

COUPLED FEATURE LEARNING VIA STRUCTURED CONVOLUTIONAL SPARSE CODING FOR MULTIMODAL IMAGE FUSION

Farshad G. Veshki and Sergiy A. Vorobyov

Department of Signal Processing and Acoustics, Aalto University, Espoo, Finland

ABSTRACT

A novel method for learning correlated features in multimodal images based on convolutional sparse coding with applications to image fusion is presented. In particular, the correlated features are captured as coupled filters in convolutional dictionaries. At the same time, the shared and independent features are approximated using separate convolutional sparse codes and a common dictionary. The resulting optimization problem is addressed using alternating direction method of multipliers. The coupled filters are fused based on a maximum-variance rule, and a maximum-absolute-value rule is used to fuse the sparse codes. The proposed method does not entail any prelearning stage. The experimental evaluations using medical and infrared-visible image datasets demonstrate the superiority of our method compared to state-of-the-art algorithms in terms of preserving the details and local intensities as well as improving objective metrics.

Index Terms— Multimodal image fusion, convolutional sparse coding, structured dictionary learning.

1. INTRODUCTION

Multimodal image fusion aims at merging the information from multiple images acquired using different imaging modalities into a single image, without introducing noise or artifacts [1, 2]. For instance, in medical image fusion, different information about the anatomies of tissues or the levels of biological activities captured using various medical imaging modalities are aggregated in a single fused image [1]. In surveillance applications, combining the visual information in optical images and the thermal information captured using infrared imaging techniques yield more informative images, and has applications, for example, in night vision [2].

A common approach for addressing the multimodal image fusion problem is to decompose the input images into multiscale or morphologically distinct components. This is usually done by employing deterministic mathematical models such as multiscale transforms [3–5]. Other techniques used for a similar purpose include subspace learning [6], dictionary learning [7, 8], and deep learning [9, 10]. An assumption made by all aforementioned fusion techniques is that the features (components) with similar structural properties convey

the same type of information. Therefore, they are appropriate for fusion. However, the multimodal images may not obey this assumption. For example, in medical imaging, computed tomography (CT) captures hard tissues and structures (e.g., bones and implants) with a higher resolution, while using magnetic resonance (MR) imaging, the details of soft tissues (e.g., fat and bone marrow) are reflected more effectively [1]. In infrared-visible images, the details in each input image provide different types of information. Thus, a fusion based on the similarity of structural properties can lead to the omission of important information.

In a recent work [11], we demonstrated that the fusion performance can be considerably improved by replacing the conventional deterministic feature-extraction techniques with a data-driven approach for extracting correlated features in multimodal images. Specifically, a method based on coupled dictionary learning [12] and a Pearson correlation constraint has been developed to decompose the multimodal images into their correlated and independent components. In particular, the correlated features have been captured as pairs of atoms in the coupled dictionaries. Then, the fusion is performed using the most significant representations of the coupled atoms. Since the information in the independent components is specific to each modality, these components are transferred to the fused image directly. This approach has shown to be superior in terms of preserving important information while yielding an improved contrast resolution [11].

In this paper, we present a coupled feature learning (CFL) method based on convolutional sparse coding (CSC) and dictionary learning. CSC incorporates a global single-valued model that, unlike standard sparse approximation, does not require patch extraction and enables shift-invariant dictionary learning [13]. In addition, instead of minimizing linear correlations between independent components (as in [11]), we incorporate a more general model that promotes statistical independence. We also propose novel schemes for fusion of correlated features and reconstruction of the final fused image. Experimental evaluations using multimodal medical and infrared-visible image datasets show that the proposed method significantly improves the performance of state-of-the-art multimodal fusion techniques. A MATLAB implementation of our fusion method is available at <https://github.com/FarshadGVeshki/ConvCFL-MMIF>.

2. CONVOLUTIONAL COUPLED FEATURE LEARNING

The proposed model decomposes n input multimodal images $\mathbf{s}^i \in \mathbb{R}^N$, $i = 1, \dots, n$, where N is the number of pixels in the images, into their correlated, shared and independent components. For simplicity of notations, we use one-dimensional arrays to represent the images. Generalization to multi-dimensional arrays is mathematically straightforward.

