# Meaning maps and saliency models based on deep convolutional neural networks are insensitive to image meaning when predicting human fixations

**Marek A. Pedziwiatr[1][*], Matthias Kümmerer[2], Thomas S.A. Wallis[2, 3], Matthias Bethge[2], Christoph Teufel[1]**

[1]Cardiff University, Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology Cardiff, United Kingdom

[2]University of Tübingen, Center for Integrative Neuroscience, Tübingen, Germany

[3]Bernstein Center for Computational Neuroscience, Tübingen, Germany

[*]Corresponding author: marek.pedziwi@gmail.com

## Abstract

Eye movements are vital for human vision, and it is therefore important to understand how observers decide where to look. Meaning maps (MMs), a technique to capture the distribution of semantic importance across an image, have recently been proposed to support the hypothesis that meaning rather than image features guide human gaze. MMs have the potential to be an important tool far beyond eye-movements research. Here, we examine central assumptions underlying MMs. First, we compared the performance of MMs in predicting fixations to saliency models, showing that DeepGaze II – a deep neural network trained to predict fixations based on high-level features rather than meaning – outperforms MMs. Second, we show that whereas human observers respond to changes in meaning induced by manipulating object-context relationships, MMs and DeepGaze II do not. Together, these findings challenge central assumptions underlying the use of MMs to measure the distribution of meaning in images.

Keywords: eye movements, natural scenes, saliency, deep neural networks, meaning maps

## Introduction

Human eyes resolve fine detail only in a small, central part of the visual field, with resolution dropping off rapidly in the periphery. To sample details, we move our eyes to orient the high-resolution part of our visual system successively to different parts of a visual scene. Information about these small scene parts is extracted during fixations – short periods in which the eyes are relatively stable. Thus, due to the structure of our visual system, human vision depends on eye movements. How the brain decides where to look in a visual scene is therefore an important question. A long-standing hypothesis suggests that semantic content of image regions is important in guiding eye movements. Recent work presented meaning maps (MMs) as a tool to test this hypothesis (Henderson & Hayes, 2017, 2018). This technique aims to index the spatial distribution of meaning across an image, which has potential applications far beyond eye-movement research. Here, we assess and challenge central assumptions of this novel tool.

A classic finding in eye-movement research shows that the specific task of an observer has an influence on where they direct their eyes (Yarbus, 1967; Hayhoe & Ballard, 2005). But in everyday life, we frequently move our eyes without any goal other than to explore the environment. In the lab, this behavior is examined in free-viewing paradigms, during which eye movements are recorded while images are viewed without an explicit task (Koehler, Guo, Zhang, & Eckstein, 2014, but see Tatler, Hayhoe, Land, & Ballard, 2011). To explain what guides eye movements during free viewing, two opposing accounts have been put forward.

According to the first account, eye movements are guided primarily by image characteristics (Borji, Sihite, & Itti, 2013; Itti & Koch, 2001; Parkhurst, Law, & Niebur, 2002). Potential support for this view comes from saliency models: algorithms, which exclusively use visual features of an image to predict human fixations. Although early models, which used only simple features such as local intensity or colors (Itti & Koch, 2000), are now deemed only moderately successful (Bylinskii et al., 2014), more recent saliency models achieve a remarkably high performance (Kümmerer, Wallis, Gatys, & Bethge, 2017). These models harness deep convolutional neural networks – biologically inspired machine learning algorithms, that somewhat resemble the human visual system (Kietzmann, McClure, & Kriegeskorte, 2019). However, even such models rely solely on visual features, albeit high-level ones.

62    In contrast to the idea underlying saliency models, several authors have argued that during

63    free viewing, eye movements are mainly guided by the semantic content of the visual scene

64    (Henderson, Malcolm, & Schandl, 2009; Nyström & Holmqvist, 2008; Onat, Açik, Schumann,

65    & König, 2014; Rider, Coutrot, Pellicano, Dakin, & Mareschal, 2018; Stoll, Thrun, Nuthmann,

66    & Einhäuser, 2015). This perspective differs fundamentally from the saliency-based

67    approach. Attributing meaning to certain parts of the scene is impossible without prior

68    knowledge of the world, i.e., a factor that is independent of the visual input (Hegde &

69    Kersten, 2010; Teufel, Dakin, & Fletcher, 2018). Consequently, the notion that semantic

70    content guides eye-movements is inconsistent with the idea that the allocation of fixations

71    is dependent solely on the distribution of image features. Given that meaning is not image-

72    computable, the notion that semantic content guides eye-movements is inconsistent with

73    the idea that the eye-movements are dependent solely on the distribution of image

74    features.

75    A string of recent studies has claimed to provide support for the role of meaning in driving

76    eye movements (Hayes & Henderson, 2019; Henderson & Hayes, 2017, 2018; Henderson,

77    Hayes, Rehrig, & Ferreira, 2018; Peacock, Hayes, & Henderson, 2018). These studies

78    (reviewed in Henderson, Hayes, Peacock, & Rehrig, 2019) are based on a novel technique

79    called meaning maps (MMs). A MM for a given image is created by breaking it down into

80    small isolated patches, which are rated for their meaningfulness independently from the

81    rest of the visual scene. These ratings are pooled together into a smooth map, which is

82    supposed to capture the distribution of meaning across the image. Compared to outputs

83    from a simple saliency model (GBVS, Harel et al., 2006), MMs were more predictive of

84    human fixations. On that basis it has been claimed that meaning guides human fixations in

85    natural scene viewing (Henderson & Hayes, 2017, 2018). Here, we examined central

86    predictions of this claim.

87    First, if MMs measure meaning and if meaning guides human eye-movements, MMs should

88    be better in predicting locations of fixations than saliency models because these models rely

89    solely on image features. Therefore, we compared MMs to a range of classic and state-of-

90    the-art models. We replicate the finding that MMs perform better than some of the most

91    basic saliency models. Contrary to the prediction, however, DeepGaze II (DGII; Kümmerer,

92  Wallis, & Bethge, 2016; Kümmerer et al., 2017), a model based on a deep convolutional

93  neural network, outperforms MMs.

94  A second prediction is that if MMs are sensitive to meaning and if meaning guides human

95  gaze, differences in eye movements that result from changes in meaning should be reflected

96  in equivalent differences in MMs. We probed this prediction experimentally using a well-

97  established effect: the same object, when presented in an atypical context (e.g., a shoe on a

98  bathroom sink) attracts more fixations than when presented in a typical context because of

99  the change in the semantic object-context relationship (Henderson, Weeks, & Hollingworth,

100  1999; Öhlschläger & Võ, 2017). Replicating previous studies, image regions attracted more

101  fixations when they contained context-inconsistent compared to context-consistent objects.

