

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

7-2017

Sparsity based reflection removal using external patch search

Renjie WAN

Boxin SHI

Ah-hwee TAN

Singapore Management University, ahtan@smu.edu.sg

Alex C. KOT

Follow this and additional works at: https://ink.library.smu.edu.sg/sis_research



Part of the [Databases and Information Systems Commons](#), and the [Graphics and Human Computer Interfaces Commons](#)

Citation

1

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylids@smu.edu.sg.

SPARSITY BASED REFLECTION REMOVAL USING EXTERNAL PATCH SEARCH

Renjie Wan^{†,§}, Boxin Shi[‡], Ah-Hwee Tan[¶], and Alex C. Kot[§]

[†]Interdisciplinary Graduate School,[§]School of Electrical and Electronic Engineering

[¶]School of Computer Science and Engineering

Nanyang Technological University, Singapore, 639798

[‡]Artificial Intelligence Research Center, National Institute of AIST, Tokyo, Japan, 1350064

rwan001@e.ntu.edu.sg, boxin.shi@aist.go.jp, {asahtan,eackot}@ntu.edu.sg

ABSTRACT

Reflection removal aims at separating the mixture of the desired background scenes and the undesired reflections, when the photos are taken through the glass. It has both aesthetic and practical applications which can largely improve the performance of many multimedia tasks. Existing reflection removal approaches heavily rely on scene priors such as separable sparse gradients brought by different levels of blur, and they easily fail when such priors are not observed in many real scenes. Sparse representation models and nonlocal image priors have shown their effectiveness in image restoration with self similarity. In this work, we propose a reflection removal method benefited from the sparsity and nonlocal image prior as a unified optimization framework. We leverage the retrieved image patch from an external database to overcome the limited prior information in the input mixture image and self similarity search. The experimental results show that our proposed model performs better than the existing state-of-the-art reflection removal method for both objective and subjective image qualities.

Index Terms— Reflection removal, image retrieval, external dataset, sparse representation

1. INTRODUCTION

The images used in most multimedia applications are often assumed to be captured under good conditions so that the desired background scene is clean and clear. However, taking images through a transparent glass is unavoidable in many scenarios such as in front of a window or for objects in a glass show case of the museum; such images are often degraded by the undesired reflections of the objects behind the camera, which may degrade the performance of multimedia applications. Reflection removal aims at removing the reflections and enhancing the visibility of the desired background scenes. It is shown to be an important pre-processing step for many computer vision approaches to largely increase the accuracy of image classification [1]. The reflection removal problem is

modeled by the following equation:

$$\mathbf{I} = \mathbf{B} + \mathbf{R}, \quad (1)$$

where \mathbf{I} is the input mixture image, \mathbf{B} is the background layer we feel interested in, and \mathbf{R} is the reflection layer we want to remove. This problem is challenging due to its ill-posed nature: the number of unknowns is twice the number of equations. To remove such ambiguity, various methods have been proposed and they can be roughly divided into two categories according to the number of images used.

One category adopts the single image as the input. Motivated by the fact that the natural image gradients have the heavy-tailed distribution, the gradient sparsity priors are widely used in many different methods [2, 3, 4]. Another prior is the GMM prior used by Shih *et al.* [1] to model the ghosting effects of the reflections. However, the special requirements for the properties of the background and reflection layers (*e.g.*, different blur levels [4, 3] or the visible ghosting effects [1]) make them difficult to handle many general scenes when such prior information is not observed.

The other category relies on multiple images as the input. With multiple images, the problem can be solved by exploring the motion field [5, 6], independent component analysis [7] or the gradient projection [8]. The patch-based sparsity prior also shows its success in solving such separation problems [9, 10]; it assumes that each layer of the input mixture image have a sparse or approximately sparse representations with respect to its corresponding dictionaries. The multiple images can make this problem less ill-posed. However, the special data capture requirements such as observing different layer motions or the demand for the polarizers largely limit these methods for practical use. For the methods based on the sparsity prior, the large structure diversity among each layer of the mixture images have to be observed to guarantee the success.