2.1. Problem Formulation

The correlated components are captured using a set of common sparse feature maps $\mathbf{\Gamma} \in \mathbb{R}^{N \times K}$ and coupled convolutional dictionaries $\mathbf{D}^i \in \mathbb{R}^{M \times K}$, $i = 1, \dots, n$. The shared and independent components are represented using a common dictionary $\mathbf{C} \in \mathbb{R}^{M \times L}$ and separate sparse feature maps $\mathbf{X}^i \in \mathbb{R}^{N \times L}$, $i = 1, \dots, n$. The decomposition problem can then be formulated as the following optimization problem

$$\begin{aligned} \underset{\{\mathbf{D}^i\}_{i=1}^n, \mathbf{C}, \{\mathbf{X}^i\}_{i=1}^n, \mathbf{\Gamma}}{\text{minimize}} \quad & \frac{1}{2} \sum_{i=1}^n \left\| \sum_{k=1}^K \mathbf{D}_k^i * \mathbf{\Gamma}_k + \sum_{l=1}^L \mathbf{C}_l * \mathbf{X}_l^i - \mathbf{s}^i \right\|_2^2 + \lambda_1 \sum_{k=1}^K \|\mathbf{\Gamma}_k\|_1 \\ & + \lambda_2 \sum_{i=1}^n \sum_{l=1}^L \|\mathbf{X}_l^i\|_1 \quad \text{s.t.} \quad \mathbf{C}_l, \mathbf{D}_k^i \in \mathcal{D} \quad \forall k, l, i, \end{aligned} \quad (1)$$

where $\mathcal{D} = \{\mathbf{d} \in \mathbb{R}^M \mid \|\mathbf{d}\|_2 \leq 1\}$ is the set of dictionary filters and $\lambda_1 > 0$ and $\lambda_2 > 0$ are regularization parameters. Subscripts are used to denote the columns of matrices.

The overlapping nonzero entries in $\{\mathbf{X}^i\}_{i=1}^n$ indicate that one of the dictionary filters $\{\mathbf{C}_l\}_{l=1}^L$ is used for approximation of multiple images at the same location, thus, it represents a shared feature. In addition, when only one of the entries in $\{\mathbf{X}^i\}_{i=1}^n$ is nonzero at one location, it means that one of the filters in $\{\mathbf{C}_l\}_{l=1}^L$ is used for only one source image. Thus, it represents an independent feature. Note that the convolutional filters are assumed to be statistically independent. As first demonstrated in [14], dictionary learning promotes statistical independence. The proof relies on the fact that accurate sparse codes preserve the information (joint entropy) in the source signal. Moreover, the sparsity regularization minimizes the entropy in each of the sparse codes (simply by maximizing the probability of one event (being zero) and minimizing the probability of all other events (being nonzero)). Therefore, by enforcing the equality of the joint entropy and the sum of the entropies of the individual sparse codes, sparse dictionary learning promotes statistical independence.

2.2. Optimization

Problem (1) is typically solved by alternating between minimization over the sparse codes and the dictionary filters. Since we address both steps in Fourier domain (using [15]), we first zero-pad all of the dictionary filters to the size of the sparse coefficient maps (\mathbb{R}^N).

2.2.1. Sparse Coding Step

Using the consensus ADMM method [16], (1) can be addressed with respect to the sparse feature maps $\{\mathbf{\Gamma}, \{\mathbf{X}^i\}_{i=1}^n\}$ by solving the following optimization problem

$$\begin{aligned} \underset{\{\mathbf{X}^i\}_{i=1}^n, \mathbf{\Gamma}}{\text{minimize}} \quad & \frac{1}{2} \sum_{i=1}^n \left\| \sum_{k=1}^K \mathbf{D}_k^i * \mathbf{\Theta}_k + \sum_{l=1}^L \mathbf{C}_l * \mathbf{Y}_l^i - \mathbf{s}^i \right\|_2^2 + \lambda_1 \sum_{k=1}^K \|\mathbf{\Gamma}_k\|_1 + \lambda_2 \sum_{i=1}^n \sum_{l=1}^L \|\mathbf{X}_l^i\|_1 \\ \text{s.t.} \quad & \mathbf{\Gamma} = \mathbf{\Theta}^i, \quad \mathbf{X}^i = \mathbf{Y}^i \quad i = 1, \dots, n. \end{aligned}$$

