# Region-based ICA Image Fusion using Textural Information

Nikolaos Mitianoudis and Sotirios-Antonios Antonopoulos
Electrical and Computer Engineering Department
Democritus University of Thrace
Xanthi, Greece
nmitiano@ee.duth.gr

Tania Stathaki
Electrical and Electronic Engineering Department
Imperial College London
London, UK
t.stathaki@imperial.ac.uk

*Abstract*—**Image Fusion is the procedure of combining useful features from multiple sensor image inputs to form a single composite image. In this work, the authors extend the previously proposed Image Fusion framework, based on self-trained Independent Component Analysis (ICA) bases, to a more sophisticated region-based Image Fusion system. The input images are segmented into three areas of different activity : edges, texture and constant background. A hierchical set of fusion rules employing textural information from the spatial-domain in the form of local variance, entropy and fourier energy is introduced. The proposed system improves the performance of our previous system.**

*Index Terms*—**Image Fusion; Independent Component Analysis; Texture Information**

## I. INTRODUCTION

Modern technology has enabled the development of low-cost, wireless sensors of various modalities that can be deployed to monitor a scene. In this paper, the case of multimodal imaging sensors of known position, that are employed to monitor a scene, will be investigated. The information provided by multimodal sensors can be quite diverse. Each image has been obtained using different instruments or acquisition techniques, allowing each image to have different characteristics, such as degradation, thermal and visual characteristics.

Let $x_1(i,j), \ldots, x_T(i,j)$ represent $T$ images of size $M_1 \times M_2$ capturing the same scene, where $i, j$ refer to the pixel coordinates in the image. The input images are assumed to have negligible registration problems. The process of combining the important features from the original $T$ images to form a single enhanced image $f(i,j)$ is referred to as *Image Fusion*. Fusion techniques can be divided into *spatial domain* and *transform domain* techniques [1], depending on the processing domain. Various transformations were proposed for image fusion, including the *Dual-Tree Wavelet Transform* [1], *Pyramid Decomposition* and self-trained Independent Component Analysis bases [2], [3]. All these transformations project the input images onto spatially localized bases, modeling sharp and abrupt transitions (edges) and therefore, transform the image into a more meaningful representation that can be used to detect and emphasize salient features, important for performing the task of image fusion.

The authors proposed a self-trained Image Fusion framework based on Independent Component Analysis, where the analysis transformation is estimated from a selection of images of similar content [2]. Several fusion rules were proposed in conjunction with this framework in [2]. The analysis framework is projecting the images into localized patches of relatively small size. The local mean value of the patches is subtracted and stored in order to reconstruct the local means of the fused image. In [2], an average of the stored means was used to reconstruct the fused image. In the case of multimodal inputs a gradient algorithm that optimises the Piella and Heijmans Fusion Quality index [4] was derived in [3] to infer an optimal means choice for the fused image. In this paper, the authors revisit the region-based rule that was proposed in [2]. The image was heuristically segmented into "active" and "non-active" regions, which were fused with the "max-abs" and the "mean" rule respectively. The proposed approach segments the image into three regions: "edges", "texture" and "background". A new set of fusion rules using the local standard deviation, entropy or fourier energy in the spatial domain were devised to fuse the class of "texture" regions. A hierchical application of fusion rules is used to construct the image in the ICA domain. The proposed system offers improved performance compared to our previous system in the case of "out-of-focus" examples.

## II. INTRODUCTION TO IMAGE FUSION USING ICA BASES

Assume an image $x(i,j)$ of size $M_1 \times M_2$. An "image patch" $x_w$ is defined as an $N \times N$ neighborhood centered around the pixel $(i_0, j_0)$. Assume that there exists a population of patches $x_w$, acquired randomly from the image $x(i,j)$. Each image patch $x_w(k,l)$ is arranged into a vector $\mathbf{x}_w(t) = \mathrm{vec}(x_w(k,l))$, using lexicographic ordering. The vectors $\mathbf{x}_w(t)$ are normalized to zero mean, producing unbiased vectors. These vectors can be expressed as linear combinations of the bases vectors $\mathbf{b}_j$ with weights $u_i(t), i = 1, \ldots, K$:

$$\mathbf{x}_w(t) = \sum_{k=1}^{K} u_k(t)\mathbf{b}_k = [\mathbf{b}_1 \ \mathbf{b}_2 \ldots \mathbf{b}_K] \begin{bmatrix} u_1(t) \\ u_2(t) \\ \ldots \\ u_K(t) \end{bmatrix} \quad (1)$$

