




## RESEARCH ARTICLE

# Improving unsupervised saliency detection by migrating from RGB to multispectral images

Miguel Á. Martínez<sup>1</sup>  | Sergi Etchebehere<sup>2</sup> | Eva M. Valero<sup>1</sup>  | Juan L. Nieves<sup>1</sup> 

<sup>1</sup>Department of Optics, Facultad de Ciencias, Universidad de Granada, Granada, Spain

<sup>2</sup>HP, Barcelona, Spain

## Correspondence

Miguel Á. Martínez, Department of Optics, Facultad de Ciencias, Universidad de Granada, Granada 18071, Spain.  
Email: martinezm@ugr.es

## Funding information

AZTI-Tecnalia, Grant/Award Number: C-3368-00; Secretaría de Estado de Investigación, Desarrollo e Innovación, Grant/Award Number: DPI2015-65471; Ministry of Economy and Competitiveness of Spain, Grant/Award Number: DPI2015-64571-R; Business-UGR Foundation; Tecnalia company

## Abstract

Saliency detection has been an important topic during the last decade. The main goal of saliency detection models is to detect the most relevant objects in a given scene. Most of these models use RGB (Red, Green, Blue) images as an input because they mainly focus on applications where features (eg, faces, textures, colors, or human silhouettes) are extracted from color images, and there are many labeled databases available for RGB-based saliency data. Nevertheless, the use of RGB inputs clearly limits the amount of information from where to extract the salient regions as spectral information is lost during the color image recording. On the contrary, multispectral systems are able to capture more than three bands in a single capture and can retrieve information from the full spectrum at a pixel. The main aim of this study is to investigate the advantages of using multispectral images instead of RGB images for saliency detection within the framework of unsupervised models. We compare the performance of several unsupervised saliency models with both RGB and multispectral images using a specific dataset of multispectral images with ground-truth data extracted from observers' fixation patterns. Our results show a general improvement when multispectral information is taken into account. The saliency maps estimated by using the multispectral features are closer to the ground-truth data, with the simplest Graph-based visual saliency and Boolean Map-based models showing good relative gain compared with other approaches.

## KEYWORDS

conspicuity, machine vision, multispectral images, spectral imaging, visual saliency

## 1 | INTRODUCTION

The human visual system is able to detect relevant or important information from all the data that enters the eye. This cognitive process, known as visual attention, is complex, and its complete understanding and simulation have been widely explored. In 1998, Itti et al<sup>1</sup> proposed the first completely functional saliency model, which tried to simulate where the human visual system would focus its attention on a given RGB image. After Itti's revolutionary work, many other models were created, which attempted to improve the

results. In order to extract salient information, most models utilize some specific features, from the more basic intensity, color, and orientation to the more advanced features such as motion, optical flow, flicker, multiple superimposed orientations (crosses or corners), and texture contrast.<sup>2</sup>

All the previously cited features use trichromatic images as an input. These are the more common types of images (RGB color images), which try to simulate how the human visual system responds to light and extracts color information. The human eye, and therefore a camera, has three kinds of channels or photoreceptors, sensitive to different parts of the

visible light spectrum. Consequently, when capturing an image, the incoming light recorded by the camera sensor is encoded with three numbers (R-, G-, and B-digital values or L-, M-, and S-cone responses), and thus, the spectral information is lost. Nevertheless, such spectral information might be useful for certain applications. In recent years, there has been a growing interest in devices (more and more affordable) able to capture all this extra information, not only with a better spectral resolution in visible light but also being able to capture light in other areas of the spectrum, such as the ultraviolet, infrared, and thermal. The increase in the availability of these multispectral and hyperspectral cameras has facilitated huge advances in fields such as robotics, remote sensing, satellite imaging, medicine, food control, and even object detection.<sup>3-5</sup>

In this study, we analyze the advantages of using multispectral images with the aim of salient object detection using some of the more known saliency models and adapting them to receive and take advantage of spectral information. The topic of saliency detection and prediction is described in general at an introductory level in the review/book by Li and Gao.<sup>6</sup> The idea is to compare the original models developed for RGB images with their adapted multispectral versions by using the most common evaluation metrics and investigate whether there is an improvement or not. Although multispectral images go beyond what the human vision can perceive, multispectral saliency detection does not imply a perfect simulation of bottom-up visual attention but rather a broader detection of objects that stand out spectrally from their neighbors, which can also be related to knowledge and task-associated visual attention, the so-called top-down visual attention. Specific visual attention models (VAMs) have been developed for spectral images<sup>7-9</sup> (Section 2), but in these studies, the comparison between RGB and multispectral images was not addressed specifically. Our study aims to tackle this issue using the least favorable situation for multispectral images, which is using models that have been specifically developed with RGB images in mind. Two specific fields of application that can benefit from the results shown in this article could be: surveillance and security field (to detect objects or events of interest in urban scenes using modified camera surveillance devices to make them multispectral), and a second one could be the active monitoring of the state of preservation of the elements present in urban scenes.

This article is organized as follows: Section 2 reviews some of the more relevant related studies modeling visual attention; Section 3 describes the methodology and the framework of the research, Section 4 analyzes the results obtained, and the conclusions are given in Section 5.

## 2 | VISUAL ATTENTION MODELING

During the last decade, it has been of great interest to determine where and why an observer aims their gaze at particular locations in a scene. When some areas in an image attract the visual attention and the point of gaze of an observer, it is said that these regions show high saliency, (ie, specific low-level visual features are attracting the observers' interest), and thus, the saliency map is a biologically plausible model for bottom-up attention as proposed by Koch and Ullman (1985).<sup>10</sup> Their definition of saliency relied on center-surround principles considering that points in the visual scene are salient if they differ from their neighbors. There are many features characterizing a visual scene, among which we could cite edges, contrast, luminance, and color as the main visual features defined at different scales. Classical bottom-up visual models obtain relatively good results when they use these features to localize the highly salient features in a scene, both for natural and artificial images. More recently, including task-dependent constraints within the saliency algorithms has been found to improve the derived salient maps.<sup>11</sup> These kinds of models, which operate at higher visual levels (ie, top-down models), use prior knowledge to gain visual attention. Eye-tracking systems are usually used to record observers' gaze paths as they view a collection of images. After discarding saccade fixation locations, the corresponding fixation map can be obtained.