In this paper, we propose a novel reflection removal approach by combining the sparsity prior and the nonlocal image prior into a unified framework. The nonlocal image prior mainly makes use of the patch recurrence among the image,

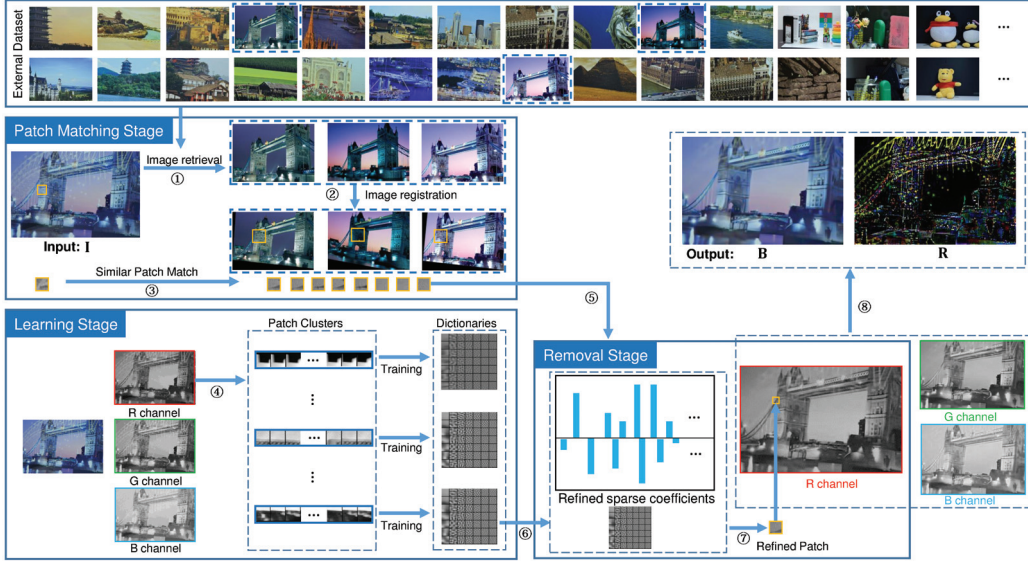


Fig. 1: The framework of our method. Our algorithm runs on RGB channel independently. For simplicity, we only show the process on R channel as an example. In the patch matching stage, we first retrieve images with similar contents from an external database (Step 1); the retrieved images are then registered to the input mixture images (Step 2); similar patches are extracted from the retrieved images based on the exemplar patches (Step 3). In the learning stage, the initial dictionary are learned from the input mixture image itself and the patch clusters are obtained using k-means and for each cluster a PCA sub-dictionary is learned (Step 4); then the nonlocal information obtained from the similar patches are used to refine the sparse codes of the exemplar patch (Step 5 and Step 6). In the removal stage, with the refined sparse codes and the dictionary, the patches are refined (Step 7) and the reflection is removed (Step 8). The reflection image \mathbf{R} is multiplied by 10 for better visual demonstration.

which is widely adopted in the patch-based image denoising [11] or super-resolution [12] methods to enhance a noisy or blurred patch from the input mixture image by reconstructing this patch with a set of similar clean patches. Most existing methods (*e.g.*, [13]) find such similar patches from the input image itself. However, the performances of these kinds of methods are often largely degraded due to the “rare patch effect” [14]. In our method, to overcome this limitation, instead of finding similar patches from the mixture image itself, the similar patches are obtained from an external database using the existing image retrieval method. The key assumption of our work is that a set of clean images that share similar contents with the background layer of the input mixture image can be retrieved from an external database and the similar patches can be extracted from the clean images. We use the nonlocal priors subtracted from the similar clean patches to regularize the sparse codes of the exemplar patches. Compared with previous methods [3, 1, 4], our method does not require special phenomenon (such as different levels of blur or ghosting effect) have to be observed on the mixture image so that we can better handle the images with general and complex structures. The framework of our complete approach is illustrated in Figure 1.

