# CANet: A Context-Aware Network for Shadow Removal

Zipei Chen[1], Chengjiang Long[2*], Ling Zhang[3], Chunxia Xiao[1*†]

[1]School of Computer Science, Wuhan University, Wuhan, Hubei, China
[2]JD Finance America Corporation, Mountain View, CA, USA
[3]Wuhan University of Science and Technology, Wuhan, Hubei, China

czpp19@whu.edu.cn, cjfykx@gmail.com, zhling@wust.edu.cn, cxxiao@whu.edu.cn

## Abstract

*In this paper, we propose a novel two-stage context-aware network named CANet for shadow removal, in which the contextual information from non-shadow regions is transferred to shadow regions at the embedded feature spaces. At Stage-I, we propose a contextual patch matching (CPM) module to generate a set of potential matching pairs of shadow and non-shadow patches. Combined with the potential contextual relationships between shadow and non-shadow regions, our well-designed contextual feature transfer (CFT) mechanism can transfer contextual information from non-shadow to shadow regions at different scales. With the reconstructed feature maps, we remove shadows at L and A/B channels separately. At Stage-II, we use an encoder-decoder to refine current results and generate the final shadow removal results. We evaluate our proposed CANet on two benchmark datasets and some real-world shadow images with complex scenes. Extensive experimental results strongly demonstrate the efficacy of our proposed CANet and exhibit superior performance to state-of-the-arts. Our source code is available at* https://github.com/Zipei-Chen/CANet.

## 1. Introduction

Shadow is a natural phenomenon appearing when the light is partially or completely blocked. As a fundamental challenge in the field of computer vision, the existence of shadow in images or videos inevitably degrades the accuracy and effectiveness of general application tasks such as intrinsic image decomposition [21, 10], visual recognition [25, 17, 24, 14], object detection and tracking [28, 1, 2], trajectory prediction [27, 33], single image super-resolution [40, 39] and image captioning [6]. Therefore, shadow removal is important and necessary to im-

prove the visual effects and avoid the performance drop on the above-mentioned computer vision tasks. However, due to the complex interactions of geometry and illumination, shadow removal remains a challenging problem.

Current shadow removal methods can be mainly divided into two categories: physical-based methods [8, 7, 12, 19, 29, 38, 43] and learning-based methods [31, 35, 15, 36, 5, 42, 44]. Compared to physical-based methods, which apply a physical model to analyze each pixel's intensities, learning-based methods analyze the image in feature maps. Recently, learning-based methods with proper model have presented potential advantages [42, 15, 30]. However, these methods mainly focus on increasing the receptive field of the model without considering other particular context-sensitive shadow-aware components, which may easily ignore the contextual matching information hidden in images.

In this paper, we propose a novel two-stage context-aware network CANet for shadow removal in an end-to-end manner. As shown in Figure 1, our CANet integrates a contextual patch matching (CPM) module and a contextual feature transfer (CFT) mechanism at Stage-I and takes Stage-II as a refinement step for shadow removal. In particular, the CPM module is designed to search for the corresponding potential relationships between shadow and non-shadow patches, which demonstrates the contextual mapping between shadow and non-shadow regions. The CFT mechanism is utilized to transfer the contextual feature at different scales from non-shadow regions to shadow regions based on the output patch matching pairs from the CPM module and the extracted contextual features.

Our CPM module is designed as a dual-head structure network with the shared patch feature extractor to predict the degree of context matching between two patches from the image, as well as determine the type of the patch pair without a shadow mask. We only focus on contextual information transfer from non-shadow regions to shadow regions. Therefore we can define three types of patch pairs, *i.e.*, (1) both from shadow or non-shadow regions, (2) the first one from the shadow region and the second one from

---

the non-shadow region, and (3) the first one from the non-shadow region and the second one from the shadow region. With these prediction types, we can filter out most irrelevant patch pairs. Unlike those traditional patch match methods, our CPM is learning-based to adaptively handle complex scenes by data-driven and can effectively avoid matching errors caused by shadows the impact of shadows by averaging the lightness. What's more, our type classification head can be used to filter out mostly patch pairs from the same shadow or non-shadow regions and only focus on other pairs with high correlation scores.

We train the CPM module with sizeable self-collected training data and apply the learned CPM module to obtain a set of patch matching pairs. Then, inspired by the idea of information transfer [32], we introduce a contextual feature transfer (CFT) mechanism to transfer the contextual feature at different scales from non-shadow patches to shadow patches, resulting in a series of feature maps without shadow information. Different from the existing information transfer strategies used in shadow removal task [38, 43, 41], which search a most relevant non-shadow patch/sub-region for each shadow patch/sub-region, our CFT mechanism performs feature transfer by applying several patch matching pairs for one shadow patch according to the similarity between two patches. With the reconstructed shadow-less feature maps, we remove shadows in the L and A/B channels separately at Stage-I. Finally, to ensure the robustness of our results, with the recovered L and A/B channel images and the shadow image as inputs, we use an encoder-decoder to predict the final shadow removal image at Stage-II.

In summary, our main contributions are three-fold as follows:

- We propose a two-stage context-aware network CANet for shadow removal in an end-to-end manner, in which the contextual information from non-shadow regions is transferred to shadow regions at the embedded feature spaces.

- We design and train a context patch matching (CPM) module to acquire the potential contextual relationships between shadow and non-shadow regions in the image, which automatically distinguishes shadow patches from non-shadow patches during the matching processing.

- The proposed contextual feature transfer (CFT) mechanism transfers the extracted contextual features from non-shadow regions to shadow regions in different scales, which remove features associated with shadows and produce superior shadow removal results.

Both quantitative and qualitative experiments demonstrate the effectiveness and efficiency of our proposed CANet, as well as its superior performance in generating realistic shadow-removal images.

## 2. Related work

**Physical-based methods for shadow removal.** Physical-based methods for shadow removal are traditional methods, which usually formulate a physical model using some prior knowledge to recover illumination in shadow regions [19, 12, 7, 38, 19, 29]. Finlayson *et al.* proposed a series of shadow removal methods [8, 7] based on gradient consistency, which reconstructs the shadow removal images based on the prior that the gradient information of the image is immutable after shadow removal. Due to the change of illumination, these methods can present obvious shadow boundary artifacts.

Another strategy is information transfer, which transfers information like color or light from one picture/region to another picture/region. It has been widely used in image process tasks. Wen *et al.* [37] proposed a user-interactive multiple local color transfer method, which sets a proper gradient-guided color transfer function for each pixel. Zhang *et al.* [45] incorporated the color-transfer techniques and gradient fusion method to altering the illumination effects of an image from another. Shor *et al.* [34] build a linear mapping model between the shadowed and non-shadow regions. Xiao *et al.* [38] conducted shadow removal task using sub-region matching illumination transfer. Zhang *et al.* [43] proposed a local-to-global shadow removal method based on illumination transfer. Although those methods using information transfer can produce pleasant shadow removal results, the effectiveness of these methods depends on the accuracy of texture matching.

