

## Aberystwyth University

### *Exploring a unified low rank representation for multi-focus image fusion*

Zhang, Qiang; Wang, Fan; Luo, Yongjiang; Han, Jungong

*Published in:*  
Pattern Recognition

*DOI:*  
[10.1016/j.patcog.2020.107752](https://doi.org/10.1016/j.patcog.2020.107752)

*Publication date:*  
2020

*Citation for published version (APA):*

Zhang, Q., Wang, F., Luo, Y., & Han, J. (2020). Exploring a unified low rank representation for multi-focus image fusion. *Pattern Recognition*, [107752]. <https://doi.org/10.1016/j.patcog.2020.107752>

**Document License**  
CC BY-NC-ND

#### **General rights**

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

#### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400  
email: [is@aber.ac.uk](mailto:is@aber.ac.uk)

# Exploring a Unified Low Rank Representation for Multi-focus Image Fusion

Qiang Zhang<sup>a, b</sup>, Fan Wang<sup>b</sup>, Yongjiang Luo<sup>c, \*</sup>, Jungong Han<sup>d, \*</sup>

<sup>a</sup> Key Laboratory of Electronic Equipment Structure Design, Ministry of Education, Xidian University, Xi'an, Shaanxi 710071, China

<sup>b</sup> Center for Complex Systems, School of Mechano-electronic Engineering, Xidian University, Xi'an, Shaanxi 710071, China

<sup>c</sup> School of Electronic Engineering, Xidian University, Xi'an, Shaanxi 710071, China

<sup>d</sup> Computer Science Department, Aberystwyth University, SY23 3FL, United Kingdom

---

**Abstract:** Recent years have witnessed a trend that uses image representation models, including sparse representation (SR), low-rank representation (LRR) and their variants for multi-focus image fusion. Despite the thrilling preliminary results, existing methods conduct the fusion patch by patch, leading to insufficient consideration of the spatial consistency among the image patches within a local region or an object. As a result, not only the spatial artifacts are easily introduced to the fused image but also the “jagged” artifacts frequently arise on the boundaries between the focused regions and the de-focused regions, which is an inherent problem in these patch-based fusion methods. Aiming to address the above problems, we propose, in this paper, a new multi-focus image fusion method integrating super-pixel clustering and a unified LRR (ULRR) model. The entire algorithm is carried out in three steps. In the first step, the source image is segmented into a few super-pixels with irregular sizes, rather than patches with regular sizes, to diminish the “jagged” artifacts and meanwhile to preserve the boundaries of objects on the fused image. Secondly, a super-pixel clustering-based fusion strategy is employed to further reduce the spatial artifacts in the fused images. This is achieved by using a proposed ULRR model, which imposes the low-rank constraints onto each super-pixel cluster. This is apparently more reasonable for those images with complicated scenes. Moreover, a Laplacian regularization term is incorporated in the proposed ULRR model to ensure the spatial consistency among the super-pixels with the same cluster. Finally, a measure of focus for each super-pixel is defined to seek the focused as well as de-focused regions in the source image via jointly using representation coefficients and sparse errors derived from the proposed ULRR model. Extensive experiments have been conducted and the results demonstrate the superiorities of the proposed fusion method in diminishing the spatial artifacts in the fused image and the “jagged” boundary artifacts between the focused and de-focused regions, compared to the state-of-the-art fusion algorithms.

**Keywords:** Multi-focus image fusion, Super-pixel clustering, Unified low-rank representation, Spatial consistency.

---

## 1. Introduction

Owing to the limited field depth of optical imaging systems, it is usually difficult, if not impossible, to acquire an image with all the objects in-focus [1]. Hence, only parts of an image have sharp appearances while the others look relatively blurring, which brings great inconvenience for human visual perception and sometimes computer processing as well. A lot of technologies are available to remedy this situation, in which multi-focus image fusion is a simple yet efficient way to combine multiple images

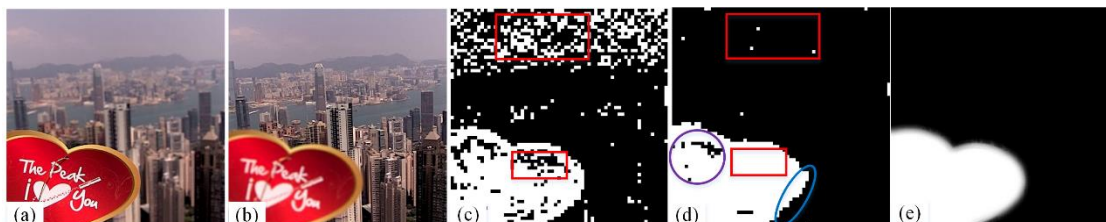
---

\* Corresponding author.  
E-mail address: jungong.han@aber.ac.uk (J. Han), yjluo@mail.xidian.edu.cn (Yongjiang Luo).

32 that shoot the same scene at different focal points into a single image, on which all objects are clearly  
33 displayed [2].

34 There are two basic requirements for a multi-focus image fusion method. One is that the focused  
35 regions should be determined and extracted from the given multi-focus input images and then preserved  
36 into the fused image, while all the defocused regions should be discarded [1]. The other is spatial artifacts  
37 or inconsistencies should be introduced to the fused image as little as possible during the fusion. Hence,  
38 how to accurately identify the focused and de-focused regions given the source images and how to  
39 combine the focused regions organically are two open questions in multi-focus image fusion. Our answer  
40 here is a new multi-focus image fusion method that employs super-pixel clustering and unified low-rank  
41 representation model.

42 So far, a number of multi-focus image fusion methods have been presented. A thorough review of  
43 these methods can be found in [2]. Among these, the fusion methods based on image representation  
44 models, e.g., sparse representation (SR) [3, 4], low rank representation (LRR) [5, 6] and their different  
45 extensions [7, 8], have attracted considerable attention in recent years attributed to their flexibilities.  
46 Usually, most image representation-based fusion methods are implemented patch by patch. Concretely,  
47 they start by dividing the input images into patches with regular shapes and the same sizes, and then  
48 carry out the fusion at the patch level.



49  
50 **Fig. 1.** Illustration of some decision maps. (a) and (b) Source images with focus on the front and the back, respectively; (c) Decision  
51 map obtained by [3] without sliding window, where the “white” and “black” points denote that the corresponding regions in the  
52 fused image are selected from Fig. 1(a) and Fig. 1(b), respectively; (d) Decision map obtained by [1] with spatial contextual  
53 information; (e) Decision map obtained by the proposed method.

54           However, most of these fusion methods just consider each image patch individually, which ignore  
55 the spatial consistency among those image patches within a local region or an object. As a result, some  
56 serious spatial artifacts appear on the focus decision maps (or the fused images), as shown in the red  
57 rectangular regions of Fig. 1(c).

58           For that, some fusion strategies have been proposed to suppress block artifacts and enhance the  
59 robustness against misregistration, among which the sliding window technology [3] is commonly  
60 employed. Despite its acceptable performance, sliding window usually leads to a huge requirement of  
61 memory storage as well as the increase of computational complexity. Alternatively, some spatial contexts  
62 or spatial consistency based strategies are presented in recent years [1, 7]. As displayed in the red  
63 rectangle regions of Fig. 1(d), these newly presented fusion strategies may diminish the spatial artifacts  
64 greatly. However, only the spatial consistency among those image patches within a local region is  
65 considered and the object area consistency among the patches within an object is ignored in these fusion  
66 strategies. Consequently, as shown in the purple circular region of Fig. 1(d), some patches may be still  
67 determined to have different focus information from those images in the same object.

68           In fact, an object in a multi-focus image is generally either wholly in-focus or out-of-focus due to  
69 the fact that the camera lens usually focuses on an object when taking a picture. Accordingly, those image  
70 patches within the same object may be similar in focus, i.e., they are all in-focus or all out-of-focus.

71           In addition to those spatial artifacts introduced in the fused image, “jagged” artifacts also arise  
72 frequently on the boundaries between the focused regions and the de-focused regions, as shown in the  
73 blue elliptical region of Fig. 1(d). This is an inherent problem in the patch-based fusion methods. Besides,  
74 it should be noted that most existing methods directly employ the intensity values as the feature for each  
75 patch, which are sensitive to the noise or illumination changes, especially, for smooth regions. As shown

76 in the rectangle regions of Fig. 1(c) and Fig. 1 (d) again, those isolated regions usually appear in those  
77 smooth regions that containing few details.

78 In order to address such problems arising in those existing image representation-based fusion  
79 methods, we present a super-pixel clustering based multi-focus image fusion method via a unified low-  
80 rank representation (ULRR) model. First, the input images are segmented into some super-pixels with  
81 irregular shapes rather than patches with fixed shapes to reduce the “jagged” artifacts between the  
82 focused and de-focused regions and meanwhile to preserve the boundaries of objects in the fused image.  
83 As well, multiple types of features, including colors, edges and textures, are extracted for each super-  
84 pixel to boost the focus discrimination.

85 Secondly, the super-pixels having similar features in each source image are first grouped into  
86 different clusters. Then these clusters are represented by using a proposed ULRR model considering the  
87 low-rankness (or correlations) of the super-pixels within a cluster. The proposed ULRR model is a sort  
88 of improved version of the traditional LRR model [9] by incorporating a Laplacian regularization term  
89 with respect to the representation coefficients. Here, the Laplacian regularization term intends to enforce  
90 the spatially adjacent super-pixels from the same cluster to be similar in representation coefficients and  
91 thus end up having similar focus information.

92 Finally, a measure of focus (MOF) is defined for each super-pixel by engaging the representation  
93 coefficients and sparse errors obtained by ULRR to compute a focus decision map, which in turn guides  
94 the fusion procedure of various source images. As displayed in Fig. 1(e), the focused and de-focused  
95 regions can be well determined by using the proposed method, on which much fewer spatial artifacts  
96 appeared on the fused image, and the “jagged” artifacts between the focused and de-focused regions  
97 disappeared. Experimental results verify the superiorities of our proposed fusion method over some state-

98 of-the-arts, even including some deep learning based methods, in diminishing the spatial artifacts in the  
99 fused image and the “jagged” boundary artifacts between the focused and de-focused regions.

100 The main contributions of this paper are highlighted as follows.

101 (1) A new multi-focus image fusion method based on clustering is proposed, where the spatial  
102 consistency among the local regions within an object is considered to reduce the spatial artifacts and to  
103 enhance the object area consistency in the fused image. This is different from that in [1] and [7], where  
104 only the local consistency among spatially-adjacent patches is considered.

105 (2) A unified low-rank representation (ULRR) model is proposed to capture the “intrinsic” low-  
106 rankness of each super-pixel cluster in our method, ensuring the spatial consistency among adjacent  
107 super-pixels within the same cluster or object. In addition, a new dictionary is constructed for ULRR.

108 (3) The proposed fusion method is implemented super-pixel by super-pixel, rather than in a patch  
109 based way as that in most existing SR and LRR based fusion methods, to reduce the “jagged” artifacts  
110 between the focused and de-focused regions and meanwhile to preserve the boundaries of objects in the  
111 fused image. Moreover, multiple types of features, including colors, edges, and textures, are extracted  
112 for each super-pixel to boost the focus discrimination. This is also beyond the traditional SR or LRR  
113 based fusion methods, where the intensity values are directly adopted as the features.

114 The rest of this paper is organized as follows. The related work is briefly introduced in Section 2,  
115 while the details of the proposed method are elaborated in Section 3. Experimental results as well as  
116 conclusions are given in Section 4 and Section 5, respectively.

## 117 **2. Related works**

118 So far, tremendous efforts have been devoted to multi-focus image fusion and numerous fusion  
119 algorithms have been presented, which fall into two groups: transform domain based methods and spatial

120 domain based methods.