where $t$ represents the $t$-th image patch selected from the original image. Equation (1) can be expressed, as follows:

$$\mathbf{x}_w(t) = B\mathbf{u}(t) \quad (2)$$

$$\mathbf{u}(t) = B^{-1}\mathbf{x}_w(t) = A\mathbf{x}_w(t) \qquad (3)$$

where $B = [\mathbf{b}_1 \ \mathbf{b}_2 \ldots \mathbf{b}_K]$ and $\mathbf{u}(t) = [u_1(t) \ u_2(t) \ldots u_K(t)]^T$. In this case, $A = B^{-1} = [\mathbf{a}_1 \ \mathbf{a}_2 \ldots \mathbf{a}_K]^T$ represents the *analysis* kernel and $B$ the *synthesis* kernel. The estimation of these basis vectors is performed using a population of training image patches $\mathbf{x}_w(t)$ and a criterion (cost function) that selects the basis vectors. Analysis/synthesis bases can be trained using *Independent Component Analysis* (ICA) and Topographic ICA, as explained in more detail in [2]. The training procedure needs to be performed only once, as the estimated transform can be used for fusing images with similar content to the training images.

A number of $N \times N$ patches (in the order of 10000 [5]) are randomly selected from similar-content training images. We perform Principal Component Analysis (PCA) on the selected patches in order to select the $K < N^2$ most important bases. Then, the ICA update rule or the topographical ICA rule in [2] for a chosen $L \times L$ neighborhood is iterated until convergence. In each iteration, the bases are orthogonalised using a symmetric decorrelation scheme. In the case of multimodal inputs, sample patches from all inputs are selected to train the ICA bases.

### A. Fusion in the ICA domain

After estimating an ICA transform, Image fusion using ICA bases is performed following the approach depicted in the generic diagram of Figure 1. Every possible $N \times N$ patch is isolated from each image $x_k(i, j)$ and is consequently re-arranged to form a vector $\mathbf{x}_k(t)$. These vectors $\mathbf{x}_k(t)$ are normalized to zero mean and the subtracted local mean $MN_k(t)$ is stored for the reconstruction process. Each of the input vectors $\mathbf{x}_k(t)$ is transformed to the ICA or Topographic ICA domain representation $\mathbf{u}_k(t)$, using equation (3). Optional denoising in the ICA representation is also possible, by applying sparse code shrinkage on the coefficients in the ICA domain [5], assuming Laplacian (generally sparse) priors for the ICA representation. The corresponding coefficients $\mathbf{u}_k(t)$ from each image are then combined to construct a composite image representation $\mathbf{u}_f(t)$ in the ICA domain. The next step is to move back to the spatial domain, using the synthesis kernel $B$. The optimal means $MN_f(t)$ are estimated using the gradient rule in [3]. In the case of images of similar contrast, one can use the average means as an optimal choice, as this is usually the answer of the gradient rule in [3]. The optimal means are then added to the corresponding image patch. The image $f(i, j)$ is synthesised by spatially averaging the image patches $\mathbf{u}_f(t)$ in the same order they were selected during the analysis step.

### B. Various fusion rules using ICA bases

Some basic rules that can be used for image fusion are described in this section. Fusion by the *absolute maximum* rule simply selects the greatest in absolute value of the corresponding coefficients in each image ("max-abs" rule). This process seems to convey all the information about the
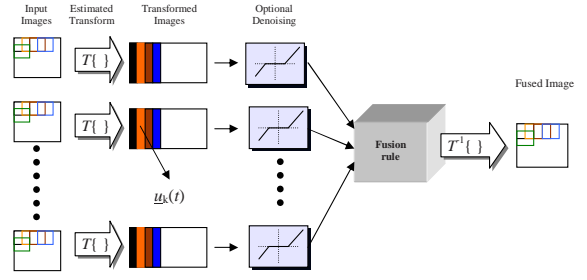


Fig. 1. The proposed fusion system using ICA / Topographical ICA bases.

edges to the fused image, however, the intensity information in constant background areas seems to be distorted. In contrast, fusion by the *averaging* rule averages the corresponding coefficients ("mean" rule). This process seems to preserve the correct contrast information, however, the edge details seem to get oversmoothed, since averaging is generally a "low-pass" filtering process.

