

Single Image Super-Resolution via Projective Dictionary Learning with Anchored Neighborhood Regression

Yihui Feng ^{1,2}, Yongbing Zhang ¹, Yulun Zhang ^{1,2}, Tao Shen ^{1,2} and Qionghai Dai ^{1,2}

¹ Graduate School at Shenzhen, Tsinghua University, Shenzhen 518055, China

² Department of Automation, Tsinghua University, Beijing 100084, China

Abstract—We propose a novel single image super-resolution (SR) algorithm based on the projective dictionary pair learning with anchored neighborhood regression. Different from previous dictionary learning methods that aim to learn only a synthesis or an analysis dictionary, our method would learn both types of dictionaries jointly for regression to achieve image SR. We first cluster the training features into K clusters in order to learn synthesis and analysis dictionaries. Moreover, we learn the regressions with the training samples at training phase and use them on reconstruction stage. As shown in our experimental results, the proposed method obtains high-quality SR results quantitatively and visually against state-of-the-art methods.

Index Terms—Analysis dictionary, anchored regression, neighbor embedding, super-resolution, synthesis dictionary.

I. INTRODUCTION

Single image super-resolution (SR) is a very attractive and important area in computer vision, because it offers the promise of generating a high-resolution (HR) image from its degraded low-resolution (LR) measurement. However, image SR is still a severely ill-posed problem with the insufficient number of LR images and the ill-conditioned registration. Various regularization methods had been proposed to further stabilize the inversion of this ill-posed problem, such as fast and robust image SR [1] and joint map registration [2]. However, these reconstruction-based SR algorithms degraded rapidly when the number of available input images is small or the desired magnification is large. In these cases, the HR images may be over-smooth and lack important high-frequency details. Then interpolation [3], [4] were proposed to solve the problem and also achieved more favorable results. But in this way, they tended to generate overly smooth images with ringing and jagged artifacts. Now more and more SR algorithms (e.g. sparse coding [5] and the neighbor embedding [6], [7]) use additional information applying machine learning (ML) techniques to obtain better performance.

Anchored neighborhood regression (ANR) [6] and Adjust anchored neighborhood regression (A+) [7] are representative ML-based fast image SR algorithms. In ANR [6], Timofte et al. found the nearest neighbor in the employed synthesis LR-HR dictionary for each input feature and then precomputed projection matrix. Soon after, Timofte et al. proposed A+ [7],

an improved variant of ANR, which used the full training material instead of the dictionary pair to learn the regressions. However, both ANR [6] and A+ [7] only employed a synthesis dictionary to learn regressions. On the other hand, Hawe et al. [8] proposed a method that learned analysis operator and applied it to image reconstruction. The image SR methods above, which only learn one type of dictionary, may lose some important details in the obtained HR images.

Recently, Gu et al. [9] proposed projective dictionary pair learning (DPL) to learn a synthesis dictionary and an analysis dictionary jointly to achieve the goal of signal representation and pattern classification. Based on the idea of DPL [9], in this paper we propose a better image SR method where we learn the synthesis and analysis dictionaries jointly. The proposed method ensures that the representation coefficients can be approximated by a simple linear projection. And the analysis dictionary is trained to generate significant coding coefficients for samples, while the synthesis dictionary is trained to reconstruct the samples from their projective coefficients. Moreover, in the reconstruction phase, we adopt the same method as that in A+ [7] to calculate the regressions.

Overall, the main contributions of this paper can be summarized in three aspects. First, we extend the conventional discriminative synthesis dictionary learning to discriminative synthesis and analysis dictionary pair learning. Second, we learn the regressions using the full training material. Then, using the learned analysis dictionary, we confirm the class label of the input feature and synthesis dictionary for computing the nearest neighbor. Third, we can obtain more sophisticated HR image features with two types of dictionaries. Experimental results demonstrate that, our method yields high-quality SR images beyond recent advanced methods.