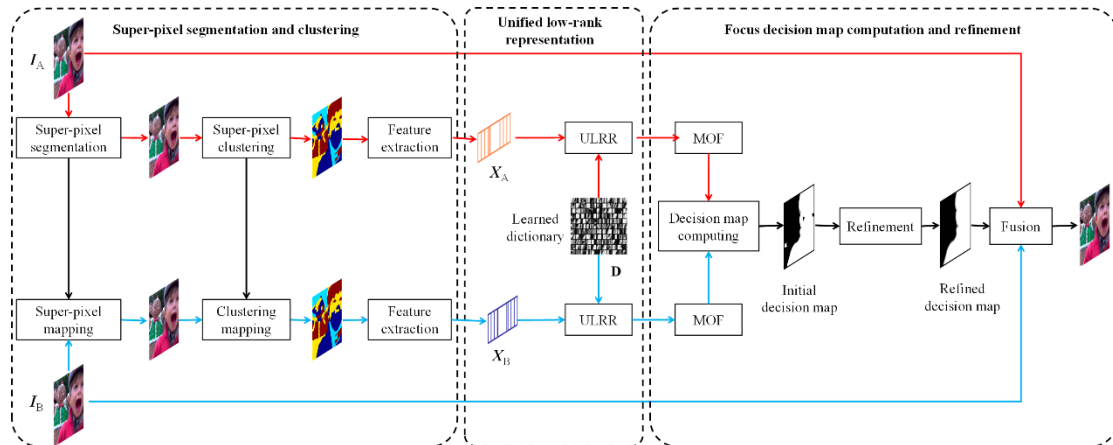
121       Among the former, multi-scale transform (MST) based fusion methods are the trends and have been  
122 discussed over the years [10, 11]. Different from transform domain based ones, spatial domain based  
123 methods directly extract the fused images (or image patches) from the source images via some measures  
124 of focus (MOFs) [12]. This usually caused many undesirable spatial artifacts. For that, some block and  
125 region based fusion methods have been presented in the past few years [13, 14, 15,]. Especially, in [15],  
126 a regional approach based on super-pixel segmentation and mean filtering was proposed. Currently, based  
127 on image matting [16], guided filtering [17], edge model [18] and conditional random field optimization  
128 [19], some new spatial domain based fusion methods have been presented to achieve state-of-the-art  
129 performance in information extraction and spatial consistency. A survey on these methods can be seen in  
130 [2, 20].

131       Recently, some new image representation models, such as sparse representation (SR) [3, 4, 21], low  
132 rank representation (LRR) [5, 6] and their variants [7, 8] have been employed to image fusion. For  
133 example, Yang *et al.* [3] took the first attempt in applying the SR theory to multi-sensor image fusion. In  
134 [21], Chen *et al.* introduced a multi-focus image fusion method based on clarity-enhanced image  
135 segmentation and regional sparse representation to strengthen its robustness against distortions that  
136 usually resulting from the pixel based coefficients selection. In our previous work [7], a robust sparse  
137 representation (RSR) based multi-focus image fusion was presented, where information from each local  
138 image patch and its spatial contextual information were jointly employed to determine the focused and  
139 de-focused regions. A multi-focus image fusion method based on dictionary learning and LRR was  
140 presented to achieve good performance in both global and local structures in [5]. The latent low-rank  
141 representation was used to extract the salient information of source images and guide the adaptive fusion

142 of low-pass sub-images in [8]. A thorough review and discussion about these fusion methods can be seen  
 143 in [22].

144 With the development of deep learning, some multi-focus image fusion methods based on deep  
 145 neural networks, e.g., convolutional neural networks (CNNs), have been proposed. Early CNN works  
 146 [23, 24] view the determination of each image patch to be in-focus or out-of-focus as a classification  
 147 problem. Later, some end-to-end networks are introduced for multi-focus image fusion [25, 26, 27].  
 148 Recently, several ensemble learning based multi-focus image fusion methods [28, 29] were presented,  
 149 where an ensemble of three CNNs were trained on three datasets to predict the decision maps without  
 150 the need of post-processing steps. Although these deep learning based methods may achieve satisfactory  
 151 performance, a massive amount of training data with labels are required to train such networks. This is a  
 152 challenging work for multi-focus image fusion.

### 153 3. Proposed method



154  
 155 **Fig. 2.** Diagram of the proposed multi-focus image fusion algorithm.

156 In this paper, only two source images are taken into account, and the images are supposed to have  
 157 been well registered in advance. Fig. 2 depicts the diagram of the proposed fusion method, which consists  
 158 of the following components: (1) Super-pixel segmentation and clustering; (2) Unified low-rank  
 159 representation; (3) Focus decision map computation and refinement. Based on the focus decision map,



160 the fused image is thus obtained. In addition, a local compact dictionary will be constructed from the  
161 average image of each pair of source images, when the source images are decomposed using the proposed  
162 ULRR model. In the following subsection, we will describe each component in detail.

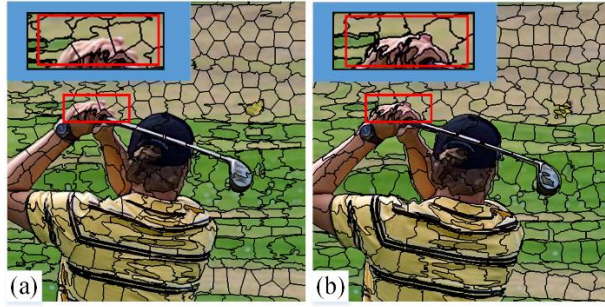
### 163 **3.1. Super-pixel segmentation and clustering**

#### 164 **3.1.1. Super-pixel segmentation**

165 Super-pixels group perceptually similar pixels to create visually meaningful entities while heavily  
166 reducing the number of primitives for subsequent processing steps [30]. Since they were first named in  
167 2003 [31], super-pixels have been widely applied to many computer vision tasks, including image fusion  
168 [15]. Compared with image patches of regular shapes, super-pixels can preferably preserve the boundary  
169 of objects in an image. Considering that, we adopt super-pixels, instead of image patches, in our proposed  
170 fusion method. So far, many super-pixel algorithms have been proposed, among which, linear spectral  
171 clustering (LSC) [32] is shown to achieve higher visual compactness and boundary adherence for natural  
172 images but with lower computational costs. Considering that, we adopt LSC for super-pixel segmentation  
173 in this paper.

174 As well, all the source images to be fused should be segmented into the same results so that their  
175 corresponding super-pixels can be properly merged in the subsequent fusion process. A commonly used  
176 way is to perform super-pixel segmentation on the average image of source images [15], and then map  
177 the super-pixel segmentation results to each source image. For multi-focus source images, some  
178 undesirable results may be obtained, especially for those transitional regions between focused and de-  
179 focused regions. For example, as shown in the rectangle region of Fig. 3(a), parts of the hands in the  
180 focused regions have been grouped into the same super-pixel with some de-focused regions. Accordingly,  
181 parts of these regions will be mistakenly determined to be focused or de-focused ones and the border of

182 the hands will be destroyed in the fused image.



183

184 **Fig. 3.** Illustration of super-pixel segmentation results by performing LSC on different input images. (a) On the average image of  
185 multi-focus source images; (b) On one of the multi-focus source images.

186 Alternatively, we will perform super-pixel segmentation on one of the source images, rather than  
187 on the average image, in order to obtain more accurate boundaries between focused and de-focused  
188 regions. As shown in the rectangle region of Fig. 3(b), each part of the hands is grouped into a super-  
189 pixel. In the subsequent fusion process, the border of the hands will be well preserved.

190 In summary, given a pair of source images, denoted by  $I_A$  and  $I_B$ , respectively, we first perform  
191 LSC on  $I_A$  to obtain a set of super-pixels  $\{sp_{A,i} | i = 1, 2, \dots, N\}$ , where  $N$  denotes the total number  
192 of super-pixels and is experimentally set to 350 in this paper. Then we map the segmentation results on  
193  $I_B$  and obtain  $\{sp_{B,i} | i = 1, 2, \dots, N\}$ .

### 194 3.1.2. Super-pixel clustering

195 In general, a super-pixel only denotes a regional atom without any perceptual meaning. Accordingly,  
196 as shown in Fig. 4 (c), each object in an image and the background may be constructed by many super-  
197 pixels with similar features. In real applications, we usually focus the lens on one object in the scene  
198 when taking a picture. As a result of that, the super-pixels within the same object in a multi-focus image  
199 may be all in-focus or all out-of-focus with a large probability. When the fusion is directly performed on  
200 super-pixels, the super-pixels within the same object may be mistakenly determined to have different  
201 focus information from the others. Some spatial artifacts may thus be easily introduced to the fused image.



202 **Fig. 4.** Super-pixel clusters. (a) and (b) Source images with focus on the front and the back, respectively; (c) Results of super-pixel  
 203 segmentation; (d) Results of super-pixel clustering.  
 204

205 In order to address such problem, we will first group the super-pixels in each source image into  
 206 different clusters and then consider the spatial consistency among the super-pixels within the same cluster  
 207 to introduce fewer spatial artifacts to the fused image. Similar to that in super-pixel segmentation, we  
 208 just group the super-pixels in one of the source images and then map the clustering results to the other  
 209 source image to ensure that the two multi-focus source images have the same clustering results.

210 In this paper, because of its popularity and simplicity, we adopt the  $k$ -means algorithm [33] to  
 211 achieve the super-pixel clustering, where the averaging RGB color values of all the pixels in each super-  
 212 pixel are employed as the super-pixel feature. Specifically, given the two source images  $I_{A/B}^1$  and their  
 213 corresponding two sets of super-pixels  $\{sp_{A/B,i} \mid i=1,2,\dots,N\}$ , two sets of super-pixel clusters  
 214  $\{C_{A/B,k} \mid k=1,2,\dots,K\}$  are obtained, where  $K$  denotes the number of clusters and will be discussed in  
 215 the experimental part. And each cluster  $C_{A/B,k}$  contains  $N_k$  super-pixels, i.e.  
 216  $C_{A/B,k} = \{sp_{A/B,k,i} \mid i=1,2,\dots,N_k\}$ . As shown in Fig. 4 (d), each object in the image is segmented into only  
 217 a fewer number of clusters, which will facilitate the consistency among the super-pixels within the same  
 218 object in the subsequent fusion process.

### 219 3.1.3. Feature extraction

220 In most SR or LRR based fusion methods, pixel intensity values are often directly employed as the

---

<sup>1</sup>The symbol  $A/B$  in  $I_{A/B}$  denotes  $A$  or  $B$ , i.e.,  $I_{A/B}$  means  $I_A$  or  $I_B$ . In the following contents, the definition of similar symbols is the same.

221 features, which are sensitive to noise or illumination changes. Some regions, especially those smooth  
 222 regions, are easily mistakenly determined to be in-focus or out-of-focus. In view of this, we extract  
 223 multiple types of features, including colors, edges and textures, rather than just the intensity, for each  
 224 super-pixel in our proposed fusion method. Specifically, feature extraction for each super-pixel and  
 225 super-pixel cluster can be described as follows.

226 (1) For each pixel  $p_{A/B,j}$  in one of the source images, construct its feature vector  $\mathbf{v}_{A/B,j} \in R^d$  of  
 227 dimension  $d = 44$ , including colors, edge and texture features. RGB color values as well as HIS (Hue,  
 228 Saturation, Intensity) components are extracted for each pixel, producing 6-dimensional color features.  
 229 For edge features, high pass filter, discrete wavelet and several edge operators (LOG, Prewitt, Sobel, *et*  
 230 *al.*) are performed onto the image, yielding 18-dimension filter responses at each location. Texture  
 231 features, which are constituted by the gray level co-occurrence matrix [34] of each super-pixel, contain  
 232 a total of 20-dimensional features, including contrast, energy, homogeneity, dissimilarity and difference  
 233 entropy.

234 (2) Construct the feature vector  $\mathbf{x}_{A/B,i} \in R^d$  for each super-pixel  $sp_{A/B,i}$  by averaging all the  
 235 feature vectors of pixels contained in the current super-pixel, i.e.,  $\mathbf{x}_{A/B,i} = \frac{1}{N_{sp_i}} \sum_{p_j \in sp_i} \mathbf{v}_{A/B,j}$ , where  $N_{sp_i}$   
 236 stands for the total number of pixels contained in the super-pixel  $sp_{A/B,i}$ .

237 (3) Construct the feature matrix  $\mathbf{X}_{A/B,k} \in R^{d \times N_k}$  for each super-pixel cluster  $C_{A/B,k}$  by using all of  
 238 the vectors of the super-pixels in the same cluster, i.e.,

$$239 \quad \mathbf{X}_{A,k} = [\mathbf{x}_{A,k,1}, \mathbf{x}_{A,k,2}, \dots, \mathbf{x}_{A,k,N_k}], \quad (1)$$

$$240 \quad \mathbf{X}_{B,k} = [\mathbf{x}_{B,k,1}, \mathbf{x}_{B,k,2}, \dots, \mathbf{x}_{B,k,N_k}], \quad (2)$$

241 where  $\mathbf{x}_{A/B,k,i}$  denotes the feature vector of the  $i$ -th super-pixel  $sp_{A/B,k,i}$  in the cluster  $C_{A/B,k}$ .

242

243 **3.2. Proposed unified low-rank representation (ULRR) for super-pixel clusters**

244 As shown in Fig. 4 (d), each super-pixel cluster represents a part of an object or a local region having  
 245 similar appearances in the scene. Therefore, the super-pixels from the same cluster are likely to be similar.  
 246 Accordingly, the feature matrix  $\mathbf{X}_k$ <sup>2</sup> for each super-pixel cluster constructed in the previous subsection  
 247 3.1.3 has “intrinsic” property of low-rankness. Therefore, a low-rank representation (LRR) model is a  
 248 natural choice for capturing the “intrinsic” low-rankness of each super-pixel cluster in our proposed  
 249 fusion method.

