
Using users' physiological responses for the estimation of websites' aesthetic judgments

Giulio Gabrieli^{1,*}, Marc H. Bornstein^{2,3} and Gianluca Esposito^{1,4}

¹Psychology Program, School of Social Sciences, Nanyang Technological University, Singapore

²Eunice Kennedy Shriver National Institute of Child Health and Human Development, Bethesda, MD, USA

³Institute for Fiscal Studies, London, UK

⁴Department of Psychology and Cognitive Science, University of Trento, Italy

Correspondence*:

Gianluca Esposito

gianluca.esposito@ntu.edu.sg

2 ABSTRACT

3 The aesthetic appearance of websites can influence the perception of their usability, reliability,
4 and trustworthiness. Several studies investigated the relationship between single aesthetic
5 features and explicit aesthetic judgments, demonstrating the existence of an attribution bias.
6 However, only a limited amount of studies focused on the interaction between multiple visual
7 properties and have considered not only explicit ratings, but also implicit judgments. In this
8 work, we employ a novel approach, based on the analysis of physiological signals (implicit
9 measures) and the application of machine learning and neural network models to predict users'
10 perceived aesthetic pleasure from the empirical analysis of web pages' advanced visual properties
11 (e.g. symmetry, visual complexity, colorfulness, ratio between visual and textual areas). Young
12 adults ($N=59$, 33 females, Mean age = 21.52 years) assessed the aesthetic appeal of websites
13 and emotional pictures while their physiological activity was recorded. Results using recursive
14 partitioning and generalized linear models demonstrate the possibility of predicting the average
15 aesthetic rating of a website using both explicit (behavioral ratings) and implicit measures
16 (physiological activities).

17 **Keywords:** web design, aesthetics, physiology, eeg, eda, emg, pupillometry, machine learning, neural networks

1 INTRODUCTION

18 People interact with websites daily for work, educational purposes, and recreation. With the advent of
19 more powerful technologies and an increasing number of active users, a new approach to the design and
20 development of web pages was born, focusing no longer exclusively on content and functionalities, but
21 also on pages' aesthetic appearance, which, as defined by Moshagen and Thielsch (2010) is "*an immediate
22 pleasurable subjective experience that is directed towards an object*".

23 Since the turn of the century, web design practices have evolved, encompassing a variety of disciplines,
24 including visual design, user interface design (*UI*), user experience design (*UX*), scripting, programming,
25 and content strategy (Robbins, 2012). User satisfaction is one of the many goals web designers aim to
26 achieve, because satisfied users are more likely to spend more time on a page, come back to the same

27 website in the future, and recommend a website to other possible users (Zhang and Von Dran, 2000).
28 Users' evaluations of interactive systems are influenced by their visual appearance, and this is especially
29 true for web pages (Karvonen, 2000; Kim et al., 2003; Zhang et al., 2001). In social psychology, the
30 influence of aesthetic factors on other attributes is called "*halo effect*" or confirmation bias (Lindgaard
31 et al., 2006). For example, individuals with more aesthetically pleasing faces are also perceived as more
32 trustworthy.

33 Currently, the evaluation of a design artifact is an iterative process that requires time, money and external
34 individuals who are asked to evaluate artifacts at different stages of the design process. Therefore,
35 researchers tried to investigate possible ways to reduce the costs associated with the evaluation using
36 heuristics or machine-based approaches.

37 Several attempts have been made to predict the perceived aesthetic perception of web pages, using a limited
38 number of visual features as well as behavioral measures (Reinecke et al., 2013). Despite this, a limited
39 number of studies have considered the role of users' exposure to different pages and their expertise in
40 designing websites on their aesthetic evaluations. Moreover, a majority of the studies focused only on
41 explicit behavioural measures.

42 In this study, we factored in multiple aesthetic features simultaneously and considered users' characteristics
43 in different neural networks and machine learning models. Perceived aesthetic appeal of websites were
44 estimated with both behavioral (self-reported; explicit measures) and physiological (ECG, EMG, EDA,
45 and Pupillometry; implicit measures) measures, to overcome possible limits of self-report measures.