II. RELATED WORK

Because the ML-based SR methods have been acquiring superior results, we limit the related work to its dictionary learning and regression. Dictionary-learning based SR approaches build upon sparse coding (SC) [10]. Yang et al. [5] adopted a SC formulation to learn HR and LR dictionaries by assuming that HR and LR share the same reconstruction coef-

ficients. Zeyde et al. [11] improved the method by performing dimensionality reduction on the patches through PCA and Orthogonal Matching Pursuit (OMP) [12] to project the features onto a low-dimensional subspace. Then Timofte et al. used the dictionary training method proposed by Zeyde et al. [11] to learn regressions anchored to the dictionary atoms within a synthesis dictionary [6] or the whole training features [7]. Yang et al. [13] clustered the input data space and learned simple functions for each subspace. Next, Dai et al. [14] used a similar approach to jointly learn the separation into cells and regressions. Recently, more sophisticated SC formulations were proposed, for example, a collaborative method was adopted for fast SR in [15] and Zhang et al. [16] used a nonparametric regression model. Very recently, deep learning shows its power in image SR by learning hierarchical representation of high dimensional data. Cui et al. [17] proposed a deep network cascade (DNC) to gradually upscale LR images layer by layer. Because independent optimization of the self-similarity search process and auto-encoder was needed in each layer of the cascade, DNC [17] failed to obtain an end-to-end solution. Then Dong et al. [18] solved the problem by proposing a model named super-resolution convolutional neural network (SRCNN) that learned structure with different mapping layers.

In this paper, our work takes advantage of the projective dictionary pair learning and adjust anchored neighborhood regression simultaneously for fast image SR, and we obtain similar or better quality while making no compromise on time.

III. PROPOSED METHOD

We branch this section into two parts, the first one is the dictionary learning which concludes the conventional dictionary learning and our method that learns a pair of synthesis and analysis dictionaries jointly. Then we combine A+ [7] that has a better quality in image SR to develop our method in order to have good performance.

A. Dictionary Learning

In this sub-section, we will introduce the conventional dictionary learning roughly and explain the new dictionary learning in detail.

Most of SR methods learn a coupled dictionary \mathbf{D}_h and \mathbf{D}_l by minimizing the objective function with a sparsity constraint and forcing the HR and LR image features to share the same coefficients:

$$\min_{\mathbf{D}_h, \mathbf{D}_l, \mathbf{Z}} \|\mathbf{X}_c - \mathbf{D}_c \mathbf{Z}\|_F^2 + \lambda \|\mathbf{Z}\|_1, \quad (1)$$

where

$$\mathbf{X}_c = \begin{bmatrix} \frac{1}{\sqrt{N}} \mathbf{X}^h \\ \frac{1}{\sqrt{M}} \mathbf{Y}^l \end{bmatrix}, \mathbf{D}_c = \begin{bmatrix} \frac{1}{\sqrt{N}} \mathbf{D}_h \\ \frac{1}{\sqrt{M}} \mathbf{D}_l \end{bmatrix}, \quad (2)$$

and $\mathbf{X}^h = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ are the set of sampled HR image patches and $\mathbf{Y}^l = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$ are the corresponding LR image patches, N and M are the dimensionality of the low and high resolution patches and \mathbf{Z} is the coefficient vector representing the sparsity constraint.

Zeyde et al. [11] used K-SVD for the LR dictionary learning and direct approach using the pseudo-inverse for HR dictionary. Moreover, they made use of PCA to reduce the image patches. A+ [7] adopted this way to learn the dictionary and had a good performance.

In the introduction, we mention that Gu et al. [9] proposed a novel idea to learn synthesis and analysis dictionaries jointly and they learned the dictionary for pattern classification. Inspired by this, we make use of the idea to learn dictionaries for image SR. In the following discussion, we denote $\mathbf{\Omega}$ as the analysis dictionary and \mathbf{D} as the synthesis dictionary:

$$\{\mathbf{\Omega}^*, \mathbf{D}^*\} = \arg \min_{\mathbf{\Omega}, \mathbf{D}} \sum_{k=1}^K \|\mathbf{X}_k - \mathbf{D}_k \mathbf{\Omega}_k \mathbf{X}_k\|_F^2 + \lambda \|\mathbf{\Omega}_k \bar{\mathbf{X}}_k\|_F^2 \quad s.t. \|\mathbf{d}_i\|_2^2 \leq 1, \quad (3)$$

where \mathbf{X} is the training feature which has been applied by PCA dimensionality reduction and projected onto a low-dimensional subspace while preserving 99.9% of the average energy. And then we cluster them to K clusters. $\bar{\mathbf{X}}_k$ denotes the complementary data matrix of \mathbf{X}_k in the whole training set \mathbf{X} , λ is a scalar constant. \mathbf{d}_i is the i_{th} atom of the synthesis dictionary \mathbf{D} .