**Learning-based methods for shadow removal.** Unlike traditional physical-based methods, learning-based methods tend to learn high-level contextual features for shadow removing [15, 41, 20, 42]. Qu *et al.* [31] proposed a multi-context embedding network DeshadowNet to integrate the information from different levels for shadow removing. Wang *et al.* [35] analyzed the relationship of shadow detection and removal and then proposed a stacked conditional generative adversarial network (ST-CGAN) model to perform shadow detection and removal jointly. Hu *et al.* [15] used direction-aware spatial context attention features for shadow detection and removal. Zhang *et al.* [42] explore relationship of the residual and inverse illumination for shadow removal and proposed a general RIS-GAN. Hieu *et al.* [20] regard shadow image as the combination of shadow-free image, shadow parameters and shadow matte, and use neural networks to predict them to remove shadows. Liu *et al.* [23] presented a LG-ShadowNet for shadow removal by training on unpaired data. Lin *et al.* [22] proposed a BEDSR-Net for document shadow removal. The BEDSR-Net is specifically designed for document image shadow re-
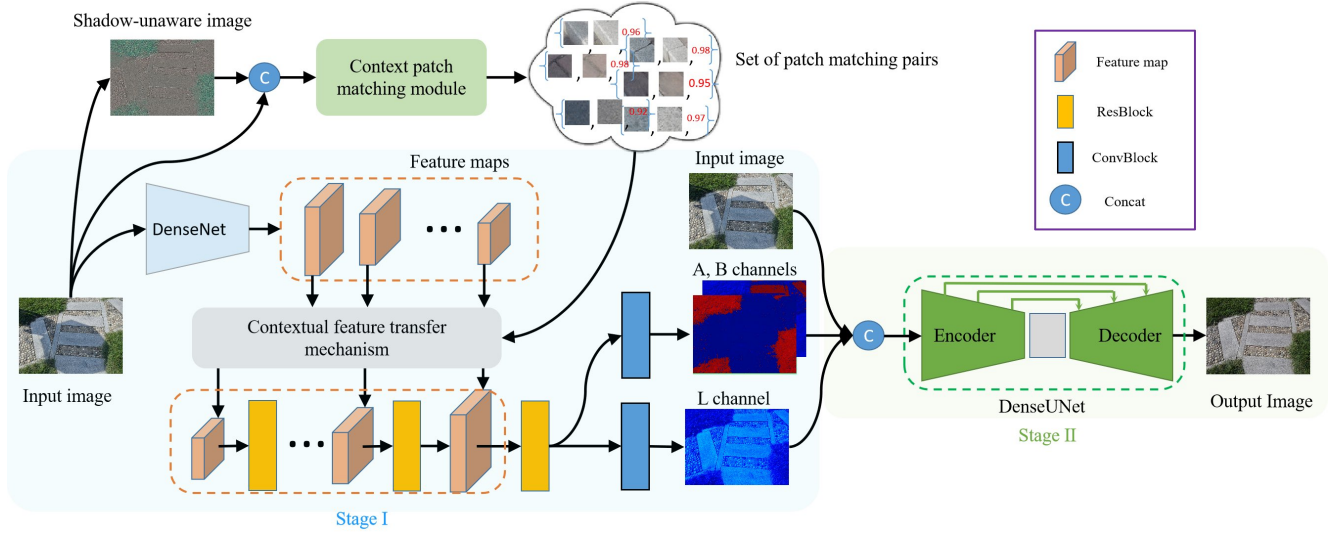
Figure 1. The overview of our proposed CANet, which takes two stages for shadow removal. At Stage-I, the contextual feature is firstly extracted via a pretrained DenseNet [18]; meanwhile, the designed contextual patch matching module (CPM) (see Figure 2) is used to acquire a set of contextual matching pairs; then, applying a contextual feature transfer mechanism (see Figure 3) to transfer contextual information from non-shadow patches to shadow patches to recover the L and A/B channels of the shadow-removal image. At Stage-II, we integrate the recovered L and A/B channel information with the input shadow image and feed them into a DenseUNet to generate the final shadow-removal result.

moval, which may be lack of expandability for other kinds of shadow images. Cun *et al.* [3] designed a network named SMGAN for shadow-removal, which can produc ghost-free shadow-removal images. Although those existing methods achieve some advances, they only focus on increasing the receptive field of the model, ignoring the paired matching information in the image. In contrast, our proposed CANet is proposed to explore the underlying contextual information between shadow and non-shadow regions.

## 3. Method

Intuitively, two patches with similar textures should have similar illumination and context under the same shadow-free environment. Based on this, we explore the idea of context-aware information transfer to conduct our shadow removal task. The proposed two-stage context-aware network (CANet) for shadow removal is illustrated in Figure 1. At Stage-I, given a shadow image, the L and A/B channels of a shadow-removal image are recovered with contextual information transferred from non-shadow patches, relying the obtained contextual patch matching information. At Stage-II, a DenseUNet is designed to integrate the recovered L and A/B channel information with the input shadow image to generate a high-quality shadow-removal image.

In the following subsections, we will introduce our contextual patch matching module, contextual feature transfer mechanism, and the two-stage CANet for shadow removal.

### 3.1. Contextual Patch Matching Module

As shown in Figure 1, the CPM module is designed to generate a set of ordered patch matching pairs of non-

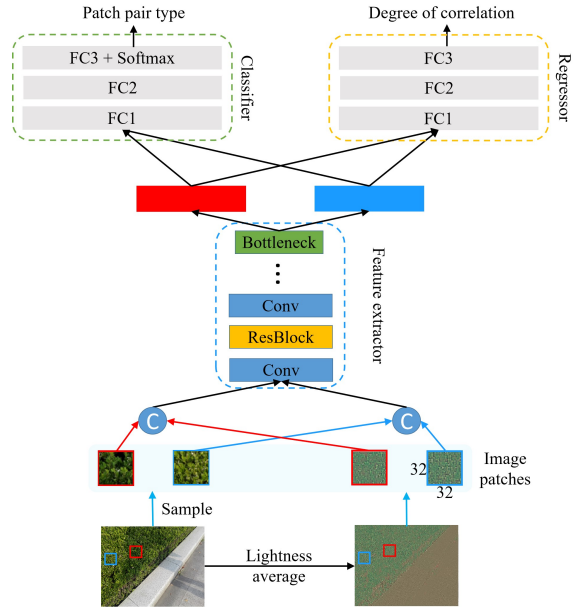shadow and shadow patches with the prediction correlation scores together.



Figure 2. The architecture of our dual-head contextual patch matching module. We first calculate the shadow-unaware image by lightness average with a mean-filter operation. Then, we extract $32 \times 32$ patches from the shadow image and the shadow-unaware image, and feed them into the shared patch descriptor to extract deep learning features. Finally, the extracted features are fed into the two heads to predict the type (denoted as -1, 0, or 1) of the input patch pair and the corresponding degree as a continuous value in the range of 0 to 1 to measure their correlation.

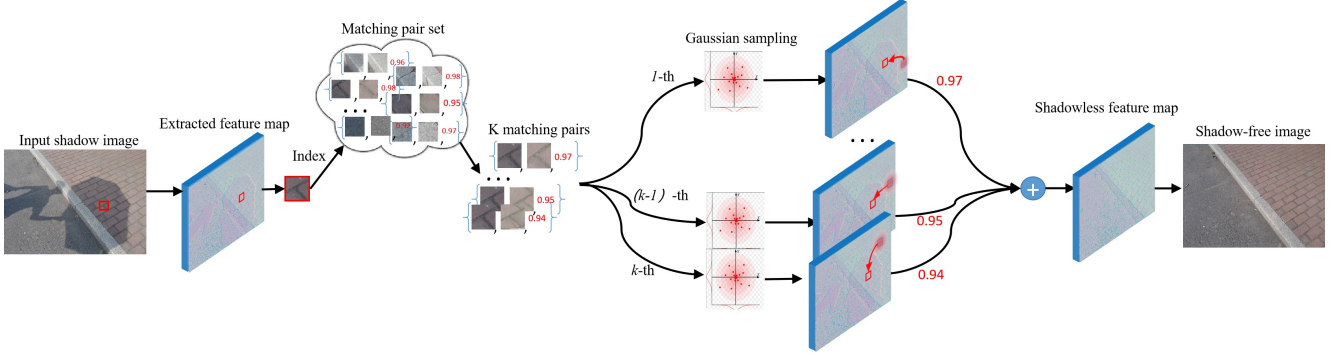To avoid matching errors caused by shadows, given an

Figure 3. The illustration of how we use our contextual feature transfer mechanism to transfer the contextual information according to captured contextual matching pairs. Note that we use the Gaussian sampling to sample contextual information from non-shadow patches and integrate them to the query shadow patches adaptively.

input shadow image, we first apply a mean-filter with kernel size 3 to get a shadow-unaware image by averaging the lightness on the shadow image. Especially, the shadow-unaware image is calculated as:

$$I_{i,j} = I_{i,j} - \frac{\sum_{(i,j)\in P} I_{i,j}}{N} + I_{avg} \qquad (1)$$

where $I_{i,j}$ denotes the lightness value at pixel $(i,j)$, $P$ denotes a $3 \times 3$ patch around pixel $(i,j)$, $N$ is the total number of pixels in the patch, and $I_{avg}$ is the global average lightness value of the image. We empirically observe that the shadow-unaware image is supplementary to the input shadow image as the input source for our CPM module, especially valuable to eliminate the effects of the shadow and distinguish shadow patches from non-shadow patches. More details can be seen in the Appendices A.