102  Crucially, however, MMs of the modified scenes did not attribute more 'meaning' to these

103  regions. DGII also failed to adjust its predictions accordingly.

104  Together, these findings suggest that semantic information contained in visual scenes is

105  critical for the control of eye movements. However, this information is captured neither by

106  MMs nor DGII. We suggest that similar to saliency models, MMs index the distribution of

107  visual features rather than meaning.

108

## Method

110  We conducted a single experiment in which human observers free-viewed natural scenes

111  while their eye-movements were being recorded. The obtained data was analyzed in two

112  complimentary ways. First, we compared how well MMs and different saliency models

113  predict locations of human fixations in natural scenes. Subsequently, we assessed the

114  sensitivity of MMs and the best-performing saliency model to manipulations of scene

115  meaning. The data, the code to create MMs, and all openly available resources used in the

116  study can be accessed via the links provided in the Supplement.

117

Fig. 1. Illustration of sample stimuli in (a) the Consistent and (b) the Inconsistent condition with the Critical Region outlined in yellow and (c, d) human fixations recorded in both conditions. In this example, a hair brush on a bathroom sink (a) – an object consistent with the scene context – has been exchanged for a shoe (b) to introduce semantic inconsistency.

**Stimuli.** We used images from two conditions of the SCEGRAM database (Öhlschläger & Võ, 2017): the Consistent and the Semantically Inconsistent conditions (called 'Inconsistent' here). In the Consistent condition (used in both analyses), scenes contain only objects that are typical for a given context. In the Inconsistent condition (used only in the second analysis), one of the objects is contextually inconsistent. For example, a hairbrush in the context of a bathroom sink from the Consistent condition is replaced with a flip-flop in the Inconsistent condition (see Figs. 1a and 1b). Such changes in object-context relationship alter the meaning attached to the manipulated object. For every scene, we indexed the location of the consistent and inconsistent objects with the superimposed bounding boxes for both objects (see Figs. 1a and 1b). We refer to this location as the Critical Region, because it is the only part of the image that changes between Consistent and Inconsistent conditions. We used 36 selected scenes in both conditions (72 photographs in total, listed in

135 the Supplement together with the selection criteria). We also replicated the main finding of
136 the first analysis in an additional set of 30, very different, images (reported in the
137 Supplement).

138

139 **Procedure.** The procedure consisted of 3 blocks, interleaved with breaks. Each participant
140 viewed all images from both conditions (Consistent and Inconsistent) and was instructed to
141 'look carefully' at each of them. Experimental blocks began with an eye tracker
142 calibration/validation. Within each block, observers free-viewed a series of 24 photographs
143 from both SCEGRAM conditions, each for 7 seconds. After image offset, observers were
144 required to press a button to view the next image. Then, a fixation point appeared centrally
145 on a screen and once observers fixate on it (as determined online by their eye-trace), the
146 actual image was displayed. Before starting the experiment, observers viewed a sample
147 image in an identical regime to familiarize themselves with the procedure. Each stimulus
148 was shown once and the order of presentation was fully randomized. The stimuli were
149 presented against a uniform grey background and had a width of 688 pixels and a height of
150 524 pixels, which subtended approximately 19.7 and 15 degrees of visual angle,
151 respectively. Our choice of task (free viewing) and stimulus parameters for size and
152 presentation time were adopted from the original study developing the SCEGRAM stimuli
153 (Öhlschläger & Võ, 2017). These design characteristics fall within the typical range used in
154 this literature (e.g. Wilming et al., 2017).

155

156 **Observers.** 20 volunteers (3 male; mean age 19.4) recruited from the Cardiff University
157 undergraduate population took part in the study. All reported normal or corrected-to-
158 normal vision, provided written consent, and received course credits in return for
159 participation. The study was approved by the Cardiff University School of Psychology
160 Research Ethics Committee. The primary units of interest in our analyses were the
161 distributions of fixations over images. The number of observers we recruited guarantees
162 that including more observers would not change these distributions significantly
163 (demonstrated in the Supplement).

164

165 **Apparatus.** The study was conducted in a dimly lit room. SCEGRAM images from both

166 conditions were presented on an LCD monitor (Iiyama ProLite B2280HS, resolution 1920 by

167 1080 pixels, 21 inches diagonal). Chin and forehead rests were used to ensure that

168 observers maintained the constant distance of 49 cm from the screen. Their eye movements

169 were recorded with the frequency of 500 Hz using an EyeLink 1000+ eye tracker placed on a

170 tower mount. The experiment was controlled by custom-written Matlab (R2017a version)

171 scripts using Psychophysics Toolbox Version 3 (Kleiner, Brainard, & Pelli, 2007).

172



173

174 Fig. 2. Illustration of the stimuli and procedure used for creating meaning maps. (**a**) Grids of

175 equally spaced circles were used to cut images into fine and coarse patches (only the latter

176 are illustrated here). The red circle indicates a sample patch in the grid. (**b**) Here, the sample

177 patch is highlighted in one of the scenes from the Consistent condition. (**c**) Patches were

178 presented in isolation and rated for their meaningfulness by three independent observers

179 on a scale from 1 to 6. The panel has illustrative purpose only – the scale presented to

180 observers included additional labels (ranging from 'Very Low' to 'Very High'). (**d**) Illustration

181 of a meaning map with greyscale values indicating 'meaningfulness'. (**e**) Simplifying

182 illustration of how meaning maps are generated from ratings. For simplicity sake, only two

183 patches are shown (step 1). Each patch is rated in isolation (step 2; here only one rating per

184 patch is shown). All pixels within an image area are then assigned average rating values,

185    taking into account all ratings for patches that overlap with this area (step 3). For the area of

186    the original patch (step 4), all pixels are then averaged and the resulting value is assigned to

187    the center of the patch (step 5). Finally, the patch centers were used as interpolation nodes

188    for thin-plate spline interpolation producing a smooth distribution of values over the image

189    (not illustrated). This procedure was conducted separately for the fine and coarse grid, and

190    the meaning map for a given image was created by averaging the two outcomes and

191    normalizing the result to a range between 0 and 1.

192

193    **Creating MMs.** To create MMs for our stimuli, we followed the procedure described by

194    Henderson & Hayes (2017, 2018; for details see Fig. 2). Each image was segmented into

195    partially overlapping patches of two sizes: fine patches had a diameter of 107 pixels (3

196    degrees of the visual angle, or 16 % of the image width), coarse patches of 247 pixels (7

197    degrees or 36% of the image width) (Fig. 2a and b). Their centers were 58 pixels (fine) and

198    97 pixels (coarse) apart from each other.

