

# Coupled Sparse Representations for Multi-Focus Image Fusion

Rui Gao, *Student Member, IEEE*, Sergiy A. Vorobyov, *Senior Member, IEEE*, and Hong Zhao

## Abstract

We address the multi-focus image fusion problem, where multiple images captured with different focal settings are to be fused into an all-in-focus image. Algorithms for this problem necessarily admit the source image characteristics along with focused and blurred feature. However, most sparsity-based approaches use a single dictionary in focused feature space to describe multi-focus images, and ignore the representations in blurred feature space. Here, we propose a novel multi-focus image fusion approach using coupled sparse representations. The approach exploits the facts that (i) the patches in given training set can be sparsely represented by a couple of overcomplete dictionaries related to the focused and blurred categories of images; and (ii) merging such representations is better than just selecting the sparsest one in the estimate of the original signal. By jointly learning the coupled dictionary, we enforce the similarity of sparse representations in the focused and blurred feature spaces, and then introduce a new fusion approach to combine these representations for generating an all-in-focus image. We also discuss the advantages of the fusion approach based on coupled sparse representations and present an efficient algorithm for learning the coupled dictionary. Extensive experimental comparisons with state-of-the-art multi-focus image fusion algorithms validate the effectiveness of the proposed approach. As a demonstration of these benefits, we present high-quality image fusion results based on this new approach.

## Index Terms

R. Gao is with Northeastern University, Dept. Computer application technology, Shenyang 110819, China. She is also with Aalto University, Dept. Signal Processing and Acoustics. E-mail: gaorui@research.neu.edu.cn

S. A. Vorobyov is with Aalto University, Dept. Signal Processing and Acoustics, FI-00076, AALTO, Finland. E-mail: svor@ieee.org

H. Zhao is with Northeastern University, Dept. Computer application technology, Shenyang 110819, China. E-mail: zhao@neusoft.com

Image fusion, coupled sparse representations, dictionary learning, multi-focus image.

## I. INTRODUCTION

Over the last several decades, much research has been devoted to multi-focus image fusion problem [1]–[5]. Multi-focus image fusion is an effective technique for combining multiple images captured with different focal distances into an all-in-focus image, without sacrificing image quality or using specialized optic sensors [6]. The problem is of high importance in various applications, such as remote sensing, defense systems, and medical imaging [7]–[9]. One important paradigm for addressing the problem is fusion processing of multi-focus images. The majority of literature on the topic is traditionally categorized into two basic approaches [10]: the *spatial frequency-based* approaches and the *transform domain-based* approaches. In the first category, the methods such as image fusion based on Laplacian pyramid (LP) [11], spatial frequency (SF) [12], multi-scale weighted gradient (MWG) [13], and variance [10] directly select the best pixels or regions to fuse multiple images. The main drawback of these methods is that they may produce blocking artifacts because of the misalignment of decision map with boundary of the focused object and wrong decision in sub-focused regions. The second category suggests to apply multiscale transforms to decompose source images, and construct an all-in-focus image in the inverse transform domain. Algorithms of this type include discrete wavelet transform (DWT) [14], curvelet transform (CVT) [7], non-subsampled contourlet transform (NSCT) [15] among others. However, all aforementioned traditional methods are still sensitive to image misregistration, and they may suffer from undesirable artifacts.

Until recently, sparsity and overcompleteness have been successfully used for image fusion [16]–[23]. The approaches exploit the fact that natural images can be compactly represented on an overcomplete dictionary as a linear combination of sparse coefficients. Formally, the basic model suggests that the signal  $\mathbf{x} \in \mathbb{R}^d$  can be described as being a linear combination of few *atoms* over an overcomplete dictionary  $\mathbf{D} \in \mathbb{R}^{d \times n}$ . The processing can be formulated as sparse representation

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_0 \quad \text{s.t. } \mathbf{x} \approx \mathbf{D}\boldsymbol{\alpha} \quad (1)$$

where  $\boldsymbol{\alpha}$  is the sparse vector of  $\mathbf{x}$  applied to *atoms* in  $\mathbf{D}$ , and  $\|\cdot\|_0$  is the  $l_0$ -norm which counts the non-zero entries. As the problem (1) is known to be NP-hard problem because of its combinational nature, a suboptimal strategy for the above problem is replacing  $l_0$ -norm

with  $l_1$ -norm and learning adaptive dictionary from the training data. Many image processing applications have benefited remarkably from learned dictionary. Representative algorithms on this topic include the K-SVD [24], the method of optimal directions (MOD) [25], online dictionary learning (OLD) [26], and others. “Good” dictionaries are expected to be highly adaptive to the observed signals and to contribute accurate sparse representations. More recently, some studies have tracked the double feature space representation problem via learning coupled dictionary [27]–[30]. These approaches explicitly learn the pair of dictionaries from a training set of double images, and formulate sparse and redundant representations according to various applications and scenarios. The combination of learned coupled dictionary and sparse approximations proved the superior in the representations of double feature spaces [31]–[33].

Current image fusion methods based on sparse representations directly present one overcomplete dictionary in a single feature space to describe multi-focus images which contain the focused and blurred categories of features. The methods ignore the sparse representations in blurred feature space, and set the limits on the sparsity of the coefficients. That consequently leads to inaccuracy in the fusing coefficients. The difficulties on working in a single feature space motivate us to perform the fusion operations on double feature spaces. Similar to the sparsity-based methods, we will rely on the sparse coefficients from multi-focus source images. The success of such approach depends on how accurately it describes the multi-focus images and how efficiently it fuses the coefficients. Instead of directly learning one overcomplete dictionary from the focused images, this paper represents the first successful attempt to use coupled sparse representations into focused and blurred feature spaces, significantly improving the performance of image fusion method. Coupled sparse representations correlate the focused and blurred feature spaces, and bridge the fusion processes of double feature spaces.

Here, we propose a novel approach towards the multi-focus image fusion problem that utilizes coupled sparse and redundant representations over learned dictionary pair. The paper presents both the algorithmic developments and simulation results in multi-focus image fusion. Compared to prior work, the main difference of the proposed approach is that the coupled dictionary is exploited for more compact and accurate representations of different focus images. Such sparse representations will be shown to lead to higher quality results for fusion problem.

### A. Contributions

Given the pair of training image sets, we seek for a couple of overcomplete dictionaries that lead to more compact representations for the focused and blurred categories of images. Based upon the coupled dictionary, we then use the plain averaging fusion criteria to find a more accurate sparse representation for reconstructing an all-in-focus image.

The main contributions of this work are threefold:

1) We cast the multi-focus image fusion problem formulation for achieving an all-in-focus image representation, based on sparsity over a couple of dictionaries. Information from different feature spaces is fused by using joint sparse coefficients. The main difficulties in the formulation are learning of coupled dictionary and fusing of corresponding sparse representations. In our conference paper [34], the first contribution has been presented.

2) We define a new fusion criteria that merging sparse coefficients over the couple of dictionaries. This criteria forces to produce a better estimated representation by plain averaging. As such, this contribution differs from previous approaches on the basis of used sparsity model. We also introduce the local sparsity in a global Bayesian reconstruction framework. The proposed criterion improves performance of multi-focus image fusion compared to some of the competing methods.

3) The K-SVD-based coupled dictionary learning algorithm is proposed such that we explicitly enforce the requirement that sparse approximations evaluated for double feature spaces can generate an all-in-focus image using their corresponding dictionaries. The proposed optimization algorithm simultaneously trains the focused and blurred dictionaries under the joint sparse coding in order to enforce collaborations between the dictionaries. A novel coupled dictionary learning algorithm alternates between the update of dictionaries atoms and sparse coding.

### B. Organization of the Paper and Notation

The remainder of this paper is organized as follows. Section II gives the problem formulation and summarizes some general assumptions. Section III gives the detailed description of our proposed fusion criteria. In Section IV, we focus on jointly learning coupled dictionary. A K-SVD-based coupled dictionary learning algorithm is presented. Simulation results are provided in Section V. Finally, we conclude the paper in Section VI with a summary of this work and a description of its future extensions.

We generally use bold letters for vectors and hold capital letters for matrices.  $\mathbb{R}^d$  and  $\mathbb{R}^{d \times n}$  denote the set of real  $d \times 1$  vectors and set of real  $d \times n$  matrices, respectively. The notation  $[\mathbf{x}]_{i,j}$  have been used to specify the element located in the  $i^{th}$  row and  $j^{th}$  column of the matrix  $\mathbf{X}$ . We use  $p(\mathbf{x})$  to denote the probability density function (pdf) of random quantity  $\mathbf{x}$ , and use  $\mathbf{X}^T$  to denote transpose. The vector norms  $\|\cdot\|_p$  for  $p \geq 1$  are the standard  $l_p$ -norms. The so-called  $l_0$ -norm, which counts the nonzero entries of its argument, is denoted by  $\|\cdot\|_0$ . The matrix norms  $\|\cdot\|$  and  $\|\cdot\|_F$  are used to denote the operator norm and the Frobenius norm, respectively. The symbol  $\odot$  represents element-wise product of matrices, and the operator  $\langle \cdot, \cdot \rangle$  means the inner product of vectors.

## II. PROBLEM FORMULATION

Consider the scenario with multi-focus source images  $\{\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_n\}$  that acquire with different focal parameters. Our objective is to recover an all-in-focus image  $\mathbf{I}_F$  from the observations  $\{\mathbf{I}_k\}_{k=1}^n$  in the same scene. Assume that the observations  $\{\mathbf{I}_k\}_{k=1}^n$  are properly aligned, and the sequence of multi-focus images and the unknown all-in-focus image are modelled as [3], [27], [35]

$$\mathbf{I}_k = \mathcal{F}(\mathbf{I}_F) \quad (2)$$

where the operator  $\mathcal{F}$  is a blurring function [36], [37]. The mathematical model describes that the physical process of capturing  $n$  multi-focus images can be represented as the convolution of the all-in-focus image.