We then extract $32 \times 32$ patches from both the input shadow image and the shadow-unaware image are concatenated and fed into our dual-head CPM module, as illustrated in Figure 2. Unlike traditional methods, which use the handcraft local descriptor such as SIFT [26] to represent image patches and apply the Euclidean distance to measure the similarity between them, our CPM is learning-based so that it can adaptively handle the complex scenes by data-driven. We make full use of convolution layers with residual blocks as the shared feature extractor to extract deep learning features for the dual heads to conduct a regression task and a classification task. To specify, one head is a correlation regressor to output a continuous value at the range from 0 to 1 as the degree of correlation, and the other one is a type classifier to predict the type of pair as one of three types denoted as -1, 0, and 1.

It is worth mentioning here we care about the order between each output patch matching pair, especially when the two patches are not from the same regions (shadow or non-shadow). In particular, 0 indicates two patches from the same shadow or non-shadow regions, -1 indicates the patch pair starting with the non-shadow patch and then the shadow patch, while 1 indicates the patch pair starting with the

shadow patch and then the non-shadow patch. Note that we use -1 and 1 to distinguish the shadow patch from the non-shadow patch, as we only transfer the contextual information from non-shadow regions to shadow regions. By doing this, our type classifier head can be used to filter out most of the patch pairs from the same shadow or non-shadow regions and only focus on other pairs with high correlation scores.

To learn a solid and robust CPM module, we make full use of the existing shadow benchmark datasets to collect a large training data set. We randomly sample $32 \times 32$ patches from shadow images. Note that we use cosine similarity as the measurement between shadow and non-shadow patches on the corresponding shadow-free images to generate ground-truth for correlation regression. For those with cosine similarity higher than 0.95, we set the ground-truth of correlation degree $\mathbf{s}_{gt}$ as 1 and 0 smaller than 0.6. For the ground-truth type $\mathbf{t}$, we utilize the ground-truth shadow mask to determine the ground-truth type to be -1, 0, or 1 for any patch pairs.

We optimize the overall loss $\mathcal{L}_{CPM}$ to train the CPM module. It contains a regression loss $\mathcal{L}_{reg}$ and a classification loss $\mathcal{L}_{cls}$, *i.e.*,

$$\mathcal{L}_{CPM} = \mathcal{L}_{reg} + \mathcal{L}_{cls} \qquad (2)$$

The regression loss $\mathcal{L}_{reg}$ is used to optimize the correlation regressor for CPM module, *i.e.*,

$$\mathcal{L}_{reg} = ||\mathbf{s}_{out} - \mathbf{s}_{gt}||_2 \qquad (3)$$

where $\mathbf{s}_{gt}$ is the correlation degree label for input pair and $\mathbf{s}_{out}$ is the output of the correlation regressor.

The classification loss $\mathcal{L}_{cls}$ is used to optimize the type classifier, which is defined as a cross-entropy:

$$\mathcal{L}_{cls} = -\sum_{i=1}^{3} t_i log(p_i) \qquad (4)$$

where $t_i$ denotes the ground truth matching type of the patch pair and $p_i$ is the output of our type classifier.

## 3.2. Contextual Feature Transfer Mechanism

The target of our contextual feature transfer mechanism is to transfer the contextual feature from non-shadow regions to shadow regions, resulting in feature maps without shadow information. Generally, directly replacing the features from non-shadow regions to shadow regions may cause sub-optimal results, such as discontinuity, artifacts. Therefore, we introduce a contextual feature transfer (CFT) mechanism to perform information transfer in feature space.

The process of our contextual feature transfer model is shown in Figure 3. Given a shadow patch from the input feature map, we first retrieve the matched non-shadow patches from the generated set of patch matching pairs. Then, for each shadow patch, we use Gaussian sampling to perform contextual feature transfer using the matched non-shadow patches. Finally, we integrate the top-$k$ transferred feature patches according to the correlation degree between each matching pair.

Let $n$ be the kernel size of Gaussian sampling and $k$ be the feature transfer times. The Gaussian sampling can be written as:

$$F'_{x,y} = \sum_{\Delta x=0}^{n} \sum_{\Delta y=0}^{n} \frac{\varphi(\Delta x, \Delta y)}{\sum_{\Delta x=0}^{n} \sum_{\Delta y=0}^{n} \varphi(\Delta x, \Delta y)} F_{x+\Delta x, y+\Delta y}$$

(5)

where $F'_{x,y}$ and $F_{x,y}$ are the feature maps after and before sampling at position $(x, y)$, respectively. $\varphi(\Delta x, \Delta y)$ is the Gaussian weight at the position $(x + \Delta x, y + \Delta y)$ and $\varphi(\Delta x, \Delta y) = exp\left(-\frac{\Delta x^2 + \Delta y^2}{2\sigma^2}\right)$, where $\sigma$ is the variance of the Gaussian distribution.

To better integrate the transferred features, we adaptively integrate the $k$ sampling results according to the correlation degree of each matching pair. The reconstructed shadowless feature $F$ can be written as:

$$F = \sum_{i=1}^{k} \frac{w_i}{\sum_{i=1}^{k} w_i} F'_i,$$

(6)

where $F'_i$ is the $i$-th sampling result, $w_i$ is the correlation degree between the matching pair. Due to Gaussian sampling has a larger receptive field and takes surrounding information into consider when sampling, the Gaussian sampling in the contextual feature transfer mechanism can help to better transfer the contextual information and get pleasant results.

## 3.3. Shadow Removal with Two-stage Strategy

Our shadow-removal network CANet adopts a two-stage strategy. At Stage-I, we first use a DenseNet [18], pretrained on ImageNet[4] as the feature extractor to extract the contextual features. Then, with the extracted contextual features as inputs, we apply our contextual feature transfer mechanism to transfer the features from non-shadow regions to shadow regions at different scales. With a series of

upsampling and residual blocks, we can remove shadows at the L and A/B channels separately. As we can see from the statistical analysis in Figure 4, the L channel is more sensitive than A/B channels to highlight the difference between shadow and non-shadow regions. The separation treatment can avoid the over-processing of A and B channels and the inadequate processing of the L channel, making it more propitious to feature transfer in the contextual feature transfer model.
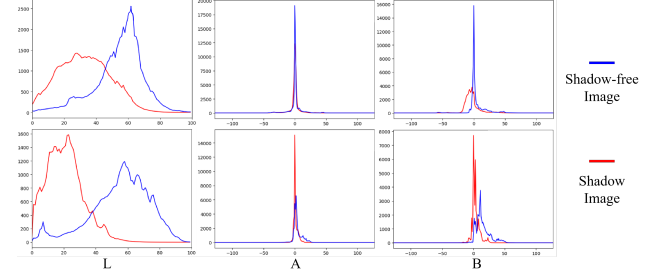


Figure 4. The difference between the input shadow image and the ground truth shadow-free image in shadow areas of each channel in LAB color space on the ISTD dataset [35](first row) and the SRD dataset [31](second row).