250 As a powerful analytical tool, LRR [9] intends to recover low-rank structures from the data  
 251 corrupted by sparse but strong noise. We may directly perform LRR on the feature matrix  $\mathbf{X}_k$ , but this  
 252 would ignore the spatial consistency among the super-pixels within the same cluster.

253 As discussed in the previous subsection 3.1.2, the spatially adjacent super-pixels residing within the  
 254 same cluster may have similar focus information, i.e., they may be all in-focus or all out-of-focus.  
 255 Therefore, these super-pixels will have similar “intrinsic” property, i.e., they may have similar  
 256 representation coefficients via LRR. Motivated by that, we think of a unified low-rank representation  
 257 (ULRR) model by incorporating a Laplacian regularization term with respect to the representation  
 258 coefficients into the traditional LRR model to capture the low-rankness of each super-pixel cluster.

259 **3.2.1. Unified low-rank representation model**

260 Given a dictionary  $\mathbf{D} \in R^{d \times M}$  with  $M$  atoms of dimension  $d$  and the feature matrix  $\mathbf{X}_k \in R^{d \times N_k}$   
 261 ( $k = 1, 2, \dots, K$ ) for each super-pixel cluster  $C_k$ , the unified low-rank representation model is  
 262 mathematically defined by

263 
$$\min_{\mathbf{Z}_1, \dots, \mathbf{Z}_k, \mathbf{E}_1, \dots, \mathbf{E}_k} \sum_{k=1}^K \|\mathbf{Z}_k\|_* + \alpha \|\mathbf{E}\|_{2,1} + \beta \text{tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^T), \text{ s.t. } \mathbf{X}_k = \mathbf{D}\mathbf{Z}_k + \mathbf{E}_k, k = 1, 2, \dots, K, \quad (3)$$

---

<sup>2</sup>In this subsection, we remove the symbol  $A/B$  from  $\mathbf{X}_{A/B,k}$  for generality.

264 where  $\mathbf{DZ}_k$  denotes the ‘‘intrinsic’’ low-rank part contained in the matrix  $\mathbf{X}_k$ .  $\mathbf{Z}_k \in R^{M \times N_k}$  denotes  
265 the representation coefficient matrix to be sought.  $\mathbf{E}_k \in R^{d \times N_k}$  represents the error or noise part.  $\|\mathbf{Z}_k\|_*$   
266 indicates the nuclear norm of the matrix  $\mathbf{Z}_k$  and is a convex relaxation of the rank function. The  
267 matrices  $\mathbf{Z} \in R^{M \times N}$  and  $\mathbf{E} \in R^{d \times N}$  are constructed by  $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_K]$  and  $\mathbf{E} = [\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_K]$ ,  
268 respectively.  $\|\mathbf{E}\|_{2,1}$  denotes the  $l_{2,1}$ -norm of the matrix  $\mathbf{E}$ . Minimizing  $\|\mathbf{E}\|_{2,1}$  enforces the matrix  $\mathbf{E}$   
269 to have column-sparsity.  $\alpha$  and  $\beta$  are two positive trade-off parameters to balance the effect of each  
270 part.

271 The Laplacian regularization term  $\text{tr}(\mathbf{ZLZ}^T)$  in Eq. (3) is defined by

$$272 \quad \text{tr}(\mathbf{ZLZ}^T) = \frac{1}{2} \sum_{i,j} \|\mathbf{z}_i - \mathbf{z}_j\|_2^2 \omega_{i,j}, \quad (4)$$

273 where  $\mathbf{z}_i$  denotes the  $i$ -th column of  $\mathbf{Z}$ . The weight  $\omega_{i,j}$  refers to the similarity between the  $i$ -th  
274 and  $j$ -th super-pixels  $sp_i$  and  $sp_j$ , and is computed by

$$275 \quad \omega_{i,j} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right), & \text{if } sp_i \text{ and } sp_j \text{ are spatially adjacent and belong to the same cluster.} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

276 Here,  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are the feature vectors of  $sp_i$  and  $sp_j$ , respectively.  $\sigma$  is a scalar parameter and  
277 we experimentally set it to  $\sqrt{0.5}$ . Given these weights, an affinity matrix  $\mathbf{W} \in R^{N \times N}$  with its  $(i, j)$ -th  
278 entry  $\mathbf{W}_{i,j} = \omega_{i,j}$  and a diagonal degree matrix  $\mathbf{\Lambda} \in R^{N \times N}$  with its  $i$ -th diagonal element  
279  $\Lambda_{i,i} = \sum_j \mathbf{W}_{i,j}$  can be constructed. The Laplacian matrix  $\mathbf{L} \in R^{N \times N}$  is then defined as  $\mathbf{L} = \mathbf{\Lambda} - \mathbf{W}$ .

280 Eq. (3) presents a convex optimization problem that can be solved by various methods. For that, we  
281 first convert it to the below equivalent one by involving some auxiliary variables in this paper:

$$282 \quad \min_{\substack{\mathbf{Z}_1, \dots, \mathbf{Z}_K, \\ \mathbf{E}_1, \dots, \mathbf{E}_K}} \sum_{k=1}^K \|\mathbf{J}_k\|_* + \alpha \|\mathbf{E}\|_{2,1} + \beta \text{tr}(\mathbf{HLH}^T), \quad \text{s.t. } \mathbf{X}_k = \mathbf{DZ}_k + \mathbf{E}_k, \mathbf{Z}_k = \mathbf{J}_k, \mathbf{Z} = \mathbf{H}. \quad (6)$$

283 To solve it, a linearized alternating direction method with adaptive penalty (LADMAP) [35] is  
284 adopted, which requires minimizing the following augmented Lagrangian function

$$L = \sum_{k=1}^K \left( \|\mathbf{J}_k\|_* + \langle \mathbf{Y}_{1,k}, \mathbf{X}_k - \mathbf{DZ}_k - \mathbf{E}_k \rangle + \langle \mathbf{Y}_{2,k}, \mathbf{Z}_k - \mathbf{J}_k \rangle + \frac{\mu}{2} \|\mathbf{X}_k - \mathbf{DZ}_k - \mathbf{E}_k\|_F^2 + \frac{\mu}{2} \|\mathbf{Z}_k - \mathbf{J}_k\|_F^2 \right) + \alpha \|\mathbf{E}\|_{2,1} + \beta \text{tr}(\mathbf{H}\mathbf{L}\mathbf{H}^T) + \langle \mathbf{Y}_3, \mathbf{Z} - \mathbf{H} \rangle + \frac{\mu}{2} \|\mathbf{Z} - \mathbf{H}\|_F^2, \quad (7)$$

where Lagrange multipliers  $\mathbf{Y}_{1,k}$ ,  $\mathbf{Y}_{2,k}$  ( $k = 1, 2, \dots, K$ ) and  $\mathbf{Y}_3$  help to remove the equality constraint in Eq. (6).  $\mu > 0$  is a penalty term.  $\langle \mathbf{A}, \mathbf{B} \rangle$  represents the Euclidean inner product of  $\mathbf{A}$  and  $\mathbf{B}$ . This problem can thus be minimized with respect to  $\mathbf{Z}_k$  (or  $\mathbf{Z}$ ),  $\mathbf{E}_k$  (or  $\mathbf{E}$ ),  $\mathbf{Y}_{1,k}$ ,  $\mathbf{Y}_{2,k}$  ( $k = 1, 2, \dots, K$ ),  $\mathbf{Y}_3$ , and  $\mathbf{H}$ , respectively. Algorithm 1 briefly summarizes how we calculate the proposed ULRR, and more details are explained in Appendix A.

---

**Algorithm 1: Solving ULRR via LADMAP**

**Input:** Observed data  $\mathbf{X}_k$  ( $k = 1, 2, \dots, K$ ), dictionary  $\mathbf{D}$ , and parameters  $\alpha$  and  $\beta$

**Output:**  $\mathbf{Z}$  and  $\mathbf{E}$

**Initialization:**  $\mathbf{Z}^0 = \mathbf{0}$ ,  $\mathbf{E}^0 = \mathbf{0}$ ,  $\mathbf{J}_k^0 = \mathbf{0}$ ,  $\mathbf{H}^0 = \mathbf{0}$ ,  $\mathbf{Y}_1^0 = \mathbf{0}$ ,  $\mathbf{Y}_2^0 = \mathbf{0}$ ,  $\mathbf{Y}_3^0 = \mathbf{0}$ ,  $\mu^0 = 10^{-6}$ ,  $\mu_{\max} = 10^6$ ,  $\varphi = 1.1$

---

**While** not converged **do**

(1) Fix the others and update  $\mathbf{J}_k$  ( $k = 1, 2, \dots, K$ ) using Eq. (A2);

(2) Fix the others and update  $\mathbf{H}$  using Eq. (A4);

(3) Fix the others and update  $\mathbf{Z}_k$  ( $k = 1, 2, \dots, K$ ) and  $\mathbf{Z}$  using Eq. (A6);

(4) Fix the others and update  $\mathbf{E}$  using Eq. (A8);

(5) Update the multipliers  $\mathbf{Y}_{1,k}$ ,  $\mathbf{Y}_{2,k}$  ( $k = 1, 2, \dots, K$ ) and  $\mathbf{Y}_3$ :

$$\mathbf{Y}_{1,k}^{i+1} = \mathbf{Y}_{1,k}^i + \mu(\mathbf{X}_k - \mathbf{DZ}_k^{i+1} - \mathbf{E}_k^{i+1}), \quad \mathbf{Y}_{2,k}^{i+1} = \mathbf{Y}_{2,k}^i + \mu(\mathbf{Z}_k^{i+1} - \mathbf{J}_k^{i+1}), \quad \mathbf{Y}_3^{i+1} = \mathbf{Y}_3^i + \mu(\mathbf{Z}^{i+1} - \mathbf{H}^{i+1});$$

(6) Update  $\mu$ :

$$\mu^{i+1} = \min(\mu^i \varphi, \mu_{\max});$$

(7) Check the convergence conditions:

$$\max_k \|\mathbf{X}_k - \mathbf{DZ}_k^{i+1} - \mathbf{E}_k^{i+1}\|_{\infty} < \varepsilon, \quad \|\mathbf{Z}^{i+1} - \mathbf{Z}^i\|_{\infty} < \varepsilon, \quad \text{and} \quad \|\mathbf{E}^{i+1} - \mathbf{E}^i\|_{\infty} < \varepsilon;$$

where  $\|\cdot\|_{\infty}$  denotes the  $l_{\infty}$ -norm of a matrix.

**end while**

---

### 3.2.2. Dictionary construction

In addition to the ULRR model, the dictionary is also crucial to fusion success. The original feature matrices (e.g.,  $\mathbf{X}_{A,k}$  or  $\mathbf{X}_{B,k}$ ) from each source image may be directly employed as the dictionary [9] for ULRR. However, it is difficult to maintain the fairness of focus measure for the corresponding super-pixels from different source images. Alternatively, an adaptive dictionary is constructed from an image

296 obtained by averaging each pair of source images in our proposed fusion method when decomposing  
 297  $\mathbf{X}_{A,k}$  and  $\mathbf{X}_{B,k}$ . Moreover, as discussed in [36], the dictionary with fairly low-rank is more desirable  
 298 for LRR. Considering that, we will perform a Gaussian filtering on the average image before constructing  
 299 the dictionary. Specifically, the dictionary for ULRR is constructed as in Algorithm 2.

---

**Algorithm2: Dictionary construction**

---

- (1) For each pair of source images  $I_A$  and  $I_B$ , an image  $\bar{I}_{AB}$  is obtained by averaging the source images;
  - (2) A blurred average image  $\bar{I}'_{AB}$  is obtained by performing a Gaussian filtering with kernel size of  $8 \times 8$  on  $\bar{I}_{AB}$ .
  - (3) The super-pixel segmentation result for  $I_A$  (or  $I_B$ ) is mapped to  $\bar{I}'_{AB}$ , obtaining a set of super-pixels  $\{sp_{AB,i} | i = 1, 2, \dots, N\}$ .
  - (4) The features  $\{\mathbf{x}_{AB,i} \in \mathbb{R}^d | i = 1, 2, \dots, N\}$  are extracted for the super-pixels  $\{sp_{AB,i} | i = 1, 2, \dots, N\}$  by using the same way as in Subsection 3.1.3. A feature matrix is thus constructed by  $\mathbf{X}_{AB} = [\mathbf{x}_{AB,1}, \mathbf{x}_{AB,2}, \dots, \mathbf{x}_{AB,N}] \in \mathbb{R}^{d \times N}$ .
  - (5) A set of eigenvalues  $\{\lambda_i | \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N \geq 0, i = 1, 2, \dots, N\}$ , sorted in descending order, and their corresponding eigenvectors  $\{\mathbf{p}_i | i = 1, 2, \dots, N\}$  are obtained by performing principal component analysis (PCA) on the matrix  $\mathbf{X}_{AB}$ .
  - (6) A compact dictionary with  $M$  atoms of dimension  $d$  is constructed by the eigenvectors  $\{\mathbf{p}_i | i = 1, 2, \dots, M\}$  corresponding to the first  $M$  largest eigenvalues, i.e.,  $\mathbf{D} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M] \in \mathbb{R}^{d \times M}$ , where  $M$  is experimentally set to 64 in this paper.
- 

300 As discussed above, a dictionary is adaptively constructed from each pair of multi-focus source  
 301 images. The adaptability will improve the representation ability of each constructed dictionary, which  
 302 will be validated in the subsequent experimental part. In addition, the number of atoms in the dictionary  
 303 gets reduced by using PCA. This also speeds up computation of the proposed fusion method.