A *Weighted Combination* (WC) pixel-based rule can be established using the ICA framework [2]. The fused image coefficients are constructed using a "weighted combination" of the input transform coefficients, i.e.

$$\mathbf{u}_f(t) = \sum_{k=1}^{T} w_k(t)\mathbf{u}_k(t) \qquad (4)$$

To estimate the contributions $w_k(t)$ of each image to the "fused" image, the mean absolute value ($\mathcal{L}_1$-norm) of each patch (arranged in a vector) in the transform domain can be employed as an activity indicator, because it fits a more general sparse profile of the ICA coefficients, denoted by a Laplacian distribution.

$$E_k(t) = ||\mathbf{u}_k(t)||_1 \qquad k = 1, \ldots, T \qquad (5)$$

The weights $w_k(t)$ should emphasize sources with more intense activity, as represented by $E_k(t)$. Consequently, the weights $w_k(t)$ for each patch $t$ can be estimated by the contribution of the $k$-th source image $\mathbf{u}_k(t)$ over the total contribution of all the $T$ source images at patch $t$, in terms of activity.

$$w_k(t) = E_k(t)/\sum_{k=1}^{T} E_k(t) \qquad (6)$$

A *regional* approach can also be established, by dividing the observed area into areas of "low" and "high" activity, using the $\mathcal{L}_1$-norm based $E_k(t)$ measurement. The areas containing salient information can be heuristically labeled as "high" activity areas, if $E_k(t) > 2\text{mean}_t\{E_k(t)\}$ and can be fused using a "max-abs" or a "weighted-combination" fusion rule. The remaining areas of "low-activity" contain background information and can be fused using the "mean" rule. Another regional approach can be to use alternative segmentations of the observed scene, based on the input sensor images and consequently fuse the different regions independently [6].

## III. An Improved Regional Fusion Rule Using Textural Information

In this section, we will describe a novel and improved regional fusion rule under the framework of ICA bases. The main aim was first to automate the procedure of selecting active and non-active regions. The second and novel aim was to attempt to identify areas of medium edge activity that can be considered to be textural information. Texture has so far been used to evaluate image fusion methods [7]. As there is no reported method to fuse textural areas in the literature (to the best of our knowledge), we will propose a relevant to method to stress the dominant textural areas in the fused image.

The first step will be to segment the input images in the ICA domain $\mathbf{u}_k(t)$ into three distinct regions: "edges", "texture" and "background". To achieve this, we will use the following activity detector $L_k(t)$, which is based the $\mathcal{L}_1$-norm based $E_k(t)$ measurement of (5).

$$L_k(t) = |E_k(t)|^p \qquad \forall\, k = 1, \ldots, T \qquad (7)$$

where the power $p \in [0.3, 0.5]$ is used to extend the value range of $E_k(t)$. The next step is to identify three distinct clusters in the values of $L_k(t)$. The cluster with larger values will correspond to "edges", the cluster with medium values to "texture" and the smaller values to "background". Clustering is performed using a typical K-Means algorithm on the values of $L_k(t)$ of each image. The segmentation result will be a map $S_k(t)$ for each input image, that will take the following form:

$$S_k(t) = \begin{cases} 3, & \text{Edge} \\ 2, & \text{Texture} \\ 1, & \text{Background} \end{cases} \qquad (8)$$

In order to define a single map for all input images, we combine the previous maps using a hierarchical approach. The main concept is that if one patch is considered to be an "edge" patch in one of the input images, then the corresponding patches should be fused using a fusion rule suitable for edges, regardless of their content. In a similar manner, a set of corresponding patches will be fused using a textural fusion rule, if at least one is tagged as "texture" and the rest are either tagged as "texture" or "background". Finally, if all corresponding patches are tagged as "background", then these patches will be fused using a rule suitable for low activity patches. This hierarchical approach will create the single map $S(t)$ simply by $S(t) = \max_k S_k(t)$. An example of the segmentation map that can be extracted using the described procedure is shown in Fig. 2. The out-of-focus fusion dataset "Bottles" is depicted in Fig. 2 (a), (b). The estimated map $S(t)$ is shown in Fig. 2(c) using white for "edge" areas, gray for "texture" areas and finally black for "background" areas.