Assuming statistically independent observations, we now define the optimal solution for  $\mathbf{I}_F$  based on maximum *a posteriori* probability (MAP)

$$\mathbf{I}_F = \underset{\mathbf{I}_F}{\operatorname{argmax}} p(\mathbf{I}_F | \{\mathbf{I}_k\}_{k=1}^n) \quad (3)$$

as the minimizer of a well-defined global penalty term.  $p(\mathbf{I}_F | \{\mathbf{I}_k\}_{k=1}^n)$  models the observation mechanism that the estimate  $\mathbf{I}_F$  is from the observations  $\{\mathbf{I}_k\}_{k=1}^n$  which are statistically coupled to the transform outputs.

This paper has at the disposal two pieces of knowledge: the first one is that the patches in the training sets can be sparsely represented by a couple of overcomplete dictionaries, and the other one is that merging multiple sparse coefficients on one signal is better than just selecting the sparsest one alone [38]. Analysing the features on multi-focus images, we objectively define two

feature spaces: the focused feature space and the blurred feature space. A couple of overcomplete dictionaries  $\{\mathbf{D}^F, \mathbf{D}^B\}$  for double spaces are jointly constructed by using the training sets from the focused and blurred feature spaces, with the same supports. With these analysis, we boost the fusion criteria based on the coupled sparse representations.

The solution of the MAP estimator (3) for multi-focus image fusion is built by

$$\begin{aligned} \left\{ \mathbf{A}, \tilde{\mathbf{I}}_F \right\} = & \underset{\mathbf{A}, \mathbf{I}_F}{\operatorname{argmin}} \| \mathbf{I}_F - \mathcal{F}(\{\mathbf{I}_k\}_{k=1}^n) \|_2^2 + \lambda \| \mathbf{A} \|_1 \\ & + \| \mathbf{D}^F \cdot \mathcal{T}(\{\mathbf{I}_k\}_{k=1}^n; \{\mathbf{D}^F, \mathbf{D}^B\}) - \mathbf{I}_F \|_2^2. \end{aligned} \quad (4)$$

where  $\mathbf{A}$  represents the sparse coefficients of  $\mathbf{I}_F$ ,  $\mathcal{T}(\cdot)$  is the fusing criteria which merges sparse coefficients from  $\{\mathbf{I}_k\}_{k=1}^n$  into  $\mathbf{A}$ , and  $\lambda$  is a regularization parameter. The first term in (4) is the log-likelihood global force that enforces the proximity between the estimated image  $\mathbf{I}_F$  and its real version. The second and third terms are the image prior that means  $\mathbf{I}_F$  has the compact sparse representations  $\mathbf{A}$  from each multi-focus image  $\mathbf{I}_k$ , using the couple of overcomplete dictionaries  $\{\mathbf{D}^F, \mathbf{D}^B\}$ . The algorithm are highly efficient, due to the natural relation to learning the coupled of dictionaries. When  $\mathbf{D}^F$  and  $\mathbf{D}^B$  are not coupled, the fusion criteria  $\mathcal{T}(\cdot)$  is unreliable.

For solving the problem (4), we divide it into two sub-problems: One problem is defining the fusion criteria  $\mathcal{T}(\cdot)$  based on coupled sparse coefficients, and the second problem is learning the couple of dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$  for describing joint sparse representation  $\mathbf{A}_k$ .

In what follows, we will present the proposed fusion criteria, and learn the couple of overcomplete dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$  via jointly sparse coding.

### III. PROPOSED FUSION CRITERIA

Assume that the couple of dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$  are known and fixed, the proposed penalty term in (4) has two kinds of unknown parameters: the compact sparse representation  $\mathbf{A}$  of all-in-focus image  $\mathbf{I}_F$ , and each sparse coefficients  $\mathbf{A}_k$  of multi-focus image  $\mathbf{I}_k$ . We start with the collection of the sparse representations  $\mathbf{A}_k$  over the coupled dictionary  $\{\mathbf{D}^F, \mathbf{D}^B\}$ , and then seek the optimal  $\mathbf{A}$  over the focused dictionary  $\mathbf{D}^F$ .

For convenience, the approach now introduces the operator  $R(\{\mathbf{A}_k\}_{k=1}^n)$  to choose the optimal sparse coefficients, and then define the fusion criteria  $\mathcal{T}$  as

$$\mathcal{T}(\{\mathbf{I}_k\}_{k=1}^n) \triangleq R(\{\mathbf{A}_k\}_{k=1}^n) \cdot \mathcal{L}(\mathbf{A}_k; \mathbf{I}_k) \quad (5)$$

where  $\mathcal{L}(\mathbf{A}_k; \mathbf{I}_k)$  is a loss function which the image  $\mathbf{I}_k$  should be sparsely represented by the coefficient  $\mathbf{A}_k$ . Using one dictionary, as in the existing approach, implies that the focused and

blurred features in each image  $\mathbf{I}_k$  are directly represented by one linear combination, and the optimal representation for blurred features is ignored. Hence, we intent to use a collection of sparse representations, with respect to the focused and blurred categories of dictionaries, in order to produce a better sparse representation.

Based on these observations, we build the averaging operator  $\mathcal{C}$  on the sparse representations via the couple of dictionaries  $\{\mathbf{D}^F, \mathbf{D}^B\}$  to seek the suitable sparse representation  $\mathbf{A}_k$ , and the representation is described by solving the following optimization function

$$\mathcal{L}(\mathbf{A}_k; \mathbf{I}_k) = \operatorname{argmin}_{\mathbf{A}_k} \left\| \mathcal{C} \cdot \begin{bmatrix} \mathbf{D}^F & \mathbf{D}^B \end{bmatrix} \mathbf{A}_k - \mathbf{I}_k \right\|_2^2 + \lambda \|\mathbf{A}_k\|_1. \quad (6)$$

In doing so, we need to work on smaller patches of size  $m \times m$  since adapting a dictionary to large signal is impractical. Using the sliding window technique, we divide the signal into small patches, from left-top to right-bottom. Visible artifacts may occur on block boundaries, and we also introduce overlapping patches for each  $m \times m$  small patch on image  $\mathbf{I}_k$ .

Put formally, we decompose the loss function (6) into two sub-functions and operate them on each small patch. The sub-functions are described as

$$[\mathbf{a}_k^F]_{ij} = \operatorname{argmin}_{[\mathbf{a}_k^F]_{ij}} \left\| [\mathbf{a}_k^F]_{ij} \right\|_1 + \left\| \mathbf{D}^F [\mathbf{a}_k^F]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2 \quad (7)$$

$$[\mathbf{a}_k^B]_{ij} = \operatorname{argmin}_{[\mathbf{a}_k^B]_{ij}} \left\| [\mathbf{a}_k^B]_{ij} \right\|_1 + \left\| \mathbf{D}^B [\mathbf{a}_k^B]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2. \quad (8)$$

Here,  $\mathbf{W}_{i,j}$  is a projection matrix that extracts the  $(i, j)^{th}$  block from the image, and the choice for  $\lambda$  dictates that the errors  $\left\| \mathbf{D}^F [\mathbf{a}_k^F]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2$  and  $\left\| \mathbf{D}^B [\mathbf{a}_k^B]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2$  go below  $\epsilon$ . We will see that the fusion quality has a close relationship with the tolerance error  $\epsilon$  in Section V. Considering the function (7) and (8) respectively, we may employ the basic orthogonal matching pursuit (OMP) algorithm [39] to estimate the sparse representations with the residual energy at each iteration. In our scheme, we intend to build a bridge in relationships between  $[\mathbf{a}_k^F]_{ij}$  and  $[\mathbf{a}_k^B]_{ij}$ , representing double kinds of feature space. Each set of sparse coefficients represents its own significant salient structures, and averaging these coefficients leads to a fusion of different features [38]. Hence, the solutions of the two functions can be jointly obtained by the randomized orthogonal matching pursuit (RandOMP) that generate two sets of competitive representations  $[\mathbf{a}_k^F]_{ij}$  and  $[\mathbf{a}_k^B]_{ij}$ . Suppose at the  $j$ -th iteration, the  $t$ -th atoms  $\left\| \mathbf{d}_t^F [\mathbf{a}_k^F]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k(t) \right\|_2^2$  and  $\left\| \mathbf{d}_t^B [\mathbf{a}_k^B]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k(t) \right\|_2^2$  need to be computed. Instead of finding the smallest errors using

OMP algorithm, we randomize the choice of the  $t$ -th *atoms* with probability linearly proportional to

$$\exp \left\{ \frac{c^2}{2\sigma^2} \cdot \frac{|\mathbf{d}_t^F \mathbf{W}_{i,j} \mathbf{I}_k^{j-1}|^2 + |\mathbf{d}_t^B \mathbf{W}_{i,j} \mathbf{I}_k^{j-1}|^2}{\|\mathbf{d}_t^F\|_2^2 + \|\mathbf{d}_t^B\|_2^2} \right\}, c^2 = \frac{\sigma_k^2}{\sigma_k^2 + \sigma^2} \quad (9)$$

where  $\sigma_k$  is the variance on the nonzero entries of the  $\mathbf{I}_k$  representation. When the errors  $\|\mathbf{D}^F [\mathbf{a}_k^F]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k\|_2^2$  and  $\|\mathbf{D}^B [\mathbf{a}_k^B]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k\|_2^2$  go below  $\epsilon$ , the randomization yields two sets of  $\{[\mathbf{a}_k^F]_{ij}, [\mathbf{a}_k^B]_{ij}\}$ . Details of this processing are given in Appendix A.