Using scaled Lagrangian multipliers $\mathbf{U}^i \in \mathbb{R}^{N \times K}$ and $\mathbf{V}^i \in \mathbb{R}^{N \times L}$, $i = 1, \dots, n$, the augmented Lagrangian is written as

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^n \left\| \sum_{k=1}^K \mathbf{D}_k^i * \mathbf{\Theta}_k + \sum_{l=1}^L \mathbf{C}_l * \mathbf{Y}_l^i - \mathbf{s}^i \right\|_2^2 + \lambda_1 \sum_{k=1}^K \|\mathbf{\Gamma}_k\|_1 + \lambda_2 \sum_{i=1}^n \sum_{l=1}^L \|\mathbf{X}_l^i\|_1 \\ + \frac{\rho}{2} \sum_{i=1}^n \left(\sum_{k=1}^K \|\mathbf{\Theta}_k^i - \mathbf{\Gamma}_k + \mathbf{U}_k^i\|_2^2 + \sum_{l=1}^L \|\mathbf{Y}_l^i - \mathbf{X}_l^i + \mathbf{V}_l^i\|_2^2 \right), \end{aligned} \quad (2)$$

where $\rho > 0$ is the penalty parameter. The ADMM iterations consist of minimizing (2) alternatively with respect to $\{\mathbf{\Theta}^i, \mathbf{Y}^i\}_{i=1}^n$, $\{\mathbf{\Gamma}, \{\mathbf{X}^i\}_{i=1}^n\}$ and $\{\mathbf{U}^i, \mathbf{V}^i\}_{i=1}^n$. The details of each subproblem are explained in the following. Denoting $\mathbf{Z}^i = \{\mathbf{\Theta}^i, \mathbf{Y}^i\}$, $\mathbf{F}^i = \{\mathbf{D}^i, \mathbf{C}\}$ and $\mathbf{W}^i = \{\mathbf{\Gamma} - \mathbf{U}^i, \mathbf{X}^i - \mathbf{V}^i\}$, we can update $\{\mathbf{\Theta}^i, \mathbf{Y}^i\}_{i=1}^n$ by solving n parallel optimization problems, which can be written as follows

$$(\mathbf{z}^i)^+ = \underset{\mathbf{z}^i}{\text{argmin}} \quad \frac{1}{2} \left\| \sum_{p=1}^P \mathbf{F}_p^i * \mathbf{z}_p^i - \mathbf{s}^i \right\|_2^2 + \frac{\rho}{2} \sum_{p=1}^P \|\mathbf{z}_p^i - \mathbf{w}_p^i\|_2^2, \quad (3)$$

where $P = K + L$ and $(\cdot)^+$ denotes the updated variables. The problem in (3) is a standard convolutional fitting problem that can be addressed using available CSC methods (e.g., [15]).

Updating $\{\mathbf{\Gamma}, \{\mathbf{X}^i\}_{i=1}^n\}$ can be efficiently addressed in an elementwise manner using the shrinkage operator $\mathcal{S}_\kappa(a) = \text{sign}(a) \max(0, |a| - \kappa)$. The updates are written as follows

$$\mathbf{\Gamma}^+ = \mathcal{S}_{\lambda_1/\rho} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{\Theta}^i + \mathbf{U}^i \right), \quad (\mathbf{X}^i)^+ = \mathcal{S}_{\lambda_2/\rho} (\mathbf{Y}^i + \mathbf{V}^i), \quad i = 1, \dots, n.$$