199    Next, we collected meaningfulness ratings from human subjects for all patches. Each patch

200    was presented in isolation and rated for its meaningfulness on a 6 point Likert scale (Fig. 2).

201    As in Henderson and Hayes (2017), we used a Qualtrics survey completed by naive

202    observers recruited via the crowdsourcing platform Amazon Mechanical Turk (see

203    Supplement for eligibility criteria). Each participant provided ratings for 305 or 303 patches

204    of both sizes (selected randomly from all images), on average spent approximately 14 min

205    on the task, and received 2.18 USD as remuneration. In total, 69 individuals were used as

206    raters, with three individuals rating each individual patch. The collected ratings were then

207    used to create MMs (see Fig. 2).

208    When creating MMs for images from both conditions, we exploited the fact that

209    photographs from the Consistent and Inconsistent conditions differ only in the Critical

210    Region (the part of the image containing the manipulated object) while the remaining parts

211    overlap. We collected meaningfulness ratings for the patches belonging to overlapping

212    areas only once, and the separate sets of ratings for Consistent and Inconsistent condition

213    were collected only for those patches that contained at least one pixel belonging to the

214    Critical Region. In total, the number of patches rated in the study amounted to 7013: 4840

215  fine patches (of which 520 belonged to the images from the Inconsistent condition) and

216  2173 coarse patches (445 Inconsistent).

217

218  **Saliency models.** In the first analysis, we compared predictive performance of MMs to four

219  saliency models of different complexity. The first two models – GBVS (Harel et al., 2006) and

220  AWS (Garcia-Diaz, Fdez-Vidal, Pardo, & Dosil, 2012) – rely on simple visual features, such as

221  local colors and edge orientations, and share the assumption that fixations land on image

222  regions distinct from their surroundings in terms of values of these features. By contrast to

223  GBVS, AWS includes a statistical whitening procedure to improve performance. Both these

224  models were previously used to estimate the influence of image features relative to

225  cognitive factors on the deployment of fixations: GBVS in the previous studies with MMs,

226  AWS elsewhere (Stoll et al., 2015).

227  Two other models that we compared to MMs – ICF and DeepGaze II (DGII) – were designed

228  in a data-driven manner (Kümmerer et al., 2017). Both have the same architecture,

229  consisting of a fixed network that extracts sets of features from images and a readout

230  network that is trained on human fixations to combine the features in a way to maximize

231  the models' predictive power. While the fixed network of ICF extracts only simple visual

232  features (local intensity and contrast), DGII is tuned to features extracted by a deep

233  convolutional neural network pre-trained for object recognition (VGG-19; Simonyan &

234  Zisserman, 2014). The key characteristic of these models that distinguishes them from

235  models such as GBVS and AWS is that they have been trained on human fixations.

236  Specifically, during the training phase, the read-out network receives its respective features

237  as an input, generates a prediction about where human observers will look in the image,

238  and gradually adjusts its parameters based on feedback comparing its prediction to human

239  fixation data to maximise the predictive power of each model. Importantly, the readout

240  network has the same architecture and number of trainable parameters for both DGII and

241  ICF. The only difference between the models is the input features, both of which are not

242  trained on human fixation data.

243  All saliency models output smooth maps that predict the probability of image regions to be

244  fixated. Human observers have the tendency to look at the center of images (Tatler, 2007),

245 and therefore this probability is usually higher in the central region of the image. This

246 'center bias' has important consequences for the evaluation of saliency models. Their

247 performance differs depending on whether they are evaluated using a metric expecting

248 some form of this bias or not (Kümmerer, Wallis, & Bethge, 2018). Here, for the sake of

249 simplicity, we do not incorporate center bias in the models or in the MMs (unlike the

250 original authors) and use an appropriate metric for this situation (see Performance metrics

251 section). Importantly, analyses addressing the issue of center bias in a more extensive way

252 (reported in the Supplement) provide only further support for our conclusions.

253

254 **Data pre-processing.** Fixation locations from the eye tracker recordings were extracted

255 using the algorithm provided by the device manufacturer operating with the default

256 parameter values. Thereby, we obtained a discrete distribution of fixations on each image

257 (see Fig. 1c and 1d). Then, in line with the previous MMs studies, we smoothed these

258 discrete distributions with a Gaussian filter with a cutoff frequency of -6 dB, using the

259 function provided by Bylinskii and colleagues (2014).

260 Next, smooth distributions from fixations, models, and MMs were separately normalized to

261 a range from 0 to 1 for each image. Finally, for each scene, histograms of all distributions

262 from both conditions were matched to histograms of smoothed fixations from Consistent

263 condition using the Matlab imhistmatch function, as in the original MMs studies. Histogram

264 matching makes distributions directly comparable as it ensures that they differ only with

265 respect to their shape, and not their total mass.

266

267 **Performance metrics.** To compare the ability of MMs and models to predict locations of

268 human fixations in Experiment 1, we use two well-established metrics (Bylinskii, Judd, Oliva,

269 Torralba, & Durand, 2016): Correlation and Shuffled Area Under ROC curve (sAUC; Zhang,

270 Marks, Tong, Shan, & Cottrell, 2007) with the implementations provided by Bylinskii and

271 colleagues (2014).