Armed with the analysis, we push the step further to define the function  $\mathcal{C}$  as the optimal value

$$\begin{aligned} [\mathbf{a}_k]_{ij} &= \mathcal{C} \left( [\mathbf{a}_k^F]_{ij}, [\mathbf{a}_k^B]_{ij} \right) \\ &= \frac{[\mathbf{a}_k^F]_{ij} + [\mathbf{a}_k^B]_{ij}}{2}. \end{aligned} \quad (10)$$

Given all  $[\mathbf{a}_k]_{ij}$ , we now turn to define the function  $R(\mathbf{A}_k)$  and seek the sparse coefficient  $\mathbf{A}$ . In a general rule, the block with a bigger defined value is chosen to construct the fused image [20]. Here, we compute

$$\mathcal{R}([\mathbf{a}_k]_{ij}) = \operatorname{argmax}_{u \neq v} \left( \left\| [\mathbf{a}_u]_{ij} \right\|_1, \left\| [\mathbf{a}_v]_{ij} \right\|_1 \right) \quad (11)$$

where  $\left\| [\mathbf{a}_u]_{ij} \right\|_1$  denotes the sum of the magnitude of its coefficients, and  $1 \leq u, v \leq n$ . Returning to (5), we can get the fusion coefficient as

$$\mathcal{T}(\{\mathbf{I}_k\}_{k=1}^n) = [\mathbf{a}]_{ij}. \quad (12)$$

In this strategy, the compact representations are built for each patches in the all-in-focus image. Gathering all  $\mathbf{a}_{ij}$ , the all-in-focus image is reconstructed as  $\mathbf{I}_F^0 = \mathbf{D}^F \mathbf{A}$ . Up to here, we have enhanced the local details (i.e., spatial edges, local textures) observed in the all-in-focus image.

In order to remove possible artifacts and improve spatially smoothness, the global reconstruction constraint as well as the consistency between the initial estimation  $\mathbf{I}_F^0$  and the final outcome can be applied to make a further improvement. In this paper, we utilize the gradient-descent based minimization [43] to solve this function as the replacement of the function (4), given by

$$\tilde{\mathbf{I}}_F = \operatorname{argmin}_{\mathbf{A}, \mathbf{I}_F} \|\mathbf{I}_F - \mathcal{F}(\{\mathbf{I}_k\}_{k=1}^n)\|_2^2 + \lambda \|\mathbf{I}_F^0 - \mathbf{I}_F\|_2^2. \quad (13)$$

Hence, the overall multi-focus image fusion algorithm is summarized as Algorithm 1.

**Algorithm 1** Image fusion from sparsity.

---

**Input:** the couple of learned dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$ ; multi-focus source images  $\{\mathbf{I}_k\}_{k=1}^n$ .

- 1: Remove the mean intensity for each source image.
- 2: **for** each  $8 \times 8$  patch  $i_k$  from  $\mathbf{I}_k$ , starting from the upper-left corner with 1 pixel overlap,
  - Solve the optimization problems (7), (8);
  - Compute the averaging coefficients using (10);
  - Solve the operator  $\mathcal{R}(\mathbf{a}_{i,j})$ , using (11);
  - Solve the fusion criteria  $\mathcal{T}(\{\mathbf{I}_k\}_{k=1}^n)$  using (4), (12);
  - Reconstruct the initial all-in-focus image  $\mathbf{I}_F^0$ .
- 3: **end for**
- 4: Build the final image using (13).
- 5: Put the mean intensity into the all-in-focus image  $\mathbf{I}_F$ .

---

**Output:** the all-in-focus image  $\mathbf{I}_F$

---

#### IV. LEARNING COUPLED DICTIONARY

The entire definition on fusion criteria so far is based on the assumption that the couple of dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$  from the focused and blurred feature spaces are known and fixed. So a key issue to achieving the fusion criteria is the couple of dictionaries to learn. In this section, we focus on the coupled dictionary learning problem and more specifically, the improvement of learning efficiency.

Assuming that the sets  $\mathbf{X}^F = [\mathbf{x}_1^F, \mathbf{x}_2^F, \dots, \mathbf{x}_n^F] \in \mathbb{R}^{d \times n}$  and  $\mathbf{X}^B \triangleq [\mathbf{x}_1^B, \mathbf{x}_2^B, \dots, \mathbf{x}_n^B] \in \mathbb{R}^{d \times n}$  as the matrix of  $n$  sampled focused image vectors emerge from focused and blurred feature spaces, respectively and  $d$  is the dimension of the sampled image vectors. The aim of coupled dictionary learning is to seek the pair of dictionaries  $\mathbf{D}^F, \mathbf{D}^B \in \mathbb{R}^{d \times N}$  adapted to two sets of focused images  $\mathbf{X}^F$  and blurred images  $\mathbf{X}^B$  with respect to the best possible representation  $\Gamma$ .

The basic coupled dictionary learning is formulated as

$$\begin{aligned} & \min_{\mathbf{D}^F, \mathbf{D}^B, \Gamma} \|\mathbf{X}^F - \mathbf{D}^F \Gamma\|_2^2 + \|\mathbf{X}^B - \mathbf{D}^B \Gamma\|_2^2 \\ & \text{s.t. } \|\Gamma\|_0 \leqslant T_0, \|\mathbf{d}_i^F\|_2^2 \leqslant 1, \|\mathbf{d}_i^B\|_2^2 \leqslant 1, \forall 1 \leqslant i \leqslant N \end{aligned} \quad (14)$$

where  $\Gamma$  is the joint spare coding of  $\mathbf{X}^F$  and  $\mathbf{X}^B$ ,  $\mathbf{d}_i^F$  and  $\mathbf{d}_i^B$  are the  $i$ -th columns of  $\mathbf{D}^F$  and  $\mathbf{D}^B$ , respectively, and  $T_0$  is the parameter controlling the sparsity penalty. Many algorithms [27]–[30], [32] have been introduced to solve the problem (14) for learning the coupled dictionary.

The coefficients are estimated via  $l_1$  minimization, keeping the two dictionaries fixed alternately, and the dictionaries are estimated through least squares, keeping the coefficients fixed. Though these algorithms benefit from having the reliability, their efficiency still need to be improved.

The K-SVD [24] is a well-known and frequently-employed tool that aims to handle dictionary learning problem. It updates the dictionary atoms and the corresponding coefficients alternately, and improves the atoms to better fit the signals efficiently. Particularly, the algorithm updates one atom and all the non-zero columns corresponding to the atom at a time in dictionary updating stage, and then sparse approximation updates the non-zero columns again in sparse coding stage. Modifying the redundant updating, the improved K-SVD algorithm [40] introduces a mask matrix of zeros and ones, and keeps all the non-zero columns intact. Armed with this algorithm, we propose an improved K-SVD-based coupled dictionary learning approach.

Formally, we add additional constraint on  $\Gamma$  for keeping all the zero atoms in  $\Gamma$  intact, given by

$$\Gamma \odot \mathbf{M} = 0 \quad (15)$$

where the mask matrix  $\mathbf{M}$  is defined as

$$\mathbf{M}_{i,j} \triangleq \begin{cases} 1, & \text{if } \Gamma_{i,j} = 0 \\ 0, & \text{if } \Gamma_{i,j} \neq 0 \end{cases}. \quad (16)$$

Thus, our coupled dictionary learning problem is expressed as

$$\begin{aligned} & \min_{\mathbf{D}^F, \mathbf{D}^B, \Gamma} \|\mathbf{X}^F - \mathbf{D}^F \Gamma\|_2^2 + \|\mathbf{X}^B - \mathbf{D}^B \Gamma\|_2^2 \\ & \text{s.t. } \|\mathbf{d}_i^F\|_2^2 \leq 1, \|\mathbf{d}_i^B\|_2^2 \leq 1, \|\Gamma\|_0 \leq T_0, \Gamma \odot \mathbf{M} = 0. \end{aligned} \quad (17)$$

We expect the learned atoms of dictionaries  $\mathbf{D}^F$ ,  $\mathbf{D}^B$  can describe the coherent structures of each individual feature space and correlation characteristics between the focused feature space and the blurred feature space, and thus facilitate the fusion criteria. The objective problem (17) has three parameters to optimize: the atoms of coupled dictionaries  $\mathbf{D}^F$ ,  $\mathbf{D}^B$ , and joint sparse coefficient  $\Gamma$  that linearly relate those atoms to the learning data sets  $\mathbf{X}^F$  and  $\mathbf{X}^B$ . Solving these parameters, the proposed approach involves two main iteration steps that (i) the atoms of the coupled dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$  are alternately updated by the sparse coefficients; and (ii) the joint sparse coefficient  $\Gamma$  is given by fixing the coupled dictionaries.

### A. Dictionary update

Let us first optimize the dictionary update steps where we assume that  $\Gamma$  is fixed, and consider the optimization problem (17) as processing the update of two dictionaries with the coefficients summarized in  $\Gamma$ , which are described as

$$\tilde{\mathbf{D}}^F = \operatorname{argmin}_{\mathbf{D}^F} \left\| \mathbf{X}^F - \sum_{i=1}^N \mathbf{d}_i^F \boldsymbol{\gamma}_i^T \right\|_2^2 \quad (18)$$

$$\tilde{\mathbf{D}}^B = \operatorname{argmin}_{\mathbf{D}^B} \left\| \mathbf{X}^B - \sum_{i=1}^N \mathbf{d}_i^B \boldsymbol{\gamma}_i^T \right\|_2^2. \quad (19)$$