Finally, the updates for the scaled Lagrangian variables $\{\mathbf{U}^i, \mathbf{V}^i\}_{i=1}^n$ are given as

$$(\mathbf{U}^i)^+ = \mathbf{\Theta}^i - \mathbf{\Gamma} + \mathbf{U}^i, \quad (\mathbf{V}^i)^+ = \mathbf{Y}^i - \mathbf{X}^i + \mathbf{V}^i, \quad i = 1, \dots, n.$$

2.2.2. Dictionary Update Step

Using the consensus ADMM, (1) can be reformulated with respect to the dictionary filters $\{\mathbf{C}, \{\mathbf{D}^i\}_{i=1}^n\}$ as

$$\begin{aligned} \underset{\{\mathbf{D}^i\}_{i=1}^n, \mathbf{C}}{\text{minimize}} \quad & \frac{1}{2} \sum_{i=1}^n \left\| \sum_{k=1}^K \mathbf{G}_k^i * \mathbf{\Gamma}_k + \sum_{l=1}^L \mathbf{H}_l^i * \mathbf{X}_l^i - \mathbf{s}^i \right\|_2^2 + \Omega(\{\mathbf{C}, \{\mathbf{D}^i\}_{i=1}^n\}) \\ \text{s.t.} \quad & \mathbf{C} = \mathbf{H}^i, \quad \mathbf{D}^i = \mathbf{G}^i, \quad i = 1, \dots, n \end{aligned}$$

where $\Omega(\cdot)$ is an indicator function of the constraint set in (1). The augmented Lagrangian is written as follows

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^n \left\| \sum_{k=1}^K \mathbf{G}_k^i * \mathbf{\Gamma}_k + \sum_{l=1}^L \mathbf{H}_l^i * \mathbf{X}_l^i - \mathbf{s}^i \right\|_2^2 + \Omega(\{\mathbf{C}, \{\mathbf{D}^i\}_{i=1}^n\}) \\ + \frac{\sigma}{2} \sum_{i=1}^n \left(\sum_{k=1}^K \|\mathbf{G}_k^i - \mathbf{D}_k^i + \mathbf{R}_k^i\|_2^2 + \sum_{l=1}^L \|\mathbf{H}_l^i - \mathbf{C}_l + \mathbf{T}_l^i\|_2^2 \right), \end{aligned} \quad (4)$$

where $\mathbf{R}^i \in \mathbb{R}^{N \times K}$ and $\mathbf{T}^i \in \mathbb{R}^{N \times L}$, $i = 1, \dots, n$, are scaled Lagrangian variables. The ADMM iterations consist of minimizing (4) alternatively with respect to $\{\mathbf{G}^i, \mathbf{H}^i\}_{i=1}^n$, $\{\mathbf{C}, \{\mathbf{D}^i\}_{i=1}^n\}$ and $\{\mathbf{R}^i, \mathbf{T}^i\}_{i=1}^n$.

Indeed, denote $\mathbf{E}^i = \{\mathbf{G}^i, \mathbf{H}^i\}$, $\mathbf{S}^i = \{\mathbf{\Gamma}, \mathbf{X}^i\}$ and $\mathbf{Q}^i = \{\mathbf{D}^i - \mathbf{R}^i, \mathbf{C} - \mathbf{T}^i\}$, $i = 1, \dots, n$. Then updating $\{\mathbf{G}^i, \mathbf{H}^i\}_{i=1}^n$ can be addressed by solving n parallel optimization problems

$$(\mathbf{E}^i)^+ = \underset{\mathbf{E}^i}{\operatorname{argmin}} \frac{1}{2} \left\| \sum_{p=1}^P \mathbf{E}_p^i * \mathbf{s}_p^i - \mathbf{s}^i \right\|_2^2 + \frac{\sigma}{2} \sum_{p=1}^P \left\| \mathbf{E}_p^i - \mathbf{Q}_p^i \right\|_2^2. \quad (5)$$

Problem (5) is similar to (3) and can be efficiently addressed using available CSC methods (e.g., [15]).

