

Video Compressive Sensing via Structured Laplacian Modelling

Chen Zhao^{#1}, Siwei Ma^{#2}, Wen Gao^{#3}

[#]*School of Electronics Engineering and Computer Science, Peking University
Beijing, China*

¹zhaochen@pku.edu.cn

²swma@pku.edu.cn

³wgao@pku.edu.cn

Abstract— Seeking a fair domain in which the signal can exhibit high sparsity is of essential significance in compressive sensing (CS). Most methods in the literature, however, use a fixed transform domain or prior information, which cannot adapt to various video contents. In this paper, we propose a video CS recovery algorithm based on the structured Laplacian model, which can effectively deal with the non-stationarity of natural videos. To build the model, structured patch groups are constructed according to the nonlocal similarity in a temporal scope. By incorporating the model into the CS paradigm, we can formulate an ℓ_1 -norm optimization problem, for which a solution based on the iterative shrinkage/thresholding algorithms (ISTA) is designed. Experimental results demonstrate that the proposed algorithm outperforms the state-of-the-art methods in both objective and subjective recovery quality.

Index Terms— Video compressive sensing, structured Laplacian sparsity, iterative shrinkage/ thresholding, nonlocal similarity, discrete concrete transform

I. INTRODUCTION

Nowadays has seen tremendous interest in compressive sensing (CS), which provides the possibility of recovering a signal at sub-Nyquist rate [1], [2]. It declares that a signal with sparse representations under some domain can be reconstructed with high probability from very few measurements, which are obtained via linearly projecting the original signal onto a random basis.

CS theory indicates that the sparser the signal is in the specified domain, the higher recovery quality could be yielded. Thus seeking a desirable domain becomes one of the main challenges for the CS recovery research. Fixed domains or a fixed basis, e.g., Discrete Cosine Transform (DCT), wavelet and contourlet, gradient domain [4], [5], [6] are lately explored in a lot of CS recovery methods. Although using these domains is intuitively comprehensible, the restoration results are far from satisfactory due to their insufficiency in adaptively representing different signals.

To deal with this problem, it is suggested incorporating additional prior knowledge into the CS recovery paradigm in the current literature. Typical works are Gaussian scale mixtures (GSM) models [6], tree-structured wavelet [7] and tree-structured DCT (TSDCT). Additionally, in [8], a projection-driven CS recovery coupled with block-based

random image sampling is presented, aiming to encourage sparsity in the domain of directional transforms. Chen et al. [9] exploited multi-hypothesis predictions to generate residuals in the domain of the CS random projections.

All the methods above are designed for single-image recovery and can be simply applied to videos by independently reconstructing each frame. This is obviously not the optimal solution, neglecting the inter-frame correlations in video sequences. Alternatively, one can exploit the temporal dependency of a video and incorporate motion estimation (ME) and motion compensation (MC) into the CS reconstruction of the video while maintaining the same frame-by-frame sampling (e.g., [10], [11], [12]). An MC prediction of the current frame is created during reconstruction such that some image CS reconstruction algorithm is applied to the residual between the current frame and its prediction. The residual is typically believed more compressible than the original frame itself, leading to the state-of-the-art results for video compressive sensing recovery [11].

In this paper, we develop a video CS recovery algorithm via the structured Laplacian modelling in the DCT domain by fully exploiting the temporal dependency of natural videos. We have three major contributions. First, a video sampling-recovery framework is designed, where video frames could be recovered collaboratively. Second, a prominent Laplacian sparsity model is achieved by structuring the video patches according to the nonlocal similarity of natural videos. This model could be casted into the CS paradigm and an ℓ_1 optimization problem is formulated. Third, for solving the optimization problem, we design a solution based on the iterative shrinkage/thresholding algorithms (ISTA).

The remainder of this paper is organized as follows. Section II is an overview of related background, providing basic ideas of the compressive sensing theory. Section III gives a detailed description of the proposed video recovery

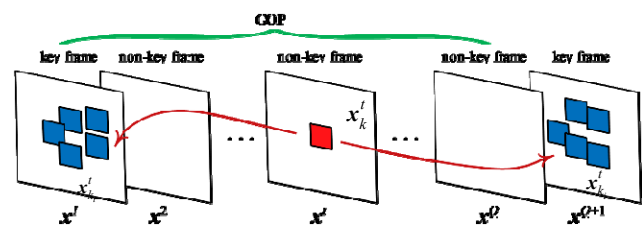


Fig. 1. Illustration of the GOP structure and the structured group construction.

algorithm and the solver design to the formulated optimization problem. Simulation results are provided in Section IV and conclusions are drawn in Section V.