Note that the recovered L and A/B channels are shadow-removal results using transfer operation, which maybe contain undesirable areas due to the inaccurate match in CPM module. Hence we turn to Stage-II to produce a fine shadow-removal result, which takes the recovered L and A/B channel as strong guidance information for shadow-removal result generation. At Stage-II, with the recovered L and A/B channel images and the shadow image as inputs, we use DenseUNet, an encoder-decoder structure, to predict the final shadow-removal image.

To get a robust parametric model for shadow removal, we use a total loss $\mathcal{L}_{CANet}$ to train our CANet. It is defined with a removal loss $\mathcal{L}_{rem}$, a perceptual loss $\mathcal{L}_{per}$, and a gradient loss $\mathcal{L}_{grad}$, i.e.,

$$\mathcal{L}_{CANet} = \lambda_1 \mathcal{L}_{rem} + \lambda_2 \mathcal{L}_{per} + \lambda_3 \mathcal{L}_{grad}$$

(7)

where $\lambda_1, \lambda_2, \lambda_3$ are the hyperparameters. In this paper, we set $\lambda_1 = 1, \lambda_2 = 25, \lambda_3 = 5$.

The removal loss $\mathcal{L}_{rem}$ is the visual-consistency loss between the shadow removal result of two stages $I_{out\_1}, I_{out\_2}$ generated by our CANet and the corresponding ground-truth $I_{gt}$, i.e.,

$$\mathcal{L}_{rem} = \|I_{gt} - I_{out}\|_2.$$

(8)

The perceptual loss $\mathcal{L}_{per}$ is perceptual-consistency loss which aims to preserve the structure of image and is defined as:

$$\mathcal{L}_{per} = \|VGG(I_{gt}) - VGG(I_{out})\|_1,$$

(9)

where $VGG(\cdot)$ is the feature extractor from the VGG19 model.

The gradient loss $\mathcal{L}_{grad}$ is used to encourage the result to be smooth and is defined as:

$$\mathcal{L}_{grad} = \|\nabla I_{gt} - \nabla I_{out}\|_1, \qquad (10)$$

where $\nabla$ is the gradient of the image at pixel-level.

## 4. Experiments

**Benchmark datasets.** We conduct various experiments on two shadow removal benchmark datasets to verify the effectiveness of our CANet. One is ISTD dataset [35], which includes 1330 training triples of shadow image, shadow mask and shadow-free image and 540 testing triplets. The other is the SRD dataset [31] with 2680 training pairs of shadow and shadow-free images and 408 testing pairs.

**Implementation details.** Our proposed method is implemented in PyTorch on two GPUs (NVIDIA GeForce 2080Ti) with the input size of the image as $400 \times 400$ and mini-batch size as 2. We empirically use the Adam Optimizer to optimize our network. In our experiments, we set the first momentum value, the second momentum value and weight decay are 0.9, 0.999 and $5 \times 10^{-4}$, respectively. We train our CPM module for 30 epochs and CANet for 50 epochs. The initial learning rate is set as 0.0001. Besides, we also partly verify our method on the Huawei MindSpore platform.

Similar to [13], we run our CPM module with two phases to avoid repeated operation on the same patch. We first extract features of all patches, then feed them respectively into two heads of our CPM module to produce the matching information of the input patch pair. Besides, to achieve a good trade-off between efficiency and accuracy, we set $k = 3$ and $n = 5$ in our experiments.

### 4.1. Comparison with State-of-the-arts

We compare our CANet with eight state-of-the-art methods, *i.e.*, Guo [12], Zhang [43] DeshadowNet [31], ST-CGAN [35], Mask-shadowGAN [16], ARGAN [5], DSC [15], and RIS-GAN [42]. Among these competing methods, the first two are traditional methods, while the last six are learning-based methods. Note that all of these learning-based methods try to explore the context information via a deep-learning model with a larger receptive field and use this information to remove shadows in the image. In particular, DSC [15] exploits the direction-aware spatial RNN, DeshadowNet [31] uses the multi-context model to capture the spatial context information.

To ensure a fair comparison, we use the same training data with the same input size (*i.e.*, $400 \times 400$) to train all the learning-based methods. We calculate the root mean square error (RMSE) in LAB color space between the generated

shadow-removal images and the shadow-free ground truth image to quantitatively evaluate the performance of shadow removal.

**Quantitative Evaluation**. We summarize the quantitative results on the test data of both SRD and ISTD in Table 1. As we can see, all the competing learning-based methods perform worse than our proposed CANet. It can be explained by the fact that these baselines ignore the potential correlation between shadow and non-shadow regions Therefore these methods may fail in handling some complex shadow scenarios, especially when the most correlative non-shadow regions are not close to the shadow regions. Conversely, by explicitly capturing the useful potential contextual matching information globally, our CANet can handle the complicated case and therefore significantly improve the results of shadow removal. Table 1 reports the quantitative evaluation results compared with state-of-the-art methods, where we can see that our CANet outperforms the other state-of-the-art methods in shadow area, non-shadow area and whole image on the two datasets, which clearly demonstrates the effectiveness of our CANet.

Table 1. The quantitative comparison results of shadow removal between our method and recent methods on ISTD and SRD datasets in terms of RMSE (where S, N, A represent the shadow area, non-shadow area and whole image respectively).

| Method | ISTD | | | SRD | | |
|---|---|---|---|---|---|---|
| | S | N | A | S | N | A |
| Guo [12] | 18.95 | 7.46 | 9.3 | 29.89 | 6.47 | 12.60 |
| Zhang [43] | 13.77 | 7.17 | 8.16 | 9.50 | 6.90 | 7.24 |
| DeshadowNet [31] | 12.76 | 7.19 | 7.83 | 17.96 | 6.53 | 8.47 |
| ST-CGAN [35] | 10.33 | 6.93 | 7.47 | 12.65 | 6.37 | 7.83 |
| Mask-shadowGAN [16] | 10.35 | 7.03 | 7.61 | 10.32 | 6.83 | 7.32 |
| ARGAN [5] | 9.21 | 6.27 | 6.63 | 8.13 | 6.05 | 6.23 |
| DSC [15] | 9.22 | 6.39 | 6.67 | 8.22 | 6.01 | 6.21 |
| RIS-GAN [15] | 9.15 | 6.31 | 6.62 | 8.09 | 6.02 | 6.17 |
| **CANet** | **8.86** | **6.07** | **6.15** | **7.82** | **5.88** | **5.98** |

**Qualitative Evaluation**. We also provide visual comparison results in Figure 15. For traditional methods, due to the local information transfer, Guo [12] can not completely remove shadows. Its results contain some artifacts, as shown in Figure 15(b). Zhang [43] also cannot handle the illumination changes at the shadow boundary well, as shown in Figure 15(c). This is due to the contextual information of the image being ignored when processing at the pixel level, causing inaccurate or wrong matching of lit block. On the contrary, with the well-designed CPM module and CFT mechanism, our CANet can better recover the illumination consistent with surroundings and solve the boundary problems such as existing artifacts, generating more realistic shadow removal results.