As explained in the previous section, the most influential attempt to create a complete saliency model was made by Itti et al,<sup>1</sup> inspired by the theoretical work of Treisman et al<sup>12</sup> in the feature integration theory, where three basic features that influence the visual attention were proposed: intensity, color, and orientation. The Itti model proposes how to extract these three features from a digital color image based on bottom-up scene-based properties by selecting preattentively computed simple features and combining all of them into a conspicuity map for each channel. Doing this to different sizes of the same image through a Gaussian blur pyramid, the center-surround difference at each feature simulates the neuronal receptive fields found in the human visual system. Finally, after obtaining the relative saliency contribution of each feature, a linear combination resulting in the final saliency map is produced. Moreover, as established by Tatler et al,<sup>13</sup> there are differences between visual features in attended and non-attended spatial locations in an image. To be more specific, these differences are determined by various contrasts, luminances, chromaticity, energy, and orientation. Nevertheless, doubt on these findings has been cast by Baddeley and Tatler,<sup>14</sup> who found that a fixation map is dominated by high-frequency edges; the authors argue that contrast does

not contribute to saliency and that the other features are “behaviorally irrelevant.”

Later on, many models appeared, improving different assets of this initial approach: the use of a log-spectrum in the input image,<sup>15</sup> using the information theory to extract salient information,<sup>16</sup> using high-level features,<sup>17</sup> and supervised learning trained by large eye-tracking datasets.<sup>18</sup> Recently, the majority of leading benchmark models has been based on convolutional networks and deep learning techniques.<sup>19</sup>

In this section, we first describe the RGB-based models used in our study and then some models developed specifically for multispectral images.

## 2.1 | RGB-based saliency prediction

Of all the existing models, we have selected five and adapted them to receive multispectral images as input. This selection was carried out taking into account their impact, their accuracy, and the feasibility of adapting them to multispectral images.

1. ITTI: Itti's model<sup>1</sup> has been selected due to its influence on saliency detection research and the many times it has been used in previous studies. As we have explained, Itti uses center-surround differentiation over three main features: intensity, color, and orientation.
2. Graph-based visual saliency (GBVS): Harel et al<sup>20</sup> proposed the graph-based visual saliency, a modification of Itti's model; whilst using the same feature extraction, it proposes new activation, normalization, and combination steps based on graph computation. Activation and normalization are achieved by implementing a Markovian approach: a fully connected graph with a weight assigned to each edge connecting one node of the feature map to all the other nodes except itself. Therefore, by adding these two graph-based approaches to the steps of activation and normalization and using the feature extraction already proposed by Itti and a linear concatenation of normalized activation maps, they were able to improve significantly both the performance and the accuracy of the other existing saliency methods.
3. RARE: Published in 2012 by Riche et al,<sup>21</sup> it proposes finding salient information by looking at the rarity of the different features. Rarity is calculated by using co-occurrence matrices of a given pixel or region, giving high values to a pixel that has values that are less frequent. It uses principal component analysis (PCA) over the RGB images in order to find higher discriminations; it also uses Gabor filters to analyze different orientations.
4. BMS: Boolean map saliency was proposed by Zhang and Sclaroff in 2013.<sup>22</sup> The idea is to binarize the different channels of the image by using random thresholding

and extract the salient information by analyzing their topological structure. This model is quite simple, and using low-cost processing, it reaches high scores when compared to other models.

5. Learning discriminative subspaces (LDS): Continuing with the same idea as RARE, learning discriminative subspaces on random contrasts,<sup>23</sup> this model tries to project the images into more discriminative subspaces that allow targets to pop out. It calculates the principal components using a big set of image patches, and by maximizing the contrast between target and background, it learns what subspaces are more suited to show this differentiation.

## 2.2 | Spectral-based saliency prediction

Although the previously cited models are able to predict salient information with a high accuracy (while lower than supervised models), they extrapolate all the information from an RGB image. The idea of using multispectral or hyperspectral images in order to predict salient information is not new, and there have been several attempts to create spectral image-based saliency models. Most of them adapted Itti's model to receive different features such as:

1. Using space transformation methods such as PCA<sup>8</sup> or Nonnegative Matrix Factorization<sup>9</sup> in order to reduce the dimensionality of the multispectral images and obtain a higher contrast of the more distinguishable objects.
2. Computing spectral differentiation metrics between the different pixels to more easily computed spectral differences between the center and surround.<sup>9,24</sup>
3. Taking advantage of the higher spectral resolution to select more accurately the blue-yellow and red-green vectors extracted from the corresponding group of spectral bands.<sup>9</sup>

All the above studies were presented as complete saliency models instead of an adaptation of previous models, so it was difficult to distinguish whether the performance of these models is related specifically to the usage of multispectral or hyperspectral information. In our case, we use models specifically developed for RGB images and adapt them to receive multispectral information as input. Our aim is to investigate if there is an improvement in the models' performance when they use a more complete source of information to obtain the saliency prediction.

We are aware that the selected models are not among the best performing since the advent of convolutional neural networks (CNN-based saliency prediction approaches<sup>25,26</sup>). However, supervised models would require a high amount of labeled spectral images to produce acceptable results

because they would, per force, have to be retrained if spectral images are to be used as input. Currently, there are no labeled spectral image databases of more than four channels for saliency detection. We think it is worth investigating whether using spectral information can provide a significant improvement in unsupervised saliency prediction before tackling the huge task of capturing and labeling a sufficient amount of spectral data to test using supervised approaches for saliency prediction. Besides, finding efficient ways to adapt existing models to receive different input data also has an intrinsic interest.

### 3 | METHODS

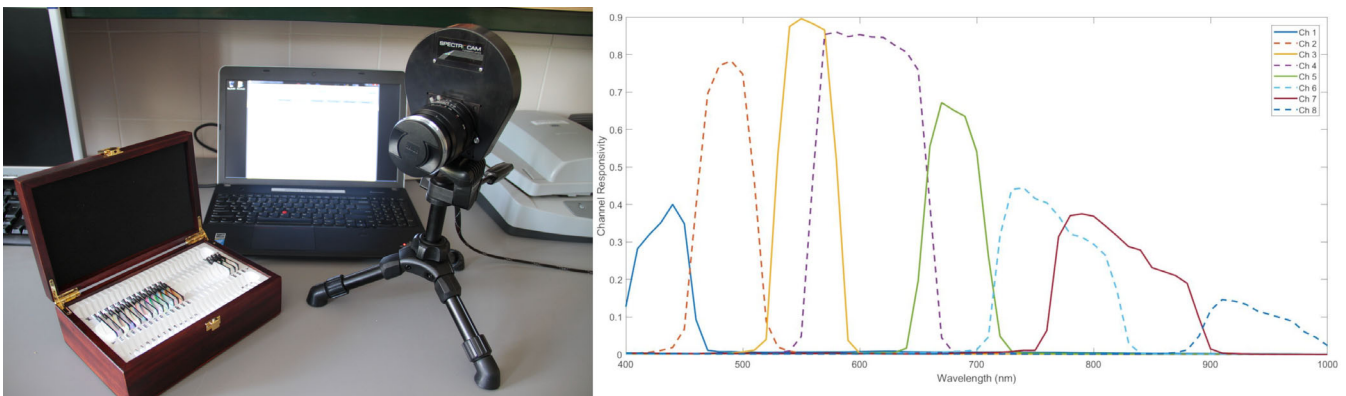
#### 3.1 | Image dataset and ground-truth data

