# 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

Pedziwiatr et al.

# Introduction

30

31 Human eyes resolve fine detail only in a small, central part of the visual field, with resolution

32 dropping off rapidly in the periphery. To sample details, we move our eyes to orient the high-

33 resolution part of our visual system successively to different parts of a visual scene.

34 Information about these small scene parts is extracted during fixations – short periods in

35 which the eyes are relatively stable. Thus, due to the structure of our visual system, human

36 vision depends on eye movements. How the brain decides where to look in a visual scene is

37 therefore an important question. A long-standing hypothesis suggests that semantic content

38 of image regions is important in guiding eye movements. Recent work presented meaning

39 maps (MMs) as a tool to test this hypothesis (Henderson & Hayes, 2017, 2018). This technique

40 aims to index the spatial distribution of meaning across an image, which has potential

41 applications far beyond eye-movement research. Here, we assess and challenge central

42 assumptions of this novel tool.

43 A classic finding in eye-movement research shows that the specific task of an observer has an

44 influence on where they direct their eyes (Yarbus, 1967; Hayhoe & Ballard, 2005). But in

45 everyday life, we frequently move our eyes without any goal other than to explore the

46 environment. In the lab, this behavior is examined in free-viewing paradigms, during which

47 eye movements are recorded while images are viewed without an explicit task (Koehler, Guo,

48 Zhang, & Eckstein, 2014, but see Tatler, Hayhoe, Land, & Ballard, 2011). To explain what

49 guides eye movements during free viewing, two opposing accounts have been put forward.

50 According to the first account, eye movements are guided primarily by image characteristics

51 (Borji, Sihite, & Itti, 2013; Itti & Koch, 2001; Parkhurst, Law, & Niebur, 2002). Potential support

52 for this view comes from saliency models: algorithms, which exclusively use visual features of

53 an image to predict human fixations. Although early models, which used only simple features

54 such as local intensity or colors (Itti & Koch, 2000), are now deemed only moderately

55 successful (Bylinskii et al., 2014), more recent saliency models achieve a remarkably high

56 performance (Kümmerer, Wallis, Gatys, & Bethge, 2017). These models harness deep

57 convolutional neural networks – biologically inspired machine learning algorithms, that

58 somewhat resemble the human visual system (Kietzmann, McClure, & Kriegeskorte, 2019).

59 However, even such models rely solely on visual features, albeit high-level ones.

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

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

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

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

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

65    Attributing meaning to certain parts of the scene is impossible without prior knowledge of

66    the world, i.e., a factor that is independent of the visual input (Hegde & Kersten, 2010; Teufel,

67    Dakin, & Fletcher, 2018). Consequently, the notion that semantic content guides eye-

68    movements is inconsistent with the idea that the allocation of fixations is dependent solely

69    on the distribution of image features. Given that meaning is not image-computable, the

70    notion that semantic content guides eye-movements is inconsistent with the idea that the

71    eye-movements are dependent solely on the distribution of image features.

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

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

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

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

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

77    isolated patches, which are rated for their meaningfulness independently from the rest of the

78    visual scene. These ratings are pooled together into a smooth map, which is supposed to

79    capture the distribution of meaning across the image. Compared to outputs from a simple

80    saliency model (GBVS, Harel et al., 2006), MMs were more predictive of human fixations. On

81    that basis it has been claimed that meaning guides human fixations in natural scene viewing

82    (Henderson & Hayes, 2017, 2018). Here, we examined central predictions of this claim.

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

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

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

86    art models. We replicate the finding that MMs perform better than some of the most basic

87    saliency models. Contrary to the prediction, however, DeepGaze II (DGII; Kümmerer, Wallis,

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

89    network, outperforms MMs.

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

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

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

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

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

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

96  1999; Öhlschläger & Võ, 2017). Replicating previous studies, image regions attracted more

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

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

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

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

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

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

103  visual features rather than meaning.

104

# Method

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

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

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

109  locations of human fixations in natural scenes. Subsequently, we assessed the sensitivity of

110  MMs and the best-performing saliency model to manipulations of scene meaning. The

111  reported experiment was not preregistered. The data, the code to create MMs, and all openly

112  available resources used in the study can be accessed via the links provided in the