2. PROPOSED APPROACH

2.1. The reflection removal model

Let \mathbf{I} , \mathbf{B} and \mathbf{R} represent the input mixture image, background and reflection, respectively. In this work, based on the model proposed in Equation (1), we define a new energy function to formulate this problem as follows:

$$L(\mathbf{B}, \mathbf{R}) = \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \lambda\rho(\mathbf{B}) + \gamma\varrho(\mathbf{R}), \quad (2)$$

where $\rho(\mathbf{B})$ and $\varrho(\mathbf{R})$ are the regularization prior terms on the background and reflection layer, respectively. Many previous methods can also be cast into such a framework where the regularization terms are crucial for obtaining better removal performances. ρ and ϱ are chosen to be the GMM priors to model the ghosting effects in the method [1]. Another category of methods [10, 15] chooses the sparsity-based prior as the regularization term. In our proposed model, we adopt the integration of the sparsity prior and nonlocal image prior to regularize \mathbf{B} and the gradient sparsity prior to regularize \mathbf{R} . Formally, given the mixture image \mathbf{I} and the set of clean images retrieved from a dataset, we want to estimate the background \mathbf{B} and reflection \mathbf{R} by

$$\{\hat{\mathbf{B}}, \hat{\mathbf{R}}\} = \underset{\mathbf{B}, \mathbf{R}}{\operatorname{argmin}} L(\mathbf{B}, \mathbf{R}), \quad (3)$$

where

$$L(\mathbf{B}, \mathbf{R}) = \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \lambda \sum_i \|\mathbf{P}_i \mathbf{B} - \mathbf{D} \alpha_i\|_2^2 + \eta \sum_i \|\alpha_i - \beta_i\|_1 + \gamma \sum_{l=1}^L |f_l * \mathbf{R}|^s. \quad (4)$$

We explain each term of the model in detail as follows:

1. The first term is the conventional constraint, which means that the mixture image \mathbf{I} should be the summation of estimated background \mathbf{B} and estimated reflection \mathbf{R} .
2. The second term means that the estimated background \mathbf{B} can be well represented with respect to its corresponding dictionaries \mathbf{D} . \mathbf{P}_i denotes the matrix extracting image patch of size $\sqrt{n} \times \sqrt{n}$. \mathbf{D} denotes the dictionary. α_i is the coefficients corresponding to the dictionary \mathbf{D} .
3. The third term is the NCSR model proposed in [11] which enforces that α_i should be as similar as β_i , where β_i is some good estimation of α_i .
4. The fourth term is a heavy tailed distribution enforced on the estimated reflection \mathbf{R} to further stabilize the solution, which is widely used in previous methods [6, 2]. Typically, the value of s is set between 0.5 to 0.8. f_l is the Laplacian filters, namely $f_1 = [-1, 1]$, $f_2 = [-1, 1]^T$, $f_3 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ as [3].

For the dictionary \mathbf{D} , previous methods usually adopt the dictionary learned from the example image patches by using the KSVD algorithm. In such case, the learned dictionaries are required to be very redundant to represent different image structures. However, the sparse coding with overcomplete dictionaries is unstable [16]. We choose the local PCA dictionaries proposed in [11], where the training patches extracted from the input mixture image itself are clustered into K ($K = 70$) clusters and a PCA sub-dictionary are learned from each cluster. For more details, please refer to [11]. Now, the task is how to estimate β_i more accurately.

2.2. The estimation of β_i

As an estimation of α_i , β_i can be estimated from the internal or external sources. In our case, we estimate β_i from an external sources where a set of similar clean images can be found for the input mixture images.