46

47 **1.1 Visual properties and aesthetic judgment**

48 Several features are known to affect website aesthetics judgments, including but not limited to visual
49 complexity, perceived colorfulness, and symmetry (Cyr et al., 2010; Karvonen, 2000; Miniukovich and
50 De Angeli, 2014). For example, Reinecke et al. (2013) assessed the impact of visual complexity and
51 colorfulness on users' first impression of 450 websites. Results of their computational models show that
52 visual complexity and colorfulness accounted for about half the variance in aesthetic judgments of web
53 pages.

54 Visual complexity, also called visual cluttering, is a widely examined factor in aesthetic decision making,
55 such that the more visually complex a design is, the higher the probability that it will be rated as more
56 aesthetically pleasing (Seckler et al., 2015a; Tuch et al., 2009, 2012).

57 Despite its wide adoption in the literature, there is no single nor standard way to estimate a website's
58 visual complexity. Some authors computed visual complexity by the weight of still-image (the weight
59 of the file, expressed in kB or MB), whereas others define it by the space taken up by text and images,
60 by calculating the number of colors, or by counting the number of images in a page (Bucy et al., 1999;
61 Ivory et al., 2001). A promising method, proposed by Zheng et al. (2009), is based on a technique called
62 Quadratic Tree Decomposition (QTD), often abbreviated as *quadtree decomposition*. The QTD recursively
63 divides (horizontally and vertically) an image into areas of smaller size, if the parent area has a complexity
64 -measured in terms of standard deviation of the area- higher than a predefined threshold. The final number
65 of obtained squares, called *leaves*, is used as an index of visual complexity. When comparing images of the
66 same size analyzed using the same complexity threshold value, the higher the number of leaves, the higher
67 the visual complexity of an image. An example of Quadratic Tree Decomposition applied to websites is
68 shown in Figure 1.

69

70 Similarly, color perception has been widely investigated in psychology, especially in relation to emotional
71 valence and arousal (Wang and Ding, 2012). In Human-Computer Interaction, colors have an influence



Figure 1. Visual representation of a Quadtree decomposition applied to a website of the AVI14 dataset. Visual complexity of an area is proportional to the number of leaves in that area.

72 on perceived trust, loyalty, and economic behavior (Cyr, 2008; Kim and Moon, 1998). The color
 73 scheme of a page can impact a user's feelings and reactions towards a page because specific colors
 74 have been demonstrated to increase—or reduce—the viewers' arousal and therefore induce excitement
 75—or relaxation—. Cooler colors are often preferred to warmer colors because they elicit relaxed feelings
 76 (Cyr et al., 2010; Hall and Hanna, 2004; Hasler and Suesstrunk, 2003; Jacobs and Hustmyer Jr, 1974). A
 77 color is composed of a hue, a level of saturation, and a value (often defined as with brightness or luminance)
 78 (HSV model). For instance, Yendrikhovskij et al. (1998) found a correlation ($r^2 = 0.91$) between users'
 79 perceived colorfulness and the sum of the average saturation value of an image and its standard deviation.
 80 Additionally, a webpage's color distribution has also been proven to affect perceived brightness and
 81 perceived colorfulness (Reinecke et al., 2013). In this study, luminance is expressed as relative luminance,
 82 as per the photometric definition (Birtolo et al., 2009).

83 Another widely adopted feature is symmetry. First introduced by Gestalt's psychologists, symmetry has
 84 been proven to be one of the most important factors in aesthetic judgment. Symmetry indicates how well
 85 one side of an image reflects the opposite side, and it can be evaluated along a horizontal axis (top versus
 86 bottom), a vertical axis (left versus right), a radial plane around the center of the image, or using QuadTree
 87 decomposition (Miniukovich and De Angeli, 2014; Reinecke et al., 2013; Wang and Li, 2016; Zheng et al.,
 88 2009). An example of symmetry estimation using QuadTree decomposition is shown in Figure 2.

89

90 Features based not only on the appearance but also on the type of content have been proposed as well.
 91 Lin et al. (2013), for example, introduced the adoption of the ratio between graphics and text. In this work,

HIGH SYMMETRY



LOW SYMMETRY

Figure 2. Example of symmetry estimation using QuadTree Decomposition (QTD) applied to a website of the AVI14 dataset. QTD is applied to the image, divided in two halves. The degree of symmetry is given by the number of overlapping rectangles between the two images. In the example above, a high degree of symmetry is present only in the upper part of the image.

