

# A feature-based approach for saliency estimation of omni-directional images

Federica Battisti, Sara Baldoni, Michele Brizzi, Marco Carli

*Università degli Studi Roma TRE  
Department of Engineering  
via Vito Volterra, 62  
00146 Rome, Italy*

---

## Abstract

Omni-directional imaging records the visual information from any direction with respect to a given view-point. It is gaining consumers' popularity due to fast spreading of low-cost devices both for acquisition and rendering. The possibility to render the whole surrounding space represents a further step towards immersivity, thus providing the user with the illusion of physically being in a virtual environment. The understanding of visual attention mechanisms for these images is a relevant topic for processing, coding, and exploiting such data. In this contribution, a saliency model for omni-directional images is presented. It is based on the combination of low-level and semantic features. The first ones account for texture, viewport saliency, hue and saturation, while the second are used to take into account the impact of the presence of human subjects on the saliency. The proposed model has been tested in the "Salient360! Visual attention modeling for 360° Images" Grand Challenge. The model, the achieved results, and finding/discussions are here presented.

*Keywords:* Saliency estimation, Human fixation, Omni-directional images, 360° images

---

## 1. Introduction

Immersivity strongly depends on successfully fooling several senses, principally sight and hearing. Thanks to this process, it provides the viewers with the feeling of physically being in the place shown by the rendering system. It

5 is generally achieved through the implementation of virtual environments, in  
6 which the viewer perceives him/herself as being surrounded by a 3D world.  
7 The surrounding world can be computer-generated or can be the rendering  
8 of real scenes. In the latter case, specific acquisition devices have to be used.  
9 The goal is to obtain an omni-directional image, or video, that allows the  
10 visual information to be seen from any direction with respect to a given  
11 view-point.

12 The imaging system may exploit mechanical or optical devices. In the  
13 first case, motorized linear or array-based cameras scan the scene resulting in  
14 very high-resolution images. In this case, the drawback is time consumption  
15 that can be very long in case of high quality scan. When the application  
16 scenario does not require very high definition, or when a real time constraint  
17 is present, optical solutions are employed. Basically, those are based on the  
18 use of mirrors or of special lenses (i.e., fish eye). Nowadays low-cost devices  
19 are available for acquiring omni-directional images, i.e. 360° cameras, thus  
20 pushing this technology and its application to the consumer market.

21 The acquired information can be rendered through 2D display, cave, or  
22 Head-Mounted Display (HMD). In particular HMD enables egocentric scene  
23 viewing, allowing the user to modify the point of view by simply moving  
24 his/her head or body, thus increasing the perceived quality of experience. In  
25 fact, the user is assumed to be placed in the center of a sphere and by moving  
26 the head, he can observe the omni-directional stimula.

27 The increasing use of this technology opens several issues such as the  
28 design of new rendering systems, new applications for exploiting the infor-  
29 mation, or new compression systems. One of the first aspects that needs  
30 to be investigated is the understanding of the modalities in which a human  
31 subject explores the omni-directional image, thus defining the salient points  
32 in the image.

33 More generally, saliency estimation refers to the localization of the areas  
34 in an image having particular clue for a human observer. This information  
35 is generally obtained by exploiting fixation points, that are the points in  
36 the visual field that are fixated by the two eyes in normal vision, and for  
37 each eye those are the points directly stimulating the fovea [1, 2]. They  
38 can be captured by means of eye-trackers, or cameras. The clustering of  
39 the fixation points, usually obtained by convolving the fixation points map  
40 with a Gaussian kernel, is used to produce the saliency map. The obtained  
41 map represents the degree of interest of an observer and can be used in  
42 many applications such as quality assessment [3], video surveillance [4], tone

43 mapping [5] or defect detection. A classification of possible applications is  
44 presented in [6].

45 In literature, many efforts have been devoted to model the saliency of  
46 images and videos [7, 8]. Proposed methods can be categorized according  
47 to the features they rely on: bottom-up approaches are based on low-level  
48 local features, like color, intensity, contrast, or orientation [9], while top-  
49 down approaches exploit high-level cues like context, semantic, knowledge,  
50 expectations or application [10].

51 However, to the best of our knowledge, very few methods are specific to  
52 omni-directional images. An application-based approach is proposed in [11]  
53 where a method for extracting visual attention-based features from panoramic  
54 (cylindrical) images is used for robot localization. A similar application-based  
55 approach is proposed in [12]. In [13] an algorithm exploiting spherical ge-  
56 ometry for reducing the geometrical distortions, that may be introduced by  
57 the plane mapping, is proposed. This method is based on the use of low-  
58 level features as proposed in [14]. In [15], the authors predict salient image  
59 regions by taking into account the image exposition time. This allows to  
60 understand the influence of this parameter in the resulting saliency map.  
61 Performed tests show that duplicating the exposition time does not modify  
62 significantly the saliency map. In [16], the authors exploit eye-tracking in a  
63 HMD system for gaze analysis. Results suggest that most eye-gaze fixations  
64 are rather far away from the center of the viewport. In [17] a method for  
65 estimating salient objects in panoramic images is presented. A draft of the  
66 saliency map is computed by background estimation and then is refined by  
67 computing the contrast only in the surrounding regions. In [15], 2D image  
68 features are estimated on a lattice of viewports and the overall saliency map  
69 is computed by considering the contribute of each viewport.

70 In this paper a novel model for saliency estimation in omni-directional  
71 images is proposed. In more details, viewports are first collected from the  
72 omni-directional content and, then, the visual attention is estimated by an-  
73 alyzing low-level and semantic features extracted from each viewport. The  
74 strategy of analyzing the viewports instead of the whole panoramic image  
75 relies on the fact that, in the exploration of the omni-directional content, the  
76 user watches only one portion of it at a time [18].