Regarding the learning-based methods, although they can handle some simple scenes well, they are still far from satisfactory for shadow images with complex scenes, resulting in some unpleasing shadow removal results. To specify,
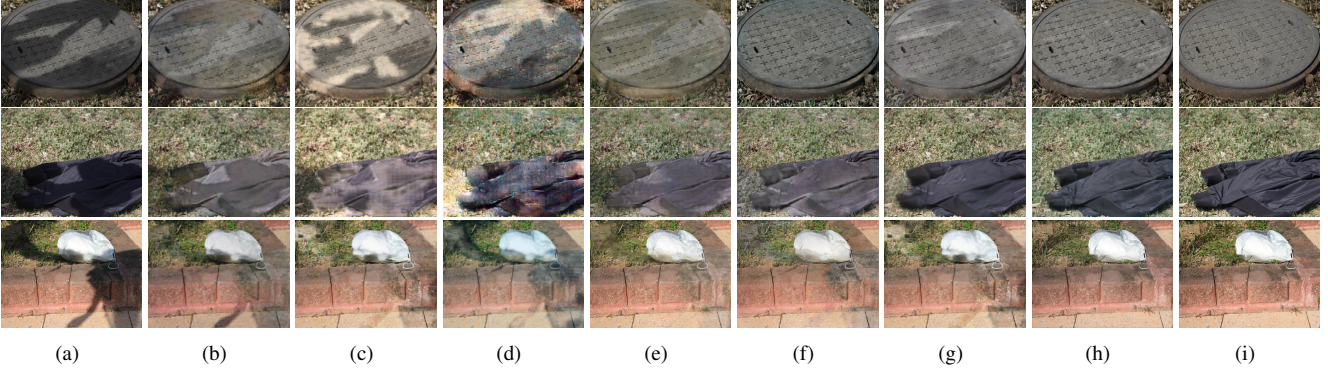
Figure 5. Shadow removal results. From left to right are: (a)input images; shadow removal results of (b) Guo, (c) Zhang, (d) ST-CGAN, (e) DSC, (f) ARGAN, (g) RIS-GAN, (h) our CANet; and (i) the corresponding shadow-free ground truth images.
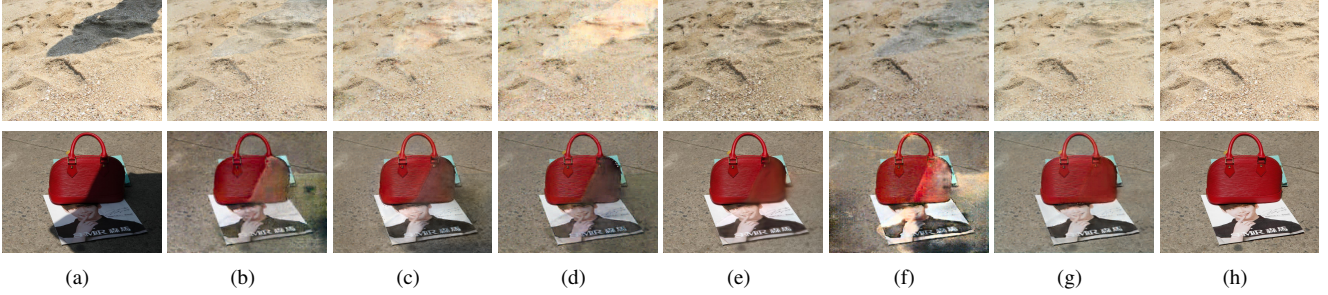


Figure 6. Shadow removal results. From left to right are: (a)input images; shadow removal results of our (b) "CANet w/ TM", (c) "CANet w/ MNet", (d) "CANet w/o CFT", (e) "CANet w/ DRCF", (f) DenseUNet, (g) our proposed CANet, and (h) the corresponding shadow-free ground truth images.

although ST-CGAN removes most of the shadows, its results still contain some artifacts, as shown in Figure 15(d). From Figure 15(e-g), we can observe that DSC severely distorts the colors around the shadows and both ARGAN and RIS-GAN have a certain degree of excessive shadow removal. The main reason for such poor results is that these methods ignore the potential correlation in images, without taking the difference between different color channels into consideration. In contrast, our proposed CANet captures the underlying correlation between shadow and non-shadow regions in image and therefore can effectively avoid the color distortion in the results. As seen in Figure 15(h), our CANet method produces more realistic and promising results than the competing methods.

## 4.2. Ablation study

To further verify the effectiveness of these components we proposed, we design a series of variants. First, we replace our CPM module with two patch matching methods. One is "traditional-match" which captures contextual matching set with a traditional hand-craft descriptor and Euclidean distance, and the other one is MatchNet [13]. We denote these two variants as "CANet w/ TM" and "CANet w/ MNet", respectively. Then we design two variants for verifying the effectiveness of the proposed CFT mechanism. One is to remove the context feature transfer mechanism completely, and the other one is to replace the contex-

tual feature directly without considering Gaussian sampling in CFT during the feature transfer. We denote these two new variants as "CANet w/o CFT" and "CANet w/ DRCF", respectively. Last but not least, we also take DenseUNet to conduct one-stage shadow removal directly.

For fairness, we train these variants on the same training data. The quantitative results are summarized in Table 2. From the table, we can observe that: (1) the contextual mapping information provided by the CFT mechanism can help to improve the accuracy of shadow removal results; (2) the CPM module is essential to ensure the quality of contextual matching information; (3) the proposed CFT mechanism ensure the best performance of our CANet.

Figure 6 presents some visual results for the different variants. As we can see, without considering the contextual matching information, there may be some shadow artifacts in the results, as shown in Figure 6(d). Replacing the contextual features directly results in unsatisfied results with discontinuous illumination and color, as shown in Figure 6(e). Besides, Figure 6(b-c) present some artifacts due to incorrect mapping information. Clearly, we can see that our CANet is the most suitable and efficient.

## 4.3. Discussion

**Robustness**. To further verify the robustness of our method, we collect some real-world shadow images with complex scenes to conduct experiments and summarize the

Table 2. The quantitative shadow removal results of ablation analysis on ISTD and SRD dataset in terms of RMSE.

| Method | ISTD | | | SRD | | |
|---|---|---|---|---|---|---|
| | S | N | A | S | N | A |
| CANet w/ TM | 9.62 | 6.33 | 6.98 | 8.44 | 6.58 | 6.89 |
| CANet w/ MNet | 9.16 | 6.20 | 6.52 | 8.17 | 6.21 | 6.35 |
| CANet w/o CFT | 10.11 | 6.88 | 7.54 | 9.28 | 6.35 | 6.96 |
| CANet w/ DRCF | 9.15 | 6.21 | 6.56 | 8.10 | 6.11 | 6.25 |
| DenseUNet | 10.22 | 7.02 | 7.58 | 10.44 | 6.71 | 7.28 |
| CANet | **8.86** | **6.07** | **6.15** | **7.82** | **5.88** | **5.98** |



(a)  (b)  (c)  (d)  (e)  (f)

Figure 7. Shadow removal results for real shadow images outside the two datasets. From left to right are: (a)input images; shadow removal results of (b)ST-CGAN, (c)DSC, (d)ARGAN, (e)RIS-GAN and (f)our CANet.

results in Figure 7. Obviously, the shadow removal results generated by our CANet look more realistic compared to the other competing algorithms. These observations clearly demonstrate the robustness of our CANet to handle the complex real-world scenes.



Figure 8. The illustration of the limitation of our method. From left to right are the input image, the result of our method and the shadow-free ground truth image. The red rectangular area in the input image does not have a matching position in the non-shadow area, which makes it difficult to recover the light in that area.

**Extension to video-level shadow removal.** We also apply our CANet to remove shadows in a video by processing each frame in the video separately. We take one frame result for every 0.5 seconds and visualize the results in Figure 9. From the results, we can see that our CANet can remove shadow well at frame-level, but the continuity of the video is not ensured, which we take as a part of our future work.



Figure 9. Shadow removal result in a video, a frame for every 0.5 seconds.

**Running time.** It firstly takes around 10 hours to train our CPM module, then takes 32 hours to train our CANet on the ISTD and SRD training set. After training, only 1.8 seconds on average is required to process an $400 \times 400$ image.