Based on the prior knowledge of sparse representation, it is not necessary to update the whole atoms in each atom updating step. For optimizing  $\mathbf{d}_i^F$  and  $\mathbf{d}_i^B$  in each  $i$ -th *atom*, we fix the remaining *atoms*  $\sum_{j \neq i} \mathbf{d}_j^F$ ,  $\sum_{j \neq i} \mathbf{d}_j^B$  and the corresponding coefficients  $\sum_{j \neq i} \boldsymbol{\gamma}_j$ , and then rewrite the functions (18) and (19) as

$$\tilde{\mathbf{d}}_i^F = \operatorname{argmin}_{\mathbf{d}_i^F} \left\| \left( \mathbf{X}^F - \sum_{j \neq i} \mathbf{d}_j^F \boldsymbol{\gamma}_j^T \right) \odot \mathbf{M}_i - \mathbf{d}_i^F \boldsymbol{\gamma}_i^T \right\|_2^2 \quad (20)$$

$$\tilde{\mathbf{d}}_i^B = \operatorname{argmin}_{\mathbf{d}_i^B} \left\| \left( \mathbf{X}^B - \sum_{j \neq i} \mathbf{d}_j^B \boldsymbol{\gamma}_j^T \right) \odot \mathbf{M}_i - \mathbf{d}_i^B \boldsymbol{\gamma}_i^T \right\|_2^2 \quad (21)$$

where  $\mathbf{M}_i \triangleq \mathbf{1}_d \cdot \mathbf{m}_i^T$ , and  $\mathbf{1}_d$  represents the  $d$  time replications of  $\mathbf{m}_i^T$ . With the mask matrix  $\mathbf{M}$ , we ignore all the columns in  $\left\{ \mathbf{X}^F - \sum_{j \neq i} \mathbf{d}_j^F \boldsymbol{\gamma}_j^T, \mathbf{X}^B - \sum_{j \neq i} \mathbf{d}_j^B \boldsymbol{\gamma}_j^T \right\}$  without using the  $i$ -th *atom*. The benefits of this update approach are two fold, not only it lowers the computational complexity, but also it avoids redundant and unnecessary updates. The error matrices  $\mathbf{E}_i^F$ ,  $\mathbf{E}_i^B$  are defined as

$$\mathbf{E}_i^F \triangleq \left( \mathbf{X}^F - \sum_{j \neq i} \mathbf{d}_j^F \boldsymbol{\gamma}_j^T \right) \odot \mathbf{M}_i \quad (22)$$

$$\mathbf{E}_i^B \triangleq \left( \mathbf{X}^B - \sum_{j \neq i} \mathbf{d}_j^B \boldsymbol{\gamma}_j^T \right) \odot \mathbf{M}_i. \quad (23)$$

Updating a new column  $\mathbf{d}_i^F$  and its coefficients  $\boldsymbol{\gamma}_i$ , we minimize the mean square errors (MSE) by

$$\tilde{\mathbf{d}}_i^F = \operatorname{argmin}_{\mathbf{d}_i^F} \|\mathbf{E}_i^F - \mathbf{d}_i^F \boldsymbol{\gamma}_i^F\|_F^2 \quad (24)$$

$$\tilde{\mathbf{d}}_i^B = \operatorname{argmin}_{\mathbf{d}_i^B} \|\mathbf{E}_i^B - \mathbf{d}_i^B \boldsymbol{\gamma}_i^B\|_F^2. \quad (25)$$

Then, we directly apply the singular value decomposition (SVD) operation on the problems (24) and (25), and get the decompositions

$$\mathbf{E}_i^F = \mathbf{U}_i^F \Delta_i^F \mathbf{V}_i^F, \quad \mathbf{E}_i^B = \mathbf{U}_i^B \Delta_i^B \mathbf{V}_i^B. \quad (26)$$

Here,  $\mathbf{d}_i^F, \mathbf{d}_i^B$  are respectively updated by  $\mathbf{U}_i^B$  and  $\mathbf{U}_i^F$ , and  $\boldsymbol{\gamma}_i^T$  is alternately valued by  $\Delta_i^F(1,1) \cdot \mathbf{V}_i^F$  and  $\Delta_i^B(1,1) \cdot \mathbf{V}_i^B$  on each iteration.

Such approach effectively updates the atoms in  $\mathbf{D}^F$  and  $\mathbf{D}^B$ , without multiplying all the non-zero representations.

### B. Coefficient update

In the sparse coding stage, one might be tempted to suggest using columns of coefficients in dictionary update stage directly, however, that does not imply a guaranteed improvement [24]. We proceed with updating  $\Gamma$  in the second stage, and search the sparsest representations for the signals  $\mathbf{X}^F$  and  $\mathbf{X}^B$  in (17).

The natural approach alternately optimizes the following sparse approximations

$$\tilde{\boldsymbol{\gamma}}_i^T = \operatorname{argmin}_{\boldsymbol{\gamma}_i^T} \left\| \mathbf{X}^F - \sum_{i=1}^N \mathbf{d}_i^F \boldsymbol{\gamma}_i^T \right\|_2^2 \quad (27)$$

$$\tilde{\boldsymbol{\gamma}}_i^T = \operatorname{argmin}_{\boldsymbol{\gamma}_i^T} \left\| \mathbf{X}^B - \sum_{i=1}^N \mathbf{d}_i^B \boldsymbol{\gamma}_i^T \right\|_2^2. \quad (28)$$

Considering computational complexity, the proposed approach ignores the formulations in (27) and (28), and try to update the coefficients using the error matrices  $\mathbf{E}_i^F$  and  $\mathbf{E}_i^B$ . Once the dictionary update is performed, we can obtain the expressions

$$\sum_{i=1}^N \left\| \mathbf{E}_i^F - \mathbf{d}_i^F \boldsymbol{\gamma}_i^F \right\|_2^2 \quad (29)$$

$$\sum_{i=1}^N \left\| \mathbf{E}_i^B - \mathbf{d}_i^B \boldsymbol{\gamma}_i^B \right\|_2^2. \quad (30)$$

The algorithm can find the suitable sparse representations during the dictionary updating stage. Then, we combine the functions (29) and (30), forcing  $\mathbf{d}_i^F$  and  $\mathbf{d}_i^B$  to share the same sparse representation  $\boldsymbol{\gamma}$ , given by

$$\begin{aligned} & \min_{\boldsymbol{\gamma}} \frac{1}{2} \left\| \mathbf{E}_i^F - \mathbf{d}_i^F \boldsymbol{\gamma} \right\|_2^2 + \frac{1}{2} \left\| \mathbf{E}_i^B - \mathbf{d}_i^B \boldsymbol{\gamma} \right\|_2^2 + \lambda \|\boldsymbol{\gamma}\|_1 \\ & \text{s.t. } \left\| \mathbf{d}_i^F \right\|_2^2 \leq 1, \quad \left\| \mathbf{d}_i^B \right\|_2^2 \leq 1, \quad \forall 1 \leq i \leq N. \end{aligned} \quad (31)$$

Note that (31) can be equivalently rewritten as the optimization function for  $\gamma$

$$\min_{\gamma} \mathcal{F}(\gamma) = \frac{1}{2} \underbrace{\left\| \tilde{\mathbf{E}}_i - \tilde{\mathbf{d}}_i \gamma \right\|_2^2}_{\mathcal{F}_1(\gamma)} + \lambda \underbrace{\|\gamma\|_1}_{\mathcal{F}_2(\gamma)} \quad (32)$$

where

$$\tilde{\mathbf{E}}_i = \begin{bmatrix} \mathbf{E}_i^F \\ \mathbf{E}_i^B \end{bmatrix}, \quad \tilde{\mathbf{d}}_i = \begin{bmatrix} \mathbf{d}_i^F \\ \mathbf{d}_i^B \end{bmatrix}. \quad (33)$$

Here, we divide the optimization function  $\mathcal{F}(\gamma)$  into two parts  $\mathcal{F}_1(\gamma)$  and  $\mathcal{F}_2(\gamma)$  and assume  $(\gamma^t, \tilde{\mathbf{d}}_j)$  be the point obtained by the function  $\mathcal{F}(\gamma)$  at iteration  $t$  point. Modified by [41], [42], our optimization approach is based on the repeated minimization of a differentiable surrogate function.

Define the function  $\mathcal{G}(\gamma, \gamma^t)$  as a surrogate function for the function  $\mathcal{F}(\gamma)$  on the point  $\gamma^t$  with the following two conditions hold.

$$\begin{aligned} \mathcal{G}(\gamma^t, \gamma^t) &= \mathcal{F}(\gamma^t) \\ \mathcal{G}(\gamma, \gamma^t) &\geq \mathcal{F}(\gamma) \end{aligned} \quad (34)$$

Then, the above optimization function can be solved by the following approximation

$$\gamma^{t+1} = \operatorname{argmin}_{\gamma} \mathcal{G}(\gamma, \gamma^t) \quad (35)$$

where

$$\mathcal{G}(\gamma, \gamma^t) \triangleq \langle \nabla_{\gamma} \mathcal{F}_1(\gamma^t), \gamma \rangle + \frac{1}{2} \|\gamma - \gamma^t\|_F^2 + \mathcal{F}_2(\gamma). \quad (36)$$

More specifically, we attempt to get the solution on the variable  $\gamma$  by this rule. We describe the well-known soft-thresholding operation  $\mathcal{S}_{\lambda}(a)$  as

$$\mathcal{S}_{\lambda}(a) \triangleq \begin{cases} a - \lambda, & \text{for } a > \lambda \\ 0, & \text{for } |a| \leq \lambda \\ a + \lambda, & \text{for } a < \lambda \end{cases}. \quad (37)$$

Combining (36) with (37), the updating formula for  $\gamma$  is obtained by the closed-form solution

$$\gamma^{t+1} = \gamma^t - \mathcal{S}_{\frac{\lambda}{\mathcal{T}^2(\tilde{\mathbf{d}}_j)}} \left( \gamma^t - \frac{1}{\mathcal{T}^2(\tilde{\mathbf{d}}_j)} \tilde{\mathbf{d}}_j^T (\tilde{\mathbf{d}}_j \gamma^t - \tilde{\mathbf{E}}_j) \right) \quad (38)$$

where  $\mathcal{T}(\tilde{\mathbf{d}}_j)$  describe the maximum singular value of  $\tilde{\mathbf{d}}_j$ . Appendix B shows the convergence analysis.

The complete implement is given as Algorithm 2.

---

**Algorithm 2** Learning the coupled dictionary.

---

**Input:** two signal sets  $\mathbf{X}^F$  and  $\mathbf{X}^B$ , initial dictionaries  $\mathbf{D}_0^F$  and  $\mathbf{D}_0^B$ , initial sparse coefficients  $\Gamma$ .