272 Correlation, used in the previous studies on MMs, is calculated as Pearson's linear

273 correlation coefficient between a smoothed distribution of observers' fixations over the
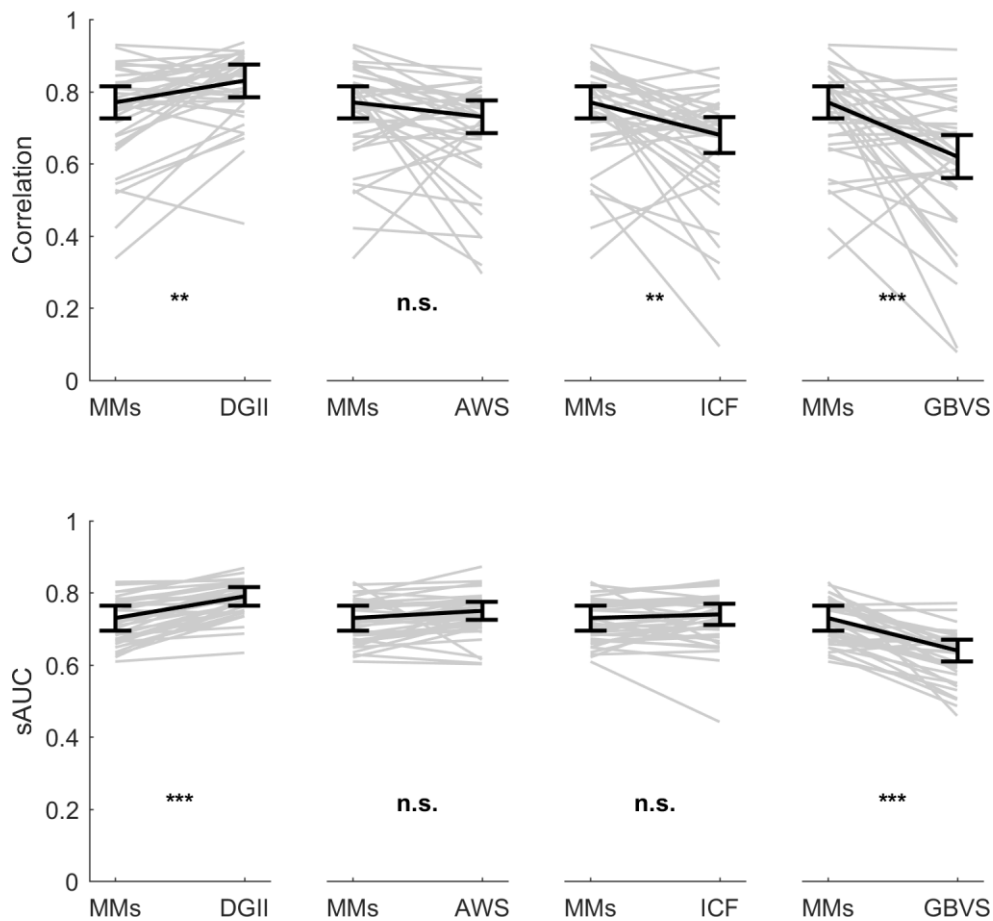
274   image and predictions of a saliency model or MMs. We additionally used sAUC (Zhang et al.,

275   2008), which, unlike Correlation, guarantees that the measured differences in performance

276   between models are driven by their sensitivity to factors guiding fixations, and not by the

277   degree to which they include human center bias in their predictions, even implicitly

278   (Kümmerer, Wallis, & Bethge, 2015; Kümmerer et al., 2018).

279

280                 **Comparing meaning maps and saliency models – results**

281   In the first analysis, we compared performance of four saliency models to MMs in predicting

282   human fixations in the Consistent condition, i.e., when viewing typical scenes with no

283   obvious object-context inconsistencies (Tab. 1, Fig. 3). If human gaze is guided by meaning,

284   and if MMs provide an index for the distribution of meaning, we would expect MMs to

285   outperform all saliency models because these models are based solely on image features.

286   Please note that for the sake of this comparison, we aggregated fixations from all observers

287   for each image and analyzed the data on a per-image basis, similarly to the original MMs

288   studies.

289

Fig. 3. Performance of MMs and saliency models in predicting human fixations according to (a) Correlation and (b) sAUC metrics. Note that according to both metrics DGII predicted human fixations better than MMs. Asterisks indicate p-values from statistical tests comparing MMs to different models (reported in Table 1.): * indicates p ≤ .05, ** p ≤ .01, *** ≤ .001 and 'n.s.' indicates the lack of statistical significance. Grey lines connect values obtained for individual images. Black vertical bars indicate 95% confidence intervals for the medians.

**Predictive power**. Correlation and sAUC values obtained for MMs and for each of the models were compared using Bonferroni-corrected paired Wilcoxon tests (Fig. 3; Tab. 1). We used non-parametric tests because for some of the distributions the assumptions of normality was not met. For the same reason we chose a median as a measure of centrality (we calculate confidence intervals for median using a bootstrapping method – see details in

305   the Supplement). Additionally, we calculated JZS Bayes Factor (Rouder, Speckman, Sun,

306   Morey, & Iverson, 2009) to quantify the evidence for (or against) the differences between

307   models and MMs (Tab. 1). While deviations from normality can be problematic for Bayes

308   factor analyses, they are most likely not an issue in the current situation: the Bayes factors

309   for the key finding are large and the deviations from normality are small.

310   As shown in Tab. 1 and on Fig. 3, according to both measures, MMs outperformed GBVS in

311   predicting human fixations, thereby replicating the results of Henderson and Hayes (2017,

312   2018) using new images and new participants. Contrary to expectations, however, both

313   metrics indicated that DGII predicted fixations better than MMs. Furthermore, performance

314   of AWS and MMs did not differ significantly irrespective of the metrics. Finally, MMs

315   outperformed ICF according to Correlation, but not sAUC. In fact, for the latter metric, JZS-

316   Bayes Factor indicated support for the null hypothesis.

317

318   Table 1. Comparison of Predictive Power of Saliency Models and MMs Using Correlation and

319   sAUC.

| Model | Median of prediction values with 95% confidence intervals | Median of differences from MMs with 95% confidence intervals | W statistic | p-value (Bonferroni-corrected) | JZS Bayes Factor |
|---|---|---|---|---|---|
| Correlation | | | | | |
| DGII | 0.83 [0.78, 0.87] | 0.07 [0.03, 0.11] | 526 | 0.00738 | 32.26 |
| MMs | 0.77 [0.72, 0.81] | – | – | – | – |
| AWS | 0.73 [0.67, 0.76] | -0.06 [-0.12, -0.01] | 192 | 0.10412 | 1.48 |
| ICF | 0.68 [0.61, 0.71] | -0.12 [-0.18, -0.06] | 144 | 0.00936 | 16.90 |
| GBVS | 0.62 [0.56, 0.68] | -0.11[-0.26, -0.05] | 94 | < .001 | 396.96 |
| sAUC | | | | | |
| DGII | 0.79 [0.77, 0.82] | 0.06 [0.05, 0.08] | 662 | < .001 | > 1000 |
| MMs | 0.73 [0.69, 0.76] | – | – | – | – |
| AWS | 0.75 [0.72, 0.77] | 0.02 [0.01, 0.04] | 490 | 0.0507 | 0.60 |
| ICF | 0.74 [0.70, 0.76] | 0.01 [-0.01, 0.02] | 383 | 1.00 | 0.19 |
| GBVS | 0.64 [0.60, 0.66] | -0.10 [-0.12, -0.08] | 13 | < .001 | > 1000 |

320

321  **Semi-partial correlations.** Because predictions of models and MMs overlap, we quantified

322  their distinct predictive power using semi-partial correlations. We conducted these analyses

323  for GBVS (used in the original MMs studies) and DGII (the only model which markedly

324  outperformed MMs).