Our framework of similar patch matching process contains three steps, which are similar image retrieval, global image registration and patch match. We adopt the image retrieval method proposed by Philbin *et al.* [17] and retrieve images from an external database. Due to the different scales and viewpoints of these retrieved images, for better patch

Algorithm 1 Sparsity prior based reflection removal

Input:

Input mixture image \mathbf{I} ;

Output:

Estimated background \mathbf{B} and reflection \mathbf{R} ;

- 1: Compute the dictionaries \mathbf{D} by k-means and PCA;
 - 2: **for** $m = 1$ to M **do**
 - 3: **for** $j = 1$ to J **do**
 - 4: Update sparse codes α_i^{j+1} by solving Equation (7);
 - 5: Update the background \mathbf{B}^{j+1} by solving Equation (10);
 - 6: Update the reflection \mathbf{R}^{j+1} by solving Equation (11);
 - 7: Set $\mathbf{B}^{m+1} = \mathbf{B}^{j+1}$ and $\mathbf{R}^{m+1} = \mathbf{R}^{j+1}$ if $j = J_{max}$
 - 8: **end for**
 - 9: If $\text{mod}(m, 5) = 0$, update the PCA dictionaries;
 - 10: **end for**
 - 11: **return** \mathbf{B}^{m+1} , \mathbf{R}^{m+1} ;
-

matching, an image registration step is needed. We use a quite standard way to register the images. We first extract SURF feature points from the mixture image and reference images, and then estimate the homographic transformation matrix by using the RANSAC algorithm. Finally, the reference images from the external database are aligned to the mixture image with the estimated transformation.

Let x_i denote the patch from the input mixture image. The nonlocal similar patches z_i that are within the first T closet to the given patch x_i are selected from a large window centered at pixel i among the registered images. Then, β_i can be computed as the weighted average of those sparse codes associated with the nonlocal similar patches as:

$$\beta_i = \sum_{t=1}^T \omega_{i,t} \alpha_{i,t}, \quad (5)$$

where $\alpha_{i,t}$ is the sparse coefficients corresponding to the patch z_i and $\omega_{i,t}$ is the weight and can be obtained as:

$$\omega_{i,t} = \frac{1}{W} \exp(-\|\hat{\mathbf{x}}_i - \hat{\mathbf{z}}_{i,t}\|_2^2/h), \quad (6)$$

where $\hat{\mathbf{x}}_i$ and $\hat{\mathbf{z}}_{i,t}$ are the estimates of the patches \mathbf{x}_i and $\mathbf{z}_{i,t}$, h is a pre-determined scalar and W is the normalization factor.

2.3. Optimization

The direct minimization of Equation (4) is difficult due to the multiple variables involved in the proposed model. Thus, we reduce the original problem into several subproblems by following the alternating minimization scheme advocated by the previous method in image deblurring and denoising works. In each step, our algorithm reduces the objective function value, and thus will converge to a local minima.