The objective function in (3) is generally non-convex. In order to solve this problem, we use a variable matrix \mathbf{A} and relax (3) to the following model:

$$\{\mathbf{\Omega}^*, \mathbf{A}^*, \mathbf{D}^*\} = \arg \min_{\mathbf{\Omega}, \mathbf{A}, \mathbf{D}} \sum_{k=1}^K \|\mathbf{X}_k - \mathbf{D}_k \mathbf{A}_k\|_F^2 + \tau \|\mathbf{\Omega}_k \mathbf{X}_k - \mathbf{A}_k\|_F^2 + \lambda \|\mathbf{\Omega}_k \bar{\mathbf{X}}_k\|_F^2 \quad s.t. \|\mathbf{d}_i\|_2^2 \leq 1, \quad (4)$$

where τ is a scalar constant. All terms in (4) are characterized by Frobenius norm and it can be easily solved. First, in order to solve this problem, we initialize the analysis dictionary $\mathbf{\Omega}$ and synthesis dictionary \mathbf{D} as random matrices with unit Frobenius norm and then alternatively update \mathbf{A} and $\{\mathbf{D}, \mathbf{\Omega}\}$. In this way, we can obtain the analysis dictionary $\mathbf{\Omega} = \{\mathbf{\Omega}_1; \mathbf{\Omega}_2; \dots; \mathbf{\Omega}_K; \dots; \mathbf{\Omega}_K\}$ and the synthesis dictionary $\mathbf{D} = \{\mathbf{D}_1; \mathbf{D}_2; \dots; \mathbf{D}_K; \dots; \mathbf{D}_K\}$, where $\{\mathbf{D}_k \in \mathbb{R}^{n \times d}, \mathbf{\Omega}_k \in \mathbb{R}^{d \times n}\}$ forms a sub-dictionary pair corresponding to cluster k . By the above algorithm, we can learn a synthesis and an analysis dictionary jointly to achieve image SR.

B. Image Reconstruction

This section we would introduce the image reconstruction based on the dictionary pair. We adopt the anchored neighborhood regression method used by A+ [7] for image SR and it had been making better results. We also use the training samples when calculating the regression during training time. Different from A+ [7], our method utilizes an analysis and a synthesis dictionary for image SR.

For each cluster, we can calculate a separate projection matrix $\mathbf{P}_{\{i,j\}}$ for each dictionary atom $\mathbf{d}_{\{i,j\}}$ representing the

j_{th} atom of the i_{th} cluster dictionary in the synthesis dictionary. The projection matrix can be calculated as:

$$P_{\{i,j\}} = X_h(Y_l^T Y_l + \lambda I)^{-1} Y_l^T, \quad (5)$$

where X_h is the set of sampled HR image patches and Y_l is the corresponding LR image patches, λ is a scalar constant.

After obtaining a separate projection matrix, the next work is to calculate the HR output patches \mathbf{x} . In order to find the nearest neighbors, we would use the analysis dictionary to confirm the cluster label for the input feature \mathbf{y} . Consequently, the reconstruction residual $\|\mathbf{y} - \mathbf{D}_k \Omega_k \mathbf{y}\|_2^2$ tends to be smaller than the residuals $\|\mathbf{y} - \mathbf{D}_i \Omega_i \mathbf{y}\|_2^2$, $i \neq k$. Based on this assumption, for each input feature we can find which cluster it belongs to so that we can find the best and nearest neighbors for reconstruction. The algorithm can be summarized as:

$$identity(\mathbf{y}) = \arg \min_i \|\mathbf{y} - \mathbf{D}_i \Omega_i \mathbf{y}\|_2. \quad (6)$$

After finding the i_{th} cluster, the SR problem can be solved by calculating for each input feature \mathbf{y}_i its nearest neighbor atom $\mathbf{d}_{\{i,j\}}$ in the dictionary, followed by the mapping to HR space using the stored projection:

$$\mathbf{x}_j = P_{\{i,j\}} \mathbf{y}_j. \quad (7)$$

Then patches calculated by equation (7) are added to the bicubically interpolated LR input image with overlapping parts averaged to create the output patch.