77 The rest of the paper is organized as follows: Section 2 details the char-  
78 acteristics of the proposed saliency model, Section 3 includes the system pa-  
79 rameters, the details of the adopted image database and a discussion on the  
80 characteristics of the proposed model. Finally, in Section 4 the conclusions

81 are drawn.

## 82 2. Proposed Method

83 The proposed saliency model is shown in Figure 1. Each input image  
 84 (i.e., the equi-rectangular image) undergoes a pre-processing step in which  
 85 the viewport extraction is performed. Then, high-level (i.e., skin color, faces,  
 86 and number of people) and low-level features (i.e., hue, saturation, intensity,  
 87 and contrast) are extracted for each viewport and averaged to obtain a first  
 88 estimation of the saliency map. This map is then refined by using an equator-  
 89 prior weighting and a smoothing operation.

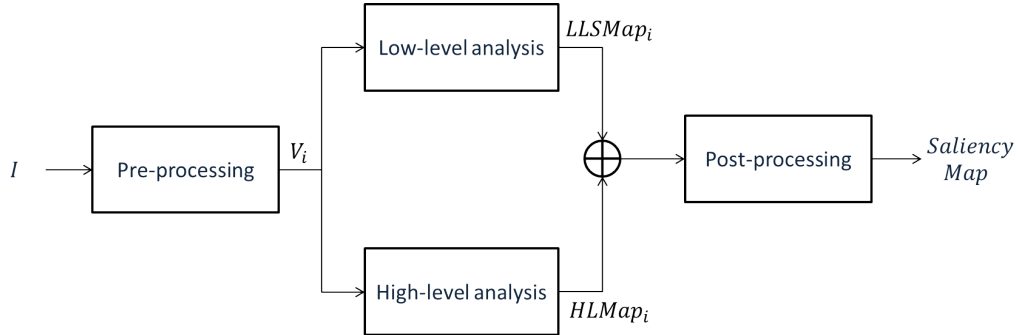


Figure 1: Proposed saliency model

90 In the following, each step of the algorithm will be detailed.

### 91 2.1. Pre-processing

92 The viewports  $V_i$ , with  $i = 1, \dots, n$ , are extracted from the input equi-  
 93 rectangular image  $I$  according to the procedure described in [19], here re-  
 94 ported for sake of clarity.

95 A non-uniform angular sampling of the sphere is performed, with  $\Delta\phi$  and  
 96  $\Delta\theta$  being the horizontal and vertical angular sampling rates and  $X_i(\phi, \theta)$  the  
 97  $i^{th}$  sampling point. Since we assume that the user will change location in the  
 98 omni-directional content by moving his head, the center of the viewport  $V_i$   
 99 will correspond to the coordinates of  $X_i$ .

100 To extract each viewport, the position of each pixel of  $V_i$  is back-projected  
 101 into the spherical reference, and then into the equi-rectangular frame. These  
 102 coordinates are used to interpolate over  $I$ .

103 Let  $(x, y)$  be the coordinates of any point  $M_v$  in the viewport  $V_i$ , whose  
 104 size is  $[V_{width}, V_{height}]$ . To represent the inverse gnomonic projection of  $V_i$   
 105 on the sphere, we define a three dimensional Cartesian coordinate system,  
 106 whose origin is surrounded by the spherical frame of unitary radius, and  
 107 place the viewport  $V_i$  on the plane tangent to the sampling point  $X_i$  (as  
 108 shown in Figure 2). Let us consider the case in which the sampling point  $X_i$   
 109 corresponds to the center of the equi-rectangular image. Then, the position  
 110 of  $M_v$  on the aforementioned plane is given by:

$$M_p(x, y, z) = \begin{bmatrix} 1 \\ pxl \cdot (x - \frac{V_{width}}{2}) \\ pxl \cdot (y - \frac{V_{height}}{2}) \end{bmatrix}$$

111 where  $pxl$  is the size of a pixel in  $V_i$ , obtained as:

$$pxl = 2 \frac{\tan(\frac{a}{2} \cdot \frac{\pi}{180})}{V_{width}}$$

112 where  $a$  is the size of the viewport in degrees.

113 The projection of  $M_p$  on the sphere,  $M_s$ , is:

$$M_s(x, y, x) = \frac{M_p(x, y, z)}{\|M_p(x, y, z)\|}$$

114 where  $\|M_p(x, y, z)\|$  denotes the  $L^2$  norm of vector  $M_p(x, y, x)$ .

115 In the case of any other sampling point,  $M_s$  needs to be multiplied by the  
 116 rotation matrix  $R_{\theta, \phi}$ :

$$R_{\theta, \phi} = \begin{pmatrix} \cos(\theta) \cos(\phi) & -\sin(\phi) & \cos(\theta) \sin(\phi) \\ \sin(\theta) \cos(\phi) & \cos(\theta) & \sin(\theta) \sin(\phi) \\ -\sin(\phi) & 0 & \cos(\phi) \end{pmatrix}$$

117 where  $\phi$  and  $\theta$  are the azimuth and elevation angles of the sampling point  $X_i$   
 118 on the sphere (Figure 2).

119 To obtain the corresponding coordinates in the equi-rectangular image,  
 120  $M_e$ , the following relation holds:

$$M_e(x, y) = \begin{bmatrix} E_{width} \cdot \left(\frac{ang}{2\pi}\right) \\ E_{height} \cdot \left(\frac{\arcsin(M_s(z))}{\pi+0.5}\right) \end{bmatrix}$$

121 where  $[E_{width}, E_{height}]$  are the sizes of the equi-rectangular image and  $ang$  is  
 122 given by:

$$ang = \tan^{-1}(M_s(y), M_s(x))$$

123 where the four-quadrant inverse tangent function is used.