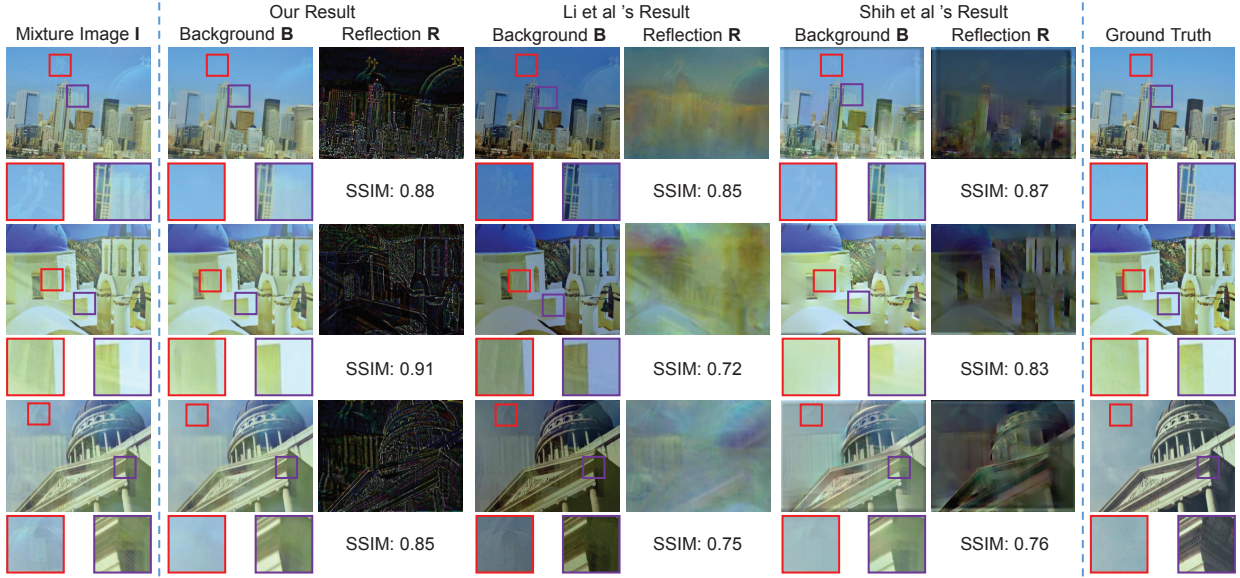


Fig. 2: Reflection removal results comparison using our method, Li *et al.*'s method [3], and Shih *et al.*'s method [1] on the postcard data. The reflection images \mathbf{R} are all multiplied by for better visual demonstration.

Solving for α_i . For a fixed \mathbf{B} and \mathbf{R} , Equation (4) reduces to a l_1 minimization problem:

$$\hat{\alpha}_i = \underset{\alpha_i}{\operatorname{argmin}} \lambda \|\mathbf{P}_i \mathbf{B} - \mathbf{D} \alpha_i\|_2^2 + \eta \|\alpha_i - \beta_i\|_1. \quad (7)$$

With fixed β_i , Equation (7) can be solved iteratively by the surrogate based algorithm [18]:

$$\alpha_i^{(t+1)} = \mathcal{S}_\tau(v_i^{(t)} - \beta_i) + \beta_i, \quad (8)$$

where $v_i^{(t)} = \mathbf{D}^\top (\mathbf{P}_i \mathbf{B} - \mathbf{D} \alpha_i^{(t)}) / c + \alpha_i^{(t)}$, $\mathcal{S}_\tau(\cdot)$ represents the soft-thresholding operator with threshold $\tau = \eta / \lambda c$, and c is a constant to guarantee the convexity. Due to the orthogonal properties of the local PCA dictionaries \mathbf{D} , the sparse coding problem of Equation (7) can be solved in just one step [12].

Solving for \mathbf{B} . When \mathbf{R} and α_i are fixed, the background \mathbf{B} can be estimated by solving the following optimization problem:

$$\hat{\mathbf{B}} = \underset{\mathbf{B}}{\operatorname{argmin}} \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \lambda \sum_i \|\mathbf{P}_i \mathbf{B} - \mathbf{D} \alpha_i\|_2^2, \quad (9)$$

where the closed-form solution can be easily obtained as follows:

$$\hat{\mathbf{B}} = (\mathbf{E} + \lambda \sum_i \mathbf{P}_i^\top \mathbf{P}_i)^{-1} (\mathbf{I} + \lambda \sum_i \mathbf{P}_i^\top \mathbf{D} \alpha_i - \mathbf{R}), \quad (10)$$

where all elements of matrix \mathbf{E} equal to one.

