

Single Image Super-Resolution via Adaptive Dictionary Pair Learning for Wireless Capsule Endoscopy Image

Y. Wang¹, C. Cai², Y. X. Zou*¹

¹ADSPLAB/ELIP, School of ECE, Peking University, Shenzhen 518055, China

²Department of Computer Science, College of Information Engineering, Northwest A&F University, Yangling 712100, China

*{zouyx@pkusz.edu.cn}

Abstract—Wireless capsule endoscopy (WCE) is an innovative solution for gastrointestinal disease detection. Limited by WCE hardware and cost of manufacture, WCE image resolution is commonly low, which creates problems for attention to image details and visual perception in medical diagnosis. Under the sparse representation framework, we propose an adaptive dictionary pair learning method to obtain more appropriate representation of each patch with more relevant atoms according to patch content. Specifically, the dictionary pair is learned from high-low resolution cluster patches based on sparse constraint of input patches. Careful examination of the WCE images show there exist unnatural block areas. In order to further improve performance, the autoregressive model is applied to enhance local structure. Intensive experiments have been conducted on WCE image dataset and natural image dataset, including comparison test between the state-of-art methods and ours, and the results validate the effectiveness of the proposed method both on visual perception effect and objective indices.

Keywords—wireless capsule endoscopy; super-resolution; sparse representation; adaptive dictionary learning; autoregressive model

I. INTRODUCTION

Nowadays, with the improvement of diet, people get increasing threat from digestive system cancers. It has been reported that almost more than 140,000 are stroked by intestinal cancer in the US each year and killed approximately 50,000 [1]. A distinguishing feature of this kind of cancer is that it can be prevented and treated with early detection. As one of the most promising detection methodologies of gastrointestinal tract examination, WCE has been widely used due to its noninvasiveness, painlessness and patient-friendly. Unfortunately, the images captured by WCE is low resolution (LR) limited by its hardware (Fig. 1). In order to offer a higher resolution images to doctors for diagnosis, it is urged to enhance WCE images' resolution, which is called super-resolution (SR).

There are three types of SR methods. The long historic one is interpolation [2, 3], but it often leads to over-smooth image and removes the details. The second is to get high resolution (HR) image from multiple LR images of the same scene [4]. Krzysztof Duda et al. had applied it to WCE image

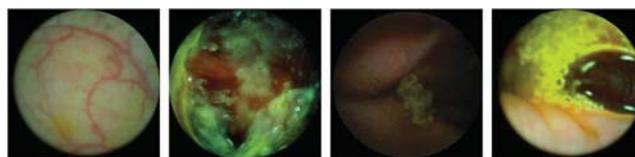


Figure 1. WCE image examples

SR for real-time applications [5]. According to [6], the essence of multi frame based SR is the registration methodology. Without better registration methods, no more improvement can be achieved by this type of SR technique. Besides, it emphasize more on speed than quality.

Besides, the more recent developed SR is learning based super-resolution approaches [7]. Compared with the first two types, it concerns on the mapping from LR patches to HR patches and the priors for constraint. Example-based [8], neighbor embedding [9] and their derivatives [10, 11] are classical methods describing the mapping. But these methods require high computational cost meanwhile they tend to produce watercolor-like artifacts. Some efforts have been made by employing the priors, such as gradient [12] or edge prior [13], to narrow down the solution domain. Their results on natural images are quite vivid but synthesized details may hinder the diagnosis for WCE image SR. Sparse representation is a technique which exploit the sparse structure of natural image region for image restoration. Yang et al. [7] had successfully applied it to SR and achieved satisfying result both in visual perception and objective indices. In their study, dictionary pair pre-trained is fixed for any input patch, which is poor for the expression of sparse representation. Inspired by the work in [14-17], we propose an adaptive dictionary pair learning technique. The HR and LR training patches are clustered by their HR features. Then it seeks the sparse representation of the input patches over a dictionary constituted by the centroids of normalized LR patches cluster, and the dictionary pair will be formed by the cluster normalized HR-LR patches based on the nonzero elements of the sparse representation. Based on the fact that the dictionary pair is produced adaptively based on the given patch content, we infer that this procedure will produce better result than Yang's. Through careful observation, there exist some block areas regions in WCE images, so we add the piecewise autoregressive (AR) constraint to enhance local image structure for refining the result further.