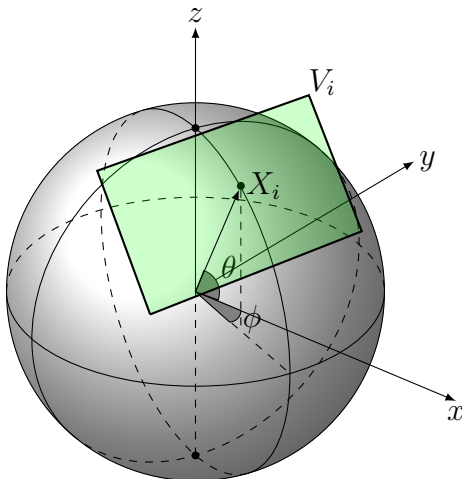


Figure 2: Viewports extraction for an arbitrary sampling point  $X_i$

124 *2.2. Low-level analysis*

125 In the low-level analysis (see Figure 3), each viewport  $V_i$  is converted from  
 126 the Red, Green, and Blue (RGB) color space to the Hue, Saturation, and  
 127 Value (HSV) one. The latter is characterized by better visual consistency  
 128 than the RGB one as suggested in [20]. For each viewport  $V_i$ , only the Hue  
 129 ( $H_i$ ) and Saturation ( $S_i$ ) components are taken into account. In more details,  
 130 for each  $V_i$ , a first map ( $LLSMap_i$ ) is obtained as combination of the result  
 131 of the texture analysis performed on  $V_i$  with the weighted sum of  $H_i$ ,  $S_i$   
 132 and the outcome of the Graph-Based Visual Saliency (GBVS) analysis [21]  
 133 of the  $H_i$  component. This model allows to estimate the human fixations  
 134 based on the creation of activation maps on specific feature channels that are  
 135 normalized to enhance the importance of the points attracting the human  
 136 attention. Moreover, it has been proven that GBVS supports a center bias,  
 137 by assigning higher saliency values in the center of the image plane. Based  
 138 on this,  $LLSMap_i$  is computed as:

$$LLSMap_i = T_i \cdot W_i \tag{1}$$

139 where:

- 140 •  $T_i$  is a binary texture map extracted from  $V_i$  by using the multi-channel
- 141 filtering approach described in [22];
- 142 •  $W_i$  is obtained as:

$$W_i = \alpha S_i + \beta H_i + \gamma G_i \quad (2)$$

143 where  $\alpha$ ,  $\beta$  and  $\gamma$  are the coefficients of the weighted sum while  $G_i$  is

144 the output of the GBVS procedure.

145 Finally, the overall low-level saliency map,  $LLSMa_{p_{tot}}$ , is obtained by

146 equi-rectangular projection of each  $LLSMa_{p_i}$ .

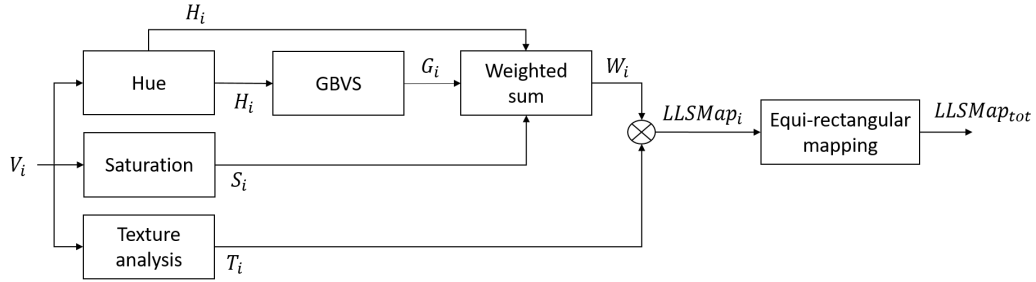


Figure 3: Low-level features analysis

### 147 2.3. High-level analysis

148 The steps performed in the high-level analysis are detailed in Figure 4.

149 The input viewport  $V_i$  undergoes two parallel processing: skin and face de-

150 tection. The former is performed based on the methods presented in [23].

151 Based on this approach, only the portions of  $V_i$  whose color components are

152 inside a pre-defined range are extracted to obtain the map  $P_i$ . Face detection

153 is achieved by using the Viola-Jones algorithm [24] to retrieve the map  $F_i$ .

154 It is useful to underline that this step is performed using not only the frontal

155 face classification model but also the upper body and profile face classifica-

156 tion models. The former detects head and shoulders areas while the latter

157 detects upright face profiles. The maps  $P_i$  and  $F_i$  are computed for each

158 viewport and they are combined through an equi-rectangular mapping to

159 respectively obtain the maps  $P_{tot}$  and  $F_{tot}$ . In order to more accurately iden-

160 tify the presence of human subjects, a weighted combination is performed to

161 obtain the fusion map  $PF_{tot}$  as:

$$PF_{tot} = \frac{2 \cdot F_{tot} + P_{tot}}{3}. \quad (3)$$

162 After the regions containing persons have been identified, a people count  
 163 is performed in order to estimate the number of persons ( $nP$ ) identified in  
 164  $PF_{tot}$ . This value, obtained through a blob analyzer, is then used to define  
 165 the weight  $w_{people}$ . Finally, the output  $HLMa_{p_{tot}}$  is obtained as:

$$HLMa_{p_{tot}} = w_{people} \cdot PF_{tot}; \quad (4)$$

166 In this work, a more relevant weight will be given to regions containing a  
 167 limited number of subjects rather than a large one. In fact, in the latter case  
 168 the human subjects in the scene are hardly distinguishable, appearing as a  
 169 texture and therefore reducing the impact of that region on the saliency.