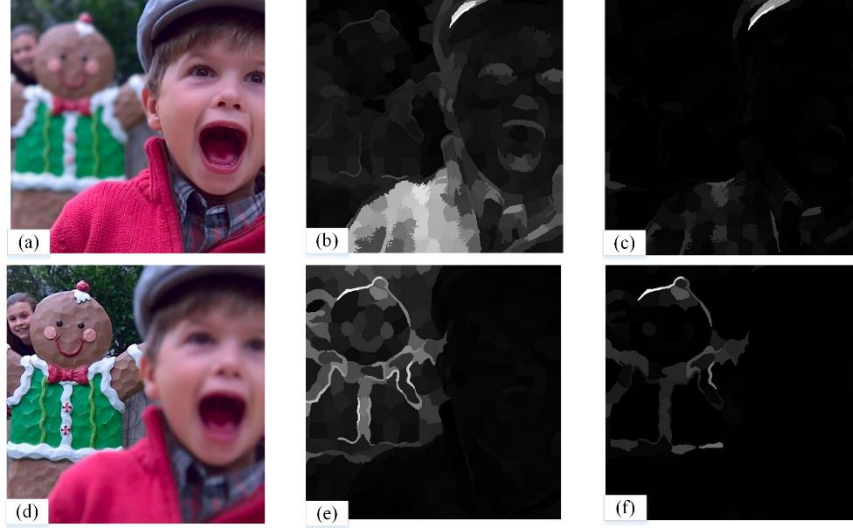
### 304 3.3. Focus decision map computation and refinement

305 Given the constructed dictionary  $\mathbf{D}$  and a set of feature matrices  $\{\mathbf{X}_{A,k} | k = 1, 2, \dots, K\}$  extracted  
 306 from a source image  $I_A$ , a representation coefficient matrix  $\mathbf{Z}_A$  and a sparse error matrix  $\mathbf{E}_A$  are  
 307 obtained by solving Eq. (3). Similarly, a representation coefficient matrix  $\mathbf{Z}_B$  and a sparse error matrix  
 308  $\mathbf{E}_B$  for the source image  $I_B$  are obtained.

309 As shown in Fig. 5, the super-pixels in the focused regions normally have larger representation  
 310 coefficient magnitudes as well as sparse errors, especially larger representation coefficient magnitudes,



311 than those super-pixels in the de-focused regions. Therefore, the focused and de-focused regions in a  
 312 multi-focus image can be determined by jointly using the representation coefficients and sparse errors.



313 **Fig. 5.** Illustration of the ULRR results on a pair of multi-focus source images. (a), (d) Source images with focus on the front and  
 314 the back, respectively; (b), (e) Representation coefficients obtained by ULRR for (a) and (d), respectively; (c), (f) Sparse errors  
 315 obtained by ULRR for (a) and (d), respectively. For better displaying, each super-pixel in the source image is replaced by the  $l_2$ -  
 316 norm of its corresponding column vector in the representation coefficient matrix and spare error matrix.  
 317

318 For that, a measure of focus (MOF) for the  $i$ -th super-pixel  $sp_{A/B,i}$  is first defined by:

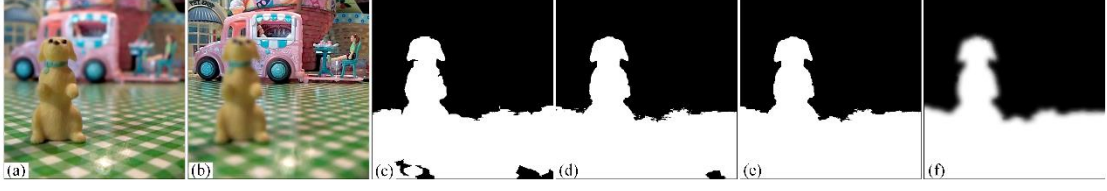
$$319 \quad MOF_{A/B,i} = \eta \|\mathbf{z}_{A/B,i}\|_2 + (1-\eta) \|\mathbf{e}_{A/B,i}\|_2, \quad (8)$$

320 where  $\mathbf{z}_{A/B,i}$  and  $\mathbf{e}_{A/B,i}$  are the  $i$ -th column of  $\mathbf{Z}_{A/B}$  and  $\mathbf{E}_{A/B}$ , respectively.  $\|\cdot\|_2$  denotes the  $l_2$ -  
 321 norm of a vector.  $\eta$  is experimentally set to 0.95 in this paper.

322 Then an initial focus decision map  $\Upsilon$  of the same size as the source images is defined and its each  
 323 element  $\Upsilon(x, y)$  is computed by

$$324 \quad \Upsilon(x, y) = \begin{cases} 1, & (x, y) \in sp_{A/B,i} \ \& \ MOF_{A,i} \geq MOF_{B,i} \\ 0, & \text{otherwise} \end{cases}. \quad (9)$$

325 Fig. 6(c) illustrates an initial focus decision map obtained by using Eq. (9). It can be obviously  
 326 found that most of the focused regions and de-focused regions can be accurately determined by using the  
 327 proposed MOF defined by Eq. (8). In addition, the boundaries between the focused regions and the de-  
 328 focused regions are naturally preserved, and few ‘‘jagged’’ artifacts are introduced because of the super-  
 329 pixel segmentation.



330

331 **Fig 6.** Illustration of the decision maps obtained by different post-processing. (a) and (b) A pair of multi-focus images with focus  
 332 on the front and the back, respectively; (c) Initial focus decision map obtained by using Eq. (9); (d) Decision map after image  
 333 matting; (e) Decision map after removing "holes"; (f) Final decision map after guided filtering.

334 In spite of that, some isolate regions still exist, as shown in Fig. 6(c). To address such problems,  
 335 some post-processing is further performed on the initial focus map, which includes: (1). Image matting  
 336 [37] to refine the boundary between the focused and de-focused region; (2). Removing holes [7] to erase  
 337 small isolate regions, i.e., a region smaller than an area threshold is reversed in the binary initial decision  
 338 map; (3). Guided filtering [17] to reduce the spatial artifacts between focused and de-focused regions.  
 339 After that, a refined focus decision map  $\Upsilon'$  is finally obtained.

340 Fig. 6 (d) indicates that the boundary accuracy between the focused and de-focused regions is  
 341 improved to some extent by using image matting. Some isolated regions are also eliminated. After  
 342 removing holes, some small isolate regions are further removed, as shown in Fig. 6(e). Finally, as shown  
 343 in Fig. 6(f), some gradual transitional regions are generated between the focused and de-focused regions,  
 344 which makes the boundaries look more natural.

### 345 3.4. Fusion

346 Given the refined focus decision map  $\Upsilon'$ , the fused image  $I_F$  can thus be obtained by using a  
 347 'weighted averaging' scheme, i.e.,

$$348 \quad I_F(x, y) = \Upsilon'(x, y)I_A(x, y) + (1 - \Upsilon'(x, y))I_B(x, y). \quad (10)$$

349 Here, the refined focus decision map  $\Upsilon'$  is used as the weighted map. In summary, the proposed fusion  
 350 method can be described in Algorithm 3.

---

**Algorithm 3: The proposed multi-focus image fusion method based on super-pixel clustering and ULRR**

---

- (1) Perform super-pixel segmentation on the source images as described in Subsection 3.1.1, and obtain two sets of super-pixels  $\{sp_{A,i} | i=1,2,\dots,N\}$  and  $\{sp_{B,i} | i=1,2,\dots,N\}$ ;
  - (2) Perform super-pixel clustering on  $\{sp_{A,i} | i=1,2,\dots,N\}$  and  $\{sp_{B,i} | i=1,2,\dots,N\}$  as described in Subsection 3.1.2, and
-

- 
- obtain two sets of super-pixel clusters  $\{C_{A,k} | k=1,2,\dots,K\}$  and  $\{C_{B,k} | k=1,2,\dots,K\}$ ;
- (3) Construct the feature matrix for each super-pixel as described in Subsection 3.1.3 and obtain two sets of feature matrices  $\{\mathbf{X}_{A,k} | k=1,2,\dots,K\}$  and  $\{\mathbf{X}_{B,k} | k=1,2,\dots,K\}$ ;
  - (4) Construct the dictionary by using Algorithm 2.
  - (5) Perform ULRR on  $\{\mathbf{X}_{A,k} | k=1,2,\dots,K\}$  and  $\{\mathbf{X}_{B,k} | k=1,2,\dots,K\}$  by using Eq. (3), and obtain representation coefficient and sparse error matrices  $\{\mathbf{Z}_A, \mathbf{E}_A\}$  and  $\{\mathbf{Z}_B, \mathbf{E}_B\}$  for source images  $I_A$  and  $I_B$ , respectively;
  - (6) Compute the focus decision map  $\Upsilon'$  by using the matrices  $\{\mathbf{Z}_A, \mathbf{Z}_B, \mathbf{E}_A, \mathbf{E}_B\}$  as described in Subsection 3.3;
  - (7) Construct the fused image  $I_F$  by using Eq. (10).
- 

## 351 4. Experiment results and analysis

352 Extensive experiments are conducted to verify the performance of the proposed multi-focus image  
 353 fusion algorithm, which are organized as: 1) the impacts of several important parameters on the proposed  
 354 method are investigated; 2) the validities of the constructed dictionaries and the proposed ULRR are  
 355 carried out; 3) comparisons against some state-of-the-art methods on two public databases; (4) extension  
 356 to the fusion of triple multi-focus images; (5) some discussions on the proposed fusion method.

### 357 4.1. Parameters setting

358 Here, we use ten pairs of multi-focus source images, which are manually generated from the images  
 359 in Fig. 7 [38], to investigate how the parameters, including the cluster number  $K$ , and the trade-off  
 360 parameters  $\alpha$  and  $\beta$  in Eq. (3), affect the proposed method. For that, image matting [37] is performed  
 361 on each image in Fig. 7 to extract the foreground object regions and the background regions, respectively.  
 362 Then the foreground regions and the background regions are blurred using a ‘Gaussian’ low-pass filter,  
 363 respectively, to obtain a pair of multi-focus images. Thus, ten pairs of manually generated multi-focus  
 364 source images are obtained. Finally, several fused images are obtained from each pair of these multi-  
 365 focus images by using the proposed method with different parameters. These fused images are compared  
 366 against their corresponding ‘ideal’ images in Fig. 7 using the metrics like mean square error ( $MSE$ ) and  
 367 difference coefficients ( $DC$ ) [1]. Smaller  $MSE$  and  $DC$  values imply higher fusion performance.

368 The experimental results in Table 1 demonstrate that the fusion performance achieves desirable  
 369 when the number of clusters  $K$  is set to 4 or 5. Similarly, the experimental results in Table 2 indicate that

370 the proposed fusion method achieves the best when  $\alpha$  is set to 0.05. The performance increases with  
 371 the decrease of  $\beta$  and keeps almost unchanged when  $\beta$  achieves 0.002. In the following experiments,  
 372 we set the parameters  $K$ ,  $\alpha$  and  $\beta$  to 4, 0.05, 0.002, respectively.



373  
 374 **Fig.7.** Original images that used to generate the multi-focus images.

375 **Table 1.** Fusion performance with different values of  $K$  on the 10 pairs of manually generated multi-focus images.

$K$	2	3	4	5	6	8
$MSE$	12.8483	12.5746	11.9715	11.9666	12.3555	13.5860
$DC$	0.0122	0.0118	0.0115	0.0114	0.0115	0.0120

376 **Table 2.** Fusion performance with different values of  $\alpha$  and  $\beta$  on the 10 pairs of manually generated multi-focus images.