113  Supplement.
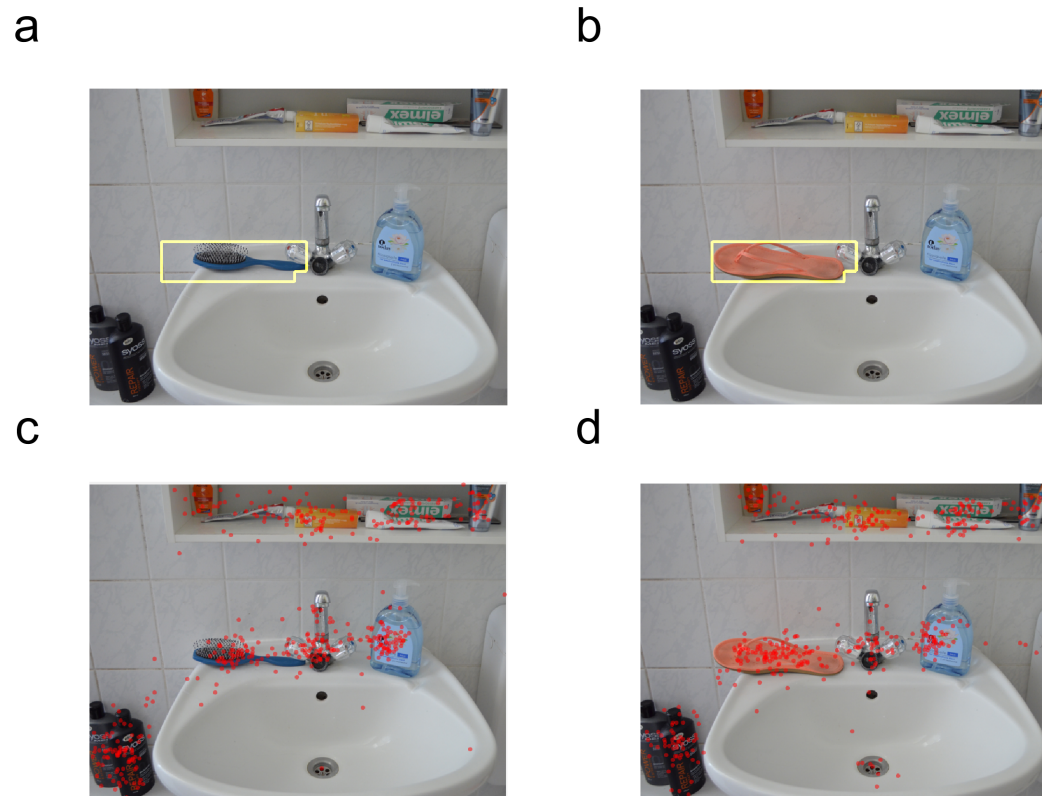
Pedziwiatr et al.

a

b

c

d

114

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.

119

**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 the Supplement together

132   with the selection criteria). We also replicated the main finding of the first analysis in an

133   additional set of 30, very different, images (reported in the Supplement).

134

135   **Procedure.** The procedure consisted of 3 blocks, interleaved with breaks. Participants were

136   instructed to 'look carefully at each' image. Experimental blocks began with an eye tracker

137   calibration/validation. Within each block, observers free-viewed a series of 24 photographs

138   from both SCEGRAM conditions, each for 7 seconds. After image offset, observers were

139   required to press a button to view the next image. Then, a fixation point appeared centrally

140   on a screen and once observers fixate on it (as determined online by their eye-trace), the

141   actual image was displayed. Before starting the experiment, observers viewed a sample image

142   in an identical regime to familiarize themselves with the procedure. Each stimulus was shown

143   once and the order of presentation was fully randomized. The stimuli were presented against

144   a uniform grey background and had a width of 688 pixels and a height of 524 pixels, which

145   subtended approximately 19.7 and 15 degrees of visual angle, respectively. Stimulus

146   presentation time and size were adopted from a previous study with the SCEGRAM database

147   (Öhlschläger & Võ, 2017).

148

149   **Observers.** 20 volunteers (3 male; mean age 19.4) recruited from the Cardiff University

150   undergraduate population took part in the study. All reported normal or corrected-to-normal

151   vision, provided written consent, and received course credits in return for participation. The

152   study was approved by the Cardiff University School of Psychology Research Ethics

153   Committee. The primary units of interest in our analyses were the distributions of fixations

154   over images. The number of observers we recruited guarantees that including more observers

155   would not change these distributions significantly (demonstrated in the Supplement).

156

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

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

159   1080 pixels, 21 inches diagonal). Chin and forehead rests were used to ensure that observers

160   maintained the constant distance of 49 cm from the screen. Their eye movements were

Pedziwiatr et al.

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

162  mount. The experiment was controlled by custom-written Matlab (R2017a version) scripts

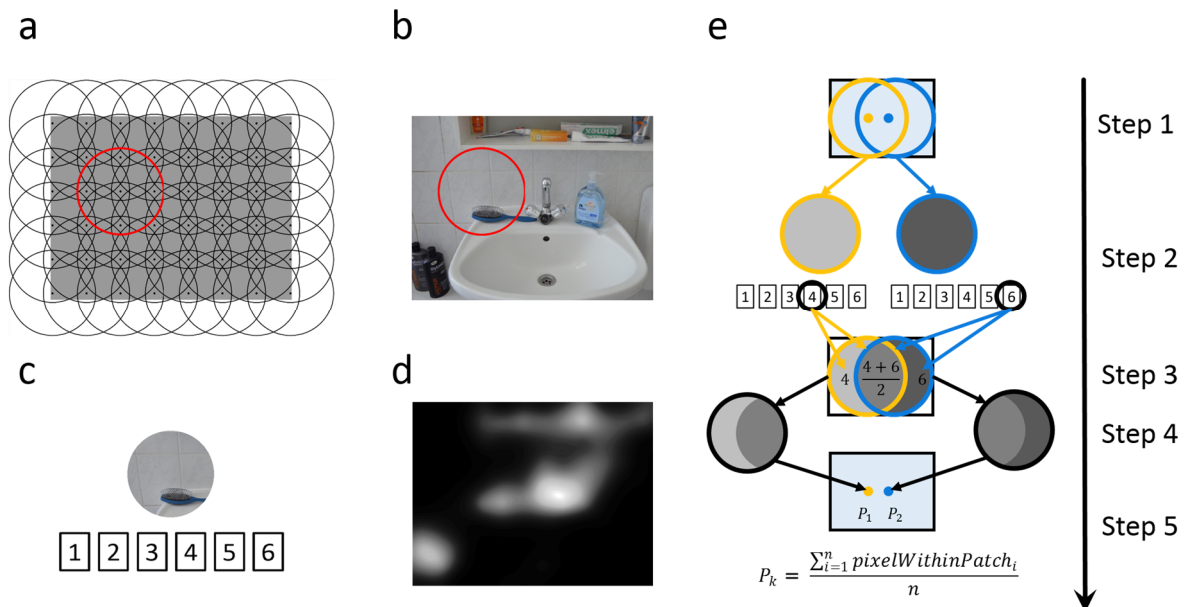163  using Psychophysics Toolbox Version 3 (Kleiner, Brainard, & Pelli, 2007).