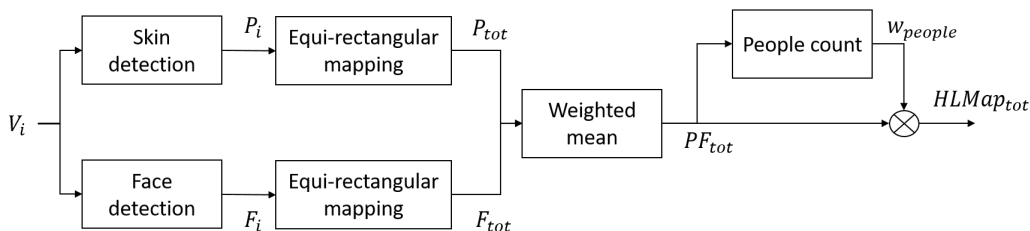


Figure 4: High-level features analysis

#### 170 2.4. Post-processing

171 After the computation of  $LLSMap_i$  and  $HLMa_{p_i}$ , they are averaged in  
 172 order to account equally for low-level and high-level features. The obtained  
 173 fusion map undergoes a post-processing step before returning the overall  
 174 Saliency Map.

175 Several studies [25, 26] show that fixations distribution tends to be strongly  
 176 biased towards the center of the screen when viewing 2D scenes on computer  
 177 monitors, to the point that a saliency map composed of a centered Gaussian  
 178 blob has good performances in predicting fixations [27]. This appears to be  
 179 independent from the distribution of the features in the images [28]. There  
 180 are different explanations for this behavior: first, objects of interest are often  
 181 placed by photographers in the center of photographs by exploiting the rule  
 182 of thirds; second, fixations might be influenced by the setup used to experi-  
 183 mentally record eye-tracking data, where users are usually placed in front of



184 the screen [29].  
 185 In the case of omni-directional images, this assumption does not hold com-  
 186 pletely, since users can explore the whole content by freely moving eyes and  
 187 head. However, even in this case, a bias towards the central area (i.e., the  
 188 equatorial area) of the omni-directional image, holds. This bias can be due  
 189 to the human posture and to the fact that moving the head for looking at a  
 190 different direction requires more intense movements with respect to the ones  
 191 required by the eyes [30, 31]. Therefore it is more likely for a subject to first  
 192 span the visible area with the eyes and then with the head, thus confirming  
 193 the results in [16]. For this reason, in the proposed method, a weighting win-  
 194 drow is applied to the estimated saliency map in the equi-rectangular format.  
 195 As can be noticed in Figure 5, the applied cost function is increasing with  
 196 the distance with respect to the equatorial line. The cost values are in the  
 197 range 1 (in the central region) to 1/4 (in the border regions).

198 Finally a low-pass filtering and a normalization step are performed for  
 199 obtaining the smoothed final saliency map.

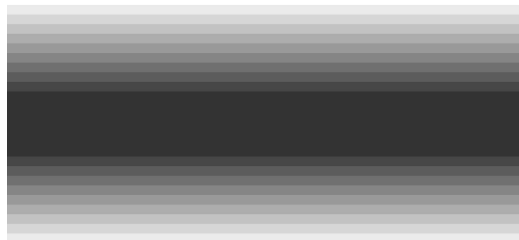


Figure 5: Weighting window

### 200 3. Evaluation tests

201 In the performed tests, the training set has been used for setting up the  
 202 system parameters while the validation set has been exploited for testing the  
 203 performances of the proposed approach. In the following, Subsection 3.1 de-  
 204 tails the procedures used for selecting the system parameters and the adopted  
 205 values, Section 3.2 describes the adopted database, and Section 3.3 presents  
 206 the performed tests and the obtained results.

#### 207 3.1. System parameters

- 208 • **Viewport extraction:** in order to perform the extraction of viewports  
 209 we used a horizontal sampling rate  $\Delta\phi = 40^\circ$  and a vertical sampling

210 rate  $\Delta\theta = 35^\circ$ . Knowing that the HMD used for the test dataset has a  
 211 resolution of 960x1080 pixels per eye and a total Field-Of-View (FOV)  
 212 of  $100^\circ$ , we extracted 1920x1080 pixels viewports in order to provide  
 213 the same FOV and set the value of the size of the viewport in degrees,  
 214  $a$ , accordingly.

215 • **Weighting constants:** the values of  $\alpha$ ,  $\beta$  and  $\gamma$  have been empir-  
 216 ically chosen for maximizing the correlation between estimated and  
 217 ground truth saliency maps in the training dataset. In more details,  
 218 the adopted normalization values are:  $\alpha = 1$ ,  $\beta = 0.1$  and  $\gamma = 0.1$ .  
 219  $LLSMap_i$  and  $HLMMap_i$  are averaged in order to equally account for  
 220 low-level and high-level features.

221 The parameter  $w_{people}$  is set according to:

$$w_{people} = \begin{cases} 0.3 & \text{if } nP \geq 10 \\ 0.6 & \text{if } 5 \leq nP < 10 \\ 0.9 & \text{if } nP < 5. \end{cases} \quad (5)$$

222 During the face detection step, the area of the smallest detectable ob-  
 223 ject is set to 100x100 pixels and during the people count step, the  
 224 minimum blob area detectable by the blob analyzer is set to 6000 pix-  
 225 els.

226 • **Gabor filters:** the adopted approach [22] is based on the use of a  
 227 bank of Gabor filters in order to almost uniformly cover the spatial-  
 228 frequency domain. In this work four orientation values were adopted  
 229  $[0^\circ, 45^\circ, 90^\circ, 135^\circ]$  and increasing values of radial frequency, raising  
 230 with step 1 octave from  $\sqrt{2}$  to the hypotenuse length of the input  
 231 image, have been taken into account.