- 1: **Initialization:** Set  $\mathbf{D}^F := \mathbf{D}_0^F$ ,  $\mathbf{D}^B := \mathbf{D}_0^B$ ,  $\Gamma := \Gamma_0$
- 2: **for**  $j = 1 \dots n$  **do**
- 3:   compute the matrices  $\mathbf{E}_i^F$  and  $\mathbf{E}_i^B$ , using (22), (23);
- 4:   update the coupled atoms  $\mathbf{d}_i^F$  and  $\mathbf{d}_i^B$  using (24), (25);
- 5:   **for**  $i = 1 \dots t$  **do**
- 6:     update sparse coefficients  $\gamma$  using (37), (38);
- 7:   **end for**
- 8: **end for**

---

**Output:** the couple of dictionaries  $\mathbf{D}^F$  and  $\mathbf{D}^B$ .

---

## V. SIMULATION AND RESULTS

This section demonstrates experiential results of our approach on simulated data. We firstly evaluate the proposed approach by visual comparisons, and then discuss the quantitative assessments of the proposed implementation. At last, we describe various the influence of factors to the proposed approach including patch size, tolerance error, and number of overlapping pixels. The objective evaluation is based on two state-of-the-art fusion performance metrics: (i)  $Q_{MI}$  [44], which measures how well the mutual information from the source images is preserved in the fused image; (ii)  $Q^{AB/F}$  [45], which evaluates how well the success of edge information transfers from the source images to the fused image. All the experiments are performed on a PC running a Inter(R) Xeon(R) 3.40GHz CPU.

### A. Experimental Setup

The proposed method is compared to the following state-of-art multi-focus image fusion algorithms:

- LP: aplacian pyramid-based image fusion approach [11];
- MWG: multi-scale weighted gradient-based spatial image fusion approach [13];
- DWT: discrete wavelet transform-based image fusion approach [14];
- NSCT: non-subsampled contourlet transform-based image fusion approach [15];

- PCA: robust principal component analysis-based image fusion approach [1];
- SR-CM: sparse representation “choose-max”-based image fusion approach [20];
- SR-KSVD: sparse representation K-SVD-based image fusion approach.

The parameter settings of these methods are as follows. For the LP and DWT methods, the source images are decomposed to 3 levels. The wavelet basis “db1” is applied to the DWT algorithm. For the NSCT approach, the direction numbers of the 4 decomposition levels from coarse to fine are selected as 4, 8, 8, and 16. For fair comparison, all algorithms based on sparsity model are implemented using the same dictionary size.

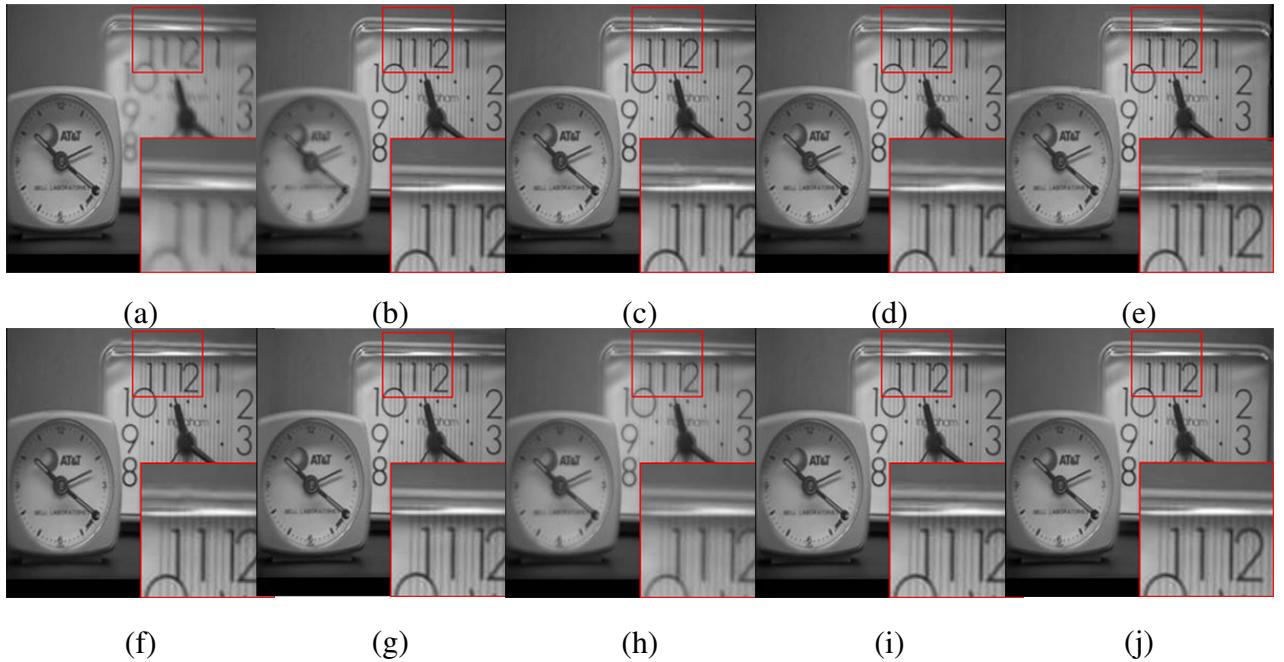


Fig. 1. Source images “Clock” and the fusion result comparisons. (a) The first source image with focus on the left. (b) The second source image with focus on the right. Fused images obtained by LP (c), MWG (d), DWT (e), NSCT (f), PCA (g), SR-CM (h), SR-KSVD (i) and the proposed method (j).

Throughout all experience, we use the following parameter setting for our method. we set the tolerance error  $\epsilon = 0.1$ , the pitch size  $m = 8$ , and the overlapping size is 1. During the coupled dictionary learning, the training data consist of 50,000  $8 \times 8$  focused patches which are randomly sampled from the database of 40 natural images and 50,000  $8 \times 8$  blurred patches created by focused patches using Gaussian blur function. The two dictionaries are produced by improved K-SVD algorithm, and initialized with samples from the training data. Larger dictionaries result in heavier computation and better quality. Considering the tradeoff between fusion quality and computation, we fix the dictionary size as  $64 \times 256$ , execute 6 multiple dictionary update cycles

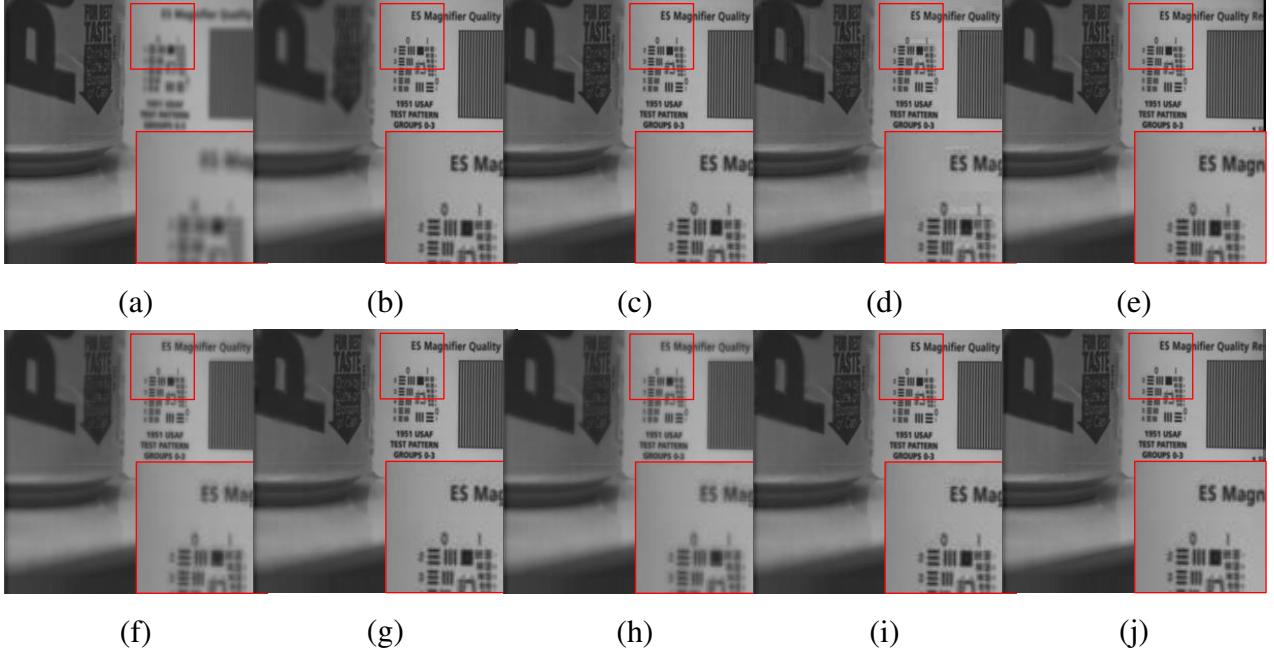


Fig. 2. Source images “Pepsi” and the fusion result comparisons. Same order in the Fig. 1.

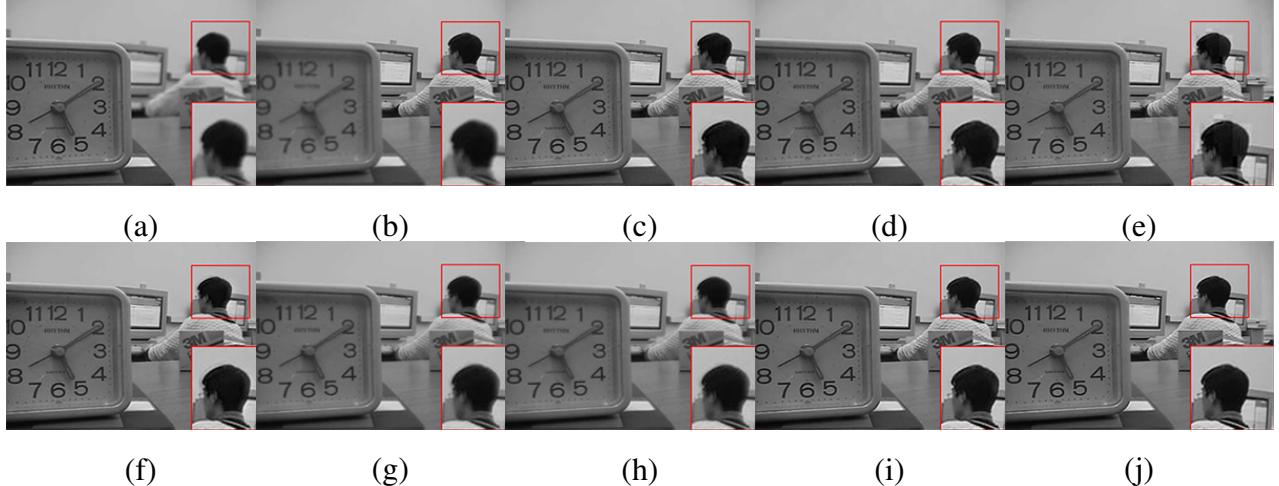


Fig. 3. Source images “Lab” and the fusion result comparisons. Same order in the Fig. 1.