325  For each scene from the Consistent condition, we calculated two semi-partial correlations

326  with the distribution from smoothed fixations: one for MMs while controlling for GBVS, and

327  one for GBVS while controlling for MMs (see Fig. 4). Consistent with findings by Henderson

328  and Hayes (2018), MMs explain more unique variance than GBVS (Fig. 6a), as indicated by

329  the significantly higher coefficients in the former than the latter case (mean difference 0.28,

330  95% confidence interval (CI) [0.17, 0.39]; paired t-test, $t(35) = 5.22$, $p < .001$). Interestingly,

331  the identical analysis with DGII revealed that DGII explained significantly more unique

332  variance than MMs (mean difference 0.15, 95% CI [0.07, 0.24]; $t(35) = 3.60$, $p < .001$, see

333  also Fig. 4b).

334



335

336  Fig. 4. Comparison of semi-partial correlations with smoothed human fixations for (a) MMs

337  and GBVS and for (b) MMs and DGII. The obtained coefficients were significantly higher

338  when assessing MMs while controlling for GBVS compared to when assessing GBVS when

339  controlling for MMs. The opposite was true for the analyses with DGII. All figure

340  characteristics are as in Fig. 3. except that means instead of medians are presented.

341

342 **Internal replication.** To demonstrate the generalizability of our conclusions beyond
343 SCEGRAM images, we replicated the main results with a different stimulus set (see the
344 Supplement).

345

346 ### Comparing meaning maps and saliency models – discussion

347 If human gaze is guided by meaning, and if MMs index the distribution of meaning across an
348 image, MMs should outperform saliency models that are exclusively based on image
349 features. Our first analysis showed that this prediction does not hold. In fact, DGII generated
350 better predictions and explained more unique variance than MMs. Therefore, at least one of
351 the two premises of our prediction is wrong: either human eye-movements are not sensitive
352 to meaning or MM do not index meaning. The second analysis allowed us to distinguish
353 between these alternatives.

354

355 ### Analyzing the effects of semantic inconsistencies within scenes – method

356 In the second analysis, we assessed how human observers, DGII, and MMs respond to
357 experimental changes in meaning induced by altered object-context relationships. We used
358 eye-movement data from both the Consistent and the Inconsistent condition. These
359 conditions differed solely in the Critical Region, an area that either contained an object that
360 was either consistent with the scene context or induce semantic conflict. For each scene, we
361 calculated the mass of the distributions of human gaze, DGII, and MMs falling into the
362 Critical Region, respectively, and divided it by the Region's area for normalization. Our
363 primary interest was the comparison between conditions: to the extent to which humans,
364 DGII, and MMs are sensitive to meaning, they should fixate more (humans) or predict more
365 fixations (DGII and MMs) on the Critical Region in the Inconsistent than the Consistent
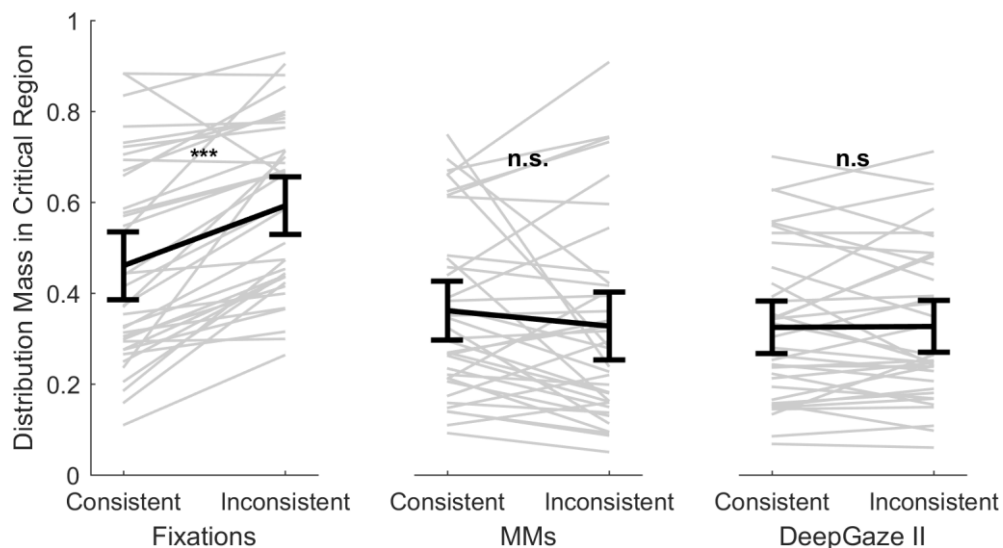366 condition.

367

368 ### Analyzing the effects of semantic inconsistencies within scenes – results

369 Our comparison indicated that, as predicted, observers fixated more on inconsistent than

370 consistent objects (Fig. 5a). By contrast, behavior of both MMs and DGII did not change

371 across conditions (Fig. 5b and c). These impressions were confirmed by a 2x3 ANOVA, with

372 condition (Consistent vs. Inconsistent) as a within-subjects factor and the distribution source

373 (human fixations vs. MMs vs. DGII) as a between-subjects factor. We found a statistically

374 significant main effect of distribution source, $F(2, 105) = 13.09$, $p < .001$, $\omega^2 = 0.16$ and

375 condition, $F(1, 105) = 7.41$ $p = 0.0076$ X, $\omega^2 = 0.005$. These main effects were qualified by a

376 significant interaction, $F(2, 105) = 16.90$, $p < .001$ X, $\omega^2 = 0.026$. Tukey post-hoc tests showed

377 that human observers looked more at the Critical Regions in the Inconsistent, than the

378 Consistent condition, $t(105) = -6.22$, $p < .001$. In contrast, no significant differences between

379 conditions were found for DGII, $t(105) = -0.09$ $p = 1.0$, and MMs, $t(105) = 1.60$ $p = 0.6028$.

380 Comparisons within conditions indicated that human fixations differed from MMs in the

381 Inconsistent condition, $t(129.91) = 5.78$ $p < .001$, but not the Consistent condition, $t(129.91)$

382 $= 2.16$ $p = 0.2662$. A significant difference between DGII and human fixations was detected

383 in both Consistent, $t(129.91) = -2.96$ $p = 0.0420$, and Inconsistent conditions, $t(129.91) = -$

384 $5.79$ $p < .001$.

385



386

387 Fig. 5. Normalized distribution mass falling within Critical Regions in both conditions for (a)

388 smoothed human fixations, (b) MMs, and (c) DGII. All figure characteristics are as in Fig. 3.

389

390 Additionally, conditions differed regarding the number of fixations per image, $t(35) = 5.67$ p

391 < .001. On average, there were 6% fewer fixations in the Inconsistent condition. This

392 excludes the possibility that higher number of fixations in this condition might drive the

393 observed increase in the distribution mass falling within the Critical Regions.