### 232 3.2. Image database

233 The test dataset [32] is composed by three sub datasets (one per model  
 234 type: head only, head+eye, scanpath). In this work, only the head motion-  
 235 based saliency model is considered, and the relevant dataset is composed by  
 236 20 training images (with size from 4000x2000 to 16000x8000 pixels) and 25  
 237 evaluation images (with size from 3000x1500 to 12000x6000 pixels), with the  
 238 corresponding ground truth. The database includes several subjects such  
 239 as vehicles (i.e. cars, public transport), urban environment (i.e. squares,

240 buildings, supermarkets, museums theaters, hotels), landscapes, animals, and  
 241 people. Moreover, for the provided images, different lighting conditions can  
 242 be found.

### 243 3.3. Experimental results

244 In order to assess the performances of the proposed method (indicated  
 245 in the following as RM3), the tools provided by the “Salient360! Visual  
 246 attention modeling for 360° Images” Grand Challenge, detailed in [32, 33],  
 247 have been used to compare the estimated saliency map with the available  
 248 ground truth. For evaluating the effectiveness of the proposed method with  
 249 respect to other methods, we computed the Correlation Coefficient (CC) and  
 250 KL Divergence (KLD) between the estimated saliency map and the ground  
 251 truth saliency map of the images in the validation set. In Figures 6-7 the  
 252 estimated saliency maps giving the best and the worst results are shown  
 253 together with the original image and the corresponding ground truth. The  
 254 corresponding CC and KLD values, compared with the ones obtained by the  
 255 best performing algorithms in the challenge, Wuhan University (WU) [34]  
 256 and Zhejiang University (ZU) [35], are reported in Tables 1-2, respectively.  
 257 As can be noticed from Table 1, the CC or KLD values are comparable with  
 258 the best performing ones. In some cases the CC value is closer to the best  
 259 one while in others the KLD value is closer. This difference is due to the  
 260 statistical measures considered by the two similarity functions.

	RM3		WU [34]		ZU [35]	
Image	CC	KLD	CC	KLD	CC	KLD
P19	0.63	0.57	0.68	0.46	0.75	0.51
P69	0.68	0.71	0.65	0.90	0.73	0.38
P73	0.63	0.51	0.84	0.20	0.86	0.19
P74	0.67	0.78	0.78	0.59	0.73	0.52
P79	0.65	0.67	0.78	0.43	0.72	0.43

Table 1: CC and KLD scores for the proposed method in the best performing cases compared with the benchmarks.

261 In order to carry out a more complete analysis of the proposed method,  
 262 we evaluated its performances in terms of strengths and weaknesses on the  
 263 entire dataset. In fact, by considering not only the validation but also the  
 264 training set, it is possible to have a better understanding of the behavior of  
 265 the proposed algorithm for different types of content.