Updating $\{\mathbf{C}, \{\mathbf{D}^i\}_{i=1}^n\}$ is performed as

$$(\mathbf{D}^i)^+ = \operatorname{proj}_{\mathcal{D}}(\mathbf{G}^i + \mathbf{R}^i), i = 1, \dots, n, \quad \mathbf{C}^+ = \operatorname{proj}_{\mathcal{D}}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{H}^i + \mathbf{T}^i\right),$$

where $\operatorname{proj}_{\mathcal{D}}(\cdot)$ denotes the orthogonal projection onto the set \mathcal{D} . This can be done by mapping the entries outside the constraint support to zero and then projecting the filters on the unit-ball.

The updates for scaled Lagrangian variables are given as

$$(\mathbf{R}^i)^+ = \mathbf{G}^i - \mathbf{D}^i + \mathbf{R}^i, \quad (\mathbf{T}^i)^+ = \mathbf{H}^i - \mathbf{C} + \mathbf{T}^i, \quad i = 1, \dots, n.$$

We perform the sparse coding and the dictionary update steps in an interleaved manner (one iteration of each step is executed before passing the variables to the next). The updated ADMM variables (auxiliary variables and scaled Lagrangian multipliers) are used to initialize the next iteration.

2.3. Projection on the Sparse Support

After the convolutional CFL stage, we can still significantly improve the approximation accuracy by orthogonalizing the residuals on the supports (the set of indices of nonzero entries) of the sparse coefficient maps. For this purpose, we use a gradient descent (GD) approach. Based on the convolution theorem, the GD iterations are found as

$$\begin{aligned} [\mathbf{\Gamma}_k^+]_{\mathcal{S}(\mathbf{\Gamma}_k)} &= [\mathbf{\Gamma}_k]_{\mathcal{S}(\mathbf{\Gamma}_k)} - \alpha \left[\operatorname{DFT}^{-1} \left(\sum_{i=1}^n \bar{\mathbf{D}}_k^i \odot \hat{\mathbf{r}}^i \right) \right]_{\mathcal{S}(\mathbf{\Gamma}_k)}, \forall k, \\ [(\mathbf{X}_l^i)^+]_{\mathcal{S}(\mathbf{X}_l^i)} &= [\mathbf{X}_l^i]_{\mathcal{S}(\mathbf{X}_l^i)} - \alpha \left[\operatorname{DFT}^{-1} \left(\bar{\mathbf{C}}_l^i \odot \hat{\mathbf{r}}^i \right) \right]_{\mathcal{S}(\mathbf{X}_l^i)}, \forall l, i, \end{aligned}$$

where $\hat{(\cdot)}$ denotes the discrete Fourier transform and $\operatorname{DFT}^{-1}(\cdot)$ represents its inverse, $\bar{(\cdot)}$ denotes the complex-conjugate, \odot is the elementwise multiplication and operator $\mathcal{S}(\cdot)$ returns the support of an array. In addition, α is the stepsize and \mathbf{r}^i represents the residuals associated with \mathbf{s}^i , that is

$$\mathbf{r}^i = \sum_{k=1}^K \mathbf{D}_k^i * \mathbf{\Gamma}_k + \sum_{l=1}^L \mathbf{C}_l * \mathbf{X}_l^i - \mathbf{s}^i, i = 1, \dots, n.$$

3. MULTIMODAL IMAGE FUSION ALGORITHM

In this section, the steps of the proposed fusion method are explained. Note that the images are considered as one-dimensional arrays. The elementwise operations are applied to all pixels.

3.1. Low-pass Filtering

The input images are first decomposed into base-layers $\{\mathbf{s}_b^i\}_{i=1}^n$ and details-layers $\{\mathbf{s}_d^i\}_{i=1}^n$ using low-pass filtering. This is done using the *lowpass* function from the SPORCO library [17] (the regularization parameter is set to 10).

3.2. Fusion of the Details-layers

The details-layers $\{\mathbf{s}_d^i\}_{i=1}^n$ are decomposed into the correlated, shared and independent components using the convolutional CFL method explained in Section 2.