	$\alpha$ with $\beta=0.002$					$\beta$ with $\alpha=0.05$				
	0.03	0.04	0.05	0.06	0.08	0.01	0.005	0.002	0.001	0.0001
$MSE$	13.5656	13.5622	11.9715	13.4647	13.6323	13.3931	12.5915	11.9715	11.9715	11.9715
$DC$	0.0119	0.0119	0.0115	0.0125	0.0125	0.0125	0.0118	0.0115	0.0115	0.0115

## 377 4.2. Validity of the constructed dictionary

378 In this subsection, we will investigate the impacts of different dictionaries on the fusion results to  
 379 verify the constructed dictionary in our proposed fusion method. To do so, six dictionaries are constructed  
 380 for fusion. The first three dictionaries ( $\mathbf{D}_{Ksvd}^{Global}$ ,  $\mathbf{D}_{Ksvd}^{Ave}$  and  $\mathbf{D}_{Ksvd}^{Blur}$ , for short) are constructed by using K-  
 381 SVD [39], and the other three dictionaries ( $\mathbf{D}_{PCA}^{Global}$ ,  $\mathbf{D}_{PCA}^{Ave}$  and  $\mathbf{D}_{PCA}^{Blur}$ , for short) are constructed by using  
 382 PCA. Especially,  $\mathbf{D}_{Ksvd}^{Global}$  and  $\mathbf{D}_{PCA}^{Global}$  are globally learned from a set of nature images with high spatial  
 383 resolutions and have 256 dictionary atoms.  $\mathbf{D}_{Ksvd}^{Ave}$  and  $\mathbf{D}_{PCA}^{Ave}$  are adaptively constructed from the  
 384 average image of each pair of source images and have 64 dictionary atoms. For that, the source images  
 385 are first averaged and the feature matrix for the average image is extracted afterwards. Then a dictionary  
 386 is learned from the feature matrix by using K-SVD or PCA.  $\mathbf{D}_{Ksvd}^{Blur}$  and  $\mathbf{D}_{PCA}^{Blur}$  (i.e., the employed

387 dictionary in our proposed fusion method) are adaptively constructed from the blurred average image of  
 388 each pair of source images and also have 64 dictionary atoms.



389 **Fig. 8.** Fusion results by using different dictionaries. (a1) and (b1) A pair of source images with focus on the front and the back,  
 390 respectively; (c1) ~ (h1) Initial focus decision maps for (a1) and (b1) obtained by using  $\mathbf{D}_{Ksvd}^{Global}$ ,  $\mathbf{D}_{Ksvd}^{Ave}$ ,  $\mathbf{D}_{Ksvd}^{Blur}$ ,  $\mathbf{D}_{PCA}^{Global}$ ,  $\mathbf{D}_{PCA}^{Ave}$   
 391 and  $\mathbf{D}_{PCA}^{Blur}$ , respectively; (a2) ~ (h2) Another pair of source images and their initial focus decision maps obtained by using different  
 392 dictionaries.  
 393

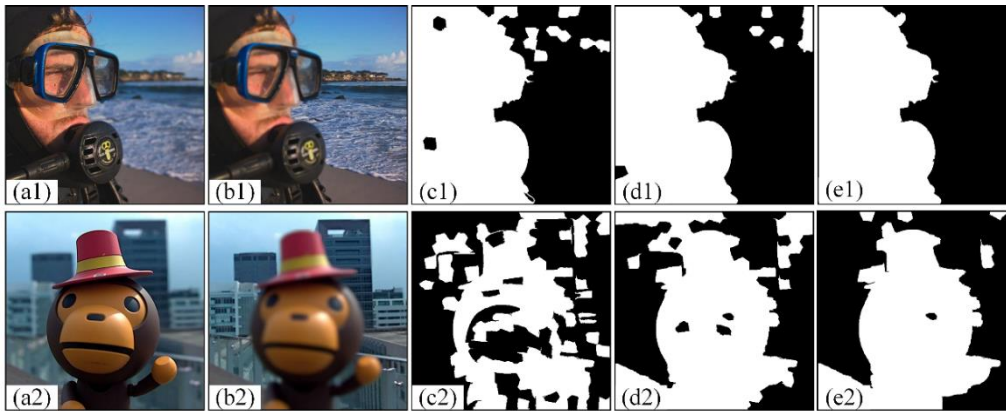
394 Fig. 8 illustrates the fusion results when using different dictionaries. It is clear that using the  
 395 dictionaries constructed by K-SVD could not lead the proposed fusion method to achieve desirable  
 396 results. Furthermore, it is also obvious that the proposed fusion method can achieve better results by  
 397 using  $\mathbf{D}_{PCA}^{Blur}$  than by using  $\mathbf{D}_{PCA}^{Ave}$  and  $\mathbf{D}_{PCA}^{Global}$ . This indicates that the representation coefficients and  
 398 the sparse errors, especially the representation coefficients, deduced from the proposed ULRR model can  
 399 better capture the “intrinsic” focus characteristics of a multi-focus image under a dictionary constructed  
 400 from blurred images than those constructed from clear images.

#### 401 4.3. Validity of the proposed ULRR model for multi-focus image fusion

402 In order to test the validity of the proposed ULRR model for multi-focus image fusion, three  
 403 versions (ULRR\_v1, ULRR\_v2, ULRR\_v3, for short, respectively) of our proposed fusion method just  
 404 with different LRR models are performed on two pairs of multi-focus source images, shown in Fig. 9. In  
 405 ULRR\_v1, the traditional LRR model [9] is employed, which is directly performed on the feature  
 406 matrices  $\mathbf{X}_{A/B}$  constructed from source image super-pixels rather than on the feature matrices  
 407  $\{\mathbf{X}_{A/B,k} \mid k = 1, 2, \dots, K\}$  constructed from the source super-pixels clusters. In ULRR\_v2, the proposed  
 408 ULRR model in Eq. (3) without the Laplacian regularization term is employed. In ULRR\_v3, i.e., the

409 proposed fusion method, the proposed ULRR model in Eq. (3) with the Laplacian regularization term is  
 410 employed.

411 As shown in Fig. 9 (c1) and Fig. 9 (c2), many isolated regions existed in the focus decision maps  
 412 obtained by using ULRR\_v1. Differently, the isolated regions are greatly reduced in the focus decision  
 413 maps obtained by using ULRR\_v2, as shown in Fig. 9 (d1) and Fig. 9 (d2). Especially, as shown in Fig.  
 414 9 (e1) and Fig. 9 (e2), the isolated regions are significantly reduced in the focus decision maps when  
 415 using ULRR\_v3. This indicates that performing ULRR on the super-pixel clusters can better capture the  
 416 “intrinsic” focus information of different regions in a multi-focus image than directly performing LRR  
 417 on super-pixels. This owes to the consideration of super-pixel clusters in the proposed method via ULRR,  
 418 especially the spatial consistency among the super-pixels within the same cluster via the Laplacian  
 419 regularization term in ULRR.



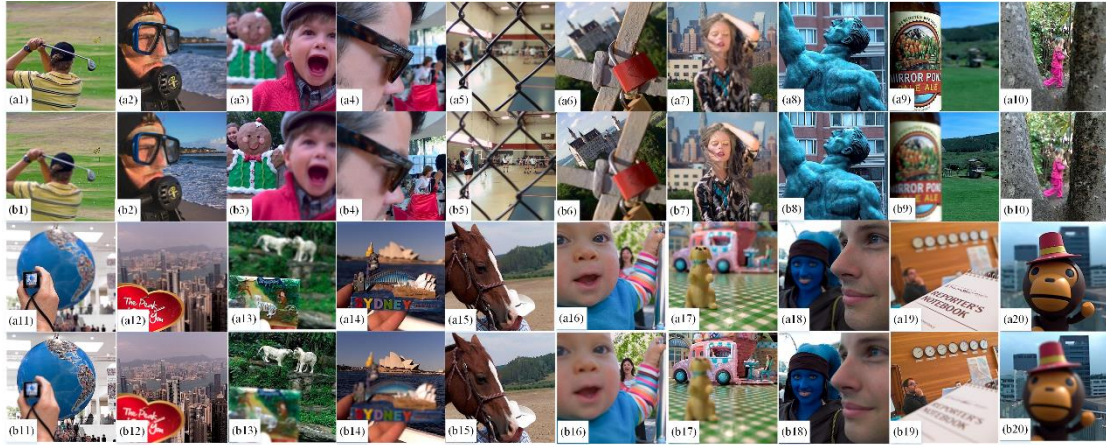
420  
 421 **Fig. 9.** Illustration of the validity of the proposed ULRR model. (a1) and (b1) A pair of source images with focus on the front and  
 422 the back, respectively; (c1) ~ (e1) Initial focus decisions maps for (a1) and (b1) obtained by ULRR\_v1, ULRR\_v2 and ULRR\_v3,  
 423 respectively; (a2) ~ (e2) Another pair of source images and their initial focus decision maps obtained by using different models.

#### 424 4.4. Comparisons with traditional fusion methods

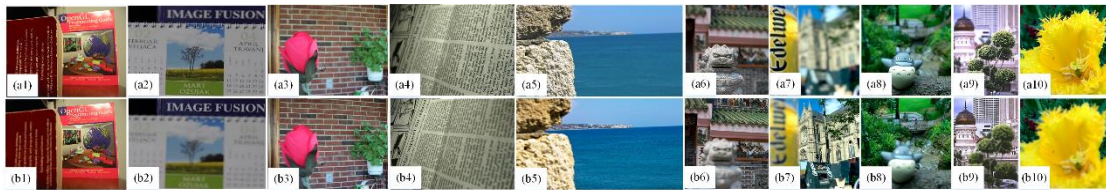
425 Here, we compare our method (ULRR, for short) with another 6 traditional state-of-the-art methods,  
 426 including GFF [17], IM [16], SPixel [15], LR\_RSR [7], SRCF [4], and DL\_LRR [5]. The public Lytro  
 427 Dataset in [4] including 20 pairs of multi-focus images and a smaller dataset (SMD, for short) including



428 10 pairs of multi-focus images that are collected from different kinds of literature are employed to test  
 429 different fusion methods, which are illustrated in Fig. 10 and Fig. 11, respectively.



430  
 431 **Fig. 10.** Lytro Dataset. (a1) ~ (a10) The first 10 input images with the focus on the front part; (b1) ~ (b10) The corresponding input  
 432 images with the focus on the back part; (a11) ~ (a20) The remaining 10 input images with the focus on the front part; (b11) ~ (b20)  
 433 The corresponding input images with the focus on the back part.

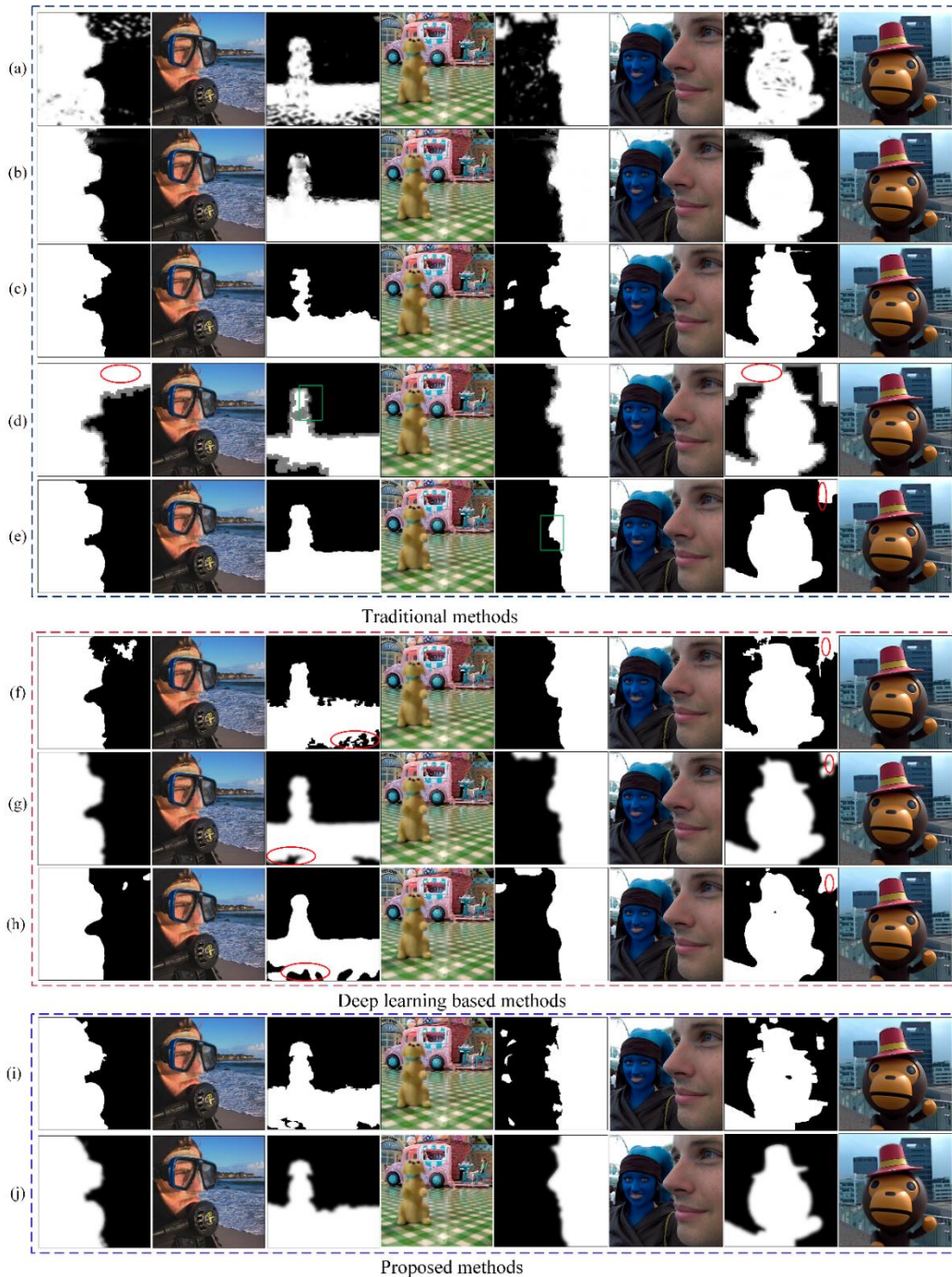


434  
 435 **Fig. 11.** SMD Dataset. (a1) ~ (a10) The 10 input images with focus on the left (or front) part; (b1) ~ (b10) The corresponding 10  
 436 input images with focus on the right (or back) part.