Solving for \mathbf{R} . Given the estimated background \mathbf{B} and sparse representation α , the estimation of reflection \mathbf{R} can be

updated. The optimization problem (5) becomes

$$\hat{\mathbf{R}} = \underset{\mathbf{R}}{\operatorname{argmin}} \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \gamma \sum_{l=1}^L |f_l * \mathbf{R}|^s. \quad (11)$$

This problem can be solved efficiently by variable substitution and Fast Fourier Transform (FFT) [19, 20]. Using the new auxiliary variables u_l ($l \in 1, 2, \dots, L$), the Equation (11) can be rewritten as:

$$\begin{aligned} \hat{\mathbf{R}} = \underset{\mathbf{R}}{\operatorname{argmin}} & \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \gamma \sum_{l=1}^L |u_l|^s \\ & + \delta \sum_{l=1}^L \|u_l - f_l * \mathbf{R}\|_2^2. \end{aligned} \quad (12)$$

It can be divided into two sub-problems: \mathbf{R} -subproblem and u -subproblem. δ is a weight value that varies during the optimization. We follow the setting in [19] to set the value of δ . In the \mathbf{R} -subproblem, the Equation (12) becomes:

$$\hat{\mathbf{R}} = \underset{\mathbf{R}}{\operatorname{argmin}} \|\mathbf{I} - \mathbf{B} - \mathbf{R}\|_2^2 + \delta \sum_{l=1}^L \|u_l - f_l * \mathbf{R}\|_2^2, \quad (13)$$

which can be solved using FFT as:

$$\hat{\mathbf{R}} = \mathcal{F}^{-1} \left(\frac{\mathcal{F}(\mathbf{I}) + \delta \sum_{l=1}^L \mathcal{F}(e_l) \star \mathcal{F}(u_l) - \mathcal{F}(\mathbf{B})}{\mathbf{E} + \delta \sum_{l=1}^L \mathcal{F}(e_l) \star \mathcal{F}(e_l)} \right), \quad (14)$$

where \mathcal{F} denotes FFT, \mathcal{F}^{-1} denotes the inverse FFT and \star is the complex conjugate.