Fig. 1. 10 test images from left to right: **home**, **bandicoot**, **flower**, **dog**, **house**, **face**, **hair**, **brow**, **kangaroo**, and **car**.

IV. EXPERIMENTAL RESULTS

In this section, we evaluate our method based on the projective dictionary pair learning with adjust anchored neighborhood regression for single image SR. First, we will outline our experimental parameters including details in datasets, evaluation metrics, and then we show the performance comparisons with other state-of-the-art methods.

A. Experimental Parameters

In this sub-section we analyze the main parameters of our proposed method. For a fair comparison, we use 61 images as the training dataset. $4\times$ magnification is conducted on 10 commonly used LR natural images shown in Fig.1 including animals, plants and objects. We employ Peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) to evaluate the performance of SR results by different methods [6], [7], [11], [14], [18]. The implementations are all from the publicly available codes provided by the authors. All of these evaluation metrics are performed between the luminance channel of the original HR and the reconstructed image.

We down-sample the original HR images to generate LR images for training and testing on the luminance channel by Bicubic interpolation. In our method, the standard settings are upscaling factor $4\times$, the cluster coefficient of 10, 5 million training samples of LR and HR patches, the dictionary size of 1024, and a neighborhood size of 2048 training samples. The HR and LR features are similar to Zeyde's methods in [11].

B. Experimental Performance

Table I shows numerical evaluations in terms of PSNR and SSIM with upscaling factor of $4\times$. As shown in the table, we can see that our method has the best results on average for each evaluation metric. An obvious average PSNR gains of our algorithm over the second best method SRCNN [18] is 0.06dB and SSIM over JOR [14] is 0.001. Although the improvement is small, it is notable because above methods are the best learning-based image SR on the current stage. This indicates that analysis and synthesis dictionary mapping LR to HR features directly contributes to improving single image SR results.

TABLE I
QUANTITATIVE COMPARISONS FOR $4\times$ MAGNIFICATION. FOR EACH IMAGE, THERE ARE TWO ROWS: PSNR (dB) AND SSIM. THE BEST RESULT FOR EACH IMAGE IS HIGHLIGHTED.

Images	Bicubic	Zeyde [11]	GR [6]	ANR [6]	SRCNN [18]	JOR [14]	A+ [7]	Ours
home	20.91	21.10	21.05	21.11	21.09	21.16	21.15	21.15
	0.483	0.514	0.508	0.515	0.512	0.525	0.524	0.522
bandicoot	29.50	29.88	29.86	29.92	29.90	29.96	29.97	29.98
	0.745	0.766	0.769	0.769	0.766	0.770	0.770	0.771
flower	31.54	32.73	32.65	32.90	33.24	33.14	33.28	33.50
	0.867	0.885	0.883	0.887	0.885	0.892	0.893	0.894
dog	28.68	29.38	29.43	29.45	29.36	29.46	29.48	29.52
	0.819	0.843	0.846	0.846	0.843	0.847	0.846	0.847
house	29.27	30.16	29.98	30.19	30.40	30.30	30.43	30.59
	0.778	0.808	0.799	0.809	0.813	0.816	0.818	0.821
face	28.13	28.66	28.65	28.70	28.79	28.64	28.53	28.67
	0.731	0.755	0.756	0.757	0.756	0.763	0.759	0.760
hair	31.11	32.09	32.06	32.15	32.34	32.37	32.14	32.39
	0.788	0.823	0.825	0.826	0.828	0.831	0.823	0.831
brow	26.86	27.28	27.26	27.33	27.33	27.43	27.35	27.38
	0.635	0.670	0.676	0.677	0.677	0.684	0.680	0.681
kangaroo	26.98	27.33	27.36	27.37	27.36	27.37	27.34	27.42
	0.679	0.707	0.715	0.711	0.707	0.708	0.705	0.713
car	21.93	22.56	22.57	22.58	22.83	22.19	22.61	22.69
	0.640	0.684	0.681	0.685	0.694	0.676	0.690	0.693
average	27.49	28.12	28.09	28.17	28.27	28.21	28.23	28.33
	0.717	0.745	0.746	0.748	0.748	0.752	0.751	0.753