We have used a set of nine multispectral images of urban scenes and their corresponding RGB versions, three of them containing people, to test RGB vs Multispectral image saliency prediction performance. The results of this study are applicable to saliency detection in any framework, although the scenes captured in this work have only urban content (buildings, vehicles, urban furniture, people, plants). The images were recorded using the PixelTeq (Halma, UK) SpectroCam VIS camera,<sup>27</sup> which is composed of a monochrome silicon sensor with a spatial resolution of  $2456 \times 2058$  pixels, sensitive to wavelengths of between 370 and 1100 nm (Figure 1 left). We are aware that more advanced sensors such as InGaAs-based ones are sensitive to spectral regions beyond this range (ie, up to 1700 or 2500 nm). These could certainly yield results for exploring different spectral regions that could provide interesting information for the saliency detection task. However, this would highly increase the cost of the imaging systems. In this regard, silicon-based sensors offer an affordable and easy-to-find solution that also demonstrates good performance for saliency detection. A filter wheel with eight slots, which is placed between the lens

and the sensor, is rotated to sequentially capture the images corresponding to each band. The exposure time for each channel was determined independently to ensure that the scene was correctly exposed for the corresponding band.

We selected a range of filters with specific transmittances to cover the visible and near-infrared (NIR) regions of the spectrum. In Figure 1 (right), the spectral responsivities of the channels are shown. Channels 1 to 5 have their responsivities within the visible range (roughly from 400 to 750 nm), channels 7 and 8 are sensitive in the NIR range (from 750 to 1000 nm), and channel 6 is both sensitive in the visible and NIR ranges. Of course, a higher number of channels with spectrally narrower sensitivities could help improve the saliency detection task by offering a larger amount of data. However, this would also increase the cost and complexity of the imaging system and the image data processing. For specific applications, an optimized filter selection could be carried out.<sup>28</sup> However, in this study, the available filters were meant for the general spectral imaging task, thus covering the whole visible and NIR range with certain overlap.

Each image has a resolution of  $2456 \times 2058$  pixels  $\times$  8 different channels corresponding to the transmittance of each filter to the scene. In order to generate the RGB images from our multispectral data, we only selected three filters that were reasonably close to the standard peak wavelengths of R, G, and B channels in a conventional RGB camera and used them as the three channels of the RGB image. These filters were those corresponding to channels 5 (680 nm), 3 (555 nm), and 1 (450 nm). At this point, one might think that the comparison is not fair as the RGB images only cover the visible range, and the multispectral system used in this article also covers the NIR range up to 1000 nm. However, this advantage is a part of the potential assets of multispectral systems, not only offering a higher number of spectral channels within the same spectral range but also extending its spectral range. Specifically, the sensor used in this study



**FIGURE 1** Left: PixelTeq SpectroCam VIS camera. Right: spectral responsivity of the eight channels used by the Spectrocam VIS camera, computed as the product of the spectral transmittance of each filter by the spectral responsivity of the monochrome sensor



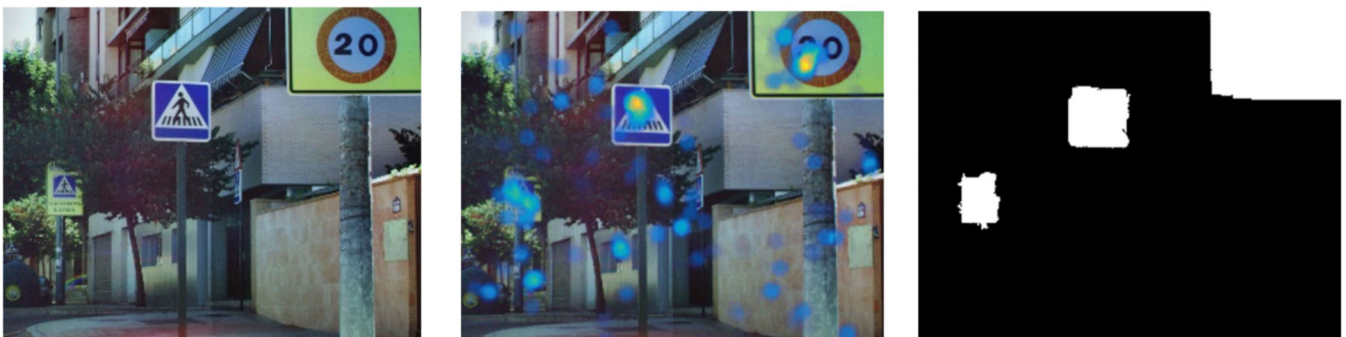
is a silicon-based sensor similar to the ones used in common RGB imaging systems. Therefore, we could extend the potential of any silicon sensor by removing the IR cut-off filter and adding the same color filters with the filter wheel. There is no need to use a more complex and costly InGaAs-based sensor for it.<sup>29</sup>

We used the RGB images to generate the ground-truth data for testing our hypothesis. Six observers, four women and two men with a mean age of 24 years, were asked to look freely at the RGB version of the images while their eye movements were being recorded with an Eye-tracker device (Tobii II, from Tobii Company, Danderyd, Sweden).<sup>30</sup> The images were presented for 6 seconds. The objects with the highest number of fixations in each image (accumulating more than 70% of the fixation time) were marked as ground-truth salient objects and manually segmented from the images to generate the ground-truth data (Figure 2).