II. BACKGROUND

Compressive sensing states that the signal can be recovered from very few samples, if it is sufficiently sparse in some domain Ψ . Concretely, supposing \mathbf{x} is the original signal of length N and \mathbf{y} is the sampling measurement of length M ($M \ll N$), the two of them satisfy $\mathbf{y} = \mathbf{A}\mathbf{x}$, in which \mathbf{A} is the random projection matrix. M/N is defined as the CS *sampling rate* or *subrate*. If the coefficients of $\boldsymbol{\alpha} = \Psi^T \mathbf{x}$ are mostly zeros or very close to zeros, where Ψ is a transform basis, \mathbf{x} could be restored through the following optimization problem

$$\min_{\mathbf{x}} \|\Psi^T \mathbf{x}\|_p, \quad \text{s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}. \quad (1)$$

In this formula, the subscript p is usually set to 0 or 1 to characterize the sparsity of the vector $\Psi^T \mathbf{x}$. When p is 0, the objective becomes an ℓ_0 norm $\|\cdot\|_0$, which counts the number of non-zero coefficients of \cdot ; when p is 1, the objective is an ℓ_1 norm $\|\cdot\|_1$, which adds up the absolute values of all coefficients of \cdot .

In practice, Eq. (1) is usually converted into an unconstrained problem by introducing a penalty parameter λ

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\Psi^T \mathbf{x}\|_p, \quad (2)$$

which can be regarded as a regularization-based paradigm for CS recovery.

Here, let us consider a group of Q consecutive frames in a video sequence, as shown in Fig. 1. We call these Q frames a group of pictures (GOP), in which the first frame is referred to as the key frame and the others the non-key frames. Each frame is sampled independently as a single image. The key frame has a relatively high *subrate* and the non-key frames has a low *subrate*.

For video recovery, the temporal dependency between frames are employed to collaboratively recover the non-key frames. The key frames are directly reconstructed using a single-image method thanks to its high subrate.

III. PROPOSED ALGORITHM

The proposed video CS recovery algorithm is composed of two steps. The first step is to apply an image CS recovery algorithm to reconstruct key frames. Then, in the second step, high-quality non-key frames are recovered via structured Laplacian modelling in the DCT domain with the help of the recovered key frames. In this section, we mainly focus on the recovery of the non-key frames.

A. Structured Laplacian Modelling in the DCT Domain

As illustrated in Fig. 1, the key frames denoted by \mathbf{x}^1 and \mathbf{x}^{Q+1} are first reconstructed by an image CS recovery method. Then the non-key frames $\mathbf{x}^2, \dots, \mathbf{x}^Q$ in between are to be recovered using their correlations with the already recovered key frames \mathbf{x}^1 and \mathbf{x}^{Q+1} . Here, we take the frame \mathbf{x}^l as an example for description of the modelling for the non-key frame recovery.

The frame \mathbf{x}^l is partitioned into D overlapped patches, which are denoted by the vector $\mathbf{x}_k^l, 1 \leq k \leq D$. Different from traditional methods that model the frame as a whole, we are dealing with each patch separately so that the non-stationarity problem of a frame is resolved and superior sparsity could be achieved. For each patch, we are constructing structured groups using similar patches from the nearest key frames \mathbf{x}^1 and \mathbf{x}^{Q+1} based on the motion consistency and nonlocal similarity [13] of video sequences.

In Fig. 1, the red square represents the patch \mathbf{x}_k^l and the blue squares are its most similar C patches in the two nearest key frames, which are notated by $\mathbf{x}_{k_i}^i, 1 \leq i \leq C$. The similarity between patches is measured by their mean squared errors (MSE). All the found similar patches and the patch \mathbf{x}_k^l itself form a group, denoted by $\mathbf{G}_k^l = [\mathbf{x}_k^l, \mathbf{x}_{k_1}^1, \dots, \mathbf{x}_{k_C}^1]$. Note that this process is conducted based an initial recovery of the frame \mathbf{x}^l and the recovered key frames from the first step.

Then, the coefficient patches are obtained via DCT as follows

$$\tilde{\mathbf{x}}_k^l = \Phi^T \mathbf{x}_k^l, \quad (3)$$

where Φ denotes the DCT basis.