164



165

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

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

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

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

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

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

172  included additional labels (ranging from 'Very Low' to 'Very High'). (**d**) Illustration of a

173  meaning map with greyscale values indicating 'meaningfulness'. (**e**) Simplifying illustration of

174  how meaning maps are generated from ratings. For simplicity sake, only two patches are

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

176  shown). All pixels within an image area are then assigned average rating values, taking into

177  account all ratings for patches that overlap with this area (step 3). For the area of the original

178  patch (step 4), all pixels are then averaged and the resulting value is assigned to the center of

179  the patch (step 5). Finally, the patch centers were used as interpolation nodes for thin-plate

180  spline interpolation producing a smooth distribution of values over the image (not illustrated).

7

Pedziwiatr et al.

181  This procedure was conducted separately for the fine and coarse grid, and the meaning map

182  for a given image was created by averaging the two outcomes and normalizing the result to a

183  range between 0 and 1.

184

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

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

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

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

189  36% of the image width) (Fig. 2a and b). Their centers were 58 pixels (fine) and 97 pixels

190  (coarse) apart from each other.

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

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

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

194  recruited via the crowdsourcing platform Amazon Mechanical Turk (see Supplement for

195  eligibility criteria). Each participant provided ratings for 305 or 303 patches of both sizes

196  (selected randomly from all images), on average spent approximately 14 min on the task, and

197  received 2.18 USD as remuneration. In total, 69 individuals were used as raters, with three

198  individuals rating each individual patch. The collected ratings were then used to create MMs

199  (see Fig. 2).

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

201  from the Consistent and Inconsistent conditions differ only in the Critical Region (the part of

202  the image containing the manipulated object) while the remaining parts overlap. We

203  collected meaningfulness ratings for the patches belonging to overlapping areas only once,

204  and the separate sets of ratings for Consistent and Inconsistent condition were collected only

205  for those patches that contained at least one pixel belonging to the Critical Region. In total,

206  the number of patches rated in the study amounted to 7013: 4840 fine patches (of which 520

207  belonged to the images from the Inconsistent condition) and 2173 coarse patches (445

208  Inconsistent).

209

**Saliency models.** In the first analysis, we compared predictive performance of MMs to four saliency models of different complexity. The first two models – GBVS (Harel et al., 2006) and AWS (Garcia-Diaz, Fdez-Vidal, Pardo, & Dosil, 2012) – rely on simple visual features, such as local colors and edge orientations, and share the assumption that fixations land on image regions distinct from their surroundings in terms of values of these features. By contrast to GBVS, AWS includes a statistical whitening procedure to improve performance. Both these models were previously used to estimate the influence of image features relative to cognitive factors on the deployment of fixations: GBVS in the previous studies with MMs, AWS elsewhere (Stoll et al., 2015).

Two other models that we compared to MMs – ICF and DeepGaze II (DGII) – were designed in a data-driven manner (Kümmerer et al., 2017). Both have the same architecture, consisting of a fixed network that extracts sets of features from images and a readout network that is trained on human fixations separately for each model to combine the features in a way to maximize the models' predictive power. While the fixed network of ICF extracts only simple visual features (local intensity and contrast), DGII is tuned to features extracted by a deep convolutional neural network pre-trained for object recognition (VGG-19; Simonyan & Zisserman, 2014).

All saliency models output smooth maps that predict the probability of image regions to be fixated. Human observers have the tendency to look at the center of images (Tatler, 2007), and therefore this probability is usually higher in the central region of the image. This 'center bias' has important consequences for the evaluation of saliency models. Their performance differs depending on whether they are evaluated using a metric expecting some form of this bias or not (Kümmerer, Wallis, & Bethge, 2018). Here, for the sake of simplicity, we do not incorporate center bias in the models or in the MMs (unlike the original authors) and use an appropriate metric for this situation (see Performance metrics section). Importantly, analyses addressing the issue of center bias in a more extensive way (reported in the Supplement) provide only further support for our conclusions.

**Data pre-processing.** Fixation locations from the eye tracker recordings were extracted using the algorithm provided by the device manufacturer operating with the default parameter

240 values. Thereby, we obtained a discrete distribution of fixations on each image (see Fig. 1c
241 and 1d). Then, in line with the previous MMs studies, we smoothed these discrete
242 distributions with a Gaussian filter with a cutoff frequency of -6 dB, using the function
243 provided by Bylinskii and colleagues (2014).

244 Next, smooth distributions from fixations, models, and MMs were separately normalized to a
245 range from 0 to 1 for each image. Finally, for each scene, histograms of all distributions from
246 both conditions were matched to histograms of smoothed fixations from Consistent condition
247 using the Matlab imhistmatch function, as in the original MMs studies. Histogram matching
248 makes distributions directly comparable as it ensures that they differ only with respect to
249 their shape, and not their total mass.

250

251 **Performance metrics.** To compare the ability of MMs and models to predict locations of
252 human fixations in Experiment 1, we use two well-established metrics (Bylinskii, Judd, Oliva,
253 Torralba, & Durand, 2016): Correlation and Shuffled Area Under ROC curve (sAUC; Zhang,
254 Marks, Tong, Shan, & Cottrell, 2007) with the implementations provided by Bylinskii and
255 colleagues (2014).