C. Visual Results

In order to have a visual comparison and demonstrate the effectiveness of our proposed method, we show the results in Figs. 2, 3, and 4 with different test images. We compare our visual results with those of state-of-the-art methods. As we can see from these visual results, Bicubic and K-SVD [11] would always produce dominant jaggy artifacts along the edges and visually displeasing blurred textural details. Figs. 2(d), 3(d), and 4(d) show the results obtained by ANR [6], which includes ringing artifacts around the over-sharped edges. SRCNN [18] produces result (e.g. Fig. 4(e)) that has unpleasing artifacts along the dominant edges. Although A+ [7] and JOR [14] have better performance in PSNR and SSIM, they can also generate blurred edges and fail to recover more detailed textures (e.g. Figs. 4(g) and 4(f)). We can observe that in Figs. 2(h), 3(h), and 4(h), our proposed method would achieve high-quality results more faithful to the original HR images with sharper edges and better details.

D. Investigation of Dictionary Size

As we have obtained the projective dictionary pair, where analysis dictionary $\{\Omega_k \in \mathbb{R}^{d \times n}\}$ is trained to produce small coefficients for samples from clusters other than k and synthesis dictionary $\{\mathbf{D}_k \in \mathbb{R}^{n \times d}\}$ is used to reconstruct the samples

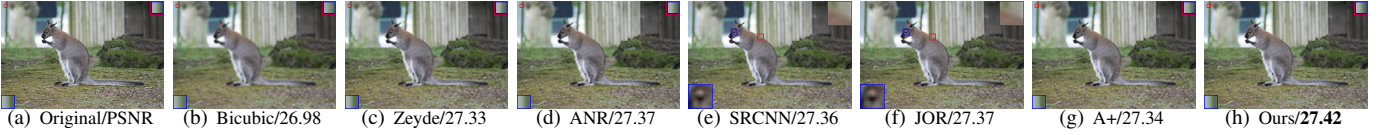


Fig. 2. Visual quality comparisons on **kangaroo** with a scaling factor 4. (Zoom in for better view.)



Fig. 3. Visual quality comparisons on **house** with a scaling factor 4. (Zoom in for better view.)

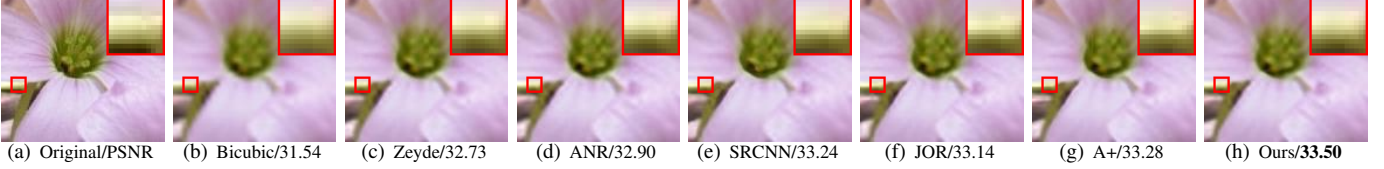


Fig. 4. Visual quality comparisons on **flower** with a scaling factor 4. (Zoom in for better view.)