Thereupon, the Laplacian distribution of the l -th coefficient in the patch $\tilde{\mathbf{x}}_k^l$ could be formulated as

$$P(\tilde{\mathbf{x}}_k^l(l)) = \frac{1}{\sqrt{2}\sigma_k^l(l)} \exp\left(-\frac{\sqrt{2}}{\sigma_k^l(l)} |\tilde{\mathbf{x}}_{i,k}^l(l) - \mu_k^l(l)|\right), \quad (4)$$

where

$$\begin{aligned} \mu_k^l(l) &= \sum_{1 \leq i \leq C} w_{k_i}^l \cdot \tilde{\mathbf{x}}_{k_i}^i(l), \\ \sigma_k^l(l) &= \sqrt{\frac{1}{C} \sum_{1 \leq i \leq C} (\tilde{\mathbf{x}}_{k_i}^i(l) - \mu_k^l(l))^2}. \end{aligned}$$

In Eq. (4), $w_{k_i}^l$ represents the weight measuring the similarity between $\tilde{\mathbf{x}}_{k_i}^i$ and $\tilde{\mathbf{x}}_k^l$, $\mu_k^l(l)$ and $\sigma_k^l(l)$ are the expectation and standard deviation of the coefficient.

Hence, assuming all coefficients in the patch are independent, their joint distribution can be described as follows

$$\begin{aligned} P(\tilde{\mathbf{x}}_k^l) &= P(\tilde{\mathbf{x}}_k^l(1))P(\tilde{\mathbf{x}}_k^l(2)) \dots P(\tilde{\mathbf{x}}_k^l(S^2)) \\ &= \prod_{1 \leq l \leq S^2} \frac{1}{\sqrt{2}\sigma_k^l(l)} \exp\left(-\frac{\sqrt{2}}{\sigma_k^l(l)} |\tilde{\mathbf{x}}_k^l(l) - \mu_k^l(l)|\right), \end{aligned} \quad (5)$$

where S^2 is the size of each patch. Then, by maximizing the probabilities of all the patches, which are assumed independent from one another, we can get the structured Laplacian sparsity prior

$$\begin{aligned} \max P(\tilde{\mathbf{x}}_1^l) P(\tilde{\mathbf{x}}_2^l) \dots P(\tilde{\mathbf{x}}_D^l) \\ = \max \prod_{1 \leq k \leq D} \prod_{1 \leq l \leq S^2} \frac{1}{\sqrt{2}\sigma_k^l(l)} \exp\left(-\frac{\sqrt{2}}{\sigma_k^l(l)} |\tilde{\mathbf{x}}_k^l(l) - \mu_k^l(l)|\right). \end{aligned} \quad (6)$$

By applying the natural logarithm, Eq. (6) is equivalent to the following expression

$$\begin{aligned} \min \sum_{1 \leq k \leq D} \sum_{1 \leq l \leq S^2} \frac{\sqrt{2}}{\sigma_k^l(l)} |\tilde{\mathbf{x}}_k^l(l) - \mu_k^l(l)| \\ = \min \sum_{1 \leq k \leq D} \|\boldsymbol{\tau}_k^l \circ (\tilde{\mathbf{x}}_k^l - \boldsymbol{\mu}_k^l)\|_1, \end{aligned} \quad (7)$$

where $\boldsymbol{\tau}_k^l$ and $\boldsymbol{\mu}_k^l$ are both vectors of length S^2 , $\boldsymbol{\tau}_k^l$ is composed of $\sqrt{2}/\sigma_k^l(l), 1 \leq l \leq S^2$ as its elements, and \circ stands for the Hadamard product of two vectors.

We now incorporate Eq. (7) as the regularization term into the optimization paradigm formulated by Eq. (2) and obtain

$$\min_x \frac{1}{2} \|y - Ax\|_2^2 + \lambda \sum_{1 \leq k \leq D} \|\tau_k' \circ (\tilde{x}_k' - \mu_k')\|. \quad (8)$$

B. Solver Design to the Optimization Problem

Having achieved the model to represent the Laplacian sparsity prior of each non-key frame, we next need to solve the optimization problem in Eq. (8). In the following description, the superscript t is omitted without confusion.

We use the iterative shrinkage/ thresholding algorithm (ISTA) [14] to solve Eq. (8) by iteratively projecting and thresholding. Referring to this method, we devise two iterative steps for Eq. (8).

$$r^{(j)} = x^{(j)} + \rho A^T (y - Ax^{(j)}), \quad (9)$$

$$x^{(j+1)} = \arg \min_x \frac{1}{2} \|x - r^{(j)}\|_2^2 + \lambda \sum_{1 \leq k \leq D} \|\tau_k' \circ (\tilde{x}_k - \mu_k)\|, \quad (10)$$

where ρ is a constant stepsize and j denotes the iteration number.