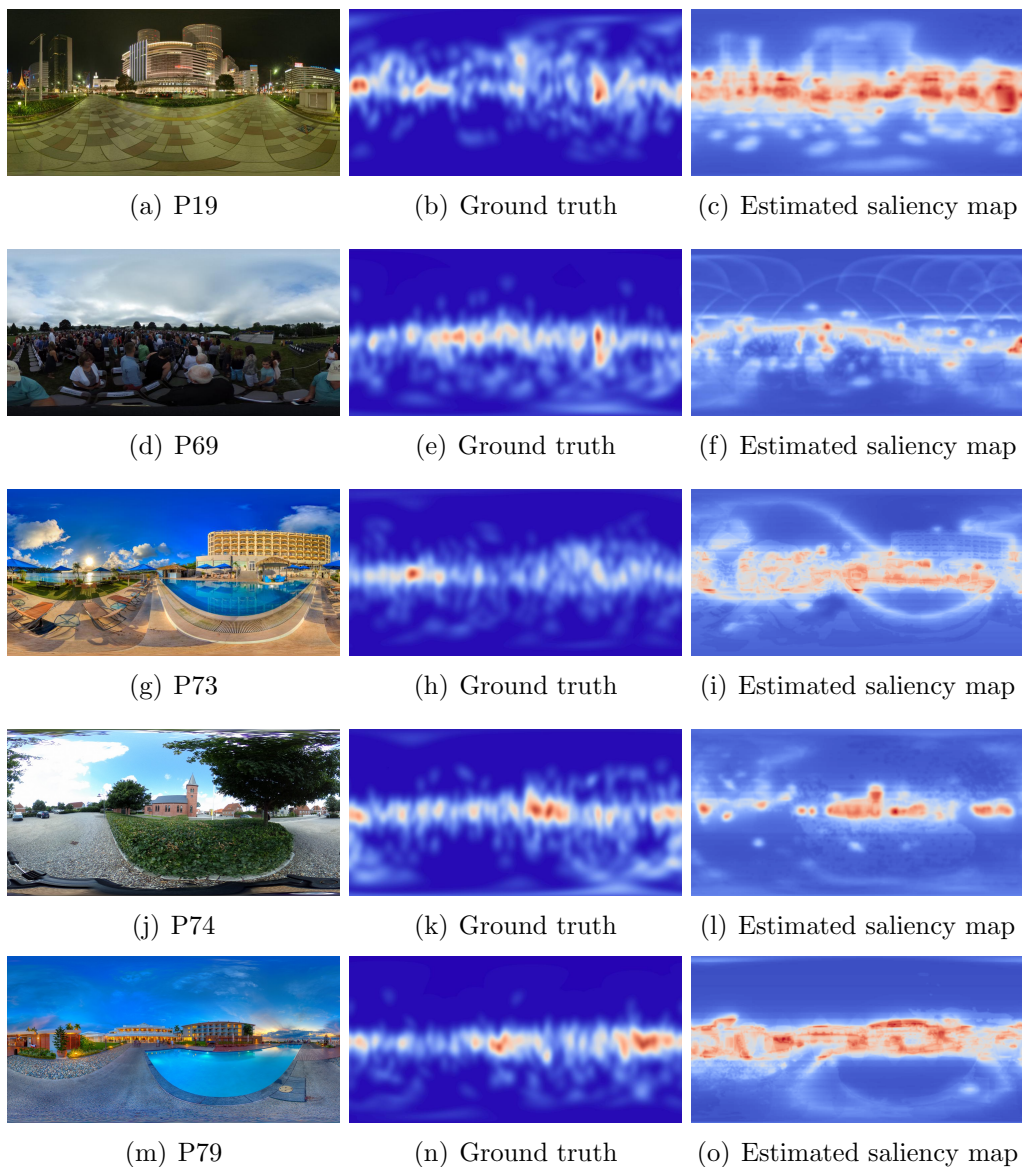


Figure 6: Left column: best performing images in the entire dataset. Center column: saliency map ground truth. Right column: saliency map estimated by the proposed algorithm

266 *3.3.1. Model strengths*

267 From the analysis performed on the images in the dataset, it is possible to  
 268 highlight some characteristics of the algorithm that allowed to obtain results

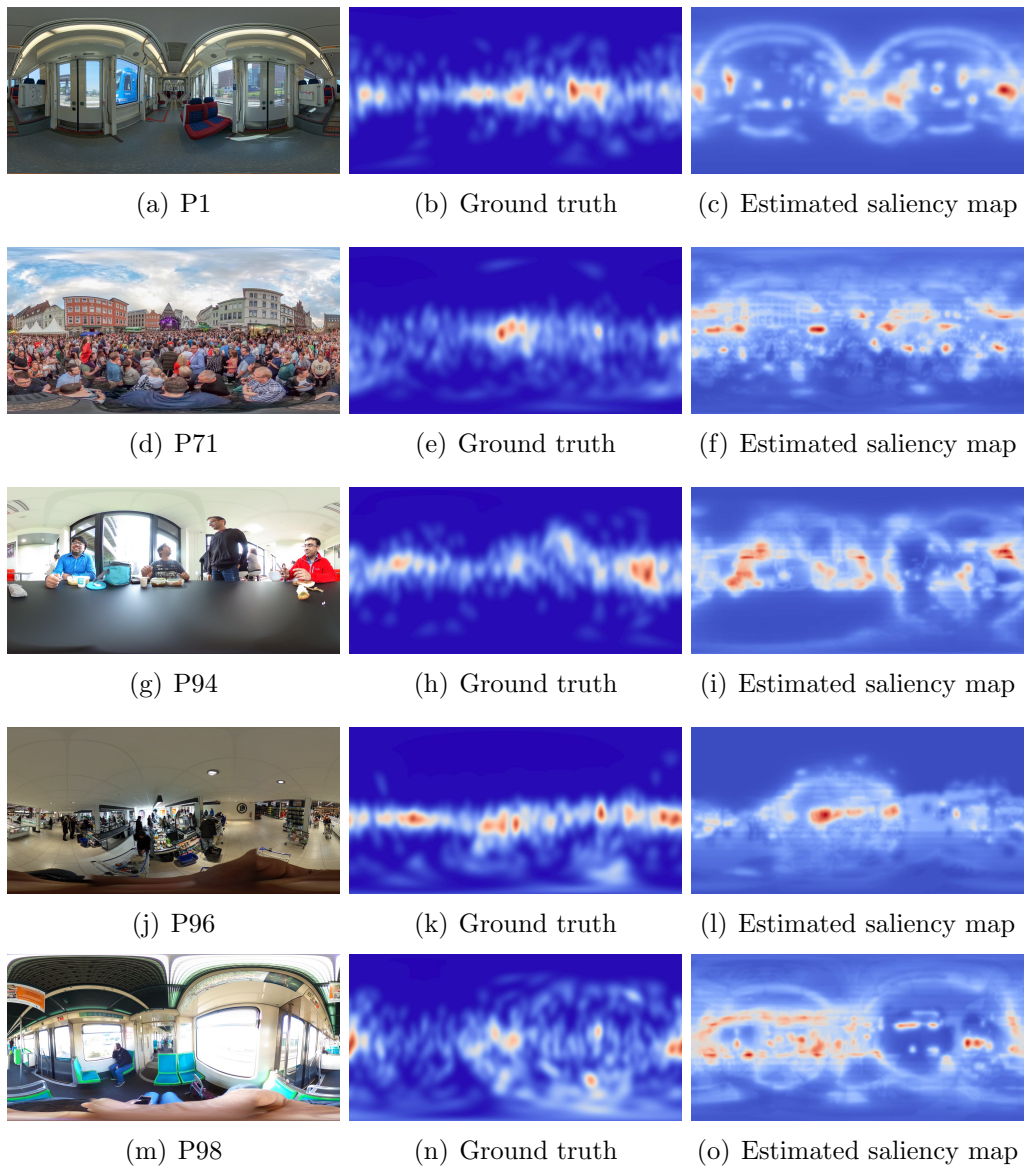


Figure 7: Left column: worst performing images in the entire dataset. Center column: saliency map ground truth. Right column: saliency map estimated by the proposed algorithm

269 similar to the best performing algorithms on the dataset:

- 270 • Use of equator-prior based weighting: the advantage of including this
- 271 step in the proposed algorithm can be noticed for stimula P24 and P73,

Image	RM3		WU [34]		ZU [35]	
	CC	KLD	CC	KLD	CC	KLD
P1	0.49	0.95	0.71	0.58	0.65	0.52
P71	0.16	1.9	0.63	0.81	0.56	0.47
P94	0.41	1.03	0.76	0.45	0.79	0.34
P96	0.52	1.24	0.68	0.87	0.80	0.40
P98	0.30	0.75	0.41	0.56	0.36	0.77

Table 2: CC and KLD scores for the proposed method in the worst performing cases compared with the benchmarks.