**Limitation.** Our proposed CANet can effectively remove shadows in images. However, it still has some limitations. (1) For some shadow images, if there is no strong contextual correlation between shadow and non-shadow regions, our CANet will fail to recover the illumination consistent with surroundings, as shown in Figure 8. (2) Since the environmental luminosity and camera exposure may vary during photo shooting, training pairs may have inconsistent colors and luminosity [15], causing our data-driven CANet produces shadow-removal results with color inconsistency.

## 5. Conclusion

In this paper, we have proposed a novel two-stage context aware network CANet for shadow removal. At stage-I, we design a contextual patch matching module (CPM) to search potential match pairs for the contextual feature transfer mechanism (CFT). At stage-II, we apply an encoder-decoder to refine the results of stage-I to generate the final high-quality shadow-removal results. The extensive experiment results have strongly confirmed the effectiveness and superiority of our method. Our framework can be extended to handle more computer vision tasks such as highlight removal [9, 11], which we take as the future work.

## Acknowledgement

# References

[1] Rita Cucchiara, Costantino Grana, Massimo Piccardi, and Andrea Prati. Detecting moving objects, ghosts, and shadows in video streams. *IEEE transactions on pattern analysis and machine intelligence*, 25(10):1337–1342, 2003. 1

[2] Rita Cucchiara, Costantino Grana, Massimo Piccardi, Andrea Prati, and Stefano Sirotti. Improving shadow suppression in moving object detection with hsv color information. In *IEEE Intelligent Transportation Systems. Proceedings*, pages 334–339. IEEE, 2001. 1

[3] Xiaodong Cun, Chi-Man Pun, and Cheng Shi. Towards ghost-free shadow removal via dual hierarchical aggregation network and shadow matting gan. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10680–10687, 2020. 3

[4] Jia Deng, Wei Dong, Richard Socher, Li Jia Li, and Fei Fei Li. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009. 5

[5] Bin Ding, Chengjiang Long, Ling Zhang, and Chunxia Xiao. Argan: Attentive recurrent generative adversarial network for shadow detection and removal. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10213–10222, 2019. 1, 6, 11

[6] Xinzhi Dong, Chengjiang Long, Wenju Xu, and Chunxia Xiao. Dual graph convolutional networks with transformer and curriculum learning for image captioning. In *Proceedings of the ACM International Conference on Multimedia*, 2021. 1

[7] Graham D Finlayson, Mark S. Drew, and Cheng Lu. Entropy minimization for shadow removal. *International Journal of Computer Vision*, 85(1):35–57, 2009. 1, 2

[8] Graham D Finlayson, Steven D Hordley, Cheng Lu, and Mark S Drew. On the removal of shadows from images. *IEEE transactions on pattern analysis and machine intelligence*, 28(1):59–68, 2005. 1, 2

[9] Gang Fu, Qing Zhang, Chengfang Song, Qifeng Lin, and Chunxia Xiao. Specular highlight removal for real-world images. *Computer Graphics Forum*, 38(7):253–263, 2019. 8

[10] Gang Fu, Qing Zhang, and Chunxia Xiao. Towards high-quality intrinsic images in the wild. In *IEEE International Conference on Multimedia and Expo*, pages 175–180. IEEE, 2019. 1

[11] Gang Fu, Qing Zhang, Lei Zhu, Ping Li, and Chunxia Xiao. A multi-task network for joint specular highlight detection and removal. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7752–7761, 2021. 8

[12] Ruiqi Guo, Qieyun Dai, and Derek Hoiem. Single-image shadow detection and removal using paired regions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2033–2040. IEEE, 2011. 1, 2, 6, 11

[13] Xufeng Han, Thomas Leung, Yangqing Jia, Rahul Sukthankar, and Alexander. C. Berg. Matchnet: Unifying feature and metric learning for patch-based matching. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 6, 7

[14] Tao Hu, Chengjiang Long, and Chunxia Xiao. A novel visual representation on text using diverse conditional gan for visual recognition. *IEEE Transactions on Image Processing*, 30:3499–3512, 2021. 1

[15] Xiaowei Hu, Chi-Wing Fu, Lei Zhu, Jing Qin, and Pheng-Ann Heng. Direction-aware spatial context features for shadow detection and removal. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. to appear. 1, 2, 6, 8, 11

[16] Xiaowei Hu, Yitong Jiang, Chi-Wing Fu, and Pheng-Ann Heng. Mask-ShadowGAN: Learning to remove shadows from unpaired data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019. to appear. 6

[17] Gang Hua, Chengjiang Long, Ming Yang, and Yan Gao. Collaborative active visual recognition from crowds: A distributed ensemble approach. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):582–594, 2018. 1

[18] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3, 5, 11

[19] Xiang Huang, Gang Hua, Jack Tumblin, and Lance Williams. What characterizes a shadow boundary under the sun and sky? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 898–905. IEEE, 2011. 1, 2

[20] Hieu Le and Dimitris Samaras. Shadow removal via shadow image decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2020. 2

[21] Zhengqi Li and Noah Snavely. Learning intrinsic image decomposition from watching the world. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2018. 1

[22] Yun-Hsuan Lin, Wen-Chin Chen, and Yung-Yu Chuang. Bedsr-net: A deep shadow removal network from a single document image. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12905–12914, 2020. 2

[23] Zhihao Liu, Hui Yin, Yang Mi, Mengyang Pu, and Song Wang. Shadow removal by a lightness-guided network with training on unpaired data. *IEEE Transactions on Image Processing*, 30:1853–1865, 2021. 2

[24] Chengjiang Long, Roddy Collins, Eran Swears, and Anthony Hoogs. Deep neural networks in fully connected crf for image labeling with social network metadata. In *IEEE Winter Conference on Applications of Computer Vision*, pages 1607–1615. IEEE, 2019. 1

[25] Chengjiang Long and Gang Hua. Correlational gaussian processes for cross-domain visual recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 118–126, 2017. 1

[26] David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1150–1157. Ieee, 1999. 4

[27] Huynh Manh and Gita Alaghband. Scene-lstm: A model for human trajectory prediction. *arXiv preprint arXiv:1808.04018*, 2018. 1

[28] Ivana Mikic, Pamela C Cosman, Greg T Kogut, and Mohan M Trivedi. Moving shadow and object detection in traffic scenes. In *Proceedings of International Conference on Pattern Recognition*, volume 1, pages 321–324. IEEE, 2000. 1

[29] Ankit Mohan, Jack Tumblin, and Prasun Choudhury. Editing soft shadows in a digital photograph. *IEEE Computer Graphics and Applications*, 27(2):23–31, 2007. 1, 2

[30] Saritha Murali and V. K. Govindan. Shadow detection and removal from a single image using lab color space. *Cybernetics and Information Technologies*, 13(1):95 – 103, 01 Mar. 2013. 1

[31] Liangqiong Qu, Jiandong Tian, Shengfeng He, Yandong Tang, and Rynson W. H. Lau. Deshadownet: A multi-context embedding deep network for shadow removal. In *IEEE Conference on Computer Vision and Pattern Recognition*, July 2017. 1, 2, 5, 6, 11

[32] Erik Reinhard, Michael Adhikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images. *IEEE Computer Graphics and Applications*, 21(5):34–41, 2001. 2

[33] Liushuai Shi, Le Wang, Chengjiang Long, Sanping Zhou, Mo Zhou, Zhenxing Niu, and Gang Hua. Sgcn: Sparse graph convolution for pedestrian trajectory prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 1

[34] Yael Shor and Dani Lischinski. The shadow meets the mask: Pyramid-based shadow removal. In *Computer Graphics Forum*, volume 27, pages 577–586. Wiley Online Library, 2008. 2

[35] Jifeng Wang, Xiang Li, and Jian Yang. Stacked conditional generative adversarial networks for jointly learning shadow detection and shadow removal. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2018. 1, 2, 5, 6, 11