92 we used a novel method for the automatic estimation of the graphics to text ratio, based on a combination
 93 of a Space-Based Decomposition (SBD) algorithm and an Optical Character recognition system (OCR).
 94 Similar to QTD, SBD uses a recursive division of the image to identify the contours of elements within
 95 an image, namely text and graphics. Once the elements have been identified, we can apply an OCR to
 96 label each element. Finally, the ratio between the area labeled as text and graphics can be automatically
 97 computed using a machine-driven approach. This same technique allows the automatic estimation of the
 98 number of images present on a page.

99 In the analysis of website visual features, it is important to rely on objective measures. Visual complexity,
 100 perceived colorfulness, graphics to text ratio, number of images, and symmetry are all automatic estimable
 101 features that can be computed in an objective algorithm.

102

103 1.2 Physiology and aesthetic judgment

104 Until recently, researchers were only able to investigate the underlying physiological correlates of
 105 aesthetic appreciation through behavioral measures of patients suffering from neurodegenerative diseases
 106 or whose brains suffered damage (Cela-Conde et al., 2011).

107 Today, researchers can investigate neurophysiological signals in a more ecological way from healthy
108 participants, using sensors applied to the surface of the body. In an electromyography (EMG) study (which
109 investigates muscles' electrical signals), Winkielman and Cacioppo (2001) demonstrated that physiological
110 measures reflect participants' affective responses to stimuli and implicit judgments of their beauty. The
111 activity of the zygomaticus major correlates with positive affective responses, and activity of a region in
112 the *corrugator supercilii* correlates with negative affective responses (Lang et al., 1993; Winkielman and
113 Cacioppo, 2001). Electrocardiography (ECG), which measures the electrical activity of the heart, also
114 shows relations between physiological responses and aesthetic judgments (de Jong, 1972; de Jong et al.,
115 1973; Ray et al., 1997).

116 In an eye-tracking study (Yanulevskaya et al., 2012), participants focused on emotionally positive parts of
117 pictures. Maughan, Gutnikov, and Stevens Maughan et al. (2007) found that positive aesthetic judgments of
118 advertisements elicited sustained attention. In addition, pupil dilation in response to pleasant images, and
119 pupil constriction in response to unpleasant images were found by Blackbourne and Schirillo Blackburn and
120 Schirillo (2012). Similarly, both ECG and EMG signals have been proven to be suitable for the empirical
121 analysis of websites' aesthetic features, as shown by Tuch et al. (2009).

122

123 **1.3 Behaviour and aesthetic judgments**

124 The analysis of participants' behavioral data (explicit ratings) has been widely adopted in previous studies
125 that investigated different aspects of websites, including their complexity and aesthetic qualities. Reinecke
126 et al. (2013), for example, employed a 9-point Likert scale to assess participants' first impressions of
127 a website's aesthetic quality, while Seckler et al. (2015b) investigated different aesthetic facets using a
128 7-point Likert scale. Those results show that by collecting self-reported measures using a Likert scale, we
129 can obtain a reliable estimate of the perceived aesthetic judgments.

130

131 **1.4 Expertise, exposure, and aesthetic judgment**

132 The mere exposure effect states that repeated exposure to a target enhances an individual's attitude towards
133 it (Zajonc, 1968; Bornstein and D'agostino, 1992). Cox and Cox (2002) found that repeated exposure to
134 a visually complex product design increased preference for it as compared to a simpler but novel design.
135 Exposure effects have also been found to evoke positive affective responses, where participants who rated
136 familiar targets as more likable than unfamiliar ones also showed more zygomatic muscle region activity
137 when viewing familiar targets (Harmon-Jones and Allen, 2001). These results suggest that individuals'
138 exposure to different websites needs to be considered when evaluating their aesthetic judgment. Since many
139 websites adopt similar designs and layouts, it is possible that not only the mere exposure to a single stimuli,
140 but also a general exposure to many different websites can play a role in shaping users' design preferences.
141 Similarly, expertise in a field affects preferences: experts and laypersons have different preferences and
142 make different aesthetic judgments (Müller et al., 2010; Orr and Ohlsson, 2005; Ulrich Kirk, 2009; Pihko
143 et al., 2011). Quispel et al. (2016) found that experts preferred familiar and novel chart designs, but
144 laypersons preferred familiar and easy-to-use designs. In addition, familiarity and perceived ease of use
145 predict the attractiveness of designs among laypersons but not experts. Bölte et al. (2017) evaluated experts'
146 and laypersons' event-related potentials to web pages: Experts more frequently rated aesthetic web pages
147 as less aesthetic than laypersons. This difference was not found in ratings of unaesthetic web pages. Given
148 the history of findings on the role of expertise in evaluating aesthetics, it is also important to consider the
149 impact of expertise in judging the aesthetic properties of a web page.