Once the single map has been created, we can fuse these regions, using a suitable fusion rule, as follows:

$$\text{if } S(t) = \begin{cases} 3, & \text{"max-abs" rule} \\ 2, & \text{"Textural" rule} \\ 1, & \text{"Mean" rule} \end{cases} \qquad (9)$$



(a) Input Image 1      (b) Input Image 2

(c) Segmentation Map      (d) Regional Fusion (entropy)

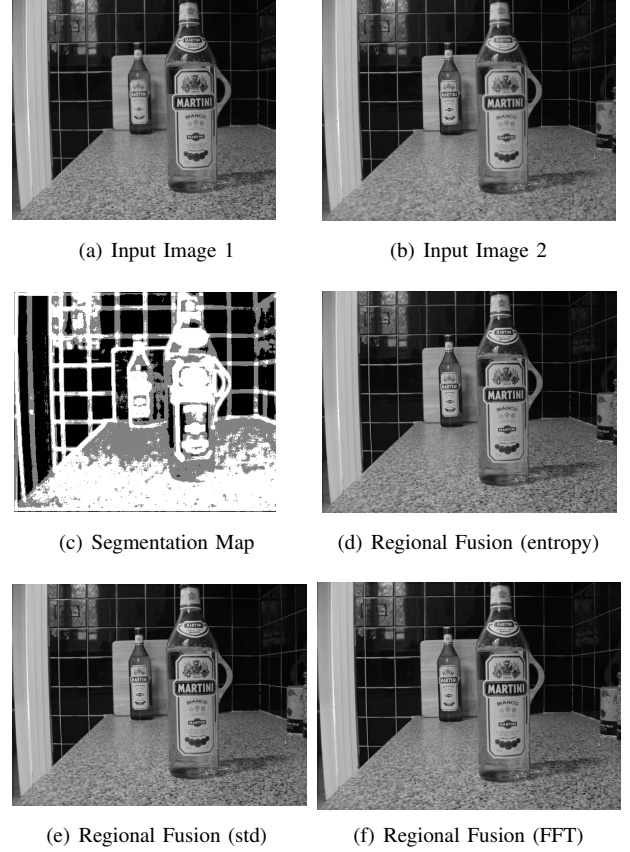(e) Regional Fusion (std)      (f) Regional Fusion (FFT)

Fig. 2. The "Bottles" fusion example: (a)-(b) Input Images, (c) Estimated three-cluster segmentation map, Regional Fusion rule using entropy (d), standard deviation (e) and FFT (f) for the "texture" areas.

The "max-abs" is suitable for fusing areas with strong edges. In contrast, the "Mean" rule is suitable for fusing areas with constant background. A novel fusion rule should be proposed for areas of medium edge activity, that can be well described as texture.

### A. Textural fusion rule

To fuse patches that contain mostly texture, we need a mechanism or features that can identify and stress the existence of texture, so that they can be highlighted in the fused image. Searching the vast literature of texture features [8], [9], we can see that texture can be identified by estimating various statistical and probabilistic measurements in the spatial domain, such as the *standard deviation* (std), *entropy* or *Fourier energy* (FFT). These measurements will be employed to calculate the weight factors $r_k(t)$ for the fusion rule.

The fused texture patch $\mathbf{u}_f(t)$ will be created using the following formula:

$$\mathbf{u}_f(t) = \sum_{k=1}^{T} r_k(t) \mathbf{u}_k(t) \qquad (10)$$

where $r_k(t)$ are weights that emphasize the most dominant texture patch. To perform this, we measure either the *standard deviation* or the *entropy* or the *Fourier energy* of the

corresponding input patches $\mathbf{x}_k(t)$ in the spatial domain, as follows:

$$r_k(t) = \frac{\text{std}\{\mathbf{x}_k(t)\}}{\sum_{k=1}^{T} \text{std}\{\mathbf{x}_k(t)\}} \tag{11}$$

$$r_k(t) = \frac{\text{entropy}\{\mathbf{x}_k(t)\}}{\sum_{k=1}^{T} \text{entropy}\{\mathbf{x}_k(t)\}} \tag{12}$$

Finally,

$$r_k(t) = \frac{\text{fft\_energ}\{\mathbf{x}_k(t)\}}{\sum_{k=1}^{T} \text{fft\_energ}\{\mathbf{x}_k(t)\}} \tag{13}$$

where fft_energ$\{\mathbf{x}_k(t)\}$ is a function that estimates the 1D-FFT of vector $\mathbf{x}_k(t)$, removes the DC component and the upper symmetrical half of the FFT and finally calculates the sum of the absolute value of the remaining coefficients. In other words, it is a measurement of periodicity using the Fourier transform. The above formulas normalise the weights to unit summation ($\sum_{k=1}^{T} r_k(t)$) to avoid inappropriate scaling of the patches. An example of using the three textural fusion rules in the aforementioned regional fusion scheme is depicted in Fig. 2 (d), (e) and (f) (entropy, standard deviation and Fourier energy respectively).