3.2.1. Fusion of Coupled Features

The coupled features are fused based on the highest visual significance, which can be measured, for example, using *variance* (denoted as $\operatorname{var}(\cdot)$). This can be formulated as follows

$$\mathbf{D}_k^F = \mathbf{D}_k^{i^*}, \quad i^* = \underset{i=1, \dots, n}{\operatorname{argmax}} \left(\operatorname{var}(\mathbf{D}_k^i) \right), \quad k = 1, \dots, K,$$

where \mathbf{D}^F is the dictionary of fused coupled features.

3.2.2. Fusion of Shared and Independent Components

The fusion of shared and independent component \mathbf{X}^F is found by combining the redundant sparse codes $\{\mathbf{X}^i\}_{i=1}^n$ using maximum-absolute-value rule. This can be written as

$$\mathbf{X}_l^F(j) = \mathbf{X}_l^{i^*}(j), \quad i^* = \underset{i=1, \dots, n}{\operatorname{argmax}} \left(|\mathbf{X}_l^i(j)| \right), \quad j = 1, \dots, N, l = 1, \dots, L.$$

This allows to transfer the independent features along with the shared features with the most significant representation coefficients into the fused image.

The fused details-layer \mathbf{s}_d^F is then reconstructed using

$$\mathbf{s}_d^F = \sum_{k=1}^K \mathbf{D}_k^F * \mathbf{\Gamma}_k + \sum_{l=1}^L \mathbf{C}_l * \mathbf{X}_l^F.$$

3.3. Fusion of the Base-layers

We form two images \mathbf{s}_b^{max} and \mathbf{s}_b^{min} representing the maximum and the minimum allowed local intensities, using

$$\mathbf{s}_b^{max} = \max_{i=1, \dots, n} (\mathbf{s}_b^i), \quad \mathbf{s}_b^{min} = \omega \left(\max_{i=1, \dots, n} (\mathbf{s}_b^i) \right) + (1 - \omega) \left(\min_{i=1, \dots, n} (\mathbf{s}_b^i) \right),$$

where $0 \leq \omega \leq 1$, and $\max(\cdot)$ and $\min(\cdot)$ are the elementwise maximum and minimum operators, respectively.

It is favorable to incorporate s_b^{max} into the final fused image. However, this can cause a loss of information due to the limited range (0 to 1) of the standard images. To achieve a compromise between contrast resolutions and local intensities, we propose the following approach. First, the difference between the local maximum and minimum intensities (local variations) of s_d^F (for example, in a 3×3 neighborhood) is stored in s_d^v . Then the fused base-layer s_b^F is computed as

$$s_b^F(j) = \begin{cases} s_b^{max}(j), & \text{if } s_d^v(j) \leq 1 - s_b^{max}(j) \\ s_b^{min}(j), & \text{if } s_d^v(j) \geq 1 - s_b^{min}(j) \\ 1 - s_b^v(j), & \text{if otherwise} \end{cases}, \quad j = 1, \dots, N.$$

A Gaussian filter may be used to smooth s_b^F so that discontinuities are not introduced.

The final fused image s^F is then reconstructed as

$$s^F = s_b^F + s_d^F.$$

4. EXPERIMENTAL RESULTS

We compare our method to four recent multimodal fusion methods both visually and using objective evaluation metrics. We use two medical image fusion methods: a method based on the non-subsampled shearlet transform (NSST) [4] and a method based on Laplacian redecomposition (LRD) [5]. We also use two infrared-visible image fusion method: a method that incorporates a hierarchical Bayesian model (Bayes) [18] and a method based on deep learning (Resnet) [9]. The multimodal medical image dataset consists of 20 pairs of images collected from [19], and the infrared-visible image dataset includes 21 pairs of images taken from https://github.com/hli1221/imagefusion_resnet50/tree/master/IV_images. Four metrics are used for objective evaluations, the objective image fusion performance measure $Q_{AB/F}$ [20], the information measure for performance of image fusion Q_{IM} [21], spatial frequency (SF) [22] and the structural similarity index (SSIM) [23]. The algorithm parameters are $\lambda_1 = \lambda_2 = 0.01$, $K=8$, $L=12$, $\rho=\sigma=10$, $\alpha=0.01$ and $\omega = 0.9$. Moreover, we use 150 ADMM iterations, 100 GD iterations and 8×8 filters (M in two-dimensional case).