REFERENCES

- [1] S. Farsiu, M. D. Robinson, M. Elad, and P. Milanfar, "Fast and robust multiframe super resolution." *IEEE Trans. Image Process.*, vol. 13, no. 10, pp. 1327–1344, 2004.
- [2] R. C. Hardie, K. J. Barnard, and E. E. Armstrong, "Joint map registration and high-resolution image estimation using a sequence of undersampled images," *IEEE Trans. Image Process.*, vol. 6, no. 12, pp. 1621–33, 1997.
- [3] H. S. Hou and H. Andrews, "Cubic splines for image interpolation and digital filtering," *IEEE Trans. Image Process.*, vol. 26, no. 6, pp. 508–517, 1979.
- [4] S. Dai, M. Han, W. Xu, Y. Wu, and Y. Gong, "Soft edge smoothness prior for alpha channel super resolution," in *CVPR*, Sep. 2014, pp. 184–199.
- [5] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE Trans. Image Process.*, vol. 19, no. 11, pp. 2861–2873, Nov. 2010.
- [6] R. Timofte, V. De Smet, and L. V. Gool, "Anchored neighborhood regression for fast example-based super-resolution," in *ICCV*, Dec. 2013, pp. 1920–1927.
- [7] R. Timofte, V. De Smet, and L. Van Gool, "A+: Adjusted anchored neighborhood regression for fast super-resolution," in *ACCV*, Nov. 2014.
- [8] K. D. Simon Hawe, Martin Kleinstaub, "Analysis operator learning and its application to image reconstruction," *IEEE Trans. Image Process.*, vol. 22, no. 6, pp. 2138–2150, 2013.
- [9] S. Gu, L. Zhang, W. Zuo, and X. Feng, "Projective dictionary pair learning for pattern classification," *Machine Learning*, vol. 1, pp. 793–801, 2014.
- [10] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed by v1?" *Vision research*, vol. 37, no. 23, pp. 3311–3325, Dec. 1997.
- [11] R. Zeyde, M. Elad, and M. Protter, "On single image scale-up using sparse-representations," in *Proc. 7th Int. Conf. Curves Surf.*, Jun. 2010, pp. 711–730.
- [12] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Trans. Inf. Theory*, vol. 53, no. 12, pp. 4655–4666, Dec. 2007.
- [13] C.-Y. Yang and M.-H. Yang, "Fast direct super-resolution by simple functions," in *ICCV*, Dec. 2013, pp. 561–568.
- [14] D. Dai, R. Timofte, and L. Van Gool, "Jointly optimized regressors for image super-resolution," in *Eurographics*, 2015.
- [15] Y. Zhang, Y. Zhang, J. Zhang, and Q. Dai, "Ccr: clustering and collaborative representation for fast single image super-resolution," *IEEE Trans. Multimedia*, vol. 18, no. 3, pp. 405–417, Mar. 2016.
- [16] Y. Zhang, Y. Zhang, J. Zhang, H. Wang, X. Wang, and Q. Dai, "Adaptive local nonparametric regression for fast single image super-resolution," in *Proc. IEEE Int. Conf. Visual Commun. Image Process.*, Dec. 2015, pp. 1–4.
- [17] Z. Cui, H. Chang, S. Shan, B. Zhong, and X. Chen, "Deep network cascade for image super-resolution," in *ECCV*, Sep. 2014, pp. 49–64.
- [18] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *ECCV*, Sep. 2014, pp. 184–199.

from cluster k . The previous mentioned experimental results show that projective dictionary pair is very effective in image SR. In those results, we fix d (the number of dictionary size used in SR) to be 1024. Intuitively, larger dictionaries should possess more expressive power and may yield more accurate results while increasing the computation cost. In this section, we present more results to investigate the influence of d on the performance of our method on the 10 test images. On the whole, in Table II, as d becomes larger, the values of PSNR, SSIM would increase. Moreover, when $d = 300$, much smaller than 1024 used in A+ [7], our proposed method outperforms all of other competing methods. Even though we set d as 100, the average performance of our method can also be better than many state-of-the-art methods. In conclusion, using the analysis and synthesis dictionaries jointly for image SR can achieve the satisfying results as well as saving computation.

TABLE II
AVERAGE PERFORMANCE OF OUR METHOD ON PSNR (DB) AND SSIM WITH DIFFERENT VALUES OF d .

d	100	200	300	400	500	600	700	800	900	1024
PSNR	28.25	28.27	28.28	28.29	28.29	28.32	28.31	28.31	28.30	28.33
SSIM	0.750	0.751	0.771	0.752	0.752	0.753	0.752	0.753	0.753	0.753

V. CONCLUSIONS

In this paper, we propose a novel method for image SR that we would cluster the training features and use them to learn synthesis and analysis dictionaries. Then we make use of the idea of A+ that uses the training samples to calculate regressions during training phase and uses these regressions to reconstruct HR patches at test phase. The results of the experiments have shown that our method obtains the best results and is superior to the state-of-art algorithms both quantitatively and visually.

ACKNOWLEDGMENT

This work was supported in part by the National Natural Science Foundation of China under Grant 61571254, Grant U1301257, and the Guangdong Natural Science Foundation 2014A030313751.