150

151 The research question on which this project is built is based on the possibility of reducing the cost, both
152 in terms of time and economical expenses, of testing the perceived aesthetic experience of a web page.
153 Results of this project may be used, in the future, to create novel technologies that will be able to support
154 designers by providing them with continuous evaluations of design artifacts, at a reduced cost.

155

156 **1.5 Purpose of the study and Hypothesis**

157 Despite an array of studies on website aesthetics (Reinecke et al., 2013; Seckler et al., 2015a; Miniukovich
158 and De Angeli, 2014; Bölte et al., 2017; Tuch et al., 2009), many have focused on limited individual
159 factors such as visual complexity or colorfulness. Few attempts have been made to study the effects of
160 multiple different visual features together on overall perceived aesthetic. Past studies have also failed to
161 consider how user exposure and expertise might affect website aesthetic judgment. Furthermore, even
162 though physiological measures have been employed previously, few have employed them to predict website
163 aesthetic judgments from multiple visual properties.

164 In this work, we apply a novel approach based on neural networks and machine learning models as well as
165 recursive partitioning and generalized linear models to estimate the perceived aesthetic appeal of a website.

166 Our proposed flow chart is illustrated in Figure 3.

167 We hypothesized that (1) the interaction between web pages' different visual properties (visual complexity,
168 colorfulness, brightness, symmetry, and text ratio) can be used to predict behavioral ratings and
169 physiologically-estimated aesthetic judgments. We also hypothesized that (2) exposure to websites
170 moderates website aesthetic judgments. Last, we hypothesized that (3) expertise on website design,
171 similarly, moderates website aesthetic judgments.