256 Correlation, used in the previous studies on MMs, is calculated as Pearson's linear correlation
257 coefficient between a smoothed distribution of observers' fixations over the image and
258 predictions of a saliency model or MMs. We additionally used sAUC (Zhang et al., 2008),
259 which, unlike Correlation, guarantees that the measured differences in performance between
260 models are driven by their sensitivity to factors guiding fixations, and not by the degree to
261 which they include human center bias in their predictions, even implicitly (Kümmerer, Wallis,
262 & Bethge, 2015; Kümmerer et al., 2018).

263

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

265 In the first analysis, we compared performance of four saliency models to MMs in predicting
266 human fixations in the Consistent condition, i.e., when viewing typical scenes with no obvious
267 object-context inconsistencies (Tab. 1, Fig. 3). If human gaze is guided by meaning, and if MMs

Pedziwiatr et al.

268    provide an index for the distribution of meaning, we would expect MMs to outperform all

269    saliency models because these models are based solely on image features.
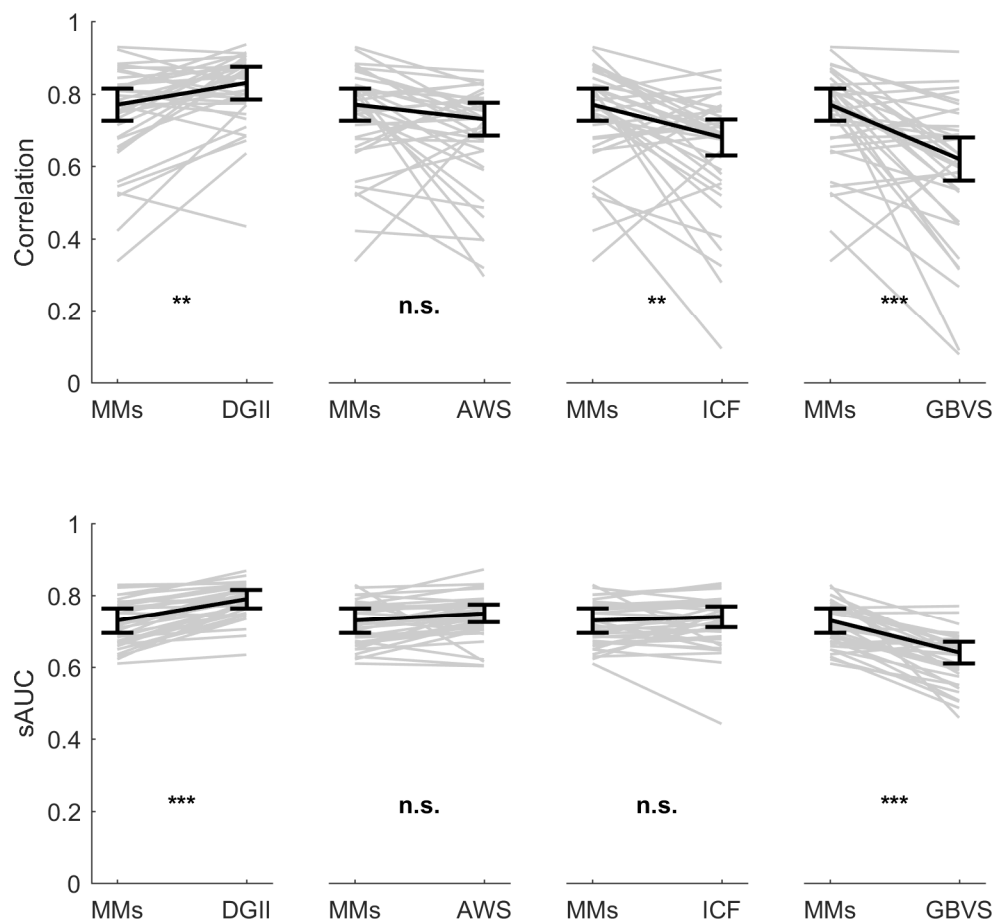
270



271

272

273    Fig. 3. Performance of MMs and saliency models in predicting human fixations according to

274    (a) Correlation and (b) sAUC metrics. Note that according to both metrics DGII predicted

275    human fixations better than MMs. Asterisks indicate p-values from statistical tests comparing

276    MMs to different models (reported in Table 1.): * indicates $p \leq .05$, ** $p \leq .01$, *** $\leq .001$ and

277    'n.s.' indicates the lack of statistical significance. Grey lines connect values obtained for

278    individual images. Black vertical bars indicate 95% confidence intervals for the medians.

279

280    **Predictive power**. Correlation and sAUC values obtained for MMs and for each of the models

281    were compared using Bonferroni-corrected paired Wilcoxon tests (Fig. 3; Tab. 1). We used

282    non-parametric tests because for some of the distributions the assumptions of normality was

11

283    not met. For the same reason we chose a median as a measure of centrality (we calculate

284    confidence intervals for median using a bootstrapping method – see details in the

285    Supplement). Additionally, we calculated JZS Bayes Factor (Rouder, Speckman, Sun, Morey, &

286    Iverson, 2009) to quantify the evidence for (or against) the differences between models and

287    MMs (Tab. 1). While deviations from normality can be problematic for Bayes factor analyses,

288    they are most likely not an issue in the current situation: the Bayes factors for the key finding

289    are large and the deviations from normality are small.