437 In order to quantitatively compare the fusion performance of different methods, six metrics are used,  
 438 including mutual information based fusion metric  $FMI$  [40], universal image quality index based metric  
 439  $Q_{uiqi}$  [41], quaternion based color image fusion quality metrics  $Q_{ssim}$  [42] and  $Q_4$  [43], and phase  
 440 consistency based metrics  $Q_{PC}$  [44] and  $ZNCC\_PC$  [45]. The former four metrics measure how  
 441 well the original information, such as entropy information and structures, from the source images, have  
 442 been preserved in the fused images. The last two metrics may evaluate different methods in spatial  
 443 consistency to some extent. Larger values of these metrics are more desirable for a fusion method.

444 Fig. 12 and Fig. 13 illustrate some fusion results on Lytro Dataset and SMD Dataset obtained by

445 some of the fusion methods mentioned above<sup>3</sup>. In addition, the initial decision maps obtained by using  
 446 our proposed ULRR method without any post-processing steps are also illustrated in Fig. 12 (i) and Fig.  
 447 13(i). Table 3 and Table 4 report the average quantitative results on the two datasets, respectively.



448  
 449 **Fig. 12.** Some fusion results on Lytro Dataset. (a) GFF; (b) IM; (c) SPixel; (d) LR\_RSR; (e) SRCF; (f) p\_CNN; (g) CNN; (h)  
 450 EN\_CNN; (i) ULRR without post-processing steps; (j) ULRR. The decision maps in (h) and (i) are initial ones without using any

<sup>3</sup>Given a pair of multi-focus images, DL\_LRR directly outputs the finally fused images without using focus decision maps. Therefore, the visual results obtained by DL\_LRR are not provided in Fig. 12 and Fig. 13.



451 post-processing steps. The remaining decision maps are the final ones after the post-processing steps.



452  
 453 **Fig. 13.** Some fusion results on SMD Dataset. (a) GFF; (b) IM; (c) SPixel; (d) LR\_RSR; (e) SRCF; (f) p\_CNN; (g) CNN; (h)  
 454 EN\_CNN; (i) ULRR without post-processing steps; (j) ULRR. The decision maps in (h) and (i) are initial ones without using any  
 455 post-processing steps. The remaining decision maps are the final ones after the post-processing steps.

456 From Fig. 12 and Fig. 13, the following facts can be easily observed. Plenty of spatial artifacts  
 457 appear on the fused images obtained by GFF. IM usually introduces some spatial artifacts on the  
 458 boundaries between focused and de-focused regions. In most cases, SPixel cannot accurately determine  
 459 the boundaries between the focused and de-focused regions, although few spatial artifacts are introduced  
 460 in their decision maps. LR\_RSR and SRCF introduce fewer spatial artifacts in their decision maps. But  
 461 some regions, especially those smooth regions as observed in the red elliptical regions in Fig. 12 (d) and  
 462 (e), are mistakenly labeled as out-of-focus (or in-focus) by the two methods. As discussed in the earlier

463 part of Section 1, this may be due to the fact that only the intensity values of the input images are  
 464 employed as the features in LR\_RSR and SRCF. Moreover, some “jagged” artifacts exist in the  
 465 boundaries between the focused and de-focused regions (e.g., the green rectangle regions) of Fig. 12(d),  
 466 Fig. 12(e), Fig. 13(d) and Fig. 13(e).

467 Differently, as shown in Fig. 12(j) and Fig. 13(j), almost no isolate regions exist in the decision  
 468 maps obtained by ULRR. This indicates that the focused and de-focused regions in the source images  
 469 are better determined by ULRR than the other methods. Accordingly, fewer spatial artifacts are involved  
 470 in the fused images by using our proposed fusion method. Moreover, no “jagged” artifacts exist in the  
 471 decision maps obtained by ULRR. The boundaries between the focused and de-focused regions in Fig.  
 472 12(j) and Fig. 13(j) look closer to the boundaries of the objects in the source images. The comparisons  
 473 between Fig. 12 (i) and (j), Fig. 13(i) and (j), indicate that the post-processing steps can partially benefit  
 474 to the improvements of our proposed method.

475 **Table 3.** Averaging performance of different traditional fusion methods on Lytro Dataset.

Methods	$FMI$	$Q_{uiqi}$	$Q_{ssim}$	$Q_4$	$Q_{pc}$	$ZNCC\_PC$
GFF	1.4153	0.9092	0.8838	0.9754	0.6817	0.9273
IM	1.4102	0.8932	0.8759	0.9736	0.6225	0.9078
SPixel	1.4194	0.9048	0.8818	0.9746	0.6736	0.9248
LR_RSR	1.4198	0.9072	0.8838	0.9751	0.6763	0.9270
SRCF	1.4195	0.9105	0.8818	0.9747	0.6790	0.9276
DL_LRR	1.4199	0.9015	0.8829	0.9748	0.6347	0.9077
ULRR	1.4197	0.9117	0.8844	0.9756	0.6819	0.9317

476 **Table 4.** Averaging performance of different traditional fusion methods on SMD Dataset.

Methods	$FMI$	$Q_{uiqi}$	$Q_{ssim}$	$Q_4$	$Q_{pc}$	$ZNCC\_PC$
GFF	1.1831	0.9053	0.7824	0.9103	0.6213	0.9189
IM	1.1931	0.9024	0.7777	0.9053	0.5948	0.9243
SPixel	1.1871	0.9070	0.7779	0.9058	0.6213	0.9287
LR_RSR	1.1852	0.9079	0.7802	0.9076	0.6297	0.9288
SRCF	1.1873	0.9141	0.7783	0.9057	0.6234	0.9283
DL_LRR	1.1860	0.8850	0.7767	0.9055	0.5223	0.8743
ULRR	1.1891	0.9119	0.7815	0.9086	0.6314	0.9356

477 The quantitative results in Table 3 and Table 4 are in line with the visual results above, which  
478 demonstrates that ULRR significantly outperforms the other methods in terms of  $Q_{PC}$  and  
479  $ZNCC_{PC}$ . This indicates that our proposed fusion method performs the best in spatial consistency,  
480 compared to those methods mentioned here, and fewer spatial artifacts have been introduced to the fused  
481 images by using ULRR than by the other methods. Table 3 also demonstrates that ULRR performs the  
482 best on Lytro Dataset in terms of  $Q_{uiqi}$ ,  $Q_{ssim}$  and  $Q_4$ . Table 4 demonstrates that ULRR always  
483 achieves the top two performance on SMD Dataset in terms of FMI,  $Q_{uiqi}$ ,  $Q_{ssim}$  and  $Q_4$ . This  
484 indicates that, in addition to spatial consistency, our proposed fusion method can also achieve better  
485 performance in information extraction than the other methods in most cases.

#### 486 4.5 Comparisons with deep learning based methods

487 In addition to those traditional methods, three deep learning (DL) based fusion methods, including  
488 CNN [23], p\_CNN [24] and EN\_CNN [28], are compared with our proposed fusion method. Some visual  
489 results on Lytro Dataset and SMD Dataset are also illustrated and provided in Fig. 12 and Fig. 13,  
490 respectively. The quality results on the two datasets are provided in Table 5 and Table 6, respectively.

491 As shown in the first columns of Fig. 12 and Fig. 13, these DL based methods can generally achieve  
492 desirable fusion results. Especially, EN\_CNN can accurately determine the focused and de-focused  
493 regions without using any post-processing steps, thanks to the strong abilities of CNNs for image  
494 representation and feature extraction. However, as shown in the red elliptical regions of Fig. 12 and Fig.  
495 13, some smooth regions are also mistakenly determined to be in-focus (or out-of-focus) by these DL  
496 based methods. For some regions with abundant textures (e.g., the blue elliptical regions in Fig. 13),  
497 these DL based fusion methods could not determine the focused and de-focused regions uniformly.  
498 Differently, our proposed fusion method can completely determine the focused and de-focused regions

499 in most cases, as illustrated in Fig. 12 and Fig. 13.

500 The experimental results in Table 5 and Table 6 indicate that ULRR achieves comparable  
501 performance with these DL based fusion methods in terms of FMI ,  $Q_{uiqi}$  ,  $Q_{ssim}$  and  $Q_4$  . In terms of  
502  $Q_{PC}$  and  $ZNCC_{PC}$  , ULRR performs competitively with CNN and outperforms p\_CNN and  
503 ES\_CNN by a clear margin. This indicates that, even in information extraction, our proposed method  
504 still performs competitively with these DL based ones, but in spatial consistency, our proposed method  
505 is clearly superior to most of these DL based ones. This further verifies the validity of our proposed  
506 fusion method in the reduction of spatial artifacts.

507 **Table 5.** Averaging performance of different deep learning based fusion methods on Lytro Dataset.

Methods	FMI	$Q_{uiqi}$	$Q_{ssim}$	$Q_4$	$Q_{PC}$	$ZNCC_{PC}$
CNN	1.4195	0.9111	0.8839	0.9753	0.6851	0.9307
p_CNN	1.4208	0.9091	0.8823	0.9747	0.6749	0.9260
ES_CNN	1.4207	0.9087	0.8816	0.9746	0.6628	0.9247
ULRR	1.4197	0.9117	0.8844	0.9756	0.6819	0.9317

508 **Table 6.** Averaging performance of different deep learning based fusion methods on SMD Dataset.

