin Huang. Learning to see through obstructions. 2020.
- [19] Daiqian Ma, Renjie Wan, Boxin Shi, Alex C Kot, and Ling-Yu Duan. Learning to jointly generate and separate reflections. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2444–2452, 2019.
- [20] Ajay Nandoriya, Mohamed Elgharib, Changil Kim, Mohamed Hefeeda, and Wojciech Matusik. Video reflection removal through spatio-temporal optimization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2411–2419, 2017.
- [21] Bernard Sarel and Michal Irani. Separating transparent layers through layer information exchange. In *European Conference on Computer Vision (ECCV)*, pages 328–341, 2004.
- [22] Yoav Y Schechner, Nahum Kiryati, and Ronen Basri. Separation of transparent layers using focus. *International Journal of Computer Vision (IJCV)*, 39(1):25–39, 2000.
- [23] Yoav Y Schechner, Joseph Shamir, and Nahum Kiryati. Polarization and statistical analysis of scenes containing a semireflector. *JOSA A*, 17(2):276–284, 2000.
- [24] YiChang Shih, Dilip Krishnan, Fredo Durand, and William T Freeman. Reflection removal using ghosting cues. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3193–3201, 2015.
- [25] Sudipta N Sinha, Johannes Kopf, Michael Goesele, Daniel Scharstein, and Richard Szeliski. Image-based rendering for scenes with reflections. *ACM Transactions on Graphics (TOG)*, 31(4):100–1, 2012.
- [26] Ofer Springer and Yair Weiss. Reflection separation using guided annotation. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages 1192–1196. IEEE, 2017.
- [27] Ruoqi Sun, Xinge Zhu, Chongruo Wu, Chen Huang, Jianping Shi, and Lizhuang Ma. Not all areas are equal: Transfer learning for semantic segmentation via hierarchical region selection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4360–4369, 2019.
- [28] Richard Szeliski, Shai Avidan, and P Anandan. Layer extraction from multiple images containing reflections and transparency. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 246–253, 2000.
- [29] Jonathan Tremblay, Aayush Prakash, David Acuna, Mark Brophy, Varun Jampani, Cem Anil, Thang To, Eric Cameracci, Shaad Boochoon, and Stan Birchfield. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 969–977, 2018.
- [30] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [31] Renjie Wan, Boxin Shi, Ling-Yu Duan, Ah-Hwee Tan, Wen Gao, and Alex C Kot. Region-aware reflection removal with unified content and gradient priors. *IEEE Transactions on Image Processing*, 27(6):2927–2941, 2018.
- [32] Renjie Wan, Boxin Shi, Ling-Yu Duan, Ah-Hwee Tan, and Alex C Kot. Benchmarking single-image reflection removal algorithms. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3922–3930, 2017.
- [33] Renjie Wan, Boxin Shi, Ling-Yu Duan, Ah-Hwee Tan, and Alex C Kot. Crnn: multi-scale guided concurrent reflection removal network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4777–4785, 2018.
- [34] Renjie Wan, Boxin Shi, Tan Ah Hwee, and Alex C Kot. Depth of field guided reflection removal. In *Image Processing (ICIP), 2016 IEEE International Conference on*, pages 21–25. IEEE, 2016.

- [35] Kaixuan Wei, Jiaolong Yang, Ying Fu, David Wipf, and Hua Huang. Single image reflection removal exploiting misaligned training data and network enhancements. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8178–8187, 2019.
- [36] Qiang Wen, Yinjie Tan, Jing Qin, Wenxi Liu, Guoqiang Han, and Shengfeng He. Single image reflection removal beyond linearity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3771–3779, 2019.
- [37] Patrick Wieschollek, Orazio Gallo, Jinwei Gu, and Jan Kautz. Separating reflection and transmission images in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 89–104, 2018.
- [38] Tianfan Xue, Michael Rubinstein, Ce Liu, and William T Freeman. A computational approach for obstruction-free photography. *ACM Transactions on Graphics (TOG)*, 34(4):79, 2015.
- [39] Jiaolong Yang, Hongdong Li, Yuchao Dai, and Robby T Tan. Robust optical flow estimation of double-layer images under transparency or reflection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1410–1419, 2016.
- [40] Jie Yang, Dong Gong, Lingqiao Liu, and Qinfeng Shi. Seeing deeply and bidirectionally: A deep learning approach for single image reflection removal. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 654–669, 2018.
- [41] Yang Yang, Wenye Ma, Yin Zheng, Jian-Feng Cai, and Weiyu Xu. Fast single image reflection suppression via convex optimization. *arXiv preprint arXiv:1903.03889*, 2019.
- [42] Xuaner Zhang, Ren Ng, and Qifeng Chen. Single image reflection separation with perceptual losses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4786–4794, 2018.