So the key to solving Eq. (8) is to find an efficient way to solve Eq. (10). As in [15], we assume that each element of $x - r^{(j)}$ is i.i.d. with an zero mean and the same variance. According to Theorem 1 in [15], $\sum_{1 \leq k \leq D} \|x_k - r_k^{(j)}\|_2^2$ and $\|x - r^{(j)}\|_2^2$ satisfy the following equation with a very large probability

$$\|x - r^{(j)}\|_2^2 = \frac{N}{K} \sum_{1 \leq k \leq D} \|x_k - r_k^{(j)}\|_2^2, \quad (11)$$

where $K = D \times S^2$ and N is the total pixel number of the frame.

Based on the Plancherel theorem, the energies in the space domain and the DCT frequency domain should be conserved. Then we have the equation below

$$\|x_k - r_k^{(j)}\|_2^2 = \|\Phi^T x_k - \Phi^T r_k^{(j)}\|_2^2 = \|\tilde{x}_k - \tilde{r}_k^{(j)}\|_2^2. \quad (12)$$

Combining Eq. (11) and Eq. (12) with Eq. (10), we get

$$\begin{aligned} \min_x \frac{1}{2} \frac{N}{K} \sum_{1 \leq k \leq D} \|\tilde{x}_k - \tilde{r}_k^{(j)}\|_2^2 + \lambda \sum_{1 \leq k \leq D} \|\tau_k' \circ (\tilde{x}_k - \mu_k)\| \\ = \min_x \sum_{1 \leq k \leq D} \left(\frac{1}{2} \|\tilde{x}_k - \tilde{r}_k^{(j)}\|_2^2 + \frac{\lambda K}{N} \|\tau_k' \circ (\tilde{x}_k - \mu_k)\| \right), \end{aligned} \quad (13)$$

which can be decomposed into D subproblems as follows

$$\tilde{x}_k^{(j+1)} = \arg \min_{\tilde{x}_k} \frac{1}{2} \|\tilde{x}_k - \tilde{r}_k^{(j)}\|_2^2 + \frac{\lambda K}{N} \|\tau_k' \circ (\tilde{x}_k - \mu_k)\|. \quad (14)$$

Obviously, Eq. (14) can be equivalently solved in an element-wise manner, i.e.,

$$\tilde{x}_k^{(j+1)}(l) = \arg \min_{\tilde{x}_k(l)} \frac{1}{2} (\tilde{x}_k(l) - \tilde{r}_k^{(j)}(l))^2 + \theta_k(l) |\tilde{x}_k(l) - \mu_k(l)|, \quad (15)$$

where $1 \leq l \leq S^2$, $\theta_k(l) = \lambda K \tau_k'(l) / N$.

By means of the soft thresholding algorithm of Lemma 2 in

[16], we can arrive at a closed-form solution for Eq. (15) as is formulated below

$$\tilde{x}_k^{(j+1)} = \text{soft}(\tilde{r}_k^{(j)} - \mu_k, \theta_k) + \mu_k. \quad (16)$$

Therefore, the corresponding patch in the space domain is

$$x_k^{(j+1)} = \Phi \tilde{x}_k^{(j+1)}. \quad (17)$$

This process is applied for all the overlapped patches, which are weightedly summed up to finally recover the entire frame.

In light of all derivations above, the complete description of video CS recovery for the non-key frames using structured Laplacian model (SLM) in the DCT domain is given below:

TABLE I NON-KEY FRAME x^t CS RECOVERY USING SLM

Input: The observed measurement y , the measurement matrix A and parameter λ ;

Initialization: set initial estimate $x^{(0)}$ using a single image recovery method;

for Iteration number $j = 0, 1, 2, \dots, \text{Max_iter}$

Get $r^{(j)}$ by computing Eq. (9);

for each patch x_k

Construct the structured group G_k by searching for similar patches in the nearest two key frames x^l and x^{l+1} ;

Calculate all μ_k and σ_k in Eq. (4);

Obtain $x_k^{(j+1)}$ by computing Eq. (16) and Eq. (17);

Recover the entire frame $x^{(j+1)}$;

end for

end for

Output: Final recovered frame \hat{x}^t .