(DUC) and 30 iterations. Experiments converge very quickly, and approximately achieve a  $\times 2$  speedup compared to the standard coupled dictionary learning [27].

### B. Simulated Data Experiments

To show the effectiveness of the proposed algorithm, we compare it with existing algorithms, including LP, MWG, DWT, NSCT, PCA, SR-CM, and SR-KSVD methods. The fusion results of the images “Clock”, “Pepsi”, and “Lab” are shown in Fig. 1, Fig. 2, and Fig. 3, including the

TABLE I  
OBJECTIVE EVALUATION OF THE IMAGE FUSION

Methods	Measures in Fig. 1			Measures in Fig. 2			Measures in Fig. 3		
	$Q_{MI}$	$Q^{AB/F}$	Time(s)	$Q_{MI}$	$Q^{AB/F}$	Time(s)	$Q_{MI}$	$Q^{AB/F}$	Time(s)
<b>LP</b>	0.9083	0.6879	<b>0.0042</b>	0.8998	0.6864	<b>0.0053</b>	0.9091	0.6901	<b>0.0039</b>
<b>MWG</b>	0.9045	0.7243	0.8383	0.8992	0.7202	0.8424	0.9083	0.7264	0.8280
<b>DWT</b>	0.8991	0.7013	0.1542	0.8876	0.6991	0.1784	0.9074	0.7083	0.1501
<b>NSCT</b>	0.9245	0.7185	1.0748	0.9188	0.7101	1.1092	0.9246	0.7189	1.0054
<b>PCA</b>	0.9381	0.6620	0.0074	0.9311	0.6601	0.0102	0.9383	0.6609	0.0070
<b>SR-CM</b>	0.9276	0.7214	3.1894	0.9241	0.7201	3.2785	0.9289	0.7215	3.1002
<b>SR-KSVD</b>	0.9386	0.7326	3.1992	0.9381	0.7314	3.2198	0.9391	0.7326	3.0784
<b>Ours</b>	<b>0.9432</b>	<b>0.7451</b>	3.2123	<b>0.9428</b>	<b>0.7439</b>	3.4009	<b>0.9445</b>	<b>0.7451</b>	3.1459

magnified details in the lower right corners. Figs. 1(a) and (b) show the source images, and the fused images obtained by different fusion methods based on the LP, MWG, DWT, NSCT, PCA, SR-CM, SR-KSVD methods, and the proposed method are presented in Figs. 1(c)-(j). The same order is in the Fig. 2 and Fig. 3.

By carefully observing the fusion results, we see clearly that the LP method does not produce continuous edges (see Fig. 1(c)) and the DWT method results in blocking artifacts (see Fig. 1(e)). Moreover, the same results can be easily observed in the magnified images. The fusion methods based on MWG (see Fig. 1(d)) and NSCT (see Fig. 1(f)) show circle blurring effect around strong boundaries. Fig. 1(f) is visually better than Fig. 1(d). In Figs. 1(g) and (i) can be seen some artificial distortions, and Fig. 1(h) obviously lacks edge information. Our approach provides the best visual appearance (see Fig. 1(j)). Similar subjective results are obtained in these experiments as shown in Fig. 2 and Fig. 3. Here again the proposed approach yields better image quality than that of the other compared methods.

Table I summarizes the evaluations on grayscale multi-focus images. The values of  $Q_{MI}$  and  $Q^{AB/F}$  range from 0 to 1, with 1 representing the ideal fusion. The bold values are the best results in the corresponding columns. As can be seen from Table I, the proposed approach generally

produces better quantitative results in terms of  $Q_{MI}$  and  $Q^{AB/F}$ . The value of  $Q_{MI}$  for the proposed approach is always larger than the values for the LP, MWG, DWT, NSCT, and SR-CM methods, and it is similar to the methods based on PCA and SR-KSVD. That means that PCA, SR-KSVD and the proposed approach well preserve the mutual information from different source images. The values of  $Q^{AB/F}$  demonstrate that our approach reduces the blocking artifacts and artificial distortions, and combines the significant edge information into the fused image. For the running time, we optimize the sparse coding using CoefROMP [40] which is initialized with the largest  $k/3$  coefficients from the previous pursuit stage unlike previous pursuit algorithm [39], so the proposed approach takes less time than the SR-KSVD approach.

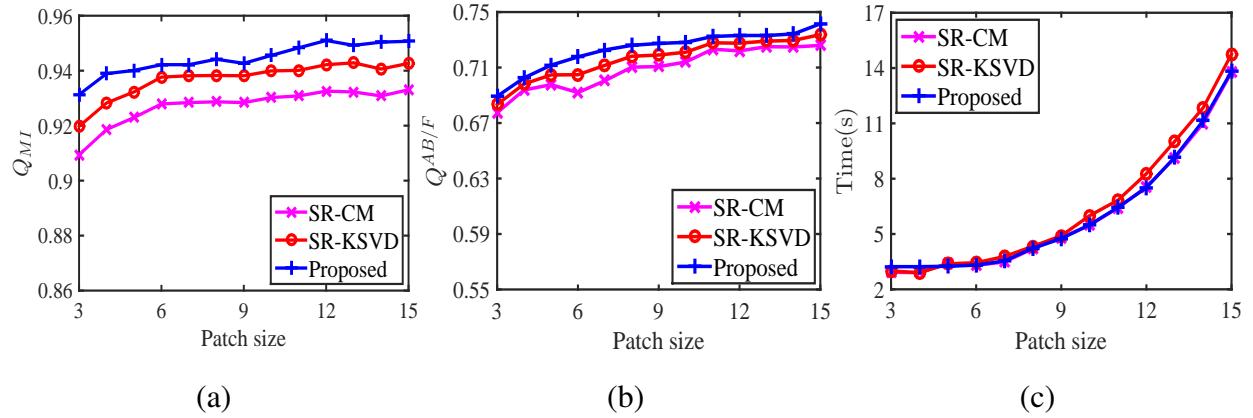


Fig. 4. Fusion performances with different patch size.

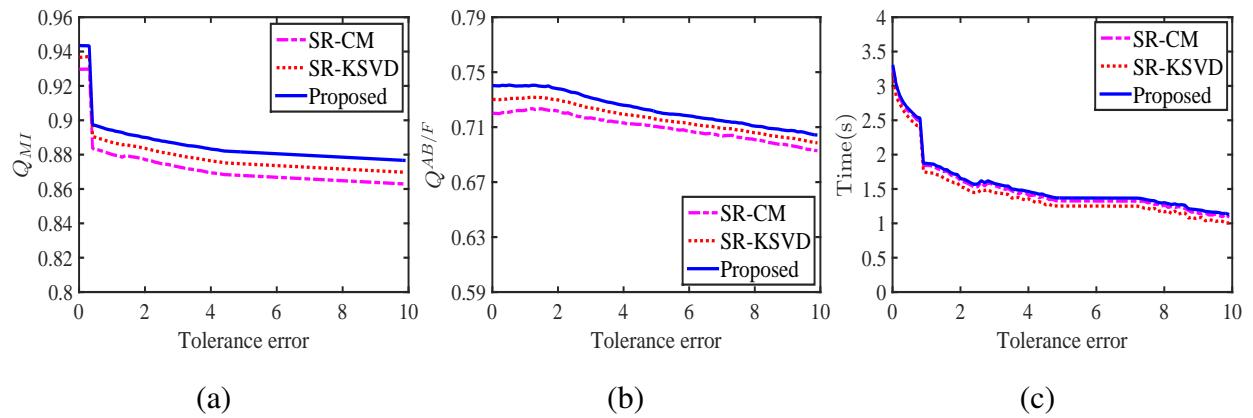


Fig. 5. Fusion performances with different tolerance error.

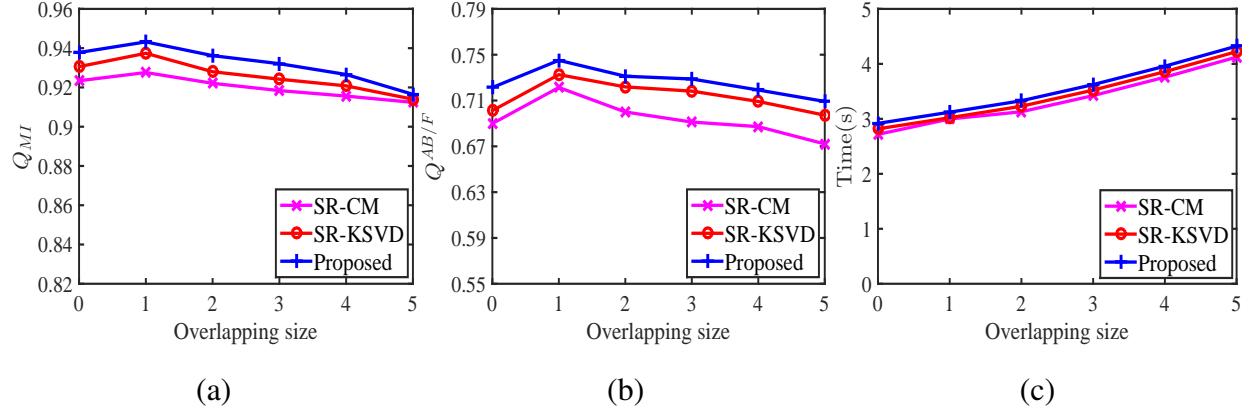


Fig. 6. Fusion performances with different overlapping size.