### 3.2 | Features analyzed

In this subsection, we describe the features extracted from the spectral images and later fed as input for the adapted version of the VAMs. We have used a range of features that can be divided into three main groups:

1. **CIELAB:** In general, color information is used in most of the saliency models that use color or RGB images as input. The raw RGB color information can be used directly by the model, or the RGB can be transformed into a different color space that better emulates human perception. The CIELAB color space<sup>31</sup> is quite widely used for this purpose. The information conveyed by the three channels of the CIELAB feature ( $L^*$ ,  $a^*$ , and  $b^*$ ) is then fed to the adapted models as a three-dimensional image. Therefore, the model processes each channel independently, and the corresponding activation maps are concatenated.
2. **PCA:** In spectral images, the information contained in each pixel is usually high dimensional. Our hypothesis is that this extra amount of information can be useful in the prediction of saliency. Nevertheless, many spectra are smooth functions, and this means that there will be some amount of correlation between adjacent spectral bands. One way to exploit this correlation and try to keep the most relevant and distinctive features of the spectra is to use a dimensionality reduction technique such as PCA, which finds the best set of orthogonal components to represent the data while capturing the highest amount of their inner variance. PCA is also used as a feature in previously developed VAM for spectral images.<sup>24</sup> The principal component basis vectors are usually ranked by variance accounted for (VAF), and the number of principal component vectors used to represent the data is selected using a threshold criterion for accumulated VAF, usually ranking from 95% to 99%. For our data, we have determined that, using three principal components, we are able to account for at least 95% of the variance, so we have decided to use the projections of our image data onto the first three principal components as an additional feature for the saliency models. The three projected images are fed independently to the models, and the activation maps are computed and then concatenated. We have computed the PCA decomposition individually for each single image to preserve its distinctive characteristics as much as possible as the images were of a size that produced a sufficiently high number of pixels to allow for this approach.
3. **Spectral angle mapper (SAM)-spectral information divergence (SID):** When comparing spectral data to analyze differences between them, it is good practice to not only compare them channel by channel but to also consider the spectrum as a whole. In the case of spectral images, as each pixel has  $N$  spectral components, the image can be considered an array of signals, and each pixel can be compared with the mean signal in the image, which could be a way to identify which are the most distinctive regions. There are different metrics used to discriminate spectral signals numerically, for instance, root mean square error computes the square root of the mean of the channel-wise differences to the square, or



**FIGURE 2** Original scene (left), fixation map (center), and segmented ground-truth image (right) for one of the scenes

Goodness-of-Fit Coefficient (GFC) is the cosine of the angle between two spectral signals (considering them as vectors on a Hilbert space<sup>32</sup>). In our case, we use the so-called SAM-SID distance,<sup>33</sup> which is a combination of both the SAM and the SID. SAM is defined as the angle between two spectral signatures  $s$  and  $s'$  (and so the  $\cos^{-1}$  of the GFC value) as expressed in the following formula:

$$\text{SAM}(s, s') = \cos^{-1} \left( \frac{\langle s, s' \rangle}{\|s\| \cdot \|s'\|} \right) \quad (1)$$

Meanwhile, SID is the discrepancy between the uncertainty of two spectral signatures,  $s$  and  $s'$ , which is computed using their respective probability density distributions  $p$  and  $q$ :

$$D(s \| s') = \sum_{j=1}^L p_j \log \left( \frac{p_j}{q_j} \right) \quad (2)$$

$$D(s' \| s) = \sum_{j=1}^L q_j \log \left( \frac{q_j}{p_j} \right) \quad (3)$$

$$\text{SID}(s, s') = D(s \| s') + D(s' \| s) \quad (4)$$

Then, the combination of both SAM and SID is performed as the sinus of the angle by the information divergence:

$$D_{\text{SAM-SID}}(c, s) = \sin[\text{SAM}(c, s)] \times \text{SID}(c, s) \quad (5)$$

The advantage of this metric is that it combines sensitivity to differences in spectral amplitude distribution (SID) with sensitivity to differences in spectral shape (SAM). By finding the product of these two measures, the spectral discriminability of the SID-SAM mixed metric is increased because it makes two similar spectral signatures even more similar and two dissimilar spectral signatures more distinct.<sup>31</sup> Therefore, a one-channel feature is introduced as input to the saliency models, showing the SAM-SID difference between each pixel and the mean spectra of the scene. This feature is activated by the model, and saliency is predicted.

Both PCA and SAM-SID are features that can only be extracted from multispectral images; nevertheless, the color information is already used as a feature in most saliency models. In Figure 3, we show one scene (original scene), its segmentation ground truth, and the corresponding feature images (PCA, SAM-SID and CIELAB). The salient objects tend to have high intensity in some of the feature images, which can be useful for improving the performance of the VAM.

### 3.3 | Model adaptations

In this section, we explain the adaptations carried out on the existing visual saliency models to enable them to receive spectral features as inputs. As the different models have a completely different architecture, we have designed different ways of adapting them to accept the spectral features as inputs.

Figure 4 summarizes the work flow of the experiment performed for each of the models selected, with the aim of establishing if the use of spectral features as input produces an increase in the performance of the models. We first used RGB images as input for the model and obtained the corresponding saliency map. Then, we used the adapted version of the model with the spectral feature images as input and obtained the spectral-based saliency map. Finally, we used the ground truth and the set of metrics described in Section 3.4 to compare the performance of the model in the two situations (RGB or spectral features as input).