## IV. EXPERIMENTS

In this section, the performance of the proposed region-based ICA scheme is evaluated using a variety of datasets that were employed by the Image Fusion community. We used the typical training procedure for the ICA framework, training 60 $8 \times 8$ ICA bases from random natural images. This is performed offline only once and the bases are used for the rest of the experiments. We used the optimal contrast correction for the multimodal examples, as described in [3]. Optimal contrast correction was used for the "out-of-focus" examples, but the algorithm attributed almost equal weights to all input images as expected, since all input images feature similar exposure and contrast. The processing of color images for the multi-modal examples is performed in a similar manner to the method described in detail in [10]. For the "out-of-focus" examples, we fused each color channel (R-G-B) independently. We performed fusion under the ICA framework, using the "max-abs", "weighted combination" and "regional" rules, as described earlier. For the novel region-based schemes, we used the activity detector in (7) with a value of $p = 0.4$ and implemented the three fusion rules based on "entropy", "standard deviation" and "Fourier Coefficients". For performance comparison, the Dual-Tree Wavelet Transform (DT-WT) method using the "max-abs" rule will also be employed[1]. The Piella Index that is calculated in this section will constantly represent the second version of the Piella Index [4].

The first task was to demonstrate the novel framework's performance on examples of "out-of-focus" fusion. We employed the commonly used datasets "Disks", "Clocks", "Books" and "Pepsi", as they were available by the ImageFusion Server [11], and two sets ("Bottles" and "Berlin") created by

[1]Code for the Dual-Tree Wavelet Transform available online by the Polytechnic University of Brooklyn, NY at http://taco.poly.edu/WaveletSoftware/

the authors. We applied the ICA-based fusion framework following previously proposed fusion rules and the novel texture-based fusion rules based on "entropy", "standard deviation" and "fourier energy". Fusion performance was measured in terms of the Piella index and is outlined in Table I. In Fig. 3, some fusion results are depicted using the "Books" dataset and in Fig. 4, some fusion results using the "Disks" dataset. The first observation is that the activity detector proposed in (7), the k-means clustering approach and hierarchical characterisation is successful at segmenting the observed scene into areas of edges, texture and constant background. The segmentation maps that are shown in 2(c), 3(c) and 4(c) and demonstrate an accurate segmentation of the observed scenes. We observed that the scheme is more efficient for higher resolution images. In comparison to the regional scheme proposed in [2], where there was a simple threshold to discriminate between active and non-active regions, we have significant improvement as now the procedure is fully automated and a more refined segmentation into three clusters is achieved. The fusion rules that are proposed to fuse the "textural" regions seem to improve the performance of the ICA fusion framework. As previously reported, the ICA-based fusion methods outperform fusion methods based on the Dual-Tree Wavelet transform (DT-WT). The novel region-based methods outperform significantly the previous regional method. In addition, the novel region-based methods based on the "standard deviation" (std) and "fourier energy" (fft) seem to outperform the method based on "entropy". From a point of view, measuring the standard deviation, i.e. the energy of a patch is equivalent to measure the energy of the Fourier coefficients due to Parseval's theorem (without the DC component which is also subtracted in the form of local mean in the calculation of standard deviation). The two methods feature similar performance, as observed in Table I. The subtle differences may be caused due to round-off errors by MATLAB. Finally, the new proposed methods seem to outperform the "max-abs" rule which is usually the best choice for "out-of-focus" fusion examples.

The second task was to explore the novel framework's performance on "multimodal" fusion examples. The two Octet image sets were employed, as they were available by the ImageFusion Server [11]. These images, captured by Octec Ltd., show men and buildings with and without a smoke screen. They were captured with a Sony Camcorder and a LWIR sensor. We also used the "Dune" and "UNcamp" datasets of surveillance images from TNO Human Factors, provided by L. Toet in the Image Fusion Server [11]. The datasets consist of two series of visual and infrared frames capturing a human subject walking through various areas. The ICA fusion system still outperforms the image fusion using DT-WT. The novel regional schemes still outperform the previous regional and the weighted combination schemes. However, the "max-abs" ICA fusion outperforms the novel regional ICA fusion schemes. This implies that texture-based rules may not be very suitable for multi-modal fusion. The different texture that exists in the input images, due to the different capture modalities may be misleading the fusion

TABLE I

AVERAGE FUSION PERFORMANCE MEASUREMENTS USING PIELLA'S INDEX FOR OF THIS EXPERIMENTAL SECTION. THE PROPOSED ICA-BASED SCHEMES WITH THE THREE TEXTURAL FUSION RULES ARE COMPARED WITH THE PREVIOUS ICA-BASED FRAMEWORK AND THE DUAL-TREE WAVELET FRAMEWORK.