The rest of this paper is organized as follows. Sec. 2 describes our proposed adaptive dictionary pair learning method and the AR constraint. Sec. 3 presents the experimental results, including natural images and WCE images. Sec. 4 summaries the paper.

II. WCE IMAGE SUPER-RESOLUTION WITH ADAPTIVE DICTIONARY PAIR LEARNING

A. An Image Super-Resolution Model

Single image super-resolution aims to recover higher resolution image X from the given LR image Y . Y is down-sampled from X , which is generally modeled as:

$$Y = SX + v \quad (1)$$

where S denotes down-sample operator and v is additive noise. If v is normally distributed, we can compute X based on maximum-likelihood estimator as:

$$\hat{X}_{ML} = \arg \min_X p(Y | X) = \arg \min_X \|Y - SX\|_2^2 \quad (2)$$

which can be solved by gradient descent or iterative back-projection algorithm [7]. However, this will lead to infinite solutions due to its ill-posedness. In order to narrow down solution domain, the appropriate prior is needed.

Based on the fact that natural images are sparse on patch scale, sparse representation model is used widely in image processing. Here the sparse constraint can be used as prior to help solve (2). Assumed that the HR image patch shares the same sparse representation with the corresponding LR one over the HR-LR dictionary pair [7]. Therefore, (2) can be reformulated as:

$$\hat{X} = \arg \min_X \left\{ \|SX - Y\|_2^2 + \lambda \sum_i \|\alpha_i\|_1 + \delta \sum_i \|D_h \alpha_i - P_i X\|_2^2 \right\} \quad (3)$$

$$\alpha_i = \arg \min_{\alpha} \|D_l \alpha - y_i\|_2^2 + \lambda \|\alpha\|_1 \quad (4)$$

where α_i is p -dim column vector which is the representation for i th patch of X , P_i is a projection matrix which extract the i th patch from X and $P_i X$ is the HR patch corresponding to the LR patch y_i (m -dim column vector). D_h (n -by- p) is an over-complete dictionary ($n << p$) trained from HR image patches and D_l (m -by- p) is the corresponding dictionary trained from corresponding LR image patches. λ and δ are the regularization coefficients [7].

The SR result mainly depends on the quality of dictionary pair applied for sparse coding. In (3), a fixed dictionary pair is learned for general expression, and it lacks adaptability to image local feature. Motivated by the idea that dictionary constructed adaptively by the content of input patch in [15], we propose an adaptive dictionary pair learning (ADPL) method to create the more appropriate dictionary pair. Different from only one dictionary used during SR process in [15], we adopt the idea of HR-LR dictionary pair from [7]. It is believed that more relevant atom signals are more helpful for reconstructing details.

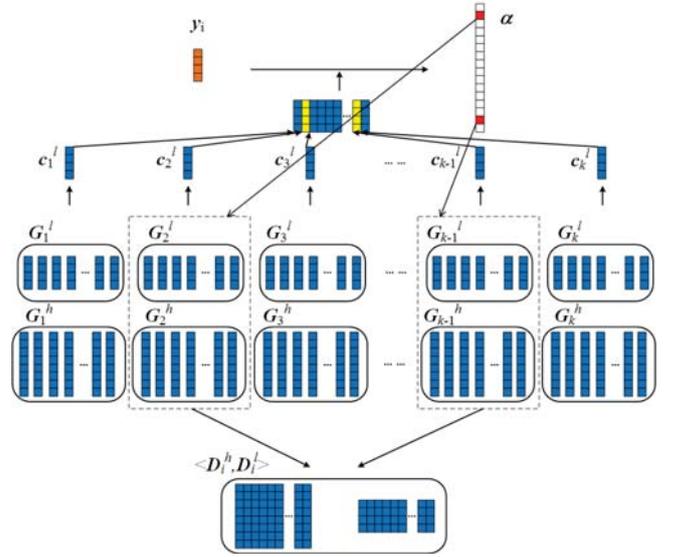


Figure 2. Adaptive dictionary pair learning process

Motivated by the idea of ADPL and (3), the ADPL representation model is formulated as follows:

$$\hat{X} = \arg \min_X \left\{ \|SX - Y\|_2^2 + \lambda \sum_i \|\alpha_i\|_1 + \delta \sum_i \|D_i^h \alpha_i - P_i X\|_2^2 \right\} \quad (5)$$

where D_i^h denotes the HR dictionary composed of auto-selected normalized HR patches for the i th patches y_i extracted from Y , and α_i is computed from y_i over D_i^l .