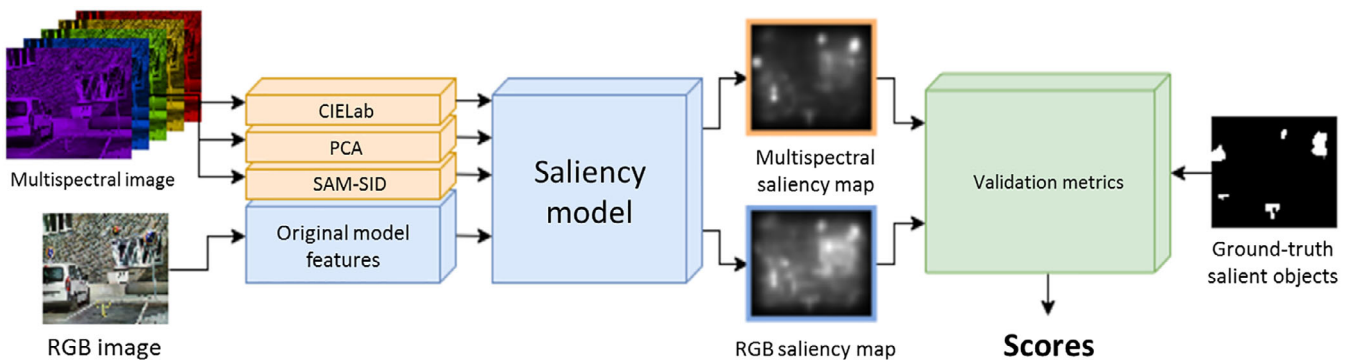
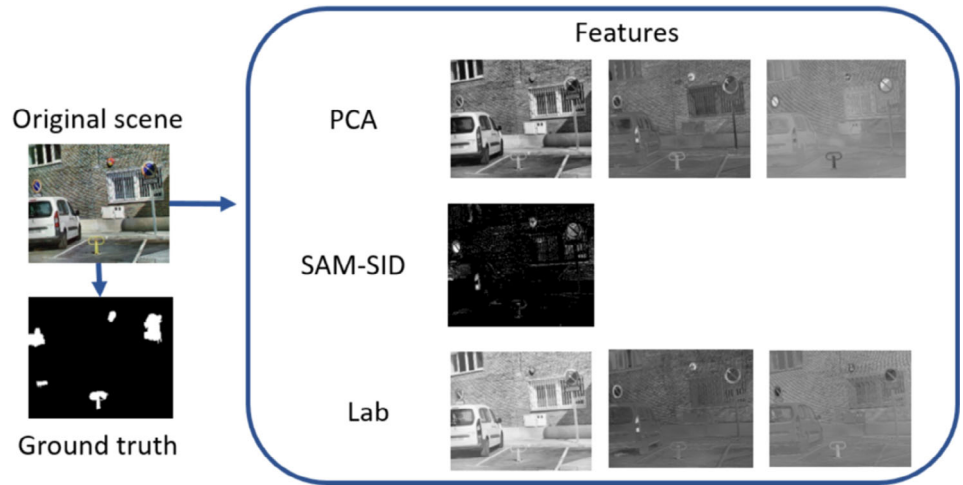
Itti and GBVS use intensity, color, and orientation as the main features, and then, the activation maps are computed. For these two models, we have substituted intensity and color for the CIELAB features, and we have used the L\* image to compute the orientation maps. Then, we added both PCAs (sequentially for each PCA component) and SAM-SID as extra features, leaving the models with a total of four feature global classes to be activated and combined. We have merged the activation maps with equal weights for all the features. Both RARE and LDS use PCAs to find a space that increases the differences between the objects. In this case, we substitute the three-dimensional input image with a seven-dimensional one, composed of the three CIELAB channels, the first three principal components, and the SAM-SID image. We then run the model with the corresponding space transformations, and the final saliency map is obtained. The BMS model applies random thresholding to the different channels of the input image. In this case, instead of applying the threshold to three different channels (RGB), we have used the random thresholding for the seven different maps (CIELAB + PCA + SAM-SID).

### 3.4 | Validation

Once a model detects the main salient regions in an image, it is necessary to validate its performance over ground-truth data. There are several metrics commonly used in this field and standardized so different models can be compared, although consistent results cannot always be obtained.<sup>34</sup> Depending on the application and the kind of data used for validation, some metrics can be more appropriate than others. We decided to use the following three metrics for our experiment:

1. Area under curve (AUC): this is computed from the receiving operator characteristic curve. For different values of

**FIGURE 3** RGB scene and corresponding feature images fed as input to the visual attention model tested



**FIGURE 4** Illustration of the work flow of our experiment. The procedure is repeated for each of the models tested

threshold in the saliency map produced by the model, true positives and false positives are computed by using the ground-truth data. Two main implementations of the AUC metrics are used: AUC-Borji<sup>35</sup> and AUC-Judd.<sup>18</sup> Another version of this metric was created in order to compensate the well-known center bias, the shuffled AUC,<sup>24</sup> which was the one we used to validate our data. The main drawback of the AUC metric is that low-valued false positives are not penalized.<sup>36</sup> This means that, if the saliency map is predicting objects as salient that are not truly salient according to the ground truth, it could still reach high values of AUC. In other words, diffuse saliency maps in which many areas are highlighted with not very extreme values of saliency are not considered poor quality.

2. Normalized scan-path saliency (NSS): This is computed as the averaged normalized saliency at the ground-truth location. Chance level is assigned a zero value, and a positive value would mean any value above the chance results. This method solves the issue of not penalizing low-valued false positives by assigning the highest score to a map that would detect all the pixels in the ground-truth salient regions as salient and would have zero values in all the other pixels in the image.<sup>37</sup>

3. Information gain (IG): This is a metric designed to compare two saliency maps taking into account the similarity of the probabilistic distribution with the ground-truth data.<sup>38</sup> Therefore, this metric is well suited for direct comparison between two different saliency methods, computing the gain or loss in information with respect to the ground-truth data for the two maps that are compared.

Although there are many more different metrics for saliency benchmarking, most of them can be highly correlated with one of the three metrics that we have chosen; these three metrics are good representatives of different strategies in the definition of the quality of saliency prediction.

## 4 | RESULTS

As we explained in the previous section, for each of the nine multispectral images, we calculated their saliency maps predicted by the five different models when using both the original features and the spectral ones. An example of these saliency maps can be seen in Figure 4.

For each of the saliency maps, the scores of the three different metrics described in Section 3.4 were calculated.

Table 1 shows the average and SD over the nine images for each of the models using both original and spectral features and each of the metrics and also the relative difference between both inputs' scores. In the case of AUC and NSS, the difference between the original and the spectral features is shown, meaning a positive, better score of the spectral features. The relative gain for the use of spectral features with respect to RGB features is also shown in the table. In the case of IG, as it already compares the two maps, only the average over the images is shown; a positive result shows better accuracy of the spectral features over the original RGB-based features.

Analyzing the results in Table 1, we can observe some differences between the different models and also between the different metrics. We can see how both ITTI and GBVS models have one of the highest scores in AUC, whilst the NSS score found is below the average across the models. One of the reasons of this noticeable difference between AUC and NSS in the ITTI and GBVS models might be the large amount of high (or salient) values in the maps; having many false positives is penalized by NSS but not by AUC. Now looking at the RARE results, this is the model scoring the highest in AUC and second highest in NSS. We can appreciate in Figure 5 how the resulting maps tend to contain high values in the salient object regions and generally low values for nonsalient regions. The BMS and LDS models are among the worst performing overall, having relatively low AUC and NSS scores both for RGB and multispectral images.