394 Any systematic differences in object size between Consistent and Inconsistent conditions

395 also could affect our results because larger objects may attract more fixations solely

396 because they occupy a larger image area. However, this factor was minimized by showing

397 each object in a consistent and an inconsistent context. Yet, the same object might be

398 shown in a slightly different position in the two conditions and might therefore occupy

399 slightly different amounts of the image. This was, however, not the case: the JZS Bayes

400 Factor of 4.26 indicated that the two conditions did not differ in the size of the bounding

401 boxes of each manipulated object (objects in the Inconsistent condition were on average

402 1562.28 pixels larger; 95% confidence interval: [-2582.74, 5707.29]).

403 Next, please note that we employed a within-subject design, which might have led to carry-

404 over effects: observer viewing a given scene in the Inconsistent condition first could be

405 biased to look at the Critical Region in the Consistent condition when they viewed the same

406 scene for a second time. Note that even if this unwanted phenomenon occurred despite a

407 randomised order of stimuli presentation, it could only decrease the magnitude of the

408 effects of interest.

409 Finally, it is possible that our observers implicitly engaged in a task. Specifically, once the

410 observers realized that the stimuli contain object-context inconsistencies, they might have

411 started actively searching for them. Engaging in this semantic oddball-search task would

412 result in very different spatial distributions of fixations compared to the ones that would be

413 obtained during free-viewing. This prediction was not supported by our findings: we

414 replicated our main experiment in a different set of observers with images that did not

415 contain semantic inconsistencies, and found that DGII still predicted fixation locations better

416 than MMs. This separate data set, therefore, suggests that observers did not engage in an

417 oddball search task and that the superiority of DGII is not specific to SCEGRAM images only

418 (details to be found in the Supplement).

419 To summarize, semantic changes induced by altering object-context relationships elicited
420 changes in distributions of human fixations, but neither MMs nor DGII could predict them.
421 These results suggest that both models might be sensitive to image features, which are
422 frequently correlated with image meaning, rather than to meaning itself.

423

# Discussion

425 A long-standing debate in visual perception concerns the extent to which visual features vs.
426 semantic content guide human eye-movements in free viewing of natural scenes. To
427 distinguish these hypotheses, indexing the distributions both of features and meaning
428 across an image is critical. While image-based saliency models have been used to index
429 features for two decades, measuring semantic importance has been difficult until meaning
430 maps (MMs) have recently been proposed. Here, we assessed the extent to which MMs
431 indeed capture the distribution of meaning across an image. First, we demonstrate that
432 despite the purported importance of meaning as measured by MMs for gaze control, MMs
433 are not better predictors of locations of human fixations than at least some saliency models,
434 which are based solely on image features. In fact, DeepGaze II (DGII), a model using deep
435 neural network features, outperformed MMs. Second, we assessed the sensitivity of human
436 eye-movements, MMs, and DGII to changes in image meaning induced by violations of
437 typical object-context relationships. Observers fixated more often on regions containing
438 objects inconsistent with scene context (thus replicating previous findings) but these regions
439 were not indexed as more meaningful by MMs, or as more salient by DGII. Together, these
440 findings challenge central assumptions of MMs, suggesting that they are insensitive to the
441 semantic information contained in the stimulus.

442 The good performance of DGII in predicting human gaze might be attributable to the high-
443 level features it extracts from images. Three other models, which use low-level features,
444 failed to decisively outperform MMs. However, unlike two of them (GBVS and AWS), DGII is
445 trained with data on human fixations to optimize performance (Kümmerer et al., 2016,
446 2017). Yet, training alone cannot explain the difference in performance. The third low-level
447 feature model (ICF) is trained in the same way (Kümmerer et al., 2017) but still achieves a
448 lower performance than DGII. These findings suggest that feature type is indeed critical for a

449 model's performance. Importantly, however, while DGII uses high-level features transferred
450 from a deep neural network trained on object recognition (Simonyan & Zisserman, 2014),
451 this is not equivalent to indexing meaning. Rather, the good performance of DGII is likely
452 due to meaning supervening on, or correlating with, some of the features indexed by this
453 model.

454 Correlation between visual features and meaning as the source of good performance in
455 saliency models has already been considered by the authors of MMs (Henderson & Hayes,
456 2017). Our findings suggest that MMs might share this characteristic with saliency models.
457 Specifically, the ratings used to construct MMs might be based on visual properties in such a
458 way that highly structured patches that contain high-level features receive high ratings.
459 These features often correlate with meaning, but in and of themselves do not amount to
460 meaning. According to this interpretation, both DGII and MMs index high-level features.
461 Their success in predicting human behavior derives from the typically strong correlation
462 between high-level features and meaning, with a higher correlation for the features
463 extracted by DGII than MMs.

464 An alternative interpretation of the finding that DGII outperforms MMs is that image
465 features rather than meaning guide human fixations. However, this interpretation is
466 inconsistent with our second analysis. Here, observers clearly exhibited sensitivity to
467 meaning, as indicated by changes in gaze-patterns elicited by introducing semantic
468 inconsistencies into the images. This experimental manipulation targets a type of meaning
469 that is based on how objects relate to the broader context in which they occur. While
470 specific, it is precisely this kind of meaning that is of high theoretical importance in eye-
471 movement research (Henderson, 2017; 429 Henderson et al., 2009). Natural scenes are
472 composed of multiple objects, and the physical and semantic relationships between these
473 objects as well as their relationship to the scene gist, determine the meaning of a scene
474 (Kaiser et al., 2019; Malcolm et al., 2016; Võ et al., 2019). Thus, the fact that MMs are not
475 sensitive to the meaning derived from object-context relationships seriously limits their
476 usefulness.

477 It is, however, possible that – as has been already suggested (Henderson et al., 2018) – MMs
478 capture some form of 'local' meaning that is important for oculomotor control. Evaluating

479   our results in this respect is complicated by the correlation between features and meaning

480   (Elazary & Itti, 2008), which we already alluded to above. Yet, at the very least, the fact that

481   MMs do not consistently outperform even simple saliency models such as AWS that by

482   design rely on low-level image features warrants caution. This finding indicates that either

483   the purported kind of meaning indexed by MMs is not of primary importance for guidance

484   of eye-movements, or that it is almost perfectly correlated with the features indexed by

485   models such as AWS. A similar issue relates to DGII: while our study shows that this model

486   does not index meaning derived from object-context relationships, one might argue that it

487   acquires sensitivity to some (local) form of meaning by virtue of being trained on human

488   data. Specifically, if eye-movements are guided by the semantic content of images, then

489   training on eye-movement data might lead to developing 'meaning-sensitivity' in the model.