B. Adaptive Dictionary Pair Learning

The intuition of this method is to reduce the reconstruction error in (4) by adding extra flexibility to the dictionaries. Inspired by the idea that the adaptive dictionary construction and the dictionary selection method in [15, 17], we form the dictionary pair by selecting relevant normalized HR-LR patches. Suppose that HR image patches set is denoted as $S_h = \{h_1, h_2, \dots, h_m\}$ and LR image patches set is denoted as $S_l = \{l_1, l_2, \dots, l_m\}$. Their corresponding features set are denoted as $\tilde{S}_h = \{\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_m\}$ and $\tilde{S}_l = \{\tilde{l}_1, \tilde{l}_2, \dots, \tilde{l}_m\}$, respectively, where m is the number of training samples. Here \tilde{S}_i is the first and second derivative features set of S_i and \tilde{S}_h is equal to S_h .

In order to achieve atoms with better expression based on input LR patches, the approximate distribution of HR features is a must for selecting the suitable atoms further. So we cluster \tilde{S}_h into k groups as G_i^h ($i = 1, 2, \dots, k$), and the corresponding LR feature patches cluster G_i^l are generated by putting the patches together according to the index set of the patches in G_i^h . The centroid of each G_i^l is computed as:

$$c_i = \text{centroid}(G_i^l) = \sum_{\tilde{l} \in G_i^l} \frac{\tilde{l}}{|G_i^l|} \quad (6)$$

where \tilde{l} is the element of G_i^l . An auxiliary dictionary \hat{D} is formed as $\hat{D} = [c_1, c_2, \dots, c_k]$ to help select dictionary pair.



Figure 3. Comparison of SR results from different methods of image monarch (from left to right: input, bicubic, Elad's [16], proposed, original HR)

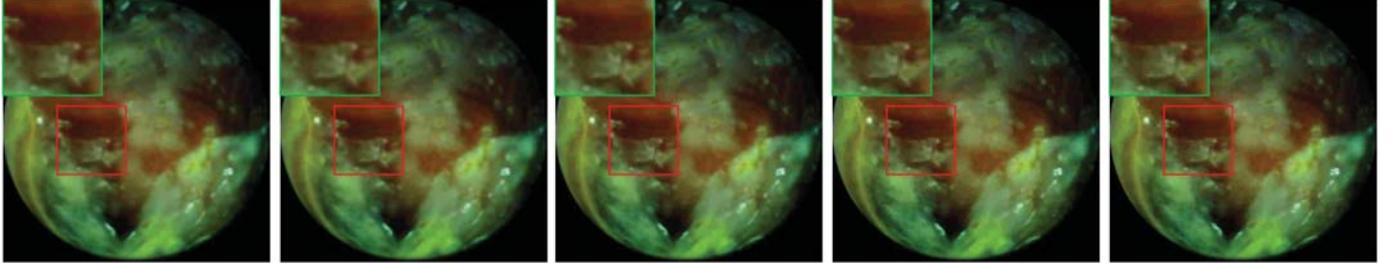


Figure 4. Comparison of SR results from different methods of WCE image J30312 (from left to right: input, bicubic, Yang's [7], proposed, original HR)

For the given LR patch \mathbf{y}_i , coarse screening is applied to $\hat{\mathbf{D}}$ based on the sparse constraint and obtain the approximate position of targeted atoms, which is denoted as

$$\alpha_i^{aux} = \arg \min_{\alpha} \|\hat{\mathbf{D}}\alpha - \mathbf{y}_i\|_2^2 + \lambda \|\alpha\|_1 \quad (7)$$

where λ is the regularization parameter.

We believe the more suitable atoms hidden in the clusters with the same indices as the non-zero elements in the computed sparse coding. In theory, it needs an exhaustive method to find the atoms, which is time-consuming. Given the tradeoff between time and preciseness, an efficient approach to locate these desired atoms is to combine these clusters into a dictionary and find the sparse coding under it. The formation of our dictionary pair is formulated as

$$[\mathbf{D}_i^h; \mathbf{D}_i^l] = \mathbf{B}\alpha_i^{aux} \quad (8)$$

where \mathbf{B} is a matrix constructed as $[\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_k]$ and $\mathbf{G}_i = [\mathbf{G}_i^h; \mathbf{G}_i^l]$.