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

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

292    2018) using new images and new participants. Contrary to expectations, however, both

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

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

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

296    Bayes Factor indicated support for the null hypothesis.

297

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

299    sAUC.

| Model | Median of prediction values with 95% confidence intervals | Median of differences from MMs with 95% confidence intervals | Z statistic | p-value (Bonferroni-corrected) | JZS Bayes Factor |
|---|---|---|---|---|---|
| Correlation | | | | | |
| DGII | 0.83 [0.78, 0.87] | 0.07 [0.03, 0.11] | -3.11 | 0.00738 | 32.26 |
| MMs | 0.77 [0.72, 0.81] | – | – | – | – |
| AWS | 0.73 [0.67, 0.76] | -0.06 [-0.12, -0.01] | -2.23 | 0.10412 | 1.48 |
| ICF | 0.68 [0.61, 0.71] | -0.12 [-0.18, -0.06] | -3.04 | 0.00936 | 16.90 |
| GBVS | 0.62 [0.56, 0.68] | -0.11 [-0.26, -0.05] | -3.97 | < .001 | 396.96 |
| sAUC | | | | | |
| DGII | 0.79 [0.77, 0.82] | 0.06 [0.05, 0.08] | -6.36 | < .001 | > 1000 |
| MMs | 0.73 [0.69, 0.76] | – | – | – | – |
| AWS | 0.75 [0.72, 0.77] | 0.02 [0.01, 0.04] | -2.49 | 0.0507 | 0.60 |
| ICF | 0.74 [0.70, 0.76] | 0.01 [-0.01, 0.02] | -0.77 | 1.00 | 0.19 |
| GBVS | 0.64 [0.60, 0.66] | -0.10 [-0.12, -0.08] | -5.96 | < .001 | > 1000 |

Pedziwiatr et al.

300

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

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

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

304 outperformed MMs).

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

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

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

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

309 significantly higher coefficients in the former than the latter case (mean difference 0.28, 95%

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

311 identical analysis with DGII revealed that DGII explained significantly more unique variance

312 than MMs (mean difference 0.15, 95% CI [0.07, 0.24]; $t(35) = 3.60$, $p < .001$, see also Fig. 4b).
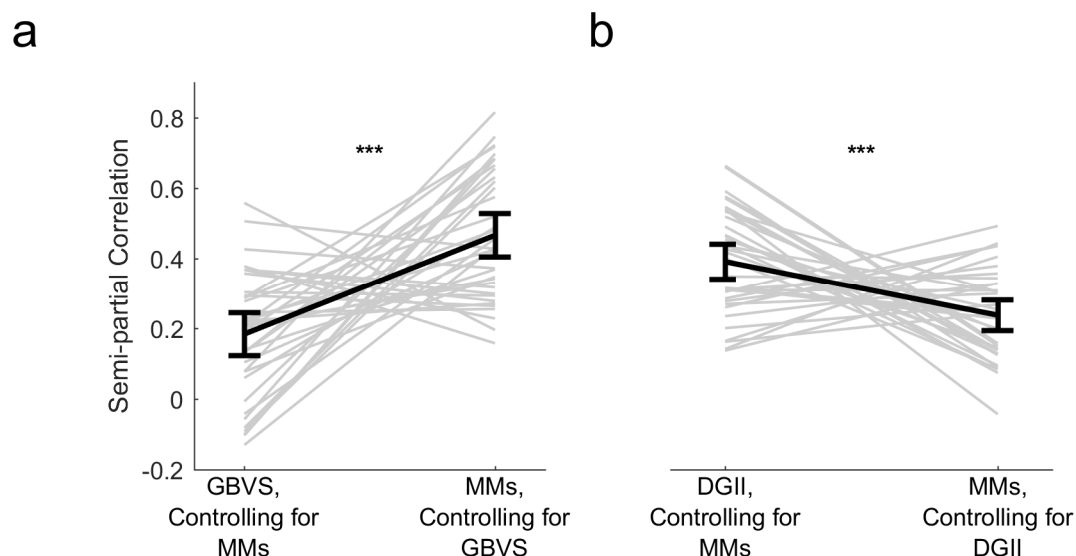
313



314

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

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

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

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

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

320

13

Pedziwiatr et al.

321 **Internal replication.** To demonstrate the generalizability of our conclusions beyond SCEGRAM

322 images, we replicated the main results with a different stimulus set (see the Supplement).

323

### Comparing meaning maps and saliency models – discussion

325 If human gaze is guided by meaning, and if MMs index the distribution of meaning across an

326 image, MMs should outperform saliency models that are exclusively based on image features.

327 Our first analysis showed that this prediction does not hold. In fact, DGII generated better

328 predictions and explained more unique variance than MMs. Therefore, at least one of the two

329 premises of our prediction is wrong: either human eye-movements are not sensitive to

330 meaning or MM do not index meaning. The second analysis allowed us to distinguish between

331 these alternatives.

332

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

334 In the second analysis, we assessed how human observers, DGII, and MMs respond to

335 experimental changes in meaning induced by altered object-context relationships. We used

336 eye-movement data from both the Consistent and the Inconsistent condition. These

337 conditions differed solely in the Critical Region, an area that either contained an object that

338 was either consistent with the scene context or induce semantic conflict. For each scene, we

339 calculated the mass of the distributions of human gaze, DGII, and MMs falling into the Critical

340 Region, respectively, and divided it by the Region's area for normalization. Our primary

341 interest was the comparison between conditions: to the extent to which humans, DGII, and

342 MMs are sensitive to meaning, they should fixate more (humans) or predict more fixations

343 (DGII and MMs) on the Critical Region in the Inconsistent than the Consistent condition.