Methods	FMI	$Q_{uiqi}$	$Q_{ssim}$	$Q_4$	$Q_{PC}$	$ZNCC_{PC}$
CNN	1.1879	0.9135	0.7815	0.9082	0.6357	0.9318
p_CNN	1.1894	0.8970	0.7808	0.9067	0.5960	0.9102
ES_CNN	1.1902	0.8949	0.7773	0.9053	0.5892	0.9152
ULRR	1.1891	0.9119	0.7815	0.9086	0.6314	0.9356

#### 509 4.6. Fusion of more than two multi-focus images

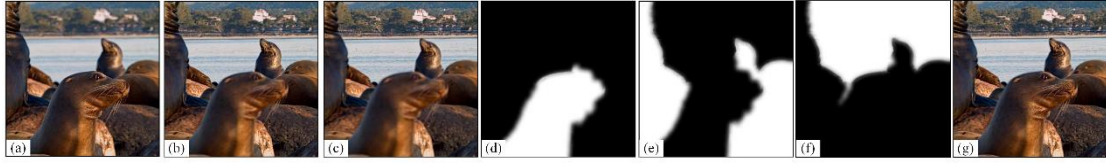
510 The proposed fusion method can be easily extend to fuse more than two multi-focus images.  
511 Suppose that there are total  $S$  images  $I_s (s=1,2,\dots,S)$  to be fused. For that, similar to Eq. (9) in the  
512 Subsection 3.3, the initial decision map  $\Upsilon_s$  for the  $s$ -th source image is determined by

$$513 \Upsilon_s(x, y) = \begin{cases} 1, & (x, y) \in sp_{s,i} \ \& \ s = \max_j MOF_{j,i} \\ 0, & otherwise \end{cases}, \quad (11)$$

514 where  $MOF_{j,i}$  denotes the measure of focus for the  $i$ -th super-pixel  $sp_{j,i}$  in the  $j$ -th image. After some  
515 post-processing, the final decision map  $\Upsilon'_s$  is obtained and the fused image  $I_F$  is obtained by

516

$$I_F(x, y) = \sum_{s=1}^S Y'_s(x, y) I_s(x, y) . \quad (12)$$

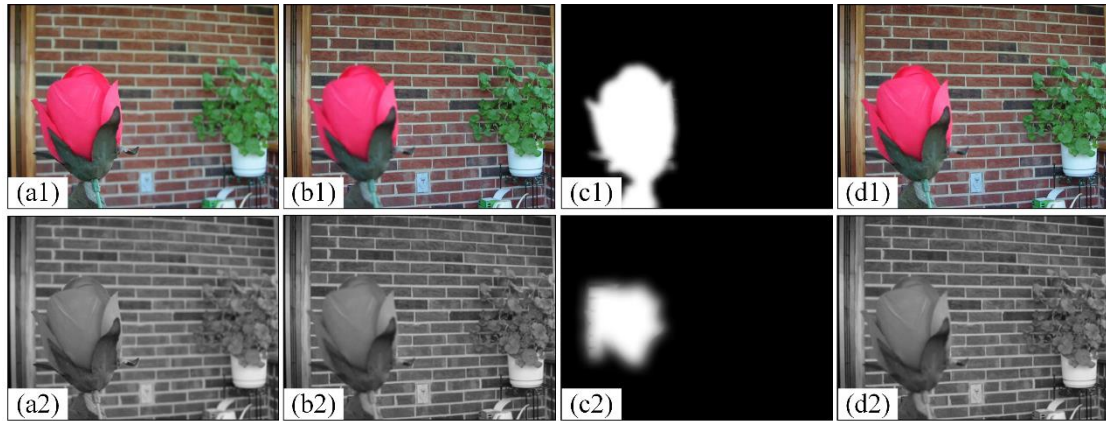


517

518 **Fig.14.** Illustration of the triple multi-focus image fusion. (a),(b) and (c) A set of triple multi-focus source images with the focus on  
519 the front, middle and back, respectively; (d), (e) and (f) The decision maps for (a),(b) and (c), respectively; (g) Fused image.

520 Fig. 14 illustrates the fusion of a set of three multi-focus images, which are also provided in the  
521 Lytro Dataset [4]. Similarly, the fusion results demonstrate that all of the focused regions within the input  
522 images can be effectively combined into the fused images without the introduction of obvious spatial  
523 artifacts.

524 **4.7 Fusion of gray-scale multi-focus images**



525

526 **Fig. 15.** Illustration of the fusion results on a pair of color multi-focus images and their gray-scale versions by using our proposed  
527 method. (a1) and (b1) A pair of color multi-focus images with the focus on the left and right parts, respectively; (c1) Focus decision  
528 map for (a1) and (b1); (d1) Fusion result on (a1) and (b1); (a2) and (b2) Gray-scale versions of (a1) and (b1), respectively; (c2)  
529 Focus decision map for (a2) and (b2); (d2) Fusion result on (a2) and (b2).

530 We have also tried to apply our proposed fusion method to fuse gray-scale multi-focus images. Fig.  
531 15 illustrates the fusion results of a pair of color multi-focus images and their gray-scale versions by  
532 using our proposed fusion method. As shown in Fig. 15, we find that the fusion results on gray-scale  
533 multi-focus images are not satisfactory although the fusion results on the color multi-focus images are  
534 desirable. This may owe to the feature extraction, super-pixel segmentation and clustering modules in  
535 our proposed method, which heavily depend on the color information of source images.

536 **4.8 Discussion**

537 In this subsection, we will discuss two issues. One is about the computational complexities of  
 538 different methods while the other is about the superiorities of our proposed methods over current DL  
 539 based fusion methods.

540 With respect to the first issue, Table 7 provides the average computational time  $T$  of different  
 541 methods on Lytro Dataset. Here, all of the traditional methods and some of the DL based methods (e.g.,  
 542 CNN and p\_CNN) are tested in Matlab R2013b environment on a PC with an Intel i7 CPU and 32 GB  
 543 of RAM. ES\_CNN is tested on an NVIDIA 1080Ti GPU with 11 G memory.

544 Table 7. Averaging computational time of different methods on Lytro Dataset.

Methods	GF	IM	SPixel	LR_RSR	SRCF	DL_LRR	CNN	p_CNN	ES_CNN	ULRR
$T(s)$	0.57	2.49	57.49	25.50	13.27	4443.15	124.38	439.41	366.37	71.00

545 As shown in Table 7, ULRR has higher computational complexity than most of the other traditional  
 546 fusion methods, such as LR\_RSR and SRCF. This may owe to the part of feature extraction for each  
 547 super-pixel in our proposed method. As shown in Table 8, for a pair of multi-focus images with size of  
 548  $520 \times 520$ , the running time of our proposed method is about 73 seconds, among which feature extraction  
 549 for each super-pixel takes about 82% of the total time. Despite that, Table 7 also demonstrates that ULRR  
 550 has higher computational efficiency than those DL based fusion methods. Especially, the average  
 551 computational time of ULRR is about half that of CNN, although the two methods perform competitively  
 552 in spatial consistency as well as in information extraction.

553 **Table 8.** Computational time of different modules in our proposed method for a pair of multi-focus images of size  $520 \times 520$ .

Module	Super-pixel segmentation	Feature extraction	Super-pixel clustering	Dictionary construction	ULRR decomposition	Post- processing	Total
Time(s)	0.47	60.56	1.51	0.12	3.32	7.47	73.45
Percentage(%)	0.64	82.45	2.06	0.16	4.52	10.17	100

554 Regarding the second issue, as discussed in Subsection 4.5, our proposed fusion method performs  
 555 competitively and even better than some DL based ones. A part of the reason might be that the training  
 556 data for these DL based fusion methods, which are manually generated by just performing different  
 557 Gaussian filters on the original images, could not fully simulate the multi-focus characteristics of an

558 image. Also, similar to most existing SR band LRR based fusion methods, these DL fusion methods  
559 usually perform fusion on image patches of fixed shapes independently, thus ignoring the spatial  
560 consistency among adjacent patches and degrading the fusion performance to some extent. Differently,  
561 the spatial consistency among adjacent super-pixels and the object area consistency among the super-  
562 pixels within an object are jointly considered in our proposed fusion method. Moreover, post-processing  
563 steps also contribute to the improvement of our fusion performance.

## 564 **5. Conclusion**

565 In this paper, we present a novel multi-focus image fusion algorithm based on super-pixel clustering  
566 and a unified low-rank representation (ULRR) model. Owing to the use of super-pixels of irregular sizes,  
567 the “jagged” artifacts between the focused and de-focused regions, which arises from the patch based  
568 fusion methods, can be effectively eliminated. Thanks to the use of multiple types of features, the focus  
569 information for the smooth regions as well as those regions with rich details can be well determined by  
570 the proposed fusion method. By further using super-pixel clustering and considering the object  
571 consistency among the super-pixels within the same cluster via the proposed ULRR model, the spatial  
572 artifacts in the fused images are greatly reduced and even eliminated by the proposed fusion method.  
573 Experimental results demonstrate that the proposed fusion method outperforms some state-of-the-arts,  
574 even including some deep learning based methods, in terms of visual and quantitative evaluations,  
575 especially in the reduction of spatial artifacts or in spatial consistency.

576 Finally, it should be also noted that the high fusion performance of our proposed method is at the  
577 cost of high computational complexity. Moreover, the proposed method works well for the color multi-  
578 focus images but it does not perform well for the gray-scale multi-focus images. In future, we will explore  
579 how to reduce the computational complexity of our proposed method and how to modify our proposed

580 method to the fusion of gray-scale multi-focus images.

### 581 Acknowledgements

582 This work is supported by the National Natural Science Foundation of China under Grant No.  
583 61773301.

### 584 Appendix A

585 (1) Update  $\mathbf{J}_k$

$$\begin{aligned}
 \mathbf{J}_k^{i+1} &= \arg \min_{\mathbf{J}_k} \|\mathbf{J}_k\|_* + \langle \mathbf{Y}_{2,k}^i, \mathbf{Z}_k^i - \mathbf{J}_k \rangle + \frac{\mu^i}{2} \|\mathbf{Z}_k^i - \mathbf{J}_k\|_F^2 \\
 &= \arg \min_{\mathbf{J}_k} \frac{1}{\mu^i} \|\mathbf{J}_k\|_* + \frac{1}{2} \left\| \mathbf{J}_k - \left( \mathbf{Z}_k^i + \frac{\mathbf{Y}_{2,k}^i}{\mu^i} \right) \right\|_F^2.
 \end{aligned} \tag{A1}$$

587 The sub-optimization has the following closed-form solution:

$$\mathbf{J}_k^{i+1} = \text{SVT}_{\frac{1}{\mu^i}} \left( \mathbf{Z}_k^i + \frac{1}{\mu^i} \mathbf{Y}_{2,k}^i \right), \tag{A2}$$

589 where  $\text{SVT}_{\delta}(\boldsymbol{\varphi})$  denotes the Singular Value Thresholding (SVT) operation on the matrix  $\boldsymbol{\varphi}$  with the  
590 threshold  $\delta$ .

591 (2) Update  $\mathbf{H}$

$$\begin{aligned}
 \mathbf{H}^{i+1} &= \arg \min_{\mathbf{H}} \beta \text{tr}(\mathbf{H}\mathbf{L}\mathbf{H}^T) + \langle \mathbf{Y}_3^i, \mathbf{Z}^i - \mathbf{H} \rangle + \frac{\mu^i}{2} \|\mathbf{Z}^i - \mathbf{H}\|_F^2 \\
 &= \arg \min_{\mathbf{H}} \beta \text{tr}(\mathbf{H}\mathbf{L}\mathbf{H}^T) + \frac{\mu^i}{2} \left\| \mathbf{Z}^i - \mathbf{H} + \frac{1}{\mu^i} \mathbf{Y}_3^i \right\|_F^2,
 \end{aligned} \tag{A3}$$

593 The optimization problem in Eq. (A3) has the flowing closed-form solution:

$$\mathbf{H}^{i+1} = \left( \mathbf{Z}^i + \frac{\mathbf{Y}_3^i}{\mu^i} \right) \left( 2 \times \frac{\beta}{\mu^i} \times \mathbf{L} + \mathbf{I} \right)^{-1}. \tag{A4}$$

595 (3) Update  $\mathbf{Z}$  ( $\mathbf{Z}_k$ )



$$\begin{aligned}
\mathbf{Z}^{i+1} &= \arg \min_{\mathbf{Z}} \sum_{k=1}^K \left( \langle \mathbf{Y}_{1,k}^i, \mathbf{X}_k - \mathbf{DZ}_k - \mathbf{E}_k^i \rangle + \langle \mathbf{Y}_{2,k}^i, \mathbf{Z}_k - \mathbf{J}_k^{i+1} \rangle + \frac{\mu^i}{2} \|\mathbf{X}_k - \mathbf{DZ}_k - \mathbf{E}_k^i\|_F^2 + \frac{\mu^i}{2} \|\mathbf{Z}_k - \mathbf{J}_k^{i+1}\|_F^2 \right) \\
&\quad + \langle \mathbf{Y}_3^i, \mathbf{Z} - \mathbf{H}^{i+1} \rangle + \frac{\mu^i}{2} \|\mathbf{Z} - \mathbf{H}^{i+1}\|_F^2 \\
596 \quad &= \arg \min_{\mathbf{Z}} \langle \mathbf{Y}_1^i, \mathbf{X} - \mathbf{DZ} - \mathbf{E}^i \rangle + \langle \mathbf{Y}_2^i, \mathbf{Z} - \mathbf{J}^{i+1} \rangle + \langle \mathbf{Y}_3^i, \mathbf{Z} - \mathbf{H}^{i+1} \rangle \\
&\quad + \frac{\mu^i}{2} \|\mathbf{X} - \mathbf{DZ} - \mathbf{E}^i\|_F^2 + \frac{\mu^i}{2} \|\mathbf{Z} - \mathbf{J}^{i+1}\|_F^2 + \frac{\mu^i}{2} \|\mathbf{Z} - \mathbf{H}^{i+1}\|_F^2 \\
&= \arg \min_{\mathbf{Z}} \frac{\mu^i}{2} \left\| \mathbf{X} - \mathbf{DZ} - \mathbf{E}^i + \frac{\mathbf{Y}_1^i}{\mu^i} \right\|_F^2 + \frac{\mu^i}{2} \left\| \mathbf{Z} - \mathbf{J}^{i+1} + \frac{\mathbf{Y}_2^i}{\mu^i} \right\|_F^2 + \frac{\mu^i}{2} \left\| \mathbf{Z} - \mathbf{H}^{i+1} + \frac{\mathbf{Y}_3^i}{\mu^i} \right\|_F^2
\end{aligned}$$