Given the fact that the selection of atoms of $[\mathbf{D}_i^h; \mathbf{D}_i^l]$ is controlled adaptively by α_i^{aux} , which is determined by the feature of \mathbf{y}_i under $\hat{\mathbf{D}}$, the proposed ADPL would produce more relevant dictionary pair than [7] at the cost of more time.

Algorithm 1

Step 1 Input: a LR image \mathbf{Y}
 Step 2 For each 3×3 patch \mathbf{y} extracted from \mathbf{Y} pixel by pixel
 Learn the dictionary pair \mathbf{D}_i and \mathbf{D}_i^l from \mathbf{y} according to (7)(8)
 Compute the sparse representation $\hat{\alpha}$ by solving the optimization problem $\min \|\mathbf{D}_i\alpha - \mathbf{y}\|_2^2 + \lambda \|\alpha\|_1$
 Produce HR patch $\mathbf{x} = \mathbf{D}_i\hat{\alpha}$ and put it into a HR image \mathbf{X}_0
 End
 Step 3 Using gradient descent to improve the SR result which can satisfies the AR model constraint $\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \{\|\mathbf{S}\mathbf{X} - \mathbf{Y}\|_2^2 + \lambda \|\mathbf{X} - \mathbf{X}_0\| + \gamma \|\mathbf{P}_i\mathbf{X} - \mathbf{a}_{index_i, \mathbf{p}_i}\|_2^2\}$

Step 4 Output: SR image $\hat{\mathbf{X}}$.

C. Autoregressive Model

According to [15], a patch from natural images can be viewed as a fixed vector, so it can be well formulated by the autoregressive model. The autoregressive model is intended for exploiting the patch structure. In part B all HR patches are clustered into k groups and for each group \mathbf{G}_i^h the autoregressive model parameters are denoted as

$$\mathbf{a}_i = \arg \min_{\mathbf{a}} \sum_{\mathbf{h}_j \in \mathbf{G}_i^h} (h_j - \mathbf{a}^T \mathbf{n}_j)^2 \quad (9)$$

where h_j is the central pixel value of patch \mathbf{h}_j and \mathbf{n}_j is a vector containing the neighboring pixels of \mathbf{h}_j . By applying (9) to all \mathbf{G}_i^h , a set of AR models $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$ would be generated for later adaptive regularization.

The selection of the AR models for the given LR patch \mathbf{y}_i is formulated as

$$index_i = \arg \min_j \|\mathbf{y}_i - \mathbf{c}_j\|_2^2 \quad (10)$$

where \mathbf{c}_j is the centroid of \mathbf{G}_i^l as (6) and $index_i$ means \mathbf{a}_{index_i} would be assigned to \mathbf{y}_i for AR regularization. It is believed that the recovered HR patch \mathbf{x}_i from \mathbf{y}_i should follow the pattern of the $index_i$ th AR model, so the autoregressive regularization can be used to calculate the minimum square error between them as

$$e = \min \|\mathbf{x}_i - \mathbf{a}_{index_i, \mathbf{p}_i}\|_2^2 \quad (11)$$

where \mathbf{p}_i is the vector containing the neighboring pixels of HR patch \mathbf{x}_i .

Adding the constraint of (9), (5) can be expanded as

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \{\|\mathbf{S}\mathbf{X} - \mathbf{Y}\|_2^2 + \lambda \sum_i \|\alpha_i\|_1 + \delta \sum_i \|\mathbf{D}_i^h \alpha_i - \mathbf{P}_i \mathbf{X}\|_2^2 + \gamma \|\mathbf{P}_i \mathbf{X} - \mathbf{a}_{index_i, \mathbf{p}_i}\|_2^2\} \quad (12)$$

where $\mathbf{P}_i \mathbf{X}$ is actually \mathbf{x}_i .