[36] Jinjiang Wei, Chengjiang Long, Hua Zou, and Chunxia Xiao. Shadow inpainting and removal using generative adversarial networks with slice convolutions. *Computer Graphics Forum*, 38(7):381–392, 2019. 1

[37] Chung Lin Wen, Chang Hsi Hsieh, Bing Yu Chen, and Ming Ouhyoung. Example-based multiple local color transfer by strokes. *Computer Graphics Forum*, 27(7):1765–1772, 2008. 2

[38] Chunxia Xiao, Donglin Xiao, Ling Zhang, and Lin Chen. Efficient shadow removal using subregion matching illumination transfer. In *Computer Graphics Forum*, volume 32, pages 421–430. Wiley Online Library, 2013. 1, 2

[39] Jiqing Zhang, Chengjiang Long, Yuxin Wang, Haiyin Piao, Haiyang Mei, Xin Yang, and Baocai Yin Yin. A two-stage attentive network for single image super resolution. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. 1

[40] Jiqing Zhang, Chengjiang Long, Yuxin Wang, Xin Yang, Haiyang Mei, and Baocai Yin. Multi-context and enhanced reconstruction network for single image super resolution. In *IEEE International Conference on Multimedia and Expo*, pages 1–6. IEEE, 2020. 1

[41] Ling Zhang, Chengjiang Long, Qingan Yan, Xiaolong Zhang, and Chunxia Xiao. Cla-gan: A context and lightness

aware generative adversarial network for shadow removal. *Computer Graphics Forum*, 39(7), 2020. 2

[42] Ling Zhang, Chengjiang Long, Xiaolong Zhang, and Chunxia Xiao. Ris-gan: Explore residual and illumination with generative adversarial networks for shadow removal. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(7):12829–12836, 2020. 1, 2, 6, 11

[43] Ling Zhang, Qing Zhang, and Chunxia Xiao. Shadow remover: Image shadow removal based on illumination recovering optimization. *IEEE Transactions on Image Processing*, 24(11):4623–4636, 2015. 1, 2, 6, 11

[44] Xuaner Zhang, Jonathan T Barron, Yun-Ta Tsai, Rohit Pandey, Xiuming Zhang, Ren Ng, and David E Jacobs. Portrait shadow manipulation. *ACM Transactions on Graphics*, 39(4):78–1, 2020. 1

[45] Yi Zhang, Tianhao Zhao, Zhipeng Mo, and Wenbo Li. A method of illumination effect transfer between images using color transfer and gradient fusion. In *Signal and Information Processing Association Annual Summit and Conference*, 2013. 2

# Appendices

## A. Details about Network Architectures

**CPM module.** The architecture of our CPM module is shown in Table 3, which divides into three parts: feature extractor, pair type classifier, and correlation degree regressor. The feature extractor is used to extract 256-dimensional feature representations for input patch pairs. It is designed with 4 convolutional layers, 3 residual blocks, and a bottleneck layer. The obtained representations are feed into the classifier and regressor separately. Both the classification head and the regression head contain 3 fully connected layers, and classification head also have a softmax layer.

Table 3. The architecture of our CPM module. It contains a feature extractor, a pair type classifier, and a correlation degree regressor.

|  | Layer | Output Size | Operation |
|---|---|---|---|
| Feature extractor | Conv1 | $16 \times 16 \times 64$ | Conv($3 \times 3$ stride 2) |
|  | Res1 | $16 \times 16 \times 64$ | Res-blocks($3 \times 3$) |
|  | Conv2 | $8 \times 8 \times 96$ | Conv($3 \times 3$ stride 2) |
|  | Res2 | $8 \times 8 \times 96$ | Res-blocks($3 \times 3$) |
|  | Conv3 | $4 \times 4 \times 96$ | Conv($3 \times 3$ stride 2) |
|  | Res3 | $4 \times 4 \times 96$ | Res-blocks($3 \times 3$) |
|  | Conv4 | $4 \times 4 \times 64$ | Conv($3 \times 3$ stride 1) |
|  | Bottleneck | 256 | FC |
| Classifier | FC1 | 256 | FC |
|  | FC2 | 128 | FC |
|  | FC3 | 3 | FC |
|  | Softmax | 3 | Softmax |
| Regressor | FC1 | 256 | FC |
|  | FC2 | 128 | FC |
|  | FC3 | 1 | FC |

**CANet.** In Figure 10, we illustrate the detailed network architecture of our proposed CANet. Each orange rectangles in the network is the feature map of the corresponding layer, and the number in the rectangles is their channel number. Note that the "DenseBlock", "Transition layer" are followed as the original version of DenseNet [18].

## B. Better Understanding for CPM Module

### B.1. Effectiveness of Light-unaware Images

As shown in Figure 11, with the supplementary light-unaware image, we can largely eliminate the difference between shadow regions and non-shadow regions, which effectively avoids the matching errors caused by shadows.

Also, from Figure 12, we can observe that there is a larger difference between shadow image and light-unaware image in shadow regions while a smaller difference in the non-shadow region. It suggests that the shadow image and light-unaware image pair can provide an indication to distinguish shadow patches from non-shadow patches, which can be used to perform our pair type classifier.

### B.2. Large-scale Training Dataset for CPM

To train our CPM module, we collect a large-scale training collection from the existing shadow benchmark datasets: ISTD [35] and SRD [31]. The collected training dataset contains more than 360,000 and 600,000 patch pairs separately (50% match pairs and 50% non-match pairs). These patch pairs are collected from two ways: (1) we select a shadow patch in the shadow image and a matched non-shadow patch in the shadow-free image, which has the same position as the shadow patch, as illustrated in Figure 13; (2) we select a shadow patch from shadow regions and find another matched patch from non-shadow regions in the shadow image. Specially, we randomly select two patches from shadow and non-shadow regions in the shadow image and calculate the cosine similarity between the two patches in the corresponding shadow-free image. We choose the pairs with cosine similarity higher than 0.95 as the matching pair and less than 0.6 as the non-match pair, as shown in Figure 13. Due to the lack of shadow mask ground-truth in SRD dataset [18], we firstly use the results of the latest shadow detection method DSD [**?**], and then manually choose the correct results as the approximate ground-truth during the process of dataset collecting.