Now, we analyze the models' performance when we use the spectral features as input, which is the main aim of our experiment. Except for the RARE model in the IG metric, we have found that there is an improvement in the models' performance when used with spectral features. This improvement is much more apparent for the NSS and IG

metrics than for the AUC metric. For the Itti and GBVS models, there is a clear improvement in NSS values, which reach a level comparable to other models for the spectral features, while the performance is much poorer if we use the RGB image as input. For the RARE model, we can see the least improvement in AUC, the second smallest in NSS, and even a decrease in the accuracy in IG.

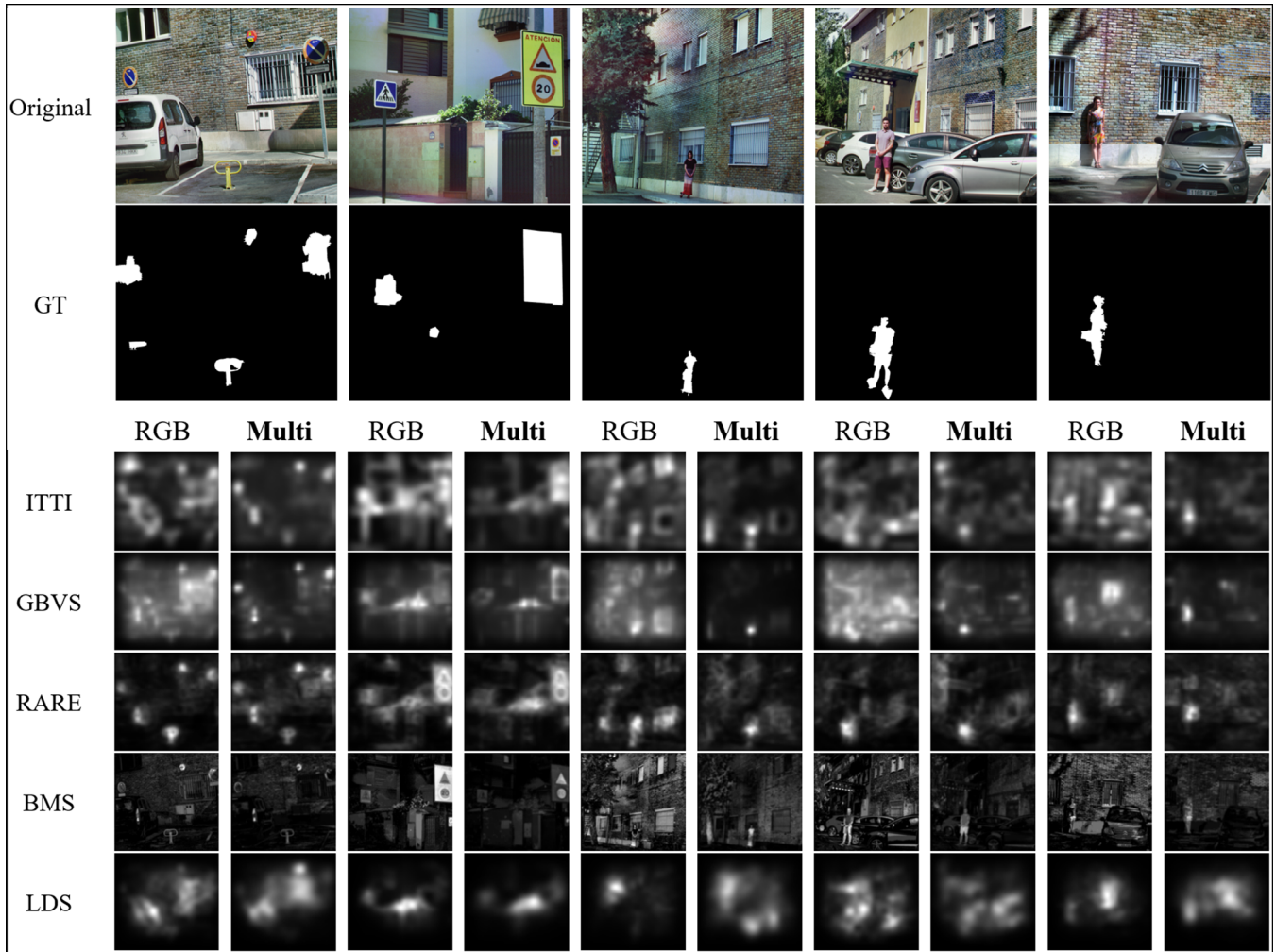
The RARE model looks for rarity instead of center-surround difference for computing the saliency map, so its strategy is markedly different from the first two models analyzed. The model is already performing quite well (compared with the others) when using the RGB image as input, and the adaptations we have introduced might not be able to add enough value to the spectral features. Regarding BMS and LDS, the accuracy of both increases when spectral information is used: around 0.6 in NSS and 0.5 in IG, with BMS reaching the highest IG score. This considerable improvement in performance might be due to a more successful adaptation strategy when introducing the spectral features. The average relative gain for all five models is 9.2% for AUC and 61.2% for NSS. Finding the precise factors that result in the observed improvement when using multispectral scenes as input for the VAMs tested is not a straightforward task. One factor is related to the new features introduced (PCA and SAM-SID), which in some instances clearly highlight the salient objects, as can be seen in Figure 3 and also in Figure 6. The remaining factors are linked to the specific way each model uses the input features to extract the saliency maps, and a detailed discussion would be excessively long considering the number and diversity of the models presented here and the fact that, for some of them, it is not easy to analyze each step sequentially and its relationship with the final saliency map delivered by the model.

**TABLE 1** Average and SDs over the nine images for (rows) each model using both original and spectral features and (columns) each of the metrics