| Method | ICA Maxabs | ICA Weight | ICA Region | ICA Textr-Entr | ICA Textr-std | ICA Textr-fft | DT-WT maxabs |
|---|---|---|---|---|---|---|---|
| *Out-of-focus fusion* | | | | | | | |
| Disks | 0.9189 | 0.9110 | 0.9059 | 0.9184 | 0.9191 | 0.9192 | 0.9095 |
| Clocks | 0.9170 | 0.9096 | 0.9041 | 0.9169 | 0.9188 | 0.9186 | 0.9105 |
| Books | 0.9162 | 0.9164 | 0.9148 | 0.9174 | 0.9175 | 0.9175 | 0.9091 |
| Berlin | 0.9665 | 0.9731 | 0.9692 | 0.9691 | 0.9693 | 0.9693 | 0.9650 |
| Bottles | 0.9594 | 0.9659 | 0.9643 | 0.9620 | 0.9614 | 0.9614 | 0.9546 |
| Pepsi | 0.9437 | 0.9344 | 0.9310 | 0.9447 | 0.9448 | 0.9448 | 0.9400 |
| *Average* | 0.9369 | 0.9351 | 0.9315 | 0.9381 | **0.9385** | **0.9385** | 0.9315 |
| *Multi-modal fusion* | | | | | | | |
| Octet1 | 0.8661 | 0.8483 | 0.8343 | 0.8485 | 0.8523 | 0.8542 | 0.8232 |
| Octet2 | 0.8442 | 0.8296 | 0.8186 | 0.8210 | 0.8251 | 0.8249 | 0.8749 |
| Dune | 0.7414 | 0.7009 | 0.6825 | 0.7149 | 0.7210 | 0.7251 | 0.7118 |
| UN camp | 0.7500 | 0.7005 | 0.6868 | 0.7294 | 0.7303 | 0.7296 | 0.7159 |
| *Average* | **0.8004** | 0.7698 | 0.7556 | 0.7785 | 0.7822 | 0.7834 | 0.7815 |



(a) Input Image 1    (b) Input Image 2    (c) Segmentation Map    (d) ICA-maxabs

(e) DTWT-maxabs    (f) ICA-Texture-entropy    (g) ICA-Texture-std    (h) ICA-Texture-fft
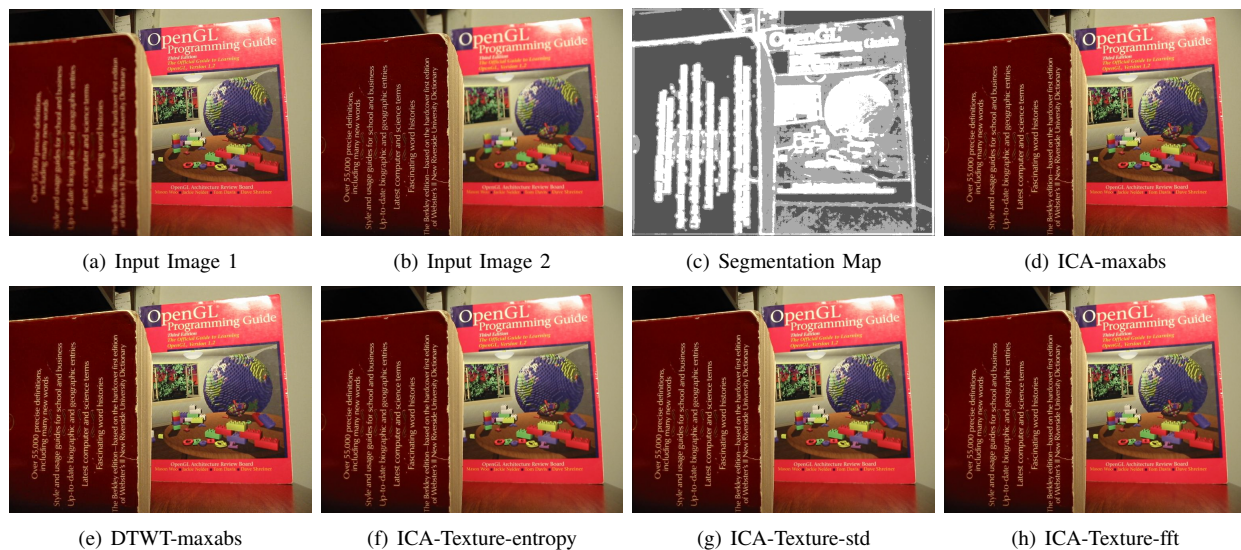
Fig. 3. The "Books" fusion example: (a)-(b) Input Images, (c) Estimated three-cluster segmentation map, Regional Fusion rule using entropy (d), standard deviation (e) an for the "texture" areas.

process, and thus the simple "edge injection" may be better in terms of image fusion performance. An example of multimodal fusion from the "Dune" sequence is shown in Fig. 5. Although the segmentation map seems to detect the correct areas of strong edges, texture and background, the Piella Index seems to favor the simple edge-injection "max-abs" rule.