272 in Figure 8(b) and (f) respectively. In these images, the horizon line  
 273 corresponds to the equatorial line and thus the proposed model can  
 274 effectively estimate the saliency of the scenes;

- 275 • Good performances in poor lighting conditions: this is achieved by  
 276 including in the salient regions only areas containing mid-level values  
 277 of saturation (S in the range [0.36-0.7]). An example is for stimula  
 278 P10 and P69, respectively shown in Figure 8(a) and (d). As can be  
 279 noticed, the method is able to detect saliency even in darker areas  
 280 without overlaying regions belonging to the border between differently  
 281 illuminated areas;
- 282 • High-level feature extraction: this step increases the weight in the  
 283 saliency map estimation of the areas containing human subjects. This  
 284 can be noted for stimulus P28 (Figure 8(c)) where the people around  
 285 the table are successfully recognized.

### 286 3.3.2. Model weaknesses

287 The analysis of the performances of the proposed algorithm on the dataset,  
 288 allowed to highlight some weak points that need a further improvement of  
 289 the system:

- 290 • Indoor scenes with presence of distributed light sources: this is the case  
 291 of Figures 8(a) and (f). A reason for such poor behavior might be the  
 292 fact that, after conversion in the HSV color space, the V component is  
 293 discarded. This operation might hinder the right handling of this type  
 294 of light source;

295 • Difficulty in people detection in crowded environments: this problem  
 296 could be solved by changing the adopted technique for discriminating  
 297 the presence of skin. An example of this problem is in Figure 8(e), in  
 298 which the presence of the crowd hinders the face detection algorithm. In  
 299 this case, the proposed algorithm detects the presence of faces, however  
 300 giving too importance to these areas. A more accurate skin selection  
 301 procedure could reduce the number of false alarm.



Figure 8: Samples of images in the dataset.

#### 302 4. Conclusions and comments

303 In this contribution, a method for visual saliency estimation for omni-  
 304 directional images is presented. The proposed method relies on the extraction  
 305 of a set of viewports from the equi-rectangular image. From each of them, a  
 306 first estimation of the saliency map is obtained by analyzing the high-level  
 307 (i.e., skin color, faces, and number of people) and low-level features (i.e.,  
 308 hue, saturation, intensity, and contrast). These intermediate maps are then  
 309 fused according to pre-defined weighting coefficients and finally refined by  
 310 using an equator-prior weighting and a smoothing operation. The proposed  
 311 model has been tested in the “Salient360! Visual attention modeling for  
 312 360° Images” Grand Challenge achieving good results especially in images  
 313 containing human subjects and in case of poor light conditions. Future work  
 314 will be devoted to a deeper investigation of the impact of image characteristics

315 (e.g., contrast, illumination) on the performances of the system thanks to the  
316 use of larger databases.

## 317 **References**

- 318 [1] I. Murakami, Fixational eye movements and motion perception,  
319 *Progress in Brain Research* 154 (2006) 193 – 209. *Visual Perception*.
- 320 [2] S. Lee, A. Bovik, Fast algorithms for foveated video processing, *IEEE*  
321 *Transactions on Circuits and Systems for Video Technology* 13 (2003)  
322 149–162.
- 323 [3] Q. Huynh-Thu, M. Barkowsky, P. L. Callet, The importance of visual  
324 attention in improving the 3D-TV viewing experience: Overview and  
325 new perspectives, *IEEE Transactions on Broadcasting* 57 (2011) 421–  
326 431.
- 327 [4] Y. Tong, F. A. Cheikh, F. F. E. Guraya, H. Konik, A. Trémeau, A  
328 Spatiotemporal Saliency Model for Video Surveillance, *Cognitive Com-*  
329 *putation* 3 (2011) 241–263.
- 330 [5] M. Narwaria, M. Perreira Da Silva, P. Le Callet, R. Pepion, Tone  
331 mapping based HDR compression: Does it affect visual experience?,  
332 *Image Commun.* 29 (2014) 257–273.
- 333 [6] M. Mancas, O. Le Meur, *Applications of Saliency Models*, Springer New  
334 York, New York, NY, pp. 331–377.
- 335 [7] A. Borji, L. Itti, State-of-the-art in visual attention modeling, *IEEE*  
336 *Transactions on Pattern Analysis and Machine Intelligence* 35 (2013)  
337 185–207.
- 338 [8] Y. Fang, W. Lin, Z. Chen, C. M. Tsai, C. W. Lin, A video saliency  
339 detection model in compressed domain, *IEEE Transactions on Circuits*  
340 *and Systems for Video Technology* 24 (2014) 27–38.
- 341 [9] L. Itti, C. Koch, E. Niebur, A model of saliency-based visual attention  
342 for rapid scene analysis, *IEEE Transactions on Pattern Analysis and*  
343 *Machine Intelligence* 20 (1998) 1254–1259.