490   While this scenario cannot be ruled out for the same reasons as in the case of MMs, recall

491   that the ICF model – which uses simpler features than DGII – is also trained on human data

492   but fails to reach the high performance of DGII. Therefore, if the high performance of DGII is

493   based on some form of 'local' meaning, then it is not training per se that leads to the

494   development of this meaning but an interaction of training and specific features.

495   If nothing else, these considerations indicate the urgent need for developing a more

496   nuanced conceptual approach and terminology to capture the intricacies of different types

497   of 'meaning', and a more appropriate language to talk about the relationship between

498   'features' and 'meaning'. Without a clearer theoretical framework, it will be difficult to

499   experimentally settle debates regarding the role of 'meaning' in natural scene perception.

500   In any case, the insensitivity to semantic inconsistencies reveals inherent limitations of both

501   MMs and DGII. The way in which MMs are constructed implicitly assumes that meaning is a

502   local image-property, which is not true for object-context (in)consistency. This limitation

503   may potentially be alleviated by 'contextualized MMs' (Peacock, Hayes, & Henderson,

504   2019), a recently suggested modification of the 'standard' MMs. These novel maps are

505   created from meaningfulness ratings by observers who see the whole scenes from which

506   the to-be-rated patches were derived. It is yet to be seen what this approach can reveal

507   about fixation selection beyond the fact that humans asked to indicate meaningful or

508   interesting regions within scenes highlight areas, which tend to be frequently fixated by

509   other observers (Nyström & Holmqvist, 2008; Onat et al., 2014). DGII, in turn, does not

510   explicitly encode semantic information, and was not trained on the relationship between

511   eye movements and semantic (in)consistency. But its failure highlights an opportunity to

512   improve saliency models by incorporating semantic relationships (Bayat, Koh, Nand, Pereira,

513   & Pomplun, 2018).

514   Taken together, our results suggest that, contrary to their core promise as a methodology,

515   meaning maps (MMs) do not offer a way to measure the spatial distribution of meaning

516   across an image. Instead of meaning per-se, they seem to index high-level features that

517   have the potential to carry meaning in typical natural scenes. They share this characteristic

518   with state-of-the-art saliency models, which are easier to use, do not require human

519   annotation, and yet predict locations of human fixations better than MMs.

520

521                                **References**

522   Bayat, A., Koh, D. H., Nand, A. K., Pereira, M., & Pomplun, M. (2018). Scene Grammar in

523        Human and Machine Recognition of Objects and Scenes. In *Proceedings of the IEEE*

524        *Conference on Computer Vision and Pattern Recognition Workshops*.

525        https://doi.org/10.1109/CVPRW.2018.00268

526   Borji, A., Sihite, D. N., & Itti, L. (2013). Objects do not predict fixations better than early

527        saliency : A re-analysis of Einhauser et al.'s data. *Journal of Vision*, *13*(2013), 1–4.

528        https://doi.org/10.1167/13.10.18

529   Bylinskii, Z., Judd, T., Borji, A., Itti, L., Durand, F., Oliva, A., & Torralba, A. (2014). MIT Saliency

530        Benchmark Results. Retrieved from http://saliency.mit.edu/

531   Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., & Durand, F. (2016). What do different

532        evaluation metrics tell us about saliency models? *ArXiv*. Retrieved from

533        http://arxiv.org/abs/1604.03605

534   Elazary, L., & Itti, L. (2008). Interesting objects are visually salient. *Journal of Vision*, *8*(3), 1–

535        15. https://doi.org/10.1167/8.3.3

536   Garcia-Diaz, A., Fdez-Vidal, X. R., Pardo, X. M., & Dosil, R. (2012). Saliency from hierarchical

537        adaptation through decorrelation and variance normalization. *Image and Vision*

538    *Computing*, *30*(1), 51–64. https://doi.org/10.1016/j.imavis.2011.11.007

539    Harel, J., Koch, C., & Perona, P. (2006). Graph-Based Visual Saliency. *Advances in Neural*

540    *Information Processing Systems 19*, *19*, 545–552. https://doi.org/10.1.1.70.2254

541    Hayes, T. R., & Henderson, J. M. (2019). Center bias outperforms image salience but not

542    semantics in accounting for attention during scene viewing. *Attention, Perception, &*

543    *Psychophysics*. https://doi.org/https://doi.org/10.3758/s13414-019-01849-7

544    Hayhoe, M., & Ballard, D. (2005). Eye movements in natural behavior. *Trends in Cognitive*

545    *Sciences*, *9*(4). https://doi.org/10.1016/j.tics.2005.02.009

546    Hegde, J., & Kersten, D. (2010). A Link between Visual Disambiguation and Visual Memory.

547    *Journal of Neuroscience*, *30*(45), 15124–15133.

548    https://doi.org/10.1523/JNEUROSCI.4415-09.2010

549    Henderson, J. M. (2017). Gaze Control as Prediction. *Trends in Cognitive Sciences*, *21*(1), 15–

550    23. https://doi.org/10.1016/j.tics.2016.11.003

551    Henderson, J. M., & Hayes, T. R. (2017). Meaning-based guidance of attention in scenes as

552    revealed by meaning maps. *Nature Human Behaviour*, *1*(October).

553    https://doi.org/10.1038/s41562-017-0208-0

554    Henderson, J. M., & Hayes, T. R. (2018). Meaning guides attention in real-world scene

555    images: Evidence from eye movements and meaning maps. *Journal of Vision*, *18*(6), 10.

556    https://doi.org/10.1167/18.6.10

557    Henderson, J. M., Hayes, T. R., Peacock, C. E., & Rehrig, G. (2019). Meaning and Attentional

558    Guidance in Scenes : A Review of the Meaning Map Approach. *Vision*, *3*(2).

559    Henderson, J. M., Hayes, T. R., Rehrig, G., & Ferreira, F. (2018). Meaning Guides Attention

560    during Real-World Scene Description. *Scientific Reports*, *8*(1), 13504.

561    https://doi.org/10.1038/s41598-018-31894-5

562    Henderson, J. M., Malcolm, G. L., & Schandl, C. (2009). Searching in the dark: Cognitive

563    relevance drives attention in real-world scenes. *Psychonomic Bulletin & Review*, *16*(5),

564    850–856. https://doi.org/10.3758/PBR.16.5.850

565    Henderson, J. M., Weeks, P. A., & Hollingworth, A. (1999). The effects of semantic

566    consistency on eye movements during complex scene viewing. *Journal of Experimental*

567    *Psychology: Human Perception and Performance*, *25*(1), 210–228.

568    https://doi.org/10.1037/0096-1523.25.1.210

569    Itti, L., & Koch, C. (2000). A saliency-based search mechanism for overt and covert shifts of

570    visual attention. *Vision Research*, *40*(10–12), 1489–1506.