TABLE I. THE RESULTS OF PSNR ON THE SET13 DATASET

Set14	Bicubic	Yang's	Elad's	ADPL	ADPL-AR
baboon	23.21	23.47	23.53	23.52	23.53
barbara	26.25	26.39	26.75	26.73	26.75
bridge	24.4	24.82	25.03	25.07	25.09
coastguard	26.55	27.02	27.15	27.13	27.14
comic	23.12	23.9	23.97	24.01	24.17
face	32.82	33.11	33.54	33.51	33.53
flowers	27.23	28.25	28.43	28.51	28.80
foreman	31.18	32.04	33.18	33.32	33.73
lenna	31.68	32.64	33.01	32.98	33.20
man	27.01	27.76	27.9	27.98	28.12
monarch	29.43	30.71	31.11	31.3	31.77
pepper	32.39	33.32	34.06	34.2	34.35
ppt3	23.71	24.98	25.23	25.49	25.68
zebra	26.63	27.95	28.5	28.68	28.82
Average	27.54	28.31	28.67	28.75	28.91

TABLE II. THE RESULTS OF PSNR ON THE WCE25 DATASET

WCE25	Bicubic	Yang's	Elad's	ADPL	ADPL-AR
J01	45.11	47.21	46.81	48.11	48.25
J02	49.29	49.49	49.54	49.45	49.52
J03	38.36	38.82	38.75	39.03	39.18
J14962	45.14	46.62	46.49	47.03	47.24
J16248	43.50	46.10	45.51	47.36	47.61
J18505	42.64	44.08	43.86	44.54	44.61
J22206	44.63	47.25	46.58	48.55	48.68
J22473	45.57	47.84	47.30	48.46	48.58
J26434	41.29	43.64	43.15	44.68	44.9
J27466	44.87	47.25	46.71	48.40	48.69
J27549	44.68	47.15	46.59	48.26	48.34
J27560	45.24	47.40	46.98	48.03	48.12
J30312	41.78	42.31	42.38	42.28	42.52
J30437	45.85	46.78	46.76	47.11	47.32
J30456	43.50	44.18	44.23	44.24	44.56
J30799	46.79	48.13	47.90	48.19	48.45
J31705	37.25	37.55	37.86	37.52	37.82
J35832	40.93	42.34	42.19	42.08	42.32
J35844	39.48	40.88	40.73	40.55	40.8
J40584	40.81	41.73	41.73	41.96	42.18
J40943	41.70	43.70	43.36	43.87	44.11
J41102	39.19	40.79	40.63	41.11	41.61
J52241	44.57	46.44	46.06	46.98	47.13
J7	41.65	43.46	43.15	43.99	44.18
J740	45.97	47.56	47.35	47.97	48.22
Average	43.19	44.75	44.50	45.19	45.40

The whole SR procedure is summarized as Algorithm 1.

III. EXPERIMENT RESULTS

In this section, several experiments are present to evaluate the performance of the proposed method. Some state-of-art methods like Yang's [7] and Elad's [16] are compared with our proposed method under the same settings, including the training data, test data and protocols. Specifically, there are 128 images for training, and it contains 91 natural images, 17 logo images, 10 license plate images and 10 randomly selected WCE images. The middle two kinds are intended for providing more edges information and WCE images are used to provide some patches for inferring their pattern. The size of patch extracted from images is set to 9×9 . The Set14 [16] is a natural image set

used to evaluate the performance by zoom factor 3x. The WCE25 consists of 25 images randomly selected from WCE images supplied from Shenzhen JiFu Technology Ltd., and is also used to evaluate the zoom factor 3x. Supplied WCE images are of size 480×480 with RGB channels. In the experiment only illuminance channel is processed, and the Cb, Cr color layers are directly bicubic interpolated.

For our method, the image dataset established ahead would be clustered into $g = 1600$ classes by using k-means. The sparsity of α^{aux} is set to 3.

A. Performance on natural images

From Table 1, it is clear to see that proposed ADPL-AR-SR outperforms other algorithms. It scores the highest PSNR value on the most images in classic Set13. Specifically, its PSNR value is higher about 0.24 dB than that of Elad's method and higher about 0.6 dB than that of Yang's method on average. Compared with bicubic, Yang's and Elad's methods on natural image from the Fig. 3, the proposed method produces more clear edges. In general, on natural images our method does recover more high frequency information (mostly edge details) and achieve better results than other methods.

B. Performance on WCE images

When processing WCE images, it can tell that the proposed method produce more clear SR result than Elad's and bicubic from Fig. 4. The magnified region marked by red of HR original image in Fig. 4 looks unnatural due to its obvious block areas and discontinuity of some edges. Through the SR algorithms, these defects seem weakened, and the effect of our method looks the most visual appeal among all, which can help diagnosis with better resolution and visual appeal WCE images. Table 2 demonstrates its effectiveness on objective evaluation. The PSNR computed based on ADPL-AR-SR is higher 0.65dB than Yang's and 0.9dB than Elad's method.