## V. CONCLUSION

In this paper, the authors extend their previous image fusion framework based on Independent Component Analysis with a novel regional fusion rule. Initially, the input images are transformed to the ICA-domain representation, where the activity of each image patch is measured. Using k-means clustering and a hierarchical grouping concept, the image patches are divided into three groups: a) high activity patches (edges), b) medium activity patches (texture) and c) low activity patches (background). A different fusion rule is used for each different group. The max-abs rule is used for edges and the mean rule is used for background. A novel fusion rule measuring several texture properties in the spatial domain is used in the texture patches. The proposed scheme provided meaningful scene segmentation into these three areas. The proposed scheme offered improved performance compared to the "max-abs" fusion rule for "out-of-focus" fusion examples. In contrast, it was not so successful in the case of multimodal examples, maybe due to the different texture properties of the various modality input images.

## REFERENCES

[1] P. Hill, N. Canagarajah, and D. Bull, "Image fusion using complex wavelets," in *Proc. 13th British Machine Vision Conference*, Cardiff, UK, 2002.
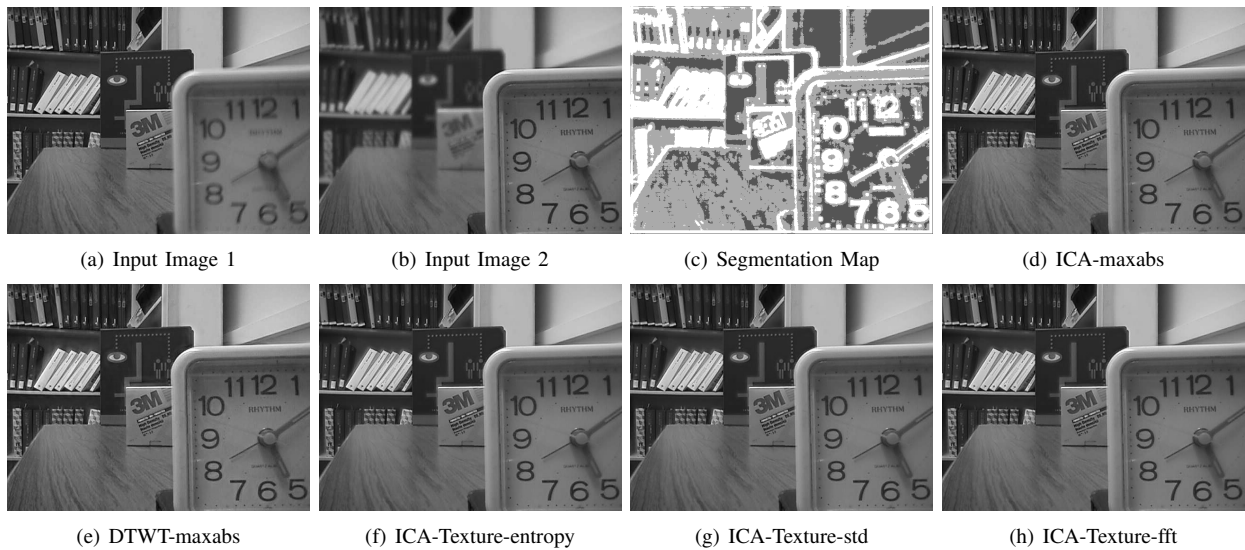
(a) Input Image 1  (b) Input Image 2  (c) Segmentation Map  (d) ICA-maxabs

(e) DTWT-maxabs  (f) ICA-Texture-entropy  (g) ICA-Texture-std  (h) ICA-Texture-fft

Fig. 4. The "Disks" fusion example: (a)-(b) Input Images, (c) Estimated three-cluster segmentation map, Regional Fusion rule using standard deviation (d) and entropy (e) for the "texture" areas.