172

2 METHODS

173 **2.1 Analytic Plan**

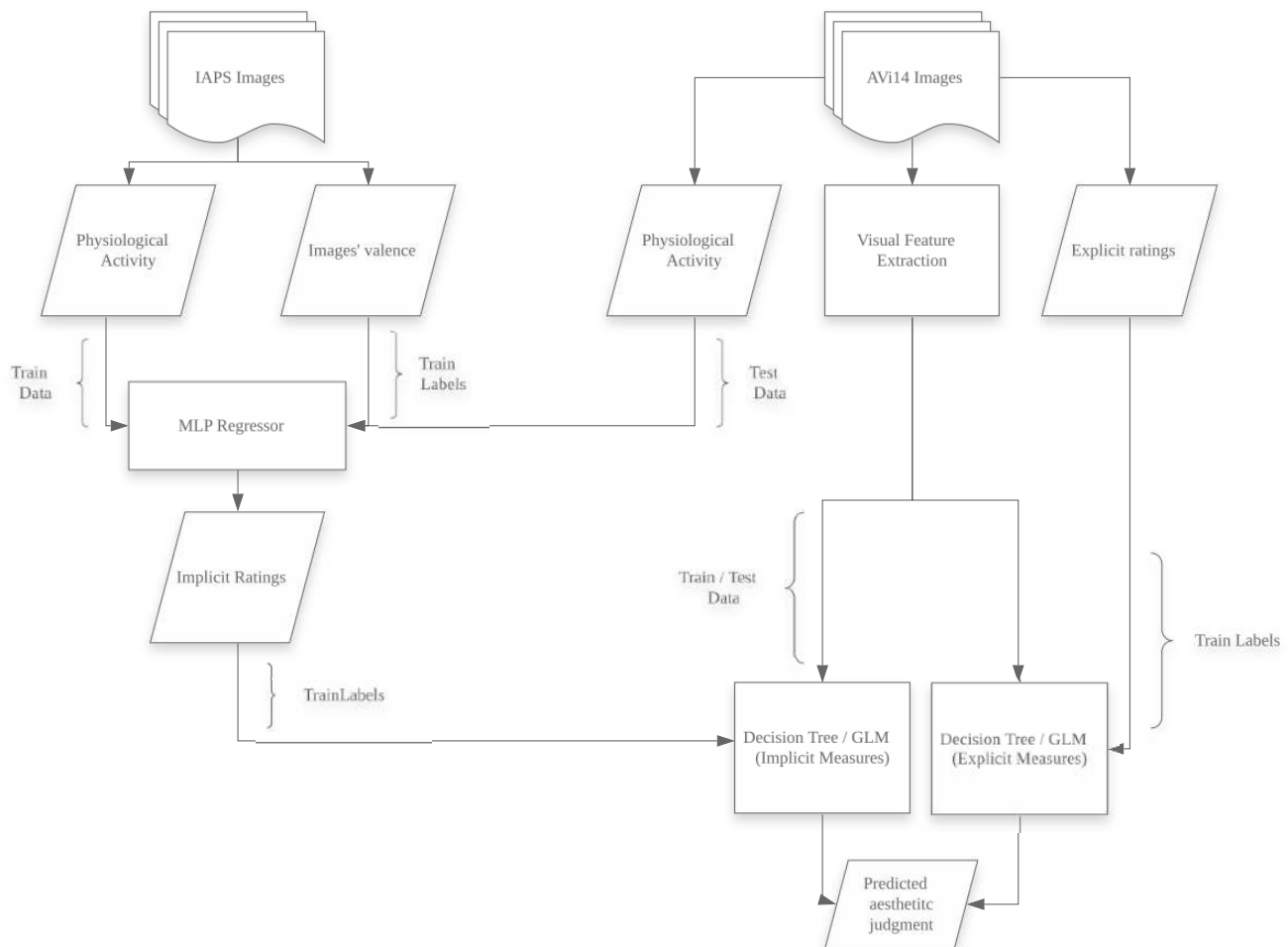
174 This work is structured as follows. First, an experimental procedure was conducted to collect behavioral
175 ratings and physiological activity of participants in an image-rating task (5-point Likert scale), where
176 participants rated screenshots of website and emotional images. Then a Multi-Layer Perceptron Neural
177 Network (MLP NN) was trained on features estimated from the physiological activity of the participants,
178 using the standardized valence values of the emotional images as training labels. The model was then
179 applied to features estimated from the physiological activity recorded while participants were exposed to
180 website images, resulting in the estimation of a valence value for the website images (implicit measures).
181 Having both the behavioral ratings and the ratings estimated from the physiological activity, we proceeded
182 with extracting a set of relevant features from the websites' images.

183 Since we know from previous studies that male and female participants may rate the aesthetic appeal of a
184 website differently, we first excluded from the analysis all the websites that received significantly different
185 ratings by male and female participants.

186 Finally, two different machine learning models —GLM and Decision Tree— were applied using websites'
187 visual features as input and the ratings —behavioral or physiological—, as labels. To reduce the influence
188 of a single participant on the overall accuracy of the model, bootstrapping is employed.

189 Performances of the models were tested not only against the 5-point Likert scale values, but also on a
190 binomial rating, obtained by clustering ratings into two groups (High ratings: 4-5, Low ratings: 1-3). This
191 was done to verify whether the models can be employed to obtain binary classifications (good/bad) of the
192 aesthetic of a webpage.

193 Then, to compare the possible differences between experts and non-expert designers and between highly



194 exposed and lowly exposed users, we assigned our participants into two groups of about the same size and
 195 **Figure 3.** Flowchart of the method employed in this project.
 196 use these covariates as factors in our models.

197 2.2 Participants

198 59 university students (33 females, *Mean age* (in years) = 21.5 ± 3.0) voluntarily enrolled for participation.
 199 Informed consent was obtained from all the participants prior to the experimental session. The study was
 200 conducted in accordance with the declaration of Helsinki.

201

202 2.3 Stimuli

203 Stimuli were selected from two different datasets: the International Affective Pictures System (IAPS)
 204 (Lang and Bradley, 2007) and the AVI14 (Miniukovich and De Angeli, 2014).

205

206 2.3.1 International Affective Picture System

207 The International Affective Picture System (IAPS) (Lang and Bradley, 2007) is a dataset of emotionally
 208 evocative pictures, developed by the NIMH Center for the Study of Emotion and Attention (University of
 209 Florida). From the 1180 pictures included in the IAPS dataset, 50 were selected for presentation in the

210 experimental procedure, 25 per block, balancing the mean valence value of each block¹. The dataset is
211 available upon request from their original authors².