### C. Effects of Main Parameters

In our experiments, we discuss three main parameters for the proposed approach to evaluate the fusion performance: patch size, tolerance error and number of overlapping pixels. We evaluate the proposed approach on the grayscale multifocus dataset which contain 10 pairs of grayscale images and 30 pairs of artifical source images with 256 level grey scales. Fig. 4, Fig. 5 and Fig. 6 respectively show the average values of  $Q_{MI}$ ,  $Q^{AB/F}$  and the running time with the the parameters pitch size, tolerance error and overlapping size.

As can be seen in Figs. 4(a) and (b), both  $Q_{MI}$  and  $Q^{AB/F}$  slightly benefit from the increasing patch size. When the patch size increases to 9, the running time begins to increase sharply (see Fig. 4(a)). Balancing the computation time and fusion quality, we set the patch size to  $8 \times 8$  in our experiments. Fig. 5 also illustrates the significant improvements in our approach compared to SR-CM and SR-KSVD. The changes of  $Q_{MI}$ ,  $Q^{AB/F}$ , and the running time are shown versus the tolerance error  $\epsilon$ . Fig. 5(a) and (b) show that the tolerance error slightly impacts on  $Q_{MI}$  and  $Q^{AB/F}$ . When  $\epsilon$  is less than 2, the tested approaches have an acceptable fusion quality (see Fig. 5(a)). When  $\epsilon$  is larger than 2,  $Q_{MI}$  is drastically decreasing. Additional experiments in Fig. 6 in show that for increased number of overlapping pixels,  $Q_{MI}$  drastically decreases and the running time sharply increases. Hence, our method has fixed overlapping size to 1.

## VI. CONCLUSION

We have presented a fusion algorithm for combing multiple images with different focal settings into an all-in-focus image. The algorithm first formulated the physical process on multi-focus image fusion, and then described the basic model based on coupled sparse representations. we

introduced the K-SVD-based coupled dictionary learning algorithm that enforced the sparse approximations for double feature spaces. Using the couple of dictionaries from the focused and blurred feature spaces, we explored the an effective and accurate fusion criteria via plain averaging. We demonstrated using simulated data that the proposed approach well preserved the edge and structural information of source images, drastically reduced the blocking artifacts, circle blurring and artifical distortions.

There are two possible extensions to the algorithm. First, while we achieved the fusion criteria in the noiseless setting, it may be possible to fuse multi-focus images in the presence of noise. This may allow us to derive theoretical analysis for fusion criteria. Second, as it was demonstrated that even a RandOMP-based averaging weighting scheme improved the performance of the algorithm, designing more advanced weighting schemes may also result in further improvement.

## APPENDIX A FUSING SPARSE COEFFICIENTS

In this appendix, we describe in details the fusion on sparse representations in Section III.

Fusing sparse coefficients aims to approximate the compact spares representations  $\mathbf{A}$  for unknown image  $\mathbf{I}_F$ . Given a pair of dictionaries  $\{\mathbf{D}^F, \mathbf{D}^B\}$ , the sparsity-based sparse coding problems are described as

$$\begin{aligned} [\mathbf{a}_k^F]_{ij} = & \arg \min_{[\mathbf{a}_k^F]_{ij}} \left\| [\mathbf{a}_k^F]_{ij} \right\|_1 \\ \text{s.t. } & \left\| \mathbf{d}^F [\mathbf{a}_k^F]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2 \leq \epsilon \end{aligned} \quad (39)$$

$$\begin{aligned} [\mathbf{a}_k^B]_{ij} = & \arg \min_{[\mathbf{a}_k^B]_{ij}} \left\| [\mathbf{a}_k^B]_{ij} \right\|_1 \\ \text{s.t. } & \left\| \mathbf{d}^B [\mathbf{a}_k^B]_{ij} - \mathbf{W}_{i,j} \mathbf{I}_k \right\|_2^2 \leq \epsilon \end{aligned} \quad (40)$$

Assume that the columns of  $\mathbf{D}^F$  and  $\mathbf{D}^B$  are normalized. The RandOMP [38] algorithm selects at each step the atoms of  $[\mathbf{a}_k^F]_{ij}$  and  $[\mathbf{a}_k^B]_{ij}$  with probability proportional to (9). When each atom is selected, the signal  $\mathbf{W}_{i,j} \mathbf{I}_k$  is orthogonally projected to the span of the elected atoms, and  $[\mathbf{a}_k^F]_{ij}$  and  $[\mathbf{a}_k^B]_{ij}$  at the  $j$ -th iteration are computed by

$$\begin{aligned} [\mathbf{a}_k^F]_{ij} &= (\mathbf{d}^F)^+ \mathbf{W}_{i,j} \mathbf{I}_k \\ &= ((\mathbf{d}^F)^T \mathbf{d}^F)^{-1} (\mathbf{d}^F)^T \mathbf{W}_{i,j} \mathbf{I}_k \end{aligned} \quad (41)$$

$$\begin{aligned} [\mathbf{a}_k^B]_{ij} &= (\mathbf{d}^B)^+ \mathbf{W}_{i,j} \mathbf{I}_k \\ &= ((\mathbf{d}^B)^T \mathbf{d}^B)^{-1} (\mathbf{d}^F)^T \mathbf{W}_{i,j} \mathbf{I}_k. \end{aligned} \quad (42)$$

When the residuals  $\mathbf{r}_j^F$  and  $\mathbf{r}_j^B$  are below  $\epsilon$ , the stopping rules are satisfied. Here, we present in Algorithm 3 the sparse coefficients fusion method.

---

**Algorithm 3** Fusing sparse coefficients.

---

**Input:** the couple of learned dictionaries  $\{\mathbf{D}^F, \mathbf{D}^B\}$ ; multi-focus source images  $\{\mathbf{I}_k\}_{k=1}^n$ ; the tolerance error  $\epsilon$ .

1: **Initialization:** Set

- $[\mathbf{a}_k^F]_{ij}^0 := 0$ ,  $[\mathbf{a}_k^B]_{ij}^0 := 0$ ;
- $\mathbf{r}_j^F := \mathbf{W}_{i,j} \mathbf{I}_k$ ,  $\mathbf{r}_j^B := \mathbf{W}_{i,j} \mathbf{I}_k$ ;
- Support  $S_F^0 = \text{Support} \left\{ [\mathbf{a}_k^F]_{ij}^0 \right\} = \emptyset$ ,
- Support  $S_B^0 = \text{Support} \left\{ [\mathbf{a}_k^B]_{ij}^0 \right\} = \emptyset$ ;

2: **for**  $j=1$  : max-iter

- Seek the random  $j_F^r$ -th atom and random  $j_B^r$ -th atom with probability proportional to (9);
- Merge the supports  $S_F^j = S_F^{j-1} \cup \{j_F^r\}$ , and  $S_B^j = S_B^{j-1} \cup \{j_B^r\}$ ;
- Compute  $[\mathbf{a}_k^F]_{ij}$  and  $[\mathbf{a}_k^B]_{ij}$ , using (41) and (42);
- Compute  $\mathbf{r}_j^F = \mathbf{W}_{i,j} \mathbf{I}_k - \mathbf{d}^F [\mathbf{a}_k^F]_{ij}$  and  
 $\mathbf{r}_j^B = \mathbf{W}_{i,j} \mathbf{I}_k - \mathbf{d}^B [\mathbf{a}_k^B]_{ij}$ ;
- When  $\mathbf{r}_j^F < \epsilon$  and  $\mathbf{r}_j^B < \epsilon$ , quit;

3: **end for**

4: Compute the averaging coefficients  $[\mathbf{a}_k]_{ij}$ , using (10).

5: Choose the compact sparse coefficients  $[\mathbf{a}]_{ij}$ , using (11).

**Output:** the fusing coefficients  $[\mathbf{a}]_{ij}$ .

---

## APPENDIX B

### CONVERGENCE ANALYSIS ON LEARNING THE COUPLED DICTIONARY

In this section, we discuss the convergence analysis on learning the coupled dictionary. The material presented here is the review of the ideas [41], [42]. The main optimization problem in (17) is not guaranteed to seek the global optimum, and we intend to discuss the local minimum. Each dictionary updating step for  $\mathbf{d}^F$ ,  $\mathbf{d}^B$  with a fixed sparse representation  $\gamma$ , we optimize the

error matrices  $\mathbf{E}_j^F$  and  $\mathbf{E}_j^B$  in (22) and (23). Empirically, the two dictionaries converge the local optimal with executing limited updating operations. In this case, the success of local convergence mainly depends on the sparse coding stage. Thus, we emphasize on discussing the minimum of joint sparse coding.

In the following convergence analysis, we assume a stationary point  $(\boldsymbol{\gamma}^*, \mathbf{d})$  generated by limited iterates during sparse coding stage. Following the definitions of  $\mathcal{F}(\cdot)$ , the following series of inequalities are described by

$$\mathcal{F}(\boldsymbol{\gamma}^1) \geq \mathcal{F}(\boldsymbol{\gamma}^2) \geq \cdots \geq \mathcal{F}(\boldsymbol{\gamma}^t) \geq \cdots \quad (43)$$

Here exists a subsequence  $(\boldsymbol{\gamma}^t, \mathbf{d})$  converging to the limit point  $(\boldsymbol{\gamma}^*, \mathbf{d})$ .  $\mathcal{G}(\cdot, \cdot)$  represents the upper bound of  $\mathcal{F}(\cdot)$ , then we obtain that

$$\mathcal{G}(\boldsymbol{\gamma}^t, \boldsymbol{\gamma}^t) = \mathcal{F}(\boldsymbol{\gamma}^t) \leq \mathcal{F}(\boldsymbol{\gamma}^{t+1}) \leq \mathcal{G}(\boldsymbol{\gamma}^{t+1}, \boldsymbol{\gamma}^t) \leq \mathcal{G}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^t). \quad (44)$$

Note that when  $t \rightarrow \infty$ , the inequalities (44) imply

$$\mathcal{G}(\boldsymbol{\gamma}^*, \boldsymbol{\gamma}^*) \leq \mathcal{G}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^*). \quad (45)$$

Thus, we conclude  $\mathcal{G}'(\boldsymbol{\gamma}, \boldsymbol{\gamma}^*) \geq 0$ , that means the update on  $\boldsymbol{\gamma}$  exists the local optimum  $(\boldsymbol{\gamma}^*, \mathbf{d})$ .