## C. More Visual Shadow Removal Results

We provide more visual shadow removal comparison results in Figure 15. Here, we compare our CANet with six state-of-the-art methods, *i.e.*, Guo [12], Zhang [43], ST-CGAN [35], ARGAN [5], DSC [15] and RIS-GAN [42].
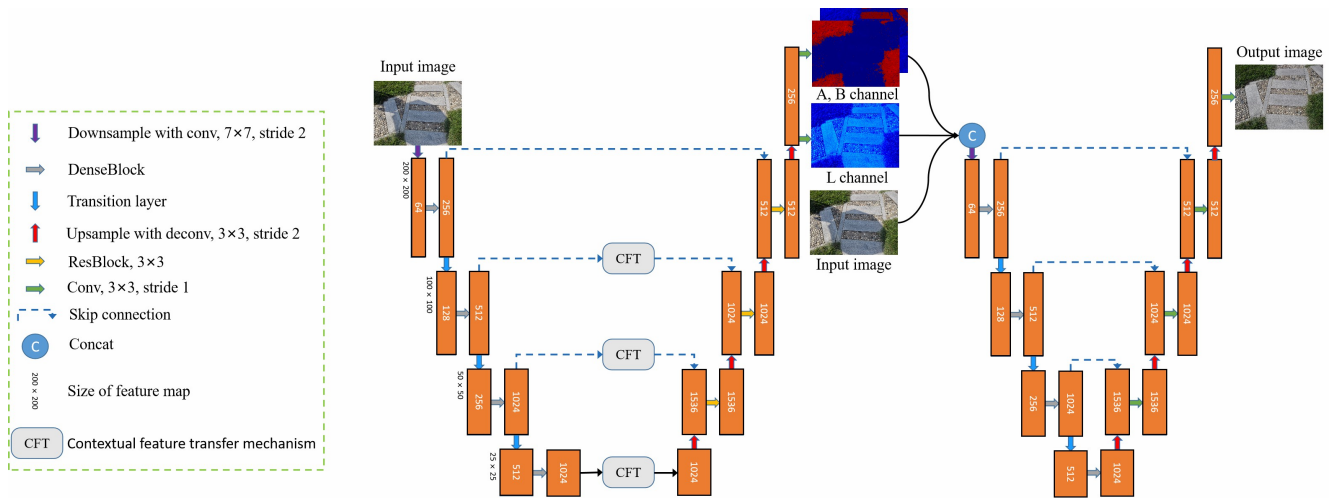
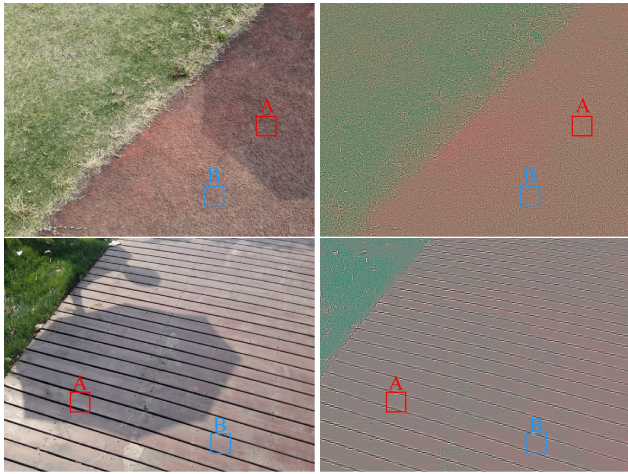Figure 10. The network architecture of our CANet.



Figure 11. From left to right are: input shadow image; and the proposed light-unaware image, which can eliminate the difference between region A and B caused by shadow.



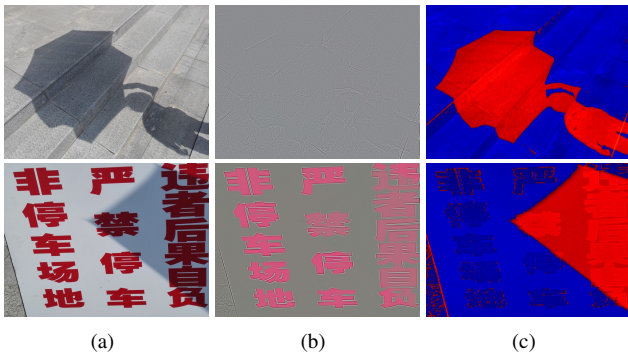Figure 13. The illustration of our first way to collect matched patch pairs.



(a) (b) (c)

Figure 12. The illustration of the difference between shadow image and light-unaware image, from left to right are: (a)input shadow image; (b)proposed light-unaware image and (c)the difference between them.
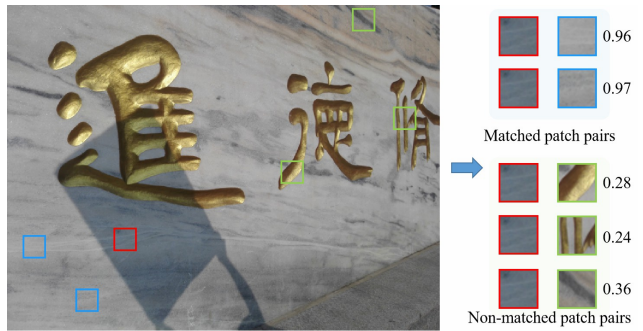


Figure 14. The illustration of our second way to collect matched and non-matched patch pairs.
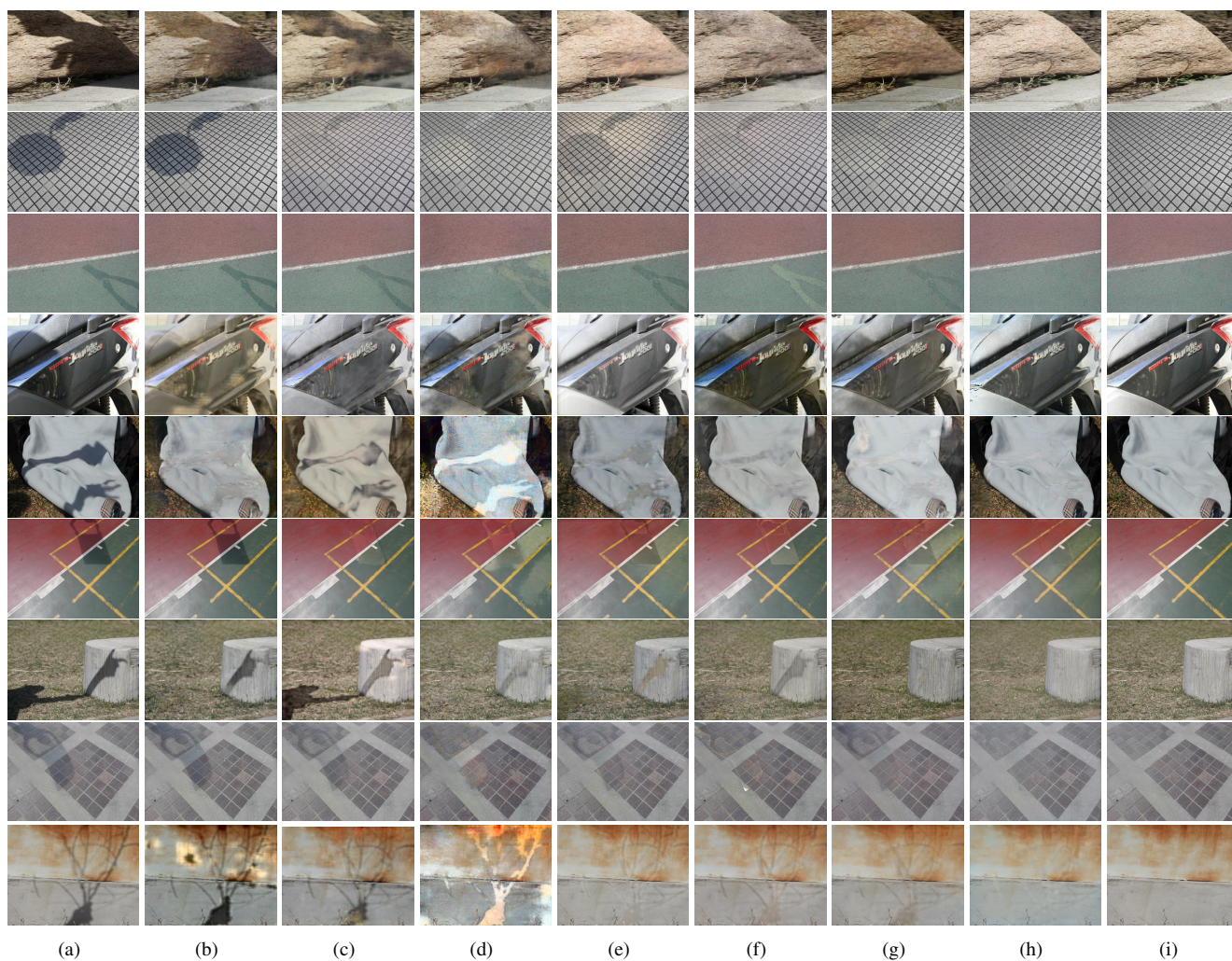
Figure 15. Shadow removal results. From left to right are: (a) input images; shadow removal results of (b) Guo, (c) Zhang, (d) ST-CGAN, (e) DSC, (f) ARGAN, (g) RIS-GAN, (h) our CANet; and (i) the corresponding shadow-free ground truth images.