212

213 2.3.2 AVI14

214 The AVI14 dataset is composed of images of 140 websites (Miniukovich and De Angeli, 2014). The
215 dataset is available online³. All the websites are in English, and the original pages have no dynamic effects.
216 Majority of the websites ($N = 115$) were selected from a public showcase of beautiful websites, and another
217 25 were selected to balance the overall aesthetic of the dataset. Used pages belong to four categories:
218 a) coffee, b) chocolate bars and shops, c) online retailers, and d) design agencies. For the purposes of
219 this work, 100 pictures were selected for presentation from the AVI14 dataset, according to their mean
220 perceived aesthetic pleasure value (Miniukovich and De Angeli, 2014) and divided semi-randomly into
221 two sets, one set per block.

222

223 2.4 Instrumentation

224 Stimuli were presented on a DELL 29" Ultrasharp Screen (U2719WM) with a fixed resolution of
225 1920x1080, (refresh rate = 60.00Hz). Pupil dilation signals were recorded using a Tobii X3-120 (sampling
226 rate: 120Hz, Tobii Technology) mounted on a tripod and placed just below the screen. ECG, EDA and
227 EMG signals were collected using a Bitalino Revolution BT board (sampling rate: 1000Hz, Wireless
228 Biosignals S.A) (Guerreiro et al., 2013; Batista et al., 2017), using disposable 36-40mm snap connector
229 foam electrodes (F9089/100, FIAB, Florence, Italy). The experimental paradigm and registration of
230 physiological measurements were implemented in Python 2.7⁴ (Oliphant, 2007; Van Rossum and Drake,
231 2011; Oliphant, 2006).

232

233 2.5 Experimental procedure

234 Participants sat approximately 50 to 70 cm away from a computer screen, in an silent and dark
235 environment.

236 Before the experimental sessions, participants were instructed on the tasks they had to perform and on the
237 physiological measures that were to be collected. The experiment consisted of two blocks, presented one
238 after the other in a semi-randomized order, with a brief pause between blocks. Each picture was presented
239 for 6 seconds, with an 8 second interval between consecutive images. A graphical representation of the
240 used procedure is shown in Figure 4A.

241

242 To record participants' physiological activity, two disposable electrodes were used to record the
243 electrodermal activity (EDA) from the left wrist, three were used to record the heart activity (ECG)
244 —one below each clavicle and one below the last rib— and three were used to record the electromyographic
245 activity (EMG) of the *corrugator supercilii* —one above the nose, one above the left eye and one on the
246 left cheek—. A graphical representation of the electrodes position is reported in Figure 4B

247 Immediately after each picture, participants rated their aesthetic appeal on a 5-point Likert scale, by clicking
248 on one of five buttons presented on the screen, with no time constraint.

¹ Pictures belonging to the following categories were removed prior to stimuli selection: "BurnVictim", "Mutilation", "DeadBody", "DeadMan", "headlessBody", "BabyTumour", "Tumor", "Accident", "SlicedHand", "Vomit", "BatteredFem". Remaining pictures were sorted by mean valence and the first and last 25 were semi-randomly selected and distributed in two sets, one per experimental block