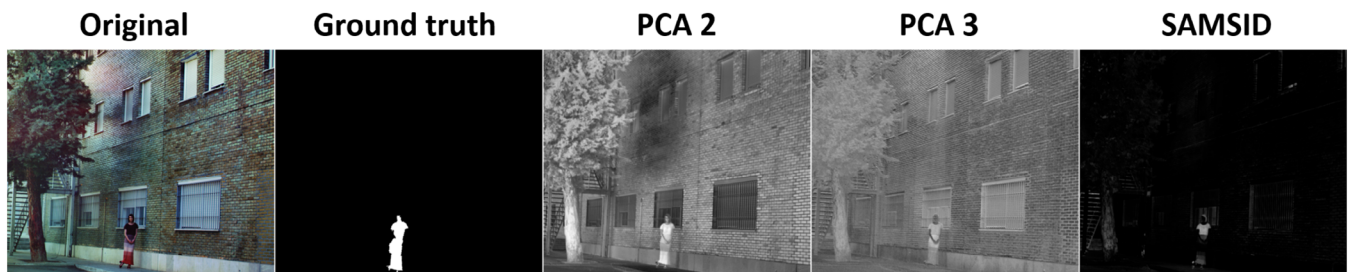
		AUC	Relative AUC variation (%)	NSS	Relative NSS variation (%)	IG
GBVS	RGB	0.792 (0.114)	9.7	1.176 (0.440)	93.2	0.526 (0.547)
	Hype	0.868 (0.064)		2.272 (1.348)		
ITTI	RGB	0.843 (0.080)	6.7	1.359 (0.535)	62.4	0.480 (0.441)
	Hype	0.904 (0.064)		2.356 (1.147)		
BMS	RGB	0.631 (0.125)	20.6	1.093 (0.809)	70.4	0.551 (1.143)
	Hype	0.761 (0.128)		1.861 (1.237)		
LDS	RGB	0.569 (0.144)	7.0	0.763 (0.604)	70.6	0.465 (0.230)
	Hype	0.609 (0.087)		1.302 (0.714)		
RARE	RGB	0.895 (0.087)	2.1	2.121 (1.081)	9.4	-0.053 (0.333)
	Hype	0.914 (0.049)		2.320 (1.007)		

Abbreviations: AUC, area under curve; BMS, Boolean map saliency; GBVS, graph-based visual saliency; IG, information gain; ITTI, XXX; LDS, learning discriminative subspace; NSS, normalized scan-path saliency; RARE, rare; RGB, Red, Green, Blue.





**FIGURE 5** An example of the saliency maps for each model using both original and spectral features of different images; ground truth (GT) is also shown for comparison



**FIGURE 6** An example of one of the images (original RGB), with its segmentation ground truth and its feature images corresponding to principal components 2 and 3 (principal component analysis 2 and 3), and spectral angle mapper-spectral information divergence

## 5 | CONCLUSIONS AND FUTURE WORK

We have used AUC, NSS, and IG metrics to assess the performance of five well-known VAMs with multispectral and conventional RGB color images. Our results suggest that the saliency maps produced by using the multispectral features are

closer to the ground-truth data. The higher gain for NSS is quite significant as this metric has advantages over AUC. In fact, NSS will be adopted as the gold standard quite soon in the VAM as the most popular benchmark.<sup>39</sup>

Saliency prediction performance has improved dramatically during the last 3 years after the outbreaks of the deep learning algorithms. Our promising results demonstrate the fact that a

CNN-based model, adequately trained using our specific spectral features, will improve the detection of the salient regions. A potential CNN-based spectral saliency detection method will carry out a prediction of the salient regions, analyzing in parallel all the spectral bands of an input image. This higher amount of information compared to RGB images will allow the Convolutional Neural Networks (CNNs) to find more complex features to detect saliency. Typically, we would need over 1000 images to obtain a decent accuracy in image classification on the cross-validation set (or even more if a transfer learning on an already trained model is not used). However, in the absence of such a number of multispectral images adapted for a saliency task, it would be difficult to hazard even a guess regarding the final spectral performance. After the results found in this study, a new multispectral image database is being built, together with its ground-truth data. It is a matter for further studies to implement a CNN-based spectral saliency model, adequately trained with this labeled multispectral image dataset.

## ACKNOWLEDGEMENTS

This research was supported by a joint agreement (reference number C-3368-00) between Tecnia company and the Business-UGR Foundation and through the Ministry of Economy and Competitiveness of Spain under research grant DPI2015-64571-R. We also thank Angela Tate for reviewing this text.

## ORCID

Miguel Á. Martínez  <https://orcid.org/0000-0003-3534-6733>

Eva M. Valero  <https://orcid.org/0000-0002-4671-5533>

Juan L. Nieves  <https://orcid.org/0000-0002-3103-8322>

## REFERENCES

- [1] Itti L, Koch C, Niebur E. A model of saliency-based visual attention for rapid scene analysis. *IEEE Trans Pattern Anal Mach Intell.* 1998;20:1254-1259.
- [2] Borji A, Itti L. State-of-the-art in visual attention modeling. *IEEE Trans Pattern Anal Mach Intell.* 2013;35(1):185-207.
- [3] Dale LM, Thewis A, Boudry C, et al. Hyperspectral imaging applications in agriculture and agro-food product quality and safety control: a review. *Appl Spectrosc Rev.* 2013;48(2):142-159.
- [4] Lu G, Fei B. Medical hyperspectral imaging: a review. *J Biomed Opt.* 2014;19(1):010901.
- [5] Liang H. Advances in multispectral and hyperspectral imaging for archaeology and art conservation. *Appl Phys A.* 2012;106(2):309-323.
- [6] Li J, Gao W, eds. *Visual Saliency Computation: A Machine Learning Perspective.* Vol 8408. Berlin, Germany: Springer; 2014.
- [7] Cornia M, Baraldi L, Serra G, Cucchiara R. Predicting human eye fixations via an LSTM-based saliency attentive model. 2016. <https://arxiv.org/pdf/1611.09571.pdf>. Accessed August 10, 2018.
- [8] Wang Q, Yan P, Yuan Y, Li X. Multi-spectral saliency detection. *Pattern Recognit Lett.* 2013;34:34-41.
- [9] Zhang J, Geng W, Zhuo L, Tian Q, Cao Y. Multiscale target extraction using a spectral saliency map for a hyperspectral image. *Appl Optics.* 2016;55:8089-8100.
- [10] Koch C, Ullman S. Shifts in selective visual attention: towards the underlying neural circuitry. *Hum Neurobiol.* 1985;4:219-227.
- [11] Canosa RL. Modelling selective perception of complex natural scenes. *Int J Artif Intell Tools.* 2005;14:233-260.
- [12] Treisman AM, Gelade G. A feature-integration theory of attention. *Cogn Psychol.* 1980;12:97-136.
- [13] Tatler BW, Baddeley RJ, Gilchrist ID. Visual correlates of fixation selection: effects of scale and time. *Vision Res.* 2005;45:643-659.
- [14] Baddeley R, Tatler B. High frequency edges (but not contrast) predict where we fixate: a Bayesian system identification analysis. *Vision Res.* 2006;46:2824-2833.
- [15] Hou X, Zhang L. Saliency detection: a spectral residual approach. Paper presented at: IEEE Conference on Computer Vision and Pattern Recognition; Minneapolis, MN; 1-8, 2007.
- [16] Bruce NDB, Tsotsos JK. Saliency, attention, visual search: an information theoretic approach. *J Vis.* 2009;9:1-24.
- [17] Sharma P, Cheikh FA, Hardeberg JY. Saliency map for human gaze prediction in images. Paper presented at: 16th Color Imaging Conf.; 2008; Portland, OR, 332-337.
- [18] Judd T, Ehinger K, Durand F, Torralba A. Learning to predict where humans look. Paper presented at: Kyoto, Japan: IEEE Int. Conf. Computer Vision; 2009; 2106-2113.
- [19] Borji A. Saliency prediction in the deep learning era: an empirical investigation. *arXiv preprint arXiv:1810.03716.* 2018.
- [20] Harel J, Koch C, Perona P. Graph-based visual saliency. Paper presented at: Advances in Neural Information Processing Systems Vancouver, Canada; 2006; 545-552.
- [21] Riche N, Mancas M, Duvinage M, Mibulumukini M, Gosselin B, Dutoit T. Rare2012: a multi-scale rarity-based saliency detection with its comparative statistical analysis. *Signal Proces Image Commun.* 2013;28:642-658.
- [22] Zhang J, Sclaroff S. Saliency detection: a boolean map approach. Paper presented at: Proceedings of the IEEE International Conference on Computer Vision Sydney, Australia; 2013; 153-160.
- [23] Fang S, Li J, Tian Y, Huang T, Chen X. Learning discriminative subspaces on random contrasts for image saliency analysis. *IEEE Trans Neural Netw Learn Syst.* 2017;28:1095-1108.
- [24] Le Moan S, Mansouri A, Hardeberg JY, Voisin Y. Saliency for spectral image analysis. *IEEE J Sel Top Appl Earth Obs Remote Sens.* 2013;6:2472-2479.
- [25] Wang L, Gao C, Jian J, Tang L, Liu J. Semantic feature based multi-spectral saliency detection. *Multimed Tools Appl.* 2018;77(3):3387-3403.
- [26] Wang T, Borji A, Zhang L, Zhang P, Lu H. A stagewise refinement model for detecting salient objects in images. Paper presented at: Proceedings of the IEEE International Conference on Computer Vision Venice, Italy; October, 2017; 4019-4028.
- [27] <http://halmapr.com/news/pixelteq/2016/07/21/new-spectrocam-swir-640-multispectral-wheel-camera-from-pixelteq>. Pixelteq Spectrocam SWIR 640. Accessed January 10, 2019.
- [28] Hardeberg JY. Filter selection for multispectral color image acquisition. *J Imaging Sci Technol.* 2004;48(2):105-110.