344

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

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

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

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

14

349 (Consistent vs. Inconsistent) as a within-subjects factor and the distribution source (human
350 fixations vs. MMs vs. DGII) as a between-subjects factor. We found a statistically significant
351 main effect of distribution source, $F(2, 105) = 13.09$, $p < .001$, $\omega^2 = 0.16$ and condition, $F(1,$
352 $105) = 7.41$ $p = 0.0076$  X, $\omega^2 = 0.005$. These main effects were qualified by a significant
353 interaction, $F(2, 105) = 16.90$, $p < .001$ X, $\omega^2 = 0.026$. Tukey post-hoc tests showed that human
354 observers looked more at the Critical Regions in the Inconsistent, than the Consistent
355 condition, $t(105) = -6.22$, $p < .001$. In contrast, no significant differences between conditions
356 were found for DGII, $t(105) = -0.09$ $p = 1.0$, and MMs, $t(105) = 1.60$ $p = 0.6028$. Comparisons
357 within conditions indicated that human fixations differed from MMs in the Inconsistent
358 condition, $t(129.91) = 5.78$ $p < .001$, but not the Consistent condition, $t(129.91) = 2.16$ $p =$
359 $0.2662$. A significant difference between DGII and human fixations was detected in both
360 Consistent, $t(129.91) = -2.96$ $p = 0.0420$, and Inconsistent conditions, $t(129.91) = -5.79$ $p <$
361 $.001$.

362

363

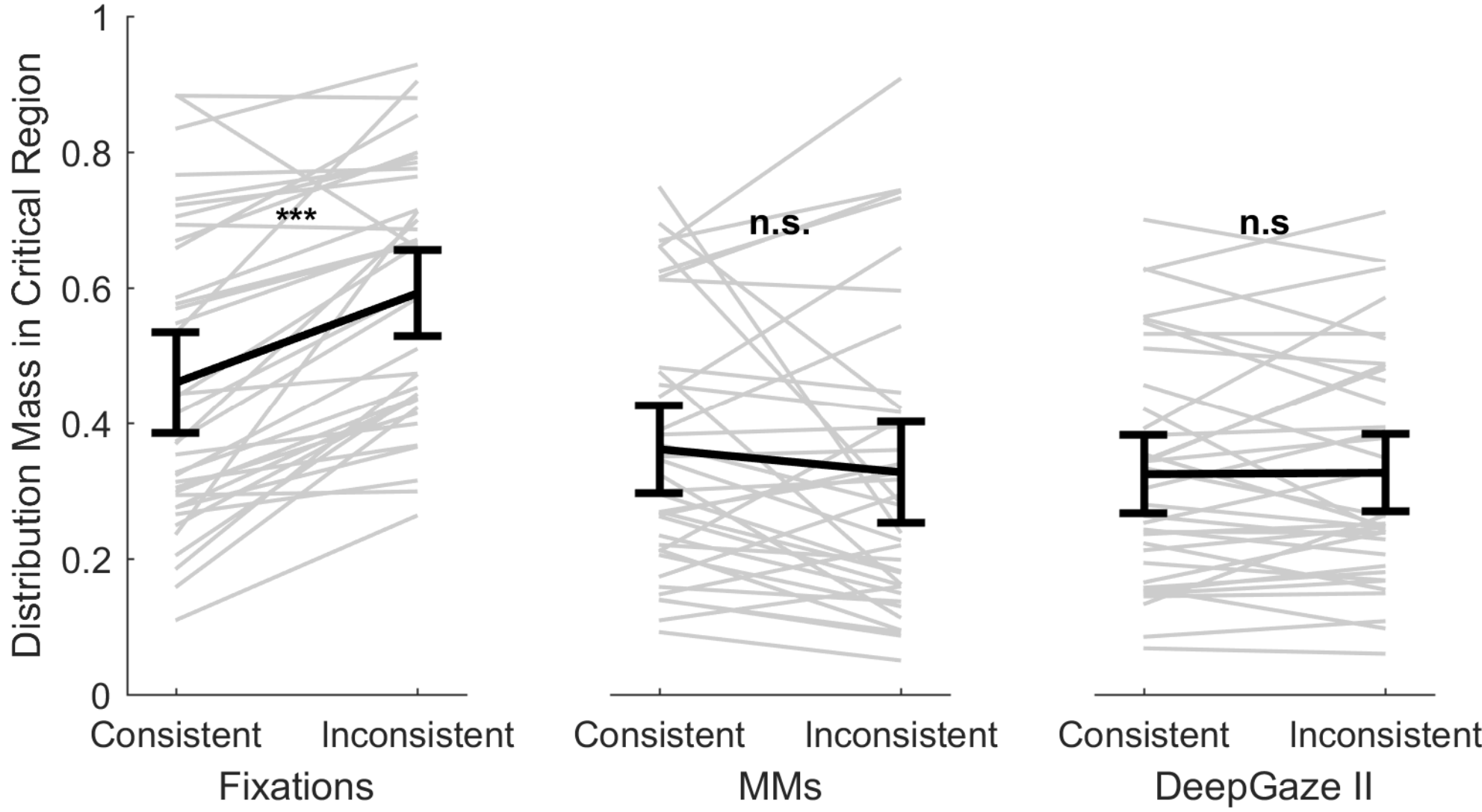


364 Fig. 5. Normalized distribution mass falling within Critical Regions in both conditions for (a)
365 smoothed human fixations, (b) MMs, and (c) DGII. All figure characteristics are as in Fig. 3.

366

367 Additionally, conditions differed regarding the number of fixations per image, $t(35) = 5.67$ $p$
368 $< .001$. On average, there were 6% fewer fixations in the Inconsistent condition. This excludes

369  the possibility that higher number of fixations in this condition might drive the observed

370  increase in the distribution mass falling within the Critical Regions.

371  Finally, systematic differences in object size between Consistent and Inconsistent conditions

372  could affect our results because larger objects may attract more fixations solely because they

373  occupy a larger image area. However, this factor was minimized by showing each object in a

374  consistent and an inconsistent context. Yet, the same object might be shown in a slightly

375  different position in the two conditions and might therefore occupy slightly different amounts

376  of the image. This was, however, not the case: the JZS Bayes Factor of 4.26 indicated that the

377  two conditions did not differ in the size of the bounding boxes of each manipulated object

378  (objects in the Inconsistent condition were on average 1562.28 pixels larger; 95% confidence

379  interval: [-2582.74, 5707.29]).