- 344 [10] A. Oliva, A. Torralba, M. Castelhana, J. M. Henderson, Top-down  
345 control of visual attention in object detection, in: IEEE International  
346 Conference on Image Processing, volume 1, pp. 253–256.
- 347 [11] A. Bur, A. Tapus, N. Ouerhani, R. Siegwar, H. Hiigli, Robot navi-  
348 gation by panoramic vision and attention guided features, in: 18th  
349 International Conference on Pattern Recognition (ICPR’06), volume 1,  
350 pp. 695–698.
- 351 [12] G. Schillaci, S. Bodiroža, V. V. Hafner, Evaluating the effect of saliency  
352 detection and attention manipulation in human-robot interaction, In-  
353 ternational Journal of Social Robotics 5 (2013) 139–152.
- 354 [13] I. Bogdanova, A. Bur, H. Hugli, Visual attention on the sphere, IEEE  
355 Transactions on Image Processing 17 (2008) 2000–2014.
- 356 [14] C. Koch, S. Ullman, Shifts in selective visual attention: Towards the  
357 underlying neural circuitry, Human Neurobiology 4 (1985) 219–227.
- 358 [15] A. De Abreu, C. Ozcinar, A. Smolic, Look around you: Saliency maps  
359 for omnidirectional images in VR applications, in: 2017 Ninth Interna-  
360 tional Conference on Quality of Multimedia Experience (QoMEX), pp.  
361 1–6.
- 362 [16] Y. Rai, P. L. Callet, P. Guillotel, Which saliency weighting for omni-  
363 directional image quality assessment?, in: 2017 Ninth International  
364 Conference on Quality of Multimedia Experience (QoMEX), pp. 1–6.
- 365 [17] B.-J. Han, J.-Y. Sim, Saliency detection for panoramic landscape im-  
366 ages of outdoor scenes, Journal of Visual Communication and Image  
367 Representation 49 (2017) 27 – 37.
- 368 [18] F. Jabar, J. Ascenso, M. P. Queluz, Perceptual analysis of perspec-  
369 tive projection for viewport rendering in 360 images, in: 2017 IEEE  
370 International Symposium on Multimedia (ISM), pp. 53–60.
- 371 [19] FOXEL, Libgnomonic, <https://github.com/FoxelSA/libgnomonic>,  
372 2017.
- 373 [20] C. Huang, Q. Liu, S. Yu, Regions of interest extraction from color image  
374 based on visual saliency, The Journal of Supercomputing 58 (2011) 20–  
375 33.

- 376 [21] J. Harel, C. Koch, P. Perona, Graph-based visual saliency, in: Pro-  
377 ceedings of the 19th International Conference on Neural Information  
378 Processing Systems, NIPS'06, MIT Press, Cambridge, MA, USA, 2006,  
379 pp. 545–552.
- 380 [22] A. K. Jain, F. Farrokhnia, Unsupervised texture segmentation using  
381 Gabor filters, in: 1990 IEEE International Conference on Systems, Man,  
382 and Cybernetics Conference Proceedings, pp. 14–19.
- 383 [23] S. Thakur, S. Paul, A. Mondal, S. Das, A. Abraham, Face detection  
384 using skin tone segmentation, in: 2011 World Congress on Information  
385 and Communication Technologies, pp. 53–60.
- 386 [24] P. Viola, M. J. Jones, Robust real-time face detection, *International*  
387 *Journal of Computer Vision* 57 (2004) 137–154.
- 388 [25] B. W. Tatler, R. J. Baddeley, I. D. Gilchrist, Visual correlates of fixation  
389 selection: effects of scale and time, *Vision Research* 45 (2005) 643 – 659.
- 390 [26] T. Judd, K. Ehinger, F. Durand, A. Torralba, Learning to predict where  
391 humans look, in: 2009 IEEE 12th International Conference on Computer  
392 Vision, pp. 2106–2113.
- 393 [27] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, G. W. Cottrell, Sun:  
394 A bayesian framework for saliency using natural statistics, *Journal of*  
395 *Vision* 8 (2008) 32.
- 396 [28] B. W. Tatler, The central fixation bias in scene viewing: Selecting  
397 an optimal viewing position independently of motor biases and image  
398 feature distributions, *Journal of Vision* 7 (2007) 4.
- 399 [29] O. L. Meur, P. L. Callet, D. Barba, Predicting visual fixations on video  
400 based on low-level visual features, *Vision Research* 47 (2007) 2483 –  
401 2498.
- 402 [30] S. M. Safavi, S. M. Sundaram, A. H. Gorji, N. S. Udaiwal, P. H. Chou,  
403 Application of infrared scanning of the neck muscles to control a cur-  
404 sor in human-computer interface, in: 2017 39th Annual International  
405 Conference of the IEEE Engineering in Medicine and Biology Society  
406 (EMBC), pp. 787–790.

- 407 [31] A. Al-Rahayfeh, M. Faezipour, Eye tracking and head movement detec-  
408 tion: A state-of-art survey, *IEEE Journal of Translational Engineering*  
409 *in Health and Medicine* 1 (2013) 2100212–2100212.
- 410 [32] Y. Rai, J. Gutiérrez, P. Le Callet, A dataset of head and eye movements  
411 for 360 degree images, in: *Proceedings of the 8th ACM on Multimedia*  
412 *Systems Conference, MMSys'17*, ACM, New York, NY, USA, 2017, pp.  
413 205–210.
- 414 [33] J. Gutiérrez, E. David, Y. Rai, P. Le Callet, Toolbox and dataset for  
415 the development of saliency and scanpath models for omnidirectional /  
416 360° still images, *Signal Processing: Image Communication* (2018).
- 417 [34] J. Ling, K. Zhang, Y. Zhang, D. Yang, Z. Chen, A saliency predic-  
418 tion model on 360 degree images using color dictionary based sparse  
419 representation, *Signal Processing: Image Communication* (2018).
- 420 [35] P. Lebreton, A. Raake, GBVS360, BMS360, ProSal: Extending existing  
421 saliency prediction models from 2D to omnidirectional images, *Signal*  
422 *Processing: Image Communication* (2018).