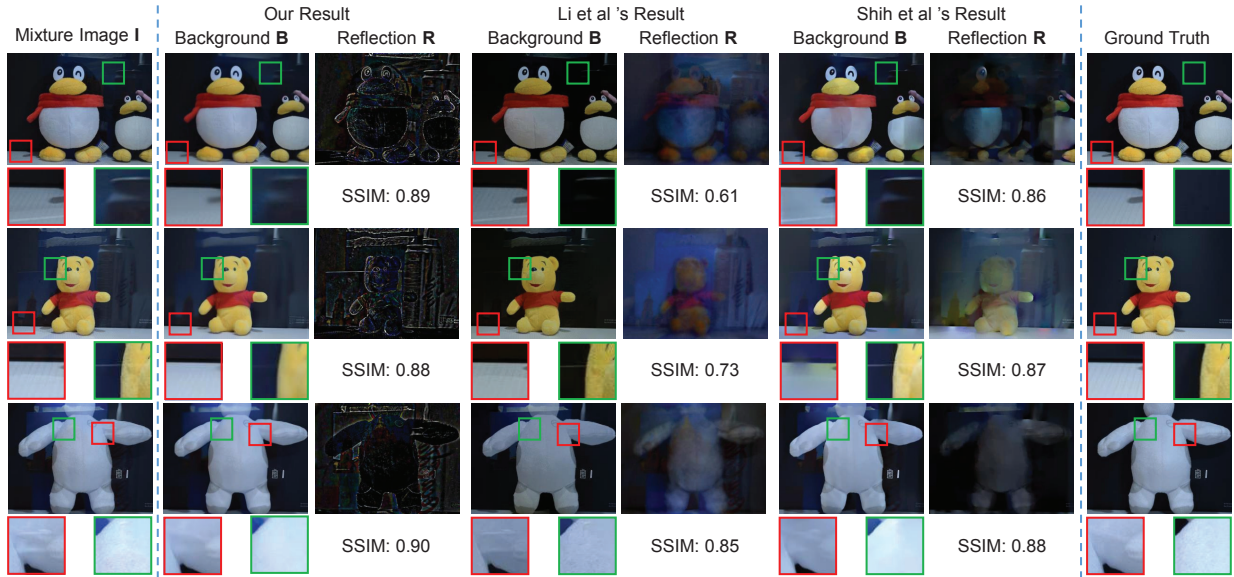


Fig. 3: Reflection removal results comparison using our method, Li *et al.*'s method [3] and Shih *et al.*'s method [1] on the solid object data. The reflection images \mathbf{R} are all multiplied by 4 for better visual demonstration.

In the u -subproblem, the u_l can be estimated by solving the following equation:

$$\hat{u}_l = \underset{u_l}{\operatorname{argmin}} \gamma \sum_{l=1}^L |u_l|^s + \delta \|u_l - f_l * \mathbf{R}\|_2^2, \quad (15)$$

which can be solved efficiently using the method in [19] over each dimension separately.

3. EXPERIMENTS

3.1. Data preparation

Existing methods usually compare the visual quality of the separated results only due to the lack of the ground truth. In order to figure out the performance of our results compared with others, we conduct evaluations for both quantitative error and visual quality. We capture the images with the ground truth following a similar way proposed in [6], where the mixture image is taken through the transparent glass and the ground truth is taken by removing the glass. We prepare two types of data capture setup: one setup uses some landmark postcards as both background and reflection objects; the other setup captures some solid objects (*e.g.*, toys of famous figures) as the background objects.

For the external database used in the patch matching stage, we collect approximately 500 images from the Internet, and three images with similar contents (the same landmark or the same toy figure captured in different environment) corresponding to each mixture image are included in the database. We then perform image retrieval [17] to find these three images before the patch matching stage. An example is shown

in Figure 1, where three images containing the Tower Bridge similar to the input mixture image are retrieved from the external dataset.

3.2. Evaluations

We show six example results in Figure 2 and Figure 3. We compare our method with two state-of-the-art methods proposed by Li *et al.* [3] and Shih *et al.* [1]. In all our experiments, the parameters are fixed as follows: T are set to 7, γ , λ and η is set to 1, 0.5 and 0.85, respectively, M is set to 15 and J is set to 10. The patch size is set to 7×7 . To quantitatively assess the algorithms, the Structural Similarity Index (SSIM) is adopted as the quality measure of the estimated background which is also use by the previous work [3, 1].

Our method shows advantage in all these results over the other two methods in terms of SSIM. Considering the visual quality of three methods, we also provide a more visually pleasing result. Li *et al.*'s method causes some color change so that the estimated background B are darker than the ground truth. For Shih *et al.*'s method, the GMM priors bring some patchy artifacts on the estimated background \mathbf{B} . Considering the results from Figure 2 and Figure 3, our algorithm can generate clearer separation and keep the details of the background. However, Li *et al.*'s method contains some residue edges from the reflection layer in the estimated background image and Shih *et al.*'s method is indeed able to suppress some of the reflections, but reflections still remain in the final estimates of \mathbf{B} . Due to the ability to reconstruct a clearer patch using external data, we recover more sharp structures in \mathbf{B} while Li and Shih's results are overly smoothed.

4. CONCLUSION

We propose a method to remove reflections based on retrieved external patch by combining the sparsity prior and the nonlocal image prior as a unified optimization. Compared with the previous methods [3, 4], we do not have special requirement for the properties of the background layer and the reflection layer, *e.g.* using different blur levels of the two layers to assist separation. Instead, we refine the sparse coefficients learned from the mixture images with the external patches to generate a more accurate sparse regularization term. Experimental results have already shown that our method outperforms the current state-of-the-art methods both from the quantitative evaluations and visual quality.