(a) Input Image 1  (b) Input Image 2  (c) Segmentation Map  (d) ICA-maxabs

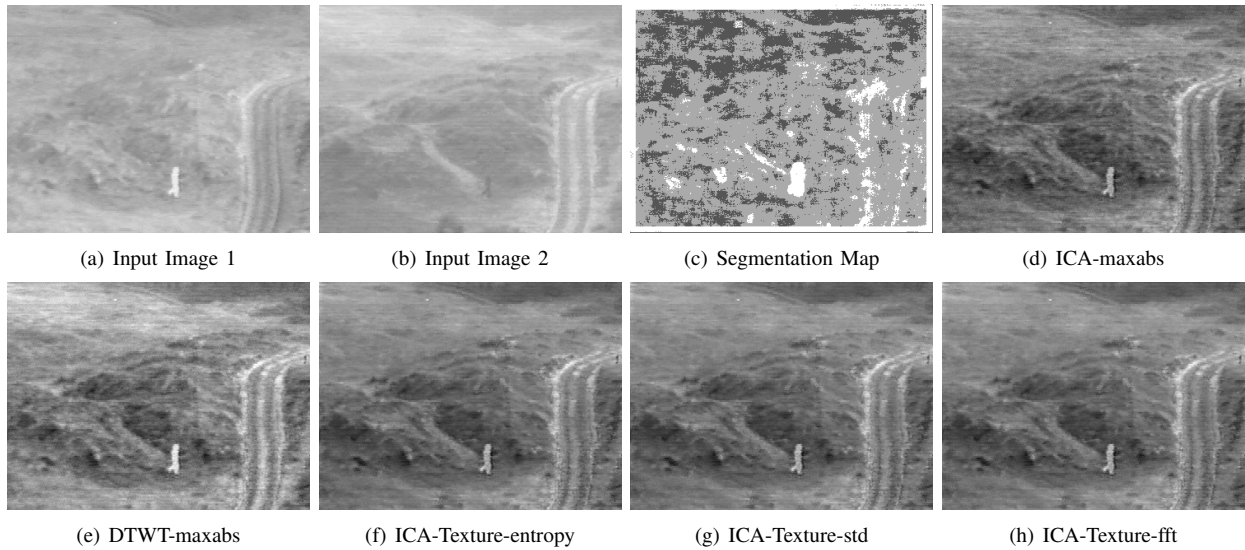(e) DTWT-maxabs  (f) ICA-Texture-entropy  (g) ICA-Texture-std  (h) ICA-Texture-fft

Fig. 5. The "Dune" fusion example: (a)-(b) Input Images, (c) Estimated three-cluster segmentation map, Regional Fusion rule using standard deviation (d) and entropy (e) for the "texture" areas.

[2] N. Mitianoudis and T. Stathaki, "Pixel-based and region-based image fusion schemes using ICA bases," *Elsevier Information Fusion*, vol. 8, no. 2, pp. 131–142, 2007.

[3] ——, "Optimal Contrast Correction for ICA-based fusion of Multimodal Images," *IEEE Sensors Journal*, vol. 8, no. 12, pp. 2016 – 2026, 2008.

[4] G. Piella, "A general framework for multiresolution image fusion: from pixels to regions," *Information Fusion*, vol. 4, pp. 259–280, 2003.

[5] A. Hyvärinen, P. O. Hoyer, and E. Oja, "Image denoising by sparse code shrinkage," in *Intelligent Signal Processing*, S. Haykin and B. Kosko, Eds. IEEE Press, 2001.

[6] N. Cvejic, D. Bull, and N. Canagarajah, "Region-based multimodal image fusion using ICA bases," *IEEE Sensors Journal*, vol. 7, no. 5, pp. 743–751, 2007.

[7] J. Majumdar and B. Patil, "A comparative analysis of image fusion methods using texture," in *Proc. of the 4th Int. Conf. on Signal and Image Processing*, Coimbatore, India, 2012, pp. 339 – 351.

[8] M. Petrou and P. Sevilla, *Image Processing: Dealing with Texture*. Wiley-Blackwell, 2006.

[9] P. Howarth and S. Ruger, "Evaluation of texture features forcontent-based image retrieval," in *Proc. of the Int. Conf. on Image and Video Retrieval*, Dublin, Ireland, 2004, pp. 326 – 334.

[10] N. Mitianoudis and T. Stathaki, "Optimal contrast for color image fusion using ICA bases," in *Proc. of 11th Int. Conf. on Information Fusion*, Cologne, Germany, July 2008.

[11] T. I. fusion server, "http://www.imagefusion.org/." [Online]. Available: http://www.imagefusion.org/