IV. CONCLUSION

This paper aims at developing a SR technique to enhance the review quality of WCE images. Motivated by the sparse representation, a sparse-based adaptive dictionary pair learning method is proposed based on input patch content. The autoregressive model has been employed to enhance the SR result. Intensive experimental results verify the effectiveness of the proposed method. It shows that the proposed ADPL-AR-SR method is able to recover texture, edge well and alleviate the blocky structure in WCE images, which makes the SR WCE images look more visual attractive and offers useful information for diagnosis.

V. ACKNOWLEDGEMENT

This work was partially supported by the Shenzhen Science & Technology Fundamental Research Program (No: JCYJ20130329175141512).

REFERENCES

- [1] R. Siegel, D. Naishadham, and A. Jemal, "Cancer statistics, 2013," *CA: A Cancer Journal for Clinicians*, vol. 63, pp. 11-30, 2013.
- [2] W. S. Dong, L. Zhang, R. Lukac, and G. M. Shi, "Sparse Representation Based Image Interpolation With Nonlocal Autoregressive Modeling," *IEEE Transactions on Image Processing*, vol. 22, pp. 1382-1394, Apr 2013.
- [3] S. C. Park, M. K. Park, and M. G. Kang, "Super-resolution image reconstruction: a technical overview," *Signal Processing Magazine, IEEE*, vol. 20, pp. 21-36, 2003.
- [4] X. Li, Y. Hu, X. Gao, D. Tao, and B. Ning, "A multi-frame image super-resolution method," *Signal Processing*, vol. 90, pp. 405-414, 2010.
- [5] K. Duda, T. Zielinski, and M. Duplaga, "Computationally simple super-resolution algorithm for video from endoscopic capsule," in *Signals and Electronic Systems, 2008. ICSES'08. International Conference on, 2008*, pp. 197-200.
- [6] S.-C. Lin and C.-T. Chen, "Reconstructing vehicle license plate image from low resolution images using nonuniform interpolation method," *International Journal of Image Processing (IJIP)*, vol. 1, pp. 21-28, 2008.
- [7] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *Image Processing, IEEE Transactions on*, vol. 19, pp. 2861-2873, 2010.
- [8] W. T. Freeman, T. R. Jones, and E. C. Pasztor, "Example-based super-resolution," *Computer Graphics and Applications, IEEE*, vol. 22, pp. 56-65, 2002.
- [9] C. Hong, Y. Dit-Yan, and X. Yimin, "Super-resolution through neighbor embedding," in *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on, 2004*, pp. I-I.
- [10] C. Xiaoxuan and Q. Chun, "Low-Rank Neighbor Embedding for Single Image Super-Resolution," *Signal Processing Letters, IEEE*, vol. 21, pp. 79-82, 2014.
- [11] M. Elad and D. Datsenko, "Example-based regularization deployed to super-resolution reconstruction of a single image," *The Computer Journal*, vol. 52, pp. 15-30, 2009.
- [12] J. Sun, Z. Xu, and H.-Y. Shum, "Image super-resolution using gradient profile prior," in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, 2008*, pp. 1-8.
- [13] T. Yu-Wing, L. Shuaicheng, M. S. Brown, and S. Lin, "Super resolution using edge prior and single image detail synthesis," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, 2010*, pp. 2400-2407.
- [14] G. Xinbo, Z. Kaibing, T. Dacheng, and L. Xuelong, "Image Super-Resolution With Sparse Neighbor Embedding," *Image Processing, IEEE Transactions on*, vol. 21, pp. 3194-3205, 2012.
- [15] D. Weisheng, D. Zhang, S. Guangming, and W. Xiaolin, "Image Deblurring and Super-Resolution by Adaptive Sparse Domain Selection and Adaptive Regularization," *Image Processing, IEEE Transactions on*, vol. 20, pp. 1838-1857, 2011.
- [16] R. Zeyde, M. Elad, and M. Protter, "On single image scale-up using sparse-representations," in *Curves and Surfaces*, ed: Springer, 2012, pp. 711-730.
- [17] D. Zhang, M. Yang, and X. Feng, "Sparse representation or collaborative representation: Which helps face recognition?," in *Computer Vision (ICCV), 2011 IEEE International Conference on, 2011*, pp. 471-478.