380  To summarize, semantic changes induced by altering object-context relationships elicited

381  changes in distributions of human fixations, but neither MMs nor DGII could predict them.

382  These results suggest that both models might be sensitive to image features, which are

383  frequently correlated with image meaning, rather than to meaning itself.

384

# Discussion

385

386  A long-standing debate in visual perception concerns the extent to which visual features vs.

387  semantic content guide human eye-movements in free viewing of natural scenes. To

388  distinguish these hypotheses, indexing the distributions both of features and meaning across

389  an image is critical. While image-based saliency models have been used to index features for

390  two decades, measuring semantic importance has been difficult until meaning maps (MMs)

391  have recently been proposed. Here, we assessed the extent to which MMs indeed capture

392  the distribution of meaning across an image. First, we demonstrate that despite the

393  purported importance of meaning as measured by MMs for gaze control, MMs are not better

394  predictors of locations of human fixations than at least some saliency models, which are based

395  solely on image features. In fact, DeepGaze II (DGII), a model using deep neural network

396  features, outperformed MMs. Second, we assessed the sensitivity of human eye-movements,

397  MMs, and DGII to changes in image meaning induced by violations of typical object-context

398 relationships. Observers fixated more often on regions containing objects inconsistent with

399 scene context (thus replicating previous findings) but these regions were not indexed as more

400 meaningful by MMs, or as more salient by DGII. Together, these findings challenge central

401 assumptions of MMs, suggesting that they are insensitive to the semantic information

402 contained in the stimulus.

403 The good performance of DGII in predicting human gaze might be attributable to the high-

404 level features it extracts from images. Three other models, which use low-level features,

405 failed to decisively outperform MMs. However, unlike two of them (GBVS and AWS), DGII is

406 trained with data on human fixations to optimize performance (Kümmerer et al., 2016, 2017).

407 Yet, training alone cannot explain the difference in performance. The third low-level feature

408 model (ICF) is trained in the same way (Kümmerer et al., 2017) but still achieves a lower

409 performance than DGII. These findings suggest that feature type is indeed critical for a

410 model's performance. Importantly, however, while DGII uses high-level features transferred

411 from a deep neural network trained on object recognition (Simonyan & Zisserman, 2014), this

412 is not equivalent to indexing meaning. Rather, the good performance of DGII is likely due to

413 meaning supervening on, or correlating with, some of the features indexed by this model.

414 Correlation between visual features and meaning as the source of good performance in

415 saliency models has already been considered by the authors of MMs (Henderson & Hayes,

416 2017). Our findings suggest that MMs might share this characteristic with saliency models.

417 Specifically, the ratings used to construct MMs might be based on visual properties in such a

418 way that highly structured patches that contain high-level features receive high ratings. These

419 features often correlate with meaning, but in and of themselves do not amount to meaning.

420 According to this interpretation, both DGII and MMs index high-level features. Their success

421 in predicting human behavior derives from the typically strong correlation between high-level

422 features and meaning, with a higher correlation for the features extracted by DGII than MMs.

423 An alternative interpretation of the finding that DGII outperforms MMs is that image features

424 rather than meaning guide human fixations. However, this interpretation is inconsistent with

425 our second analysis. Here, observers clearly exhibited sensitivity to meaning, as indicated by

426 changed gaze patterns after introducing semantic inconsistencies into the scenes. This

427 experimental manipulation targets a type of meaning that is based on how objects relate to

428 the broader context in which they occur. While specific, it is precisely this kind of meaning

17

429     that is of high theoretical importance in eye-movement research (Henderson, 2017;

430     Henderson et al., 2009). Thus, even if MMs were to measure other types of meaning, as has

431     been suggested (Henderson et al., 2018), the fact that they are not sensitive to meaning

432     derived from object-context relationships seriously limits their usefulness. Moreover, the idea

433     that MMs indeed index other kinds of meaning that are important for guidance of fixations is

434     not consistent with our findings. If this were the case, then we would expect MMs to predict

435     human fixations better than saliency models that solely rely on image features, which is not

436     the case.

437     The insensitivity to semantic inconsistencies reveals inherent limitations of both MMs and

438     DGII. The way in which MMs are constructed implicitly assumes that meaning is a local image-

439     property, which is not true for object-context (in)consistency. This limitation may potentially

440     be alleviated by 'contextualized MMs' (Peacock, Hayes, & Henderson, 2019), a recently

441     suggested modification of the 'standard' MMs. These novel maps are created from

442     meaningfulness ratings by observers who see the whole scenes from which the to-be-rated

443     patches were derived. It is yet to be seen what this approach can reveal about fixation

444     selection beyond the fact that humans asked to indicate meaningful or interesting regions

445     within scenes highlight areas, which tend to be frequently fixated by other observers

446     (Nyström & Holmqvist, 2008; Onat et al., 2014). DGII, in turn, does not explicitly encode

447     semantic information, and was not trained on the relationship between eye movements and

448     semantic (in)consistency. But its failure highlights an opportunity to improve saliency models

449     by incorporating semantic relationships (Bayat, Koh, Nand, Pereira, & Pomplun, 2018).

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

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

452     an image. Instead of meaning per-se, they seem to index high-level features that have the

453     potential to carry meaning in typical natural scenes. They share this characteristic with state-

454     of-the-art saliency models, which are easier to use, do not require human annotation, and yet

455     predict locations of human fixations better than MMs.