#### ACKNOWLEDGMENT

The authors would like to thank...

#### REFERENCES

- [1] T. Wan, C. Zhu, and Z. Qin, "Multifocus image fusion based on robust principal component analysis," *Pattern Recognit. lett.*, vol. 34, no. 9, pp. 1001–1008, Jul. 2013.
- [2] Y. Liu, S. Liu, and Z. Wang, "Multi-focus image fusion with dense SIFT," *Inf. Fusion*, vol. 23, pp. 139–155, May 2015.
- [3] S. Pertuz, D. Puig, M. A. Garcia, and A. Fusielo, "Generation of all-in-focus images by noise-robust selective fusion of limited depth-of-field images," in *IEEE Trans. Image Process.*, vol. 22, no. 3, pp. 1242–1251, Mar. 2013.
- [4] J. Tian, and L. Chen, "Multi-focus image fusion using wavelet-domain statistics," in *Proc. IEEE Int. Conf. Image Process.*, Hong Kong, 2010, pp. 1205–1208.
- [5] J. Tian, L. Chen, L. Ma, and W. Yu, "Multi-focus image fusion using a bilateral gradient-based sharpness criterion," in *Opt. Commun.*, vol. 284, no. 1, pp. 80–87, Jan. 2011.
- [6] Q. Zhang, and B. L. Guo, "Multifocus image fusion using the nonsubsampled contourlet transform," *Signal Process.*, vol. 89, pp. 1334–1346, Jul. 2009.
- [7] F. Nencini, A. Garzelli, S. Baronti, and L. Alparone, "Remote sensing image fusion using the curvelet transform," *Inf. Fusion*, vol. 8, no. 2, pp. 143–156, Apr. 2007.

- [8] G. Pajares and J. Cruz, "A wavelet-based image fusion tutorial," *Pattern Recognit.*, vol. 37, no. 9, pp. 1855–1872, Sep. 2004.
- [9] V. D. Calhoun and T. Adali, "Feature-based fusion of medical imaging data," *IEEE Trans. Inf. Technol. Biomedicine*, vol. 13, no. 5, pp. 711–720, Sep. 2009.
- [10] V. Aslantas and R. Kurban, "Fusion of multi-focus images using differential evolution algorithm," *Expert Syst. Appl.*, vol. 37, no. 12, pp. 8861–8870, Dec. 2010.
- [11] P. Burt and E. Adelson, "The laplacian pyramid as a compact image code," *IEEE Trans. Commun.*, vol. 31, no. 4, pp. 532–540, Apr. 1983.
- [12] S. Li, B. Yang, "Multifocus image fusion using region segmentation and spatial frequency," *Inf. Fusion*, vol. 26, no. 7, pp. 971–979, Jul. 2008.
- [13] Z. Zhou, S. Li, and B. Wang, "Multi-scale weighted gradient-based fusion for multi-focus images," *Image Vision Comput.*, vol. 20, pp. 60–72, Nov. 2014.
- [14] J. Tian and L. Chen, "Adaptive multi-focus image fusion using a wavelet-based statistical sharpness measure," *Signal Process.*, vol. 92, no. 9, pp. 2137–2146, Sep. 2012.
- [15] Q. Zhang and B. Guo, "Multifocus image fusion using the nonsubsampled contourlet transform," *Signal Process.*, vol. 89, no. 7, pp. 1334–1346, Jul. 2009.
- [16] T. Wan, N. Canagarajah, and A. Achim, "Compressive image fusion," in *Proc. IEEE 15th Int. Conf. Image Process.*, San Diego, CA, 2008, pp. 1308–1311.
- [17] S. Ambat, S. Chatterjee, and K. Hari, "Fusion of algorithms for compressed sensing," *IEEE Trans. Signal Process.*, vol. 61, no. 14, pp. 3699–3704, May. 2010.
- [18] T. Wan, Z. Qin, C. Zhu, and R. Liao, "A robust fusion scheme for multifocus images using sparse features," in *Proc. IEEE 38th Int. Conf. Acoustics, Speech and Signal Process.*, British Columbia, Canada, 2013, pp. 1957–1961.
- [19] H. Li, L. Li, and J. Zhang, "Multi-focus image fusion based on sparse feature matrix decomposition and morphological filtering," *Opt. Commun.*, vol. 342, pp. 1–11, May. 2015.
- [20] B. Yang and S. Li, "Multifocus image fusion and restoration with sparse representation," *IEEE Trans. Instrum. Meas.*, vol. 59, no. 4, pp. 884–892, Apr. 2010.
- [21] M. Nejati, S. Samavi, and S. hirani, "Multi-focus image fusion using dictionary-based sparse representation," *Inf. Fusion*, vol. 25, pp. 72–84, Sep. 2015.
- [22] Q. Zhang, and M. D. Levine, "Robust Multi-Focus Image Fusion Using Multi-Task Sparse Representation and Spatial Context," in *IEEE Trans. Image Process.*, vol. 25, no. 5, pp. 2045–2058, Mar. 2016.
- [23] L. Cao, L. Jin, H. Tao, G. Li, Z. Zhuang and Y. Zhang, "Multi-focus image fusion based on spatial frequency in discrete cosine transform domain," in *IEEE Signal Process. Lett.*, vol. 22, no. 2, pp. 220–224, Sep. 2015
- [24] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 4311–4322, Nov. 2006.
- [25] K. Engan, S. O. Aase, and J. H. Husoy, "Method of optimal directions for frame design," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Process.*, Phoenix, AZ, USA, 1999, pp. 2443–2446.
- [26] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online dictionary learning for sparse coding," in *Proc. 26th ACM Int. Conf. Mach. Learn.*, Montreal, QC, Canada, 2009, pp. 689–696.
- [27] 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, May. 2010.
- [28] J. Yang, Z. Wang, Z. Lin, S. Cohen, and T. Huang, "Coupled dictionary training for image super-resolution," *IEEE Trans. Image Process.*, vol. 21, no. 8, pp. 3467–3478, Aug. 2012.

- [29] J. Sadasivan, S. Mukherjee, and C. S. Seelamantula, "Joint dictionary training for bandwidth extension of speech signals," in *Proc. IEEE 41st Int. Conf. Acoustics, Speech and Signal Process.*, Shanghai, China, 2016, pp. 5925–5929.
- [30] S. Wang, L. Zhang, Y. Liang, and Q. Pan, "Semi-coupled dictionary learning with applications to image super-resolution and photo-sketch synthesis," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit.*, Rhode Island, USA, 2012, pp. 2216–2223.
- [31] R. Rubinstein, M. Zibulevsky, and M. Elad, "Double sparsity: Learning sparse dictionaries for sparse signal approximation," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1553–1564, Mar. 2010.
- [32] T. Peleg, and M. Elad, "A Statistical Prediction Model Based on Sparse Representations for Single Image Super-Resolution," *IEEE Trans. Image Process.*, vol. 23, no. 6, pp. 2569–2582, Jun. 2014.
- [33] A. Mesaros, T. Heittola, O. Dikmen, and T. Virtanen, "Sound event detection in real life recordings using coupled matrix factorization of spectral representations and class activity annotations," in *Proc. IEEE 40st Int. Conf. Acoustics, Speech and Signal Process.*, South Brisbane, QLD, 2015, pp. 151–155.
- [34] R. Gao, S. A. Vorobyov, and H. Zhao, "Multi-focus image fusion via coupled dictionary training," in *Proc. IEEE 41st Int. Conf. Acoustics, Speech and Signal Process.*, Shanghai, China, 2016, pp. 1666–1670.
- [35] T. Khler, X. Huang, F. Schebesch, A. Aichert, A. Maier, and J. Horngger , "Robust Multiframe Super-Resolution Employing Iteratively Re-Weighted Minimization," *IEEE Trans. Comput. Imag.*, vol. 2, no. 1, pp. 42–58, Mar. 2016.
- [36] M. Subbarao, T. Choi, and A. Nikzad, "Focusing techniques," *Opt. Eng.*, vol. 32, pp. 2824–2836, Mar. 1993.
- [37] M. Born, and E. Wolf, "Principles of Optics," in *Cambridge Univ. Press.*, 1999.
- [38] M. Elad and I. Yavneh, "A plurality of sparse representations is better than the sparsest one alone," *IEEE Trans. Inf. Theory*, vol. 55, no. 10, pp. 4701–4714, Oct. 2009.
- [39] R. Rubinstein, M. Zibulevsky, and M. Elad, "Efficient implementation of the K-SVD algorithm using batch orthogonal matching pursuit," *CS Technion*, vol. 40, no. 8, pp. 1–15, Apr. 2008.
- [40] L. N. Smith and M. Elad, "Improving dictionary learning: multiple dictionary updates and coefficient reuse," *IEEE Signal Process. Lett.*, vol. 20, no. 1, pp. 79–82, Jan. 2013.
- [41] M. Razaviyayn, H. W. Tseng and Z. Q. Luo, "Computational Intractability of Dictionary Learning for Sparse Representation," *arXiv preprint arXiv:1511.01776*, 2015.
- [42] M. Razaviyayn, M. Hong and Z. Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM J. Optim.*, vol. 23, no. 2, pp. 1126–1153, Jan. 2013.
- [43] M. Nikolova, S. Esedoglu, and T. F. Chan, "Algorithms for finding global minimizers of image segmentation and denoising models," *SIAM J. Appl. Math.*, vol. 66, no. 5, pp. 1632–1648, Jun. 2006.
- [44] M. Hossny, S. Nahavandi, and D. Creighton, "Comments on information measure for performance of image fusion," *Electron. Lett.*, vol. 44, no. 18, pp. 1066–1067, Aug. 2008.
- [45] C. Xydeas and V. Petrović, "Objective image fusion performance measure," *Electron. Lett.*, vol. 36, no. 4, pp. 308–309, Feb. 2000.