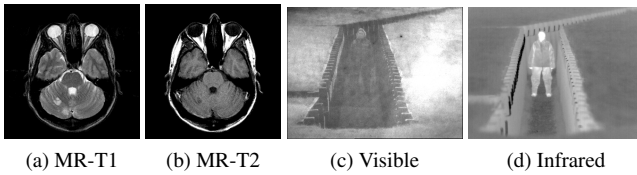


Fig. 1: Examples of multimodal images.

Fig. 1 shows a pair of images from each dataset used. The results obtained using different methods are shown in Fig. 2. Table 1 compares the average results for objective evaluation metrics for each dataset. The results show that the LRD

method leads to low contrast-resolutions, which is reflected in very low SF and $Q_{AB/F}$ values for this method. NSST also loses/blurs high-resolution information, while this information is well preserved using our method (see Figs. 2a and 2b, for example). Results obtained using Resnet and Bayes show an inferior fusion of local intensities, which results in low visibility of the details in the fused images (see Figs. 2d and 2e, for example). Overall, the proposed method results in the best performance in terms of the fusion of the high-resolution information as well as the local intensities (for example, see Figs. 2c and 2f). These observations can be validated by the objective evaluation results in Table 1, where our method obtains the best results in all cases.

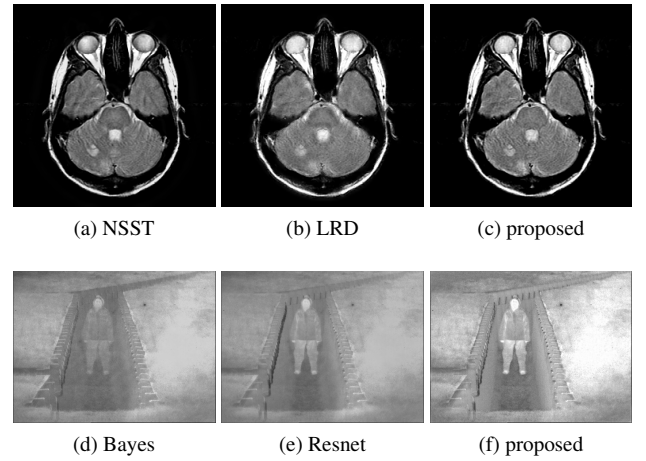


Fig. 2: The fusion results for the multimodal images in Fig. 1 using different methods.

Metrics	Medical			Infrared-Visible		
	NSST	LRD	proposed	Bayes	Resnet	proposed
$Q_{AB/F}$	0.5646	0.5278	0.5667	0.4561	0.3520	0.4941
Q_{IM}	0.6778	0.7341	0.7424	0.4044	0.3167	0.4311
SF	30.02	28.93	31.93	7.63	6.13	11.11
SSIM	0.5501	0.6601	0.7047	0.5040	0.5072	0.5121

Table 1: Average objective evaluation results for each dataset using different methods. The Best results are shown in bold.

5. CONCLUSION

A novel multimodal image fusion method based on convolutional sparse coding has been developed. A convolutional coupled feature learning algorithm has been proposed for the decomposition of multimodal images into correlated, shared, and independent features. Appropriate schemes have been proposed for the fusion of extracted features and reconstruction of the final image. The experimental results show significant improvements by the proposed method compared to the state-of-the-art multimodal image fusion methods.