Limitations. Currently, our method can only handle the landmark scenes or some well-known objects, that can be efficiently retrieved. It is still difficult for our method to deal with the general objects or scenes, for which similar contents cannot be retrieved from the external database. Furthermore, we manually prepare the external database which already includes the similar images for the mixture image as a proof-of-concept for reflection removal. We plan to perform quantitative experiments to verify the influence from the accuracy of image retrieval to our problem and build a complete system by searching from Internet-scale database with millions of image as our future work.

5. ACKNOWLEDGEMENT

This research was carried out at the Rapid-Rich Object Search (ROSE) Lab at the Nanyang Technological University, Singapore. The ROSE Lab is supported by the National Research Foundation, Singapore, under its Interactive Digital Media (IDM) Strategic Research Programme. Boxin Shi is supported by a project commissioned by the New Energy and Industrial Technology Development Organization (NEDO).

6. REFERENCES

- [1] Y. Shih, D. Krishnan, F. Durand, and W. T. Freeman, "Reflection removal using ghosting cues," in *Proc. CVPR*, 2015.
- [2] A. Levin and Y. Weiss, "User assisted separation of reflections from a single image using a sparsity prior," *IEEE TPAMI*, 2007.
- [3] Y. Li and M. S. Brown, "Single image layer separation using relative smoothness," in *Proc. CVPR*, 2014.
- [4] R. Wan, B. Shi, A. H. Tan, and A. C. Kot, "Depth of field guided reflection removal," in *Proc. ICIP*, 2016.
- [5] K. Gai, Z. Shi, and C. Zhang, "Blind separation of superimposed moving images using image statistics," *IEEE TPAMI*, 2012.
- [6] T. Xue, M. Rubinstein, C. Liu, and W. T. Freeman, "A computational approach for obstruction-free photography," *ACM TOG*, 2015.
- [7] H. Farid and E. H. Adelson, "Separating reflections and lighting using independent components analysis," in *Proc. CVPR*, 1999.
- [8] A. Agrawal, R. Raskar, S. K. Nayar, and Y. Li, "Removing photography artifacts using gradient projection and flash-exposure sampling," *ACM TOG*, 2005.
- [9] J. Bobin, Y. Moudden, J. L. Starck, and M. Elad, "Morphological diversity and source separation," *IEEE SPL*, 2006.
- [10] V. Abolghasemi, S. Ferdowsi, and S. Sanei, "Blind separation of image sources via adaptive dictionary learning," *IEEE TIP*, 2012.
- [11] W. Dong, L. Zhang, G. Shi, and X. Li, "Nonlocally centralized sparse representation for image restoration," *IEEE TIP*, 2013.
- [12] Y. Li, W. Dong, G. Shi, and X. Xie, "Learning parametric distributions for image super-resolution: Where patch matching meets sparse coding," in *Proc. ICCV*, 2015.
- [13] W. Zuo, L. Zhang, C. Song, D. Zhang, and H. Gao, "Gradient histogram estimation and preservation for texture enhanced image denoising," *IEEE TIP*, 2014.
- [14] C. Deledalle, V. Duval, and J. Salmon, "Non-local methods with shape-adaptive patches (nlm-sap)," *Springer JMIV*, 2012.
- [15] G. Peng and W. Hwang, "Reweighted and adaptive morphology separation," *SIAM SIIMS*, 2014.
- [16] M. Elad and I. Yavneh, "A plurality of sparse representations is better than the sparsest one alone," *IEEE TIT*, 2009.
- [17] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *Proc. CVPR*, 2007.
- [18] X. Zhang, "Matrix analysis and applications," *Tsinghua and Springer Publishing house, Beijing*, 2004.
- [19] D. Krishnan and R. Fergus, "Fast image deconvolution using hyper-laplacian priors," in *Proc. NIPS*, 2009.
- [20] H. Zhang, J. Yang, Y. Zhang, N. M. Nasrabadi, and T. S. Huang, "Close the loop: Joint blind image restoration and recognition with sparse representation prior," in *Proc. ICCV*, 2011.