456

457

458

459 **Author contributions**

460 C.T. and M.P. conceived of the study. M.P., T.W., and C.T. designed the experiment. M.P.

461 collected and analyzed the data under the supervision of C.T. with support from M.K., T.W.,

462 and M.B. The paper was drafted by M.P. and C.T.; T.W., M.K., and M.B. provided detailed

463 comments.

464

465 **References**

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

486  Hayhoe, M., & Ballard, D. (2005). Eye movements in natural behavior. *Trends in Cognitive*
487  *Sciences*, *9*(4). https://doi.org/10.1016/j.tics.2005.02.009

488  Hegde, J., & Kersten, D. (2010). A Link between Visual Disambiguation and Visual Memory.
489  *Journal of Neuroscience*, *30*(45), 15124–15133.
490  https://doi.org/10.1523/JNEUROSCI.4415-09.2010

491  Henderson, J. M. (2017). Gaze Control as Prediction. *Trends in Cognitive Sciences*, *21*(1), 15–
492  23. https://doi.org/10.1016/j.tics.2016.11.003

493  Henderson, J. M., & Hayes, T. R. (2017). Meaning-based guidance of attention in scenes as
494  revealed by meaning maps. *Nature Human Behaviour*, *1*(October).
495  https://doi.org/10.1038/s41562-017-0208-0

496  Henderson, J. M., & Hayes, T. R. (2018). Meaning guides attention in real-world scene
497  images: Evidence from eye movements and meaning maps. *Journal of Vision*, *18*(6), 10.
498  https://doi.org/10.1167/18.6.10

499  Henderson, J. M., Hayes, T. R., Peacock, C. E., & Rehrig, G. (2019). Meaning and Attentional
500  Guidance in Scenes : A Review of the Meaning Map Approach. *Vision*, *3*(2).

501  Henderson, J. M., Hayes, T. R., Rehrig, G., & Ferreira, F. (2018). Meaning Guides Attention
502  during Real-World Scene Description. *Scientific Reports*, *8*(1), 13504.
503  https://doi.org/10.1038/s41598-018-31894-5

504  Henderson, J. M., Malcolm, G. L., & Schandl, C. (2009). Searching in the dark: Cognitive
505  relevance drives attention in real-world scenes. *Psychonomic Bulletin & Review*, *16*(5),
506  850–856. https://doi.org/10.3758/PBR.16.5.850

507  Henderson, J. M., Weeks, P. A., & Hollingworth, A. (1999). The effects of semantic
508  consistency on eye movements during complex scene viewing. *Journal of Experimental*
509  *Psychology: Human Perception and Performance*, *25*(1), 210–228.
510  https://doi.org/10.1037/0096-1523.25.1.210

511  Itti, L., & Koch, C. (2000). A saliency-based search mechanism for overt and covert shifts of
512  visual attention. *Vision Research*, *40*(10–12), 1489–1506.
513  https://doi.org/10.1016/S0042-6989(99)00163-7

514    Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews*
515        *Neuroscience*, *2*(3), 194–203. https://doi.org/10.1038/35058500

516    Kietzmann, T. C., McClure, P., & Kriegeskorte, N. (2019). Deep Neural Networks in
517        Computational Neuroscience. In *Oxford Research Encyclopedia of Neuroscience*.

518    Kleiner, M., Brainard, D., & Pelli, D. G. (2007). What's new in psychtoolbox-3? *Perception*,
519        *36*(1).

520    Koehler, K., Guo, F., Zhang, S., & Eckstein, M. P. (2014). What do saliency models predict?
521        *Journal of Vision*, *14*(3). https://doi.org/10.1167/14.3.14

522    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2015). Information-theoretic model
523        comparison unifies saliency metrics. *Proceedings of the National Academy of Sciences*,
524        *112*(52), 16054–16059. https://doi.org/10.1073/pnas.1510393112

525    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2016). DeepGaze II: Reading fixations from
526        deep features trained on object recognition, 1–16. Retrieved from
527        http://arxiv.org/abs/1610.01563

528    Kümmerer, M., Wallis, T. S. A., & Bethge, M. (2018). Saliency Benchmarking Made Easy:
529        Separating Models, Maps and Metrics. In V. Ferrari, M. Hebert, C. Sminchisescu, & Y.
530        Weiss (Eds.), *Computer Vision – ECCV 2018. ECCV 2018. Lecture Notes in Computer*
531        *Science* (Vol. 11220, pp. 798–814). Springer. https://doi.org/10.1007/978-3-030-01270-
532        0_47

533    Kümmerer, M., Wallis, T. S. A., Gatys, L. A., & Bethge, M. (2017). Understanding Low- and
534        High-Level Contributions to Fixation Prediction. In *The IEEE International Conference on*
535        *Computer Vision (ICCV)*. https://doi.org/10.1109/ICCV.2017.513

536    Nyström, M., & Holmqvist, K. (2008). Semantic override of low-level features in image
537        viewing–both initially and overall. *Journal of Eye Movement Research*, *2*(2), 1–11.
538        https://doi.org/10.16910/jemr.2.2.2

539    Öhlschläger, S., & Võ, M. L. H. (2017). SCEGRAM: An image database for semantic and
540        syntactic inconsistencies in scenes. *Behavior Research Methods*, *49*(5).
541        https://doi.org/10.3758/s13428-016-0820-3

Pedziwiatr et al.

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

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

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

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

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

547        6989(01)00250-4

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

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

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

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

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

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

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

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

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

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

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

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

560        225–237. https://doi.org/10.3758/PBR.16.2.225

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

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

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

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

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

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

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

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

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

Pedziwiatr et al.

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

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

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

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

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

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

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

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

578

23