IV. EXPERIMENTAL RESULTS

In this section, experimental results are presented to evaluate the performance of the proposed video CS recovery algorithm SLM. Four standard Cif (352×288) video sequences are tested, *Foreman*, *Mobile*, *Bus*, and *Akiyo*. In our experiments, the CS measurements are obtained by applying a Gaussian random projection matrix to the original video frames at the block level, i.e., block-based CS with block size of 16×16 [16], with key-frame *subrate* 0.7 and non-key frame *subrate* 0.2. The proposed SLM is compared with three representative CS recovery methods in literature, i.e., wavelet method (DWT) [8], multi-hypothesis (MH) method [9] and motion compensation/motion estimation (MC/ME) method [11]. It is worth emphasizing that MC/ME is known as the state-of-the-art algorithm for video CS recovery.

In our implementation, all the parameters of SLM are set empirically for all test sequences. Concretely, the size of each patch S^2 is set 8×8 , the area for searching similar patches is 20×20 , the number of selected similar patches C is 10 and the value for λ is 2.24. It is necessary to stress that the choice for all the parameters can be generalized to other natural videos, which has been verified in our experiments. In this paper, the

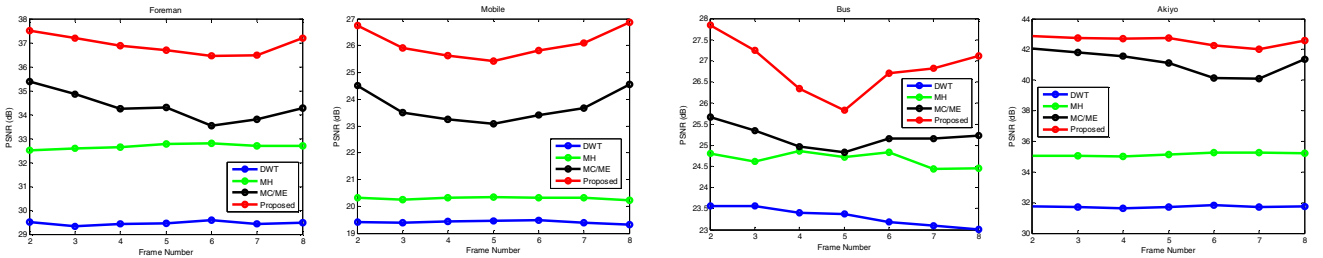


Fig.2. CS recovery Performance of non-key frames in the first GOP for various algorithms as respect to four sequences.

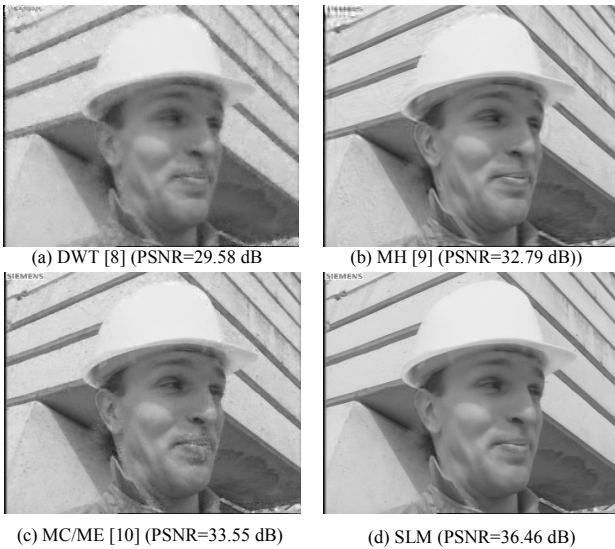


Fig. 3. Visual comparison of CS recovered results for 6th frame in *Foreman* by different methods.

results of MC/ME are exploited to recover the key frames and initialize the non-key frames for the proposed SLM. Hence, MC/ME and SLM have the same CS recovery results for key frames. Therefore, we only provide the CS recovery results of non-key frames below.

Figure 2 illustrates the CS recovery performance of non-key frames in the first GOP for the four algorithms as respect to four sequences. It is obvious to see that the proposed algorithm (red) and MC/ME (black), which exploit temporal dependency of videos, achieve much better results than single image CS recovery algorithms, i.e., DWT (blue) and MH (green). Furthermore, the proposed SLM obtains about 2 dB gains over MC/ME in PSNR on average, which is attributed to the structured Laplacian modelling based on video nonlocal similarity, which offers a powerful mechanism of characterizing the structured sparsity of natural video signals.

The visual results of the recovered 6th frame of *Foreman* by the four algorithms are presented in Fig. 3. Obviously, our proposed SLM algorithm not only yields the highest objective score in PSNR, but also preserves the fine details in the frames and shows much clearer and better visual results than the other comparative methods.