- [29] Martin T, Brubaker R, Dixon P, Gagliardi MA, Sudol T. 640x512 InGaAs focal plane array camera for visible and SWIR imaging. *Infrared Technology and Applications XXXI*. Vol 5783. Bellingham, Washington, USA: International Society for Optics and Photonics; 2005:12-20.
- [30] <https://www.tobii.com/>. Tobii eyetracker. Accessed January 10, 2019
- [31] Schanda J, ed. *Colorimetry: Understanding the CIE System*. Hoboken, Nueva Jersey, USA: John Wiley & Sons; 2007.
- [32] Hernández-Andrés J, Romero J, García-Beltrán A, Nieves JL. Testing linear models on spectral daylight measurements. *Appl Optics*. 1998;37:971-977.
- [33] Du Y, Chang CI, Ren H, Chang CC, Jensen JO, D'Amico FM. New hyperspectral discrimination measure for spectral characterization. *Opt Eng*. 2004;43:1777-1778.
- [34] Kümmerer M, Wallis TS, Bethge M. Saliency Benchmarking: Separating Models, Maps and Metrics. *arXiv preprint arXiv:1704.08615*. 2017.
- [35] Borji A, Sihite DN, Itti L. Quantitative analysis of human-model agreement in visual saliency modelling: a comparative study. *IEEE Trans Image Process*. 2012;22:55-69.
- [36] Bylinskii Z, Judd T, Oliva A, Torralba A, Durand F. What do different evaluation metrics tell us about saliency models? *arXiv:1604.03605*. 2016.
- [37] Peters RJ, Iyer A, Itti L, Koch C. Components of bottom-up gaze allocation in natural images. *Vision Res*. 2005;45:2397-2416.
- [38] Kummerer M, Wallis T, Bethge M. Information-theoretic model comparison unifies saliency metrics. *Proc Natl Acad Sci*. 2015; 112:16054-16059.
- [39] [http://saliency.mit.edu/results\\_mit300.html](http://saliency.mit.edu/results_mit300.html). MIT Saliency benchmark. Accessed January 10, 2019

## AUTHOR BIOGRAPHIES

**Miguel Á. Martínez** is a post-doctoral researcher working on multi- and hyperspectral high-dynamic range image capture and processing in the visible and near infrared. His research interests are in the field of high-dynamic range imaging, spectral imaging, digital image processing and analysis, color and spectral sciences, machine/deep learning, saliency detection, color vision, etc. He has a BS in telecommunications engineering, specializing in image and sound; an MS in color in informatics and media technology; and a PhD in physics and space sciences.

**Sergi Etchebehere** is a researcher who worked on the detection and classification of salient objects in hyperspectral images. His research interests are spectral image processing and analysis and use of machine learning techniques in computer vision. He currently works at Hewlett Packard.

**Eva Valero** is an associate professor in the Department of Optics at the University of Granada. Her recent research interests include hyperspectral imaging, spectral estimation, HDR Imaging, and color and spectral image processing. She is involved in teaching several subjects in the BS courses of Physics, Optics, and Optometry and the CIMET/COSI international master program.

**Juan L. Nieves** received M.S. and Ph.D. degrees in physics from the University of Granada, Granada, Spain, in 1991 and 1996, respectively. He is currently a Full Professor with the Department of Optics, Science Faculty at the University of Granada, where he conducts research in the Color Imaging Laboratory. His current research interests include computational color vision (color constancy, human visual system processing of spatiochromatic information) and spectral analysis of color images. Dr Nieves is the President of the Spanish Color Committee and a representative of this Committee in the International Color Association (AIC) and is currently the coordinator of the Erasmus I Joint Master Degree "Color in Science and Industry (COSI)."

**How to cite this article:** Martínez MÁ, Etchebehere S, Valero EM, Nieves JL. Improving unsupervised saliency detection by migrating from RGB to multispectral images. *Color Res Appl*. 2019; 1-11. <https://doi.org/10.1002/col.22421>