6. REFERENCES

- [1] B. Huang, F. Yang, M. Yin, X. Mo, and C. Zhong, "A review of multimodal medical image fusion techniques," *Comput. Math. Methods. Med.*, vol. 2020, 2020.
- [2] X. Zhang, P. Ye, H. Leung, K. Gong, and G. Xiao, "Object fusion tracking based on visible and infrared images: A comprehensive review," *Inf. Fusion*, vol. 63, pp. 166–187, 2020.
- [3] G. Li, Y. Lin, and X. Qu, "An infrared and visible image fusion method based on multi-scale transformation and norm optimization," *Inf. Fusion*, vol. 71, pp. 109–129, 2021.
- [4] W. Tan, P. Tiwari, H. M. Pandey, C. Moreira, and A. K. Jaiswal, "Multimodal medical image fusion algorithm in the era of big data," *Neural Comput. Appl.*, 2020.
- [5] X. Li, X. Guo, P. Han, X. Wang, H. Li, and T. Luo, "Laplacian re-decomposition for multimodal medical image fusion," *IEEE Trans. Instrum. Meas.*, vol. 69, no. 9, pp. 6880–6890, 2020.
- [6] Y. Yang, S. Cao, S. Huang, and W. Wan, "Multimodal medical image fusion based on weighted local energy matching measurement and improved spatial frequency," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–16, 2021.
- [7] Y. Liu, X. Chen, R. K. Ward, and Z. J. Wang, "Medical image fusion via convolutional sparsity based morphological component analysis," *IEEE Signal Process. Lett.*, vol. 26, no. 3, pp. 485–489, 2019.
- [8] H. Li, X. He, D. Tao, Y. Tang, and R. Wang, "Joint medical image fusion, denoising and enhancement via discriminative low-rank sparse dictionaries learning," *Pattern Recognit.*, vol. 79, pp. 130–146, 2018.
- [9] H. Li, X. Wu, and T. S. Durrani, "Infrared and visible image fusion with resnet and zero-phase component analysis," *Infrared Physics and Technology*, vol. 102, 2019.
- [10] Z. Wang, X. Li, H. Duan, Y. Su, X. Zhang, and X. Guan, "Medical image fusion based on convolutional neural networks and non-subsampled contourlet transform," *Expert Systems with Applications*, vol. 171, 2021.
- [11] F. G. Veshki, N. Ouzir, S. A. Vorobyov, and E. Ollila, "Coupled feature learning for multimodal medical image fusion," *arXiv:2102.08641*, 2021.
- [12] F. G. Veshki and S. A. Vorobyov, "An efficient coupled dictionary learning method," *IEEE Signal Process. Lett.*, vol. 26, no. 10, pp. 1441–1445, 2019.
- [13] C. Garcia-Cardona and B. Wohlberg, "Convolutional dictionary learning: A comparative review and new algorithms," *IEEE Trans. Comput. Imaging*, vol. 4, no. 3, pp. 366–381, 2018.
- [14] B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, pp. 607–609, 1996.
- [15] F. G. Veshki and S. A. Vorobyov, "Efficient ADMM-based algorithms for convolutional sparse coding," *IEEE Signal Process. Lett.*, vol. 29, pp. 389–393, 2022.
- [16] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [17] B. Wohlberg, "SParse Optimization Research COde (SPORCO)," Software library available from <http://purl.org/brendt/software/sporco>, 2017.
- [18] Z. Zhao, S. Xu, C. Zhang, J. Liu, and J. Zhang, "Bayesian fusion for infrared and visible images," *Signal Process.*, vol. 177, pp. 1–12, 2020.
- [19] Harvard Medical School, "The Whole Brain Atlas," <http://www.med.harvard.edu/AANLIB/>, [Online; accessed 16-sep-2021].
- [20] C. Xydeas and V. Petrovic, "Objective image fusion performance measure," *Electron. Lett.*, vol. 36, no. 4, pp. 308–309, 2000.
- [21] G. Qu, D. Zhang, and P. Yan, "Information measure for performance of image fusion," *Electronics Letters*, vol. 38, no. 7, pp. 313–315, 2002.
- [22] A.M. Eskicioglu and P.S. Fisher, "Image quality measures and their performance," *IEEE Transactions on Communications*, vol. 43, no. 12, pp. 2959–2965, 1995.
- [23] Z. Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.