597 (A5)

598 where  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$ ,  $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_K]$ ,  $\mathbf{Y}_1^i = [\mathbf{Y}_{1,1}^i, \mathbf{Y}_{1,2}^i, \dots, \mathbf{Y}_{1,K}^i]$ ,  $\mathbf{Y}_2^i = [\mathbf{Y}_{2,1}^i, \mathbf{Y}_{2,2}^i, \dots, \mathbf{Y}_{2,K}^i]$ ,

599  $\mathbf{J}^{i+1} = [\mathbf{J}_1^{i+1}, \mathbf{J}_2^{i+1}, \dots, \mathbf{J}_K^{i+1}]$  and  $\mathbf{E}^i = [\mathbf{E}_1^i, \mathbf{E}_2^i, \dots, \mathbf{E}_K^i]$ . Eq. (A5) is a convex function and has the

600 following optimal solutions:

$$601 \quad \mathbf{Z}^{i+1} = (\mathbf{D}^T \mathbf{D} + 2\mathbf{I})^{-1} \times \left[ \mathbf{D}^T \left( \mathbf{X} - \mathbf{E}^i + \frac{\mathbf{Y}_1^i}{\mu^i} \right) + \mathbf{J}^{i+1} + \mathbf{H}^{i+1} - \frac{\mathbf{Y}_2^i}{\mu^i} - \frac{\mathbf{Y}_3^i}{\mu^i} \right], \quad (\text{A6})$$

602 (4) Update  $\mathbf{E}$

$$\begin{aligned}
603 \quad \mathbf{E}^{i+1} &= \arg \min_{\mathbf{E}} \alpha \|\mathbf{E}\|_{2,1} + \sum_{k=1}^K \left( \langle \mathbf{Y}_{1,k}^i, \mathbf{X}_k - \mathbf{DZ}_k^{i+1} - \mathbf{E}_k \rangle + \frac{\mu^i}{2} \|\mathbf{X}_k - \mathbf{DZ}_k^{i+1} - \mathbf{E}_k\|_F^2 \right) \\
&= \arg \min_{\mathbf{E}} \alpha \|\mathbf{E}\|_{2,1} + \frac{\mu^i}{2} \left\| \mathbf{X} - \mathbf{DZ}^{i+1} - \mathbf{E} + \frac{\mathbf{Y}_1^i}{\mu^i} \right\|_F^2, \quad (\text{A7})
\end{aligned}$$

604 where  $\mathbf{E} = [\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_K]$ . This sub-optimization problem has the following closed-form solution [35]:

$$605 \quad \mathbf{E}^{i+1}(:, i) = \begin{cases} \left( \frac{\|\mathbf{G}(:, i)\|_2 - \frac{\alpha}{\mu^i}}{\|\mathbf{G}(:, i)\|_2} \right) \mathbf{G}(:, i), & \text{if } \|\mathbf{G}(:, i)\|_2 \geq \frac{\alpha}{\mu^i}, \\ 0, & \text{otherwise} \end{cases} \quad (\text{A8})$$

606 where  $\mathbf{G} = \mathbf{X} - \mathbf{DZ}^{i+1} + \frac{\mathbf{Y}_1^i}{\mu^i}$ .  $\mathbf{E}(:, i)$  and  $\mathbf{G}(:, i)$  denote the  $i$ -th column of  $\mathbf{E}$  and  $\mathbf{G}$ , respectively.

## 607 References

608 [1]. Q. Zhang, M. D. Levine, Robust multi-focus image fusion using multi-task sparse representation and spatial context, IEEE

609 Transactions on Image Processing 25 (5) (2016) 2045-2058.

610 [2]. S. Li, X. Kang, L. Fang, J. Hu, H. Yin, Pixel-level image fusion: A survey of the state of the art, Information Fusion 33

611 (2017) 100-112.

- 612 [3]. B. Yang, S. Li, Multifocus image fusion and restoration with sparse representation, *IEEE Transactions on Instrumentation*  
613 *and Measurement* 59 (4) (2010) 884-892.
- 614 [4]. M. Nejati, S. Samavi, S. Shirani, Multi-focus image fusion using dictionary-based sparse representation, *Information Fusion*  
615 25 (2015) 72-84.
- 616 [5]. H. Li, X. Wu, Multi-focus image fusion using dictionary learning and low-rank representation, *International Conference on*  
617 *Image and Graphics* 10666 (2017) 675-686.
- 618 [6]. H. Li, X. He, D. Tao, Y. Tang, R. Wang, Joint medical image fusion, denoising and enhancement via discriminative low-  
619 rank sparse dictionaries learning, *Pattern Recognition* 79 (2018) 130-146.
- 620 [7]. Q. Zhang, T. Shi, F. Wang, R. S. Blum, J. Han, Robust sparse representation based multi-focus image fusion with dictionary  
621 construction and local spatial consistency, *Pattern Recognition* 83 (2018) 299-313.
- 622 [8]. B. Cheng, L. Jin, G. Li, General fusion method for infrared and visual images via latent low-rank representation and local  
623 non-subsampled shearlet transform, *Infrared Physics & Technology* 92 (2018) 68-77.
- 624 [9]. G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, Y. Ma, Robust recovery of subspace structures by low-rank representation, *IEEE*  
625 *Transaction on Pattern Analysis and Machine Intelligence* 35 (1) (2013) 171-184.
- 626 [10]. F. Kou, Z. Li, C. Wen, W. Chen, Edge-preserving smoothing pyramid based multi-scale exposure fusion, *Journal of Visual*  
627 *Communication and Image Representation* 53 (2018) 235-244.
- 628 [11]. H. Zhao, Z. Shang, Y. Tang, B. Fang, Multi-focus image fusion based on the neighbor distance, *Pattern Recognition* 46 (3)  
629 (2013) 1002-1011.
- 630 [12]. W. Huang, Z. Jing, Evaluation of focus measures in multi-focus image fusion, *Pattern Recognition Letters* 28 (2007) 493-  
631 500.
- 632 [13]. S. Li, B. Yang, Multifocus image fusion using region segmentation and spatial frequency, *Image and Vision Computing* 26  
633 (2008) 971-979.
- 634 [14]. Y. Liu, J. Jin, Q. Wang, Y. Shen, X. Dong, Region level based multi-focus image fusion using quaternion wavelet and  
635 normalized cut, *Signal Processing*, 97 (2014) 9-30.
- 636 [15]. J. Duan, L. Chen, C. L. P. Chen, Multifocus image fusion using superpixel segmentation and superpixel-based mean filtering,  
637 *Applied Optics* 55 (36) (2016) 10352-10362.
- 638 [16]. S. LI, X. Kang, J. Hu, B. Yang, Image matting for fusion of multi-focus images in dynamic scenes, *Information Fusion* 14  
639 (2013) 147-162.
- 640 [17]. S. Li, X. Kang, J. Hu, Image fusion with guided filtering, *IEEE Transactions on Image Processing* 22 (7) (2013) 2864-2875.
- 641 [18]. Y. Chen, J. Guan, W. Cham, Robust multi-focus image fusion using edge model and multi-matting, *IEEE Transactions on*

- 642 Image Processing 27 (3) (2018) 1526-1541.
- 643 [19]. O. Bouzos, I. Andreadis, N. Mitianoudis, Conditional random field model for robust multi-focus image fusion, IEEE  
644 Transactions on Image Processing 28 (11) (2019) 5636-5648.
- 645 [20]. B. Meher, S. Agrawal, R. Panda, A. Abraham, A survey on region based image fusion methods, Information Fusion 48 (2019)  
646 119-132.
- 647 [21]. L. Chen, J. Li, C. L. Chen, Regional multifocus image fusion using sparse representation, Optics Express 21 (4) 2013 5182-  
648 5197.
- 649 [22]. Q. Zhang, Y. Liu, R.S. Blum, J. Han, D. Tao, Sparse representation based multi-sensor image fusion for multi-focus and  
650 multi-modality images: a review, Information Fusion 40 (2018) 57-75.
- 651 [23]. Y. Liu, X. Chen, H. Peng, Z. Wang, Multi-focus image fusion with a deep convolutional neural network, Information Fusion  
652 36 (2017) 191-207.
- 653 [24]. H. Tang, B. Xiao, W. Li, G. Wang, Pixel convolutional neural networks for multi-focus image fusion, Information Science  
654 433 (2018) 125-141.
- 655 [25]. W. Zhao, D. Wang, H. Lu, Multi-focus image fusion with a natural enhancement via joint multi-level deeply supervised  
656 convolutional neural network, IEEE Transactions on Circuits and Systems for Video Technology 29 (4) (2019) 1102-1115.
- 657 [26]. H. T. Mustafa, F. Liu, J. Yang, Z. Khan, Q. Huang, Dense multi-focus fusion net: A deep unsupervised convolutional network  
658 for multi-focus image fusion, In: Internal Conference on Artificial Intelligence and Soft Computing, 2019, pp. 153-163.
- 659 [27]. H. Ma, J. Zhang, S. Liu, Q. Liao, Boundary aware multi-focus image fusion using deep neural network, In: IEEE  
660 International Conference on Multimedia and Expo, 2019, pp. 1150-1155.
- 661 [28]. M. Amin-Naji, A. Aghagolzadeh, M. Ezoji, Ensemble of CNN for multi-focus image fusion, Information Fusion 51 (2019)  
662 201-214.
- 663 [29]. M. Amin-Naji, A. Aghagolzadeh, M. Ezoji, CNNs hard voting for multi-focus image fusion, Journal of Ambient Intelligence  
664 and Humanized Computing 11 (2020) 1749-7969.
- 665 [30]. D. Stutz, A. Hermans, B. Leibe, Superpixels: An evaluation of the state-of-the-art, Computer Vision and Image  
666 Understanding 166 (2018) 1-27.
- 667 [31]. X. Ren, J. Malik, Learning a classification model for segmentation, In: International Conference on Computer Vision, 2003,  
668 pp. 10-17.
- 669 [32]. J. Chen, Z. Li, B. Huang, Linear spectral clustering superpixel, IEEE Transactions on Image Processing 26 (7) (2017) 3317-  
670 3329.
- 671 [33]. S. P. Lloyd, Least squares quantization in PCM, IEEE Transaction on Information Theory 28 (2) (1982) 129-137.

- 672 [34]. F. R. Siqueira, W. R. Schwartz, H. Pedrini, Multi-scale gray level co-occurrence matrices for texture description,  
673 Neurocomputing 120 (2013) 336-345.
- 674 [35]. Z. C. Lin, R. S. Liu, Z. X. Su, Linearized alternating direction method with adaptive penalty for low-rank representation, In:  
675 Advances in neural information processing systems, 2011, pp. 612-620.
- 676 [36]. G. Liu, Q. Liu, P. Li, Blessing of dimensionality: recovering mixture data via dictionary pursuit, IEEE Transaction on Pattern  
677 Analysis and Machine Intelligence 39 (1) (2017) 47-60.
- 678 [37]. A. Levin, D. Lischinski, Y. Weiss, A closed-form solution to natural image matting, IEEE Transaction on Pattern Analysis  
679 and Machine Intelligence 30 (2) (2008) 228-242.
- 680 [38]. <http://r0k.us/graphics/kodak>.
- 681 [39]. M. Aharon, M. Elad, A. Bruckstein, K-SVD: an algorithm for designing over-complete dictionaries for sparse representation,  
682 IEEE Transaction on Signal Processing 54 (11) (2006) 4311-4322.
- 683 [40]. M.B.A. Haghghat, A. Aghagolzadeh, H. Seyedarabi, A non-reference image fusion metric based on mutual information of  
684 image features, Computers & Electrical Engineering 37 (2011) 744-756.
- 685 [41]. G. Piella, H. J. Heijmans, A new quality metric for image fusion, In: International Conference on Image Processing, 2003.
- 686 [42]. A. Kolaman, O. Yadidpecht, Quaternion structural similarity: A new quality index for color images, IEEE Transactions on  
687 Image Processing 21(4) (2012) 1526-1536.
- 688 [43]. L. Alparone, S. Baronti, A. Garzelli, F. Nencini, A global quality measurement of pan-sharpened multispectral imagery, IEEE  
689 Geoscience and Remote Sensing Letters 1 (4) (2004) 313-317.
- 690 [44]. J. Zhao, R. Laganieri, Z. Liu, Performance assessment of combinative pixel-level image fusion based on an absolute feature  
691 measure, International Journal of Innovative Computing, Information and Control 3(6A) (2007) 1433-1447.
- 692 [45]. Z. Liu, D. S. Forsyth, R. Laganière, A feature-based metric for the quantitative evaluation of pixel-level image fusion,  
693 Computer Vision and Image Understanding 109 (1) (2008) 56-68.