The complexity of SLM is provided as follows. Assume for each patch the average time to search for similar patches is T and its DCT operation is $O(S^2 \log(S))$. Hence, the total complexity of SLM is $O(N(S^2 \log(S) + T))$. For a frame with size 352×288 , the proposed CS recovery algorithm requires about 3~4 minutes for on an Intel Core2 Duo 2.96G PC under Matlab R2011a environment.

V. CONCLUSIONS

In this paper, we design an algorithm for video CS recovery by the structured Laplacian modelling in the DCT domain. The nonlocal similarity of natural videos is exploited during structuring the video patches, and superior sparsity could be achieved in this way. This prior information is incorporated into the CS paradigm and an ℓ_1 optimization problem is

formulated. Additionally, for solving the optimization problem generated from the techniques above, we design an efficient solution based on the iterative shrinkage/thresholding algorithms (ISTA).

Experimental results prove that the proposed algorithm can beat the other methods with very high objective gain as well as generating better visual video frames. This work can be used in scenarios that require high-quality video reconstruction, e.g., distributed video coding, Magnetic resonance image CS recovery.

ACKNOWLEDGMENT

This work is supported in part by the National High-tech R&D Program of China (863 Program, 2012AA010805), National Science Foundation (61322106, 61390515, and 61103088).

REFERENCES

- [1] E. J. Candes et al., "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Image Process.*, vol. 52, no. 2, pp. 489–509, 2006.
- [2] D. L. Donoho, "Compressed sensing," *IEEE Trans. Image Process.*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [3] M. Davenport, M. F. Duarte, Y. C. Eldar and G. Kutyniok, *Compressed Sensing: Theory and Applications*, Cambridge University Press, 2011.
- [4] l1-Magic Toolbox available at <http://users.ece.gatech.edu/~justin/l1magic/>.
- [5] C. Li, W. Yin, and Y. Zhang, "TVAL3: TV Minimization by Augmented Lagrangian and Alternating Direction Algorithm," 2009.
- [6] Y. Kim, M. S. Nadar and A. Bilgin, "Compressive sensing using a Gaussian scale mixtures model in wavelet do-main," *Proc. of IEEE Int. Conf. Image Process.*, pp. 3365–3368, 2010.
- [7] L. He and L. Carin, "Exploiting structure in wavelet-based Bayesian compressive sensing," *IEEE Trans. Signal Process.*, vol. 57, no. 9, pp. 3488–3497, 2009.
- [8] S. Mun and J. E. Fowler, "Block compressive sensing of images using directional transforms," *Proc. of IEEE Int. Conf. Image Process.*, pp. 3021–3024, 2009.
- [9] C. Chen, E. W. Tramel, and J. E. Fowler, "Compressed-Sensing Recovery of Images and Video Using Multihypo-thesis Predictions," *Proc. of the 45th Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, CA, pp. 1193–1198, Nov. 2011.
- [10] T. T. Do, Y. Chen, D. T. Nguyen, N. Nguyen, L. Gan, and T. D. Tran, "Distributed compressed video sensing," *Proc. of International Conference on Image Processing*, Cairo, Egypt, November 2009, pp. 1393–1396.
- [11] S. Mun and J. E. Fowler, "Residual reconstruction for block-based compressed sensing of video," *Proc. of IEEE Data Compression Conference*, Snowbird, UT, March 2011, pp. 183–192.
- [12] J. Prades-Nebot, Y. Ma, and T. Huang, "Distributed video coding using compressive sampling," in *Proceedings of the Picture Coding Symposium*, Chicago, IL, May 2009.
- [13] A. Buades, B. Coll, and J. M. Morel, "A non-local algorithm for image denoising," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2005, pp. 60–65.
- [14] M. Afonso, J. Bioucas-Dias and M. Figueiredo, "Fast image recovery using variable splitting and constrained op-timization," *IEEE Trans. on Image Process.*, vol. 19, no. 9, pp. 2345–2356, 2010.
- [15] J. Zhang, C. Zhao, D. Zhao, and W. Gao, "Image compressive sensing recovery using adaptively learned sparsifying basis via L0 minimization," *Signal Processing* (2013), <http://dx.doi.org/10.1016/j.sigpro.2013.09.025>.
- [16] J. Zhang, D. Zhao, C. Zhao, R. Xiong, S. Ma and W. Gao, "Image compressive sensing recovery via collaborative sparsity," *IEEE J. on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 3, pp. 380–391, Sep. 2012.