571    https://doi.org/10.1016/S0042-6989(99)00163-7

572    Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews*

573    *Neuroscience*, *2*(3), 194–203. https://doi.org/10.1038/35058500

574    Kaiser, D., Quek, G. L., Cichy, R. M., & Peelen, M. V. (2019). Object Vision in a Structured

575    World. *Trends in Cognitive Sciences*, *23*(8), 672–685.

576    https://doi.org/10.1016/j.tics.2019.04.013

577    Kietzmann, T. C., McClure, P., & Kriegeskorte, N. (2019). Deep Neural Networks in

578    Computational Neuroscience. In *Oxford Research Encyclopedia of Neuroscience*.

579    Kleiner, M., Brainard, D., & Pelli, D. G. (2007). What's new in psychtoolbox-3? *Perception*,

580    *36*(1).

581    Koehler, K., Guo, F., Zhang, S., & Eckstein, M. P. (2014). What do saliency models predict?

582    *Journal of Vision*, *14*(3). https://doi.org/10.1167/14.3.14

583    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2015). Information-theoretic model

584    comparison unifies saliency metrics. *Proceedings of the National Academy of Sciences*,

585    *112*(52), 16054–16059. https://doi.org/10.1073/pnas.1510393112

586    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2016). DeepGaze II: Reading fixations from

587    deep features trained on object recognition, 1–16. Retrieved from

588    http://arxiv.org/abs/1610.01563

589    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2018). Saliency Benchmarking Made Easy:

590    Separating Models, Maps and Metrics. In V. Ferrari, M. Hebert, C. Sminchisescu, & Y.

591    Weiss (Eds.), *Computer Vision – ECCV 2018. ECCV 2018. Lecture Notes in Computer*

592    *Science* (Vol. 11220, pp. 798–814). Springer. https://doi.org/10.1007/978-3-030-01270-

593    0_47

594    Kümmerer, M., Wallis, T. S. A., Gatys, L. A., & Bethge, M. (2017). Understanding Low- and

595          High-Level Contributions to Fixation Prediction. In *The IEEE International Conference on*

596          *Computer Vision (ICCV)*. https://doi.org/10.1109/ICCV.2017.513

597    Malcolm, G. L., Groen, I. I. A., & Baker, C. I. (2016). Making Sense of Real-World Scenes.

598          *Trends in Cognitive Sciences*, *20*(11), 843–856.

599          https://doi.org/10.1016/j.tics.2016.09.003

600    Nyström, M., & Holmqvist, K. (2008). Semantic override of low-level features in image

601          viewing–both initially and overall. *Journal of Eye Movement Research*, *2*(2), 1–11.

602          https://doi.org/10.16910/jemr.2.2.2

603    Öhlschläger, S., & Võ, M. L. H. (2017). SCEGRAM: An image database for semantic and

604          syntactic inconsistencies in scenes. *Behavior Research Methods*, *49*(5).

605          https://doi.org/10.3758/s13428-016-0820-3

606    Onat, S., Açik, A., Schumann, F., & König, P. (2014). The contributions of image content and

607          behavioral relevancy to overt attention. *PLoS ONE*, *9*(4).

608          https://doi.org/10.1371/journal.pone.0093254

609    Parkhurst, D., Law, K., & Niebur, E. (2002). Modeling the role of salience in the allocation of

610          overt visual attention. *Vision Research*, *42*(1), 107–123. https://doi.org/10.1016/S0042-

611          6989(01)00250-4

612    Peacock, C. E., Hayes, T. R., & Henderson, J. M. (2018). Meaning guides attention during

613          scene viewing, even when it is irrelevant. *Attention, Perception, and Psychophysics*, 20–

614          34. https://doi.org/10.3758/s13414-018-1607-7

615    Peacock, C. E., Hayes, T. R., & Henderson, J. M. (2019). The role of meaning in attentional

616          guidance during free viewing of real-world scenes. *Acta Psychologica*, *198*(June).

617          https://doi.org/10.1016/j.actpsy.2019.102889

618    Rider, A. T., Coutrot, A., Pellicano, E., Dakin, S. C., & Mareschal, I. (2018). Semantic content

619          outweighs low-level saliency in determining children's and adults' fixation of movies.

620          *Journal of Experimental Child Psychology*, *166*, 293–309.

621          https://doi.org/10.1016/j.jecp.2017.09.002

622 Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for

623     accepting and rejecting the null hypothesis. *Psychonomic Bulletin and Review*, *16*(2),

624     225–237. https://doi.org/10.3758/PBR.16.2.225

625 Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale

626     Image Recognition. *CoRR, Abs/1409.1556*. Retrieved from

627     http://arxiv.org/abs/1409.1556

628 Stoll, J., Thrun, M., Nuthmann, A., & Einhäuser, W. (2015). Overt attention in natural scenes:

629     Objects dominate features. *Vision Research*, *107*, 36–48.

630     https://doi.org/10.1016/j.visres.2014.11.006

631 Tatler, B. (2007). The central fixation bias in scene viewing: Selecting an optimal viewing

632     position independently of motor biases and image feature distributions. *Journal of*

633     *Vision*, *7*(4), 1–17. https://doi.org/10.1167/7.14.4

634 Tatler, B. W., Hayhoe, M. M., Land, M. F., & Ballard, D. H. (2011). Eye guidance in natural

635     vision: Reinterpreting salience. *Journal of Vision*, *11*(5), 5–5.

636     https://doi.org/10.1167/11.5.5

637 Teufel, C., Dakin, S. C., & Fletcher, P. C. (2018). Prior object-knowledge sharpens properties

638     of early visual feature- detectors. *Scientific Reports*, (June), 1–12.

639     https://doi.org/10.1038/s41598-018-28845-5

640 Võ, M. L. H., Boettcher, S. E., & Draschkow, D. (2019). Reading scenes: how scene grammar

641     guides attention and aids perception in real-world environments. *Current Opinion in*

642     *Psychology*, *29*, 205–210. https://doi.org/10.1016/j.copsyc.2019.03.009

643 Wilming, N., Onat, S., Ossandón, J. P., Açik, A., Kietzmann, T. C., Kaspar, K., Gameiro, R. R.,

644 Vormberg, A., & König, P. (2017). An extensive dataset of eye movements during viewing of

645 complex images. *Scientific Data*, *4*, 1–11. https://doi.org/10.1038/sdata.2016.126

646 Zhang, L., Tong, M. H., Marks, T. K., & Cottrell, G. W. (2008). SUN: A Bayesian framework for

647     saliency using natural statistics. *Journal of Vision*, *8*(32). https://doi.org/10.1167/8.7.32