² <https://csea.php.ufl.edu/Media.html>

³ <https://github.com/aliko-str/avi14dataset>

⁴ v. 2.7.12

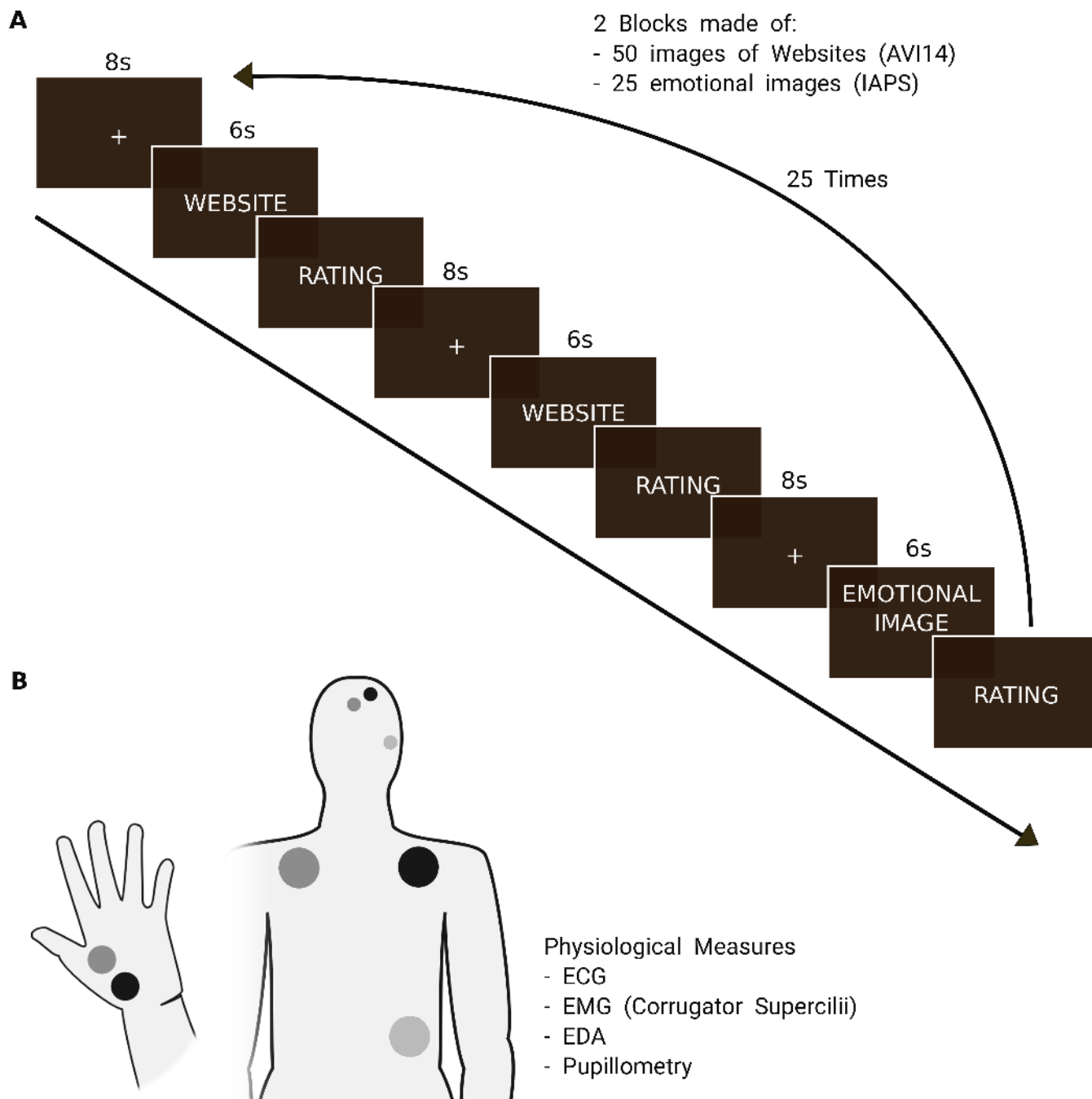


Figure 4. Graphical representation of the (A) experimental paradigm and (B) electrodes position for physiological signals recording

249 At the end of the experimental procedure, participants completed a 7-item survey on their browsing habits
 250 (exposure) and expertise in design, development, and management of websites. Finally, participants were
 251 debriefed.
 252

253 **2.6 Data modeling and analysis**

254 2.6.1 Website Feature Extraction

255 Website features were extracted using a self-developed tool released under the name "PrettyWebsite"
256 (Gabrieli, 2019a). The package is available through the Python Package manager (Pypi) and the project
257 repository ⁵.

258

259 2.6.2 Physiological Feature Extraction

260 Physiological features were extracted from collected signals using "Pysiology (Gabrieli et al., 2020;
261 Gabrieli, 2019b), a Python package designed for physiological signal processing.

262 For each stimulus, physiological measures were computed in epochs of 8 seconds. For ECG signals, the 20
263 seconds of recording before the first stimulus of each block was used as a baseline. For pupil diameter,
264 signals 6 seconds preceding each stimulus served as a baseline. Detailed information about the parameters
265 used to clean the signals and estimate features are reported in Supplementary Materials.

266

267 2.6.3 Estimation of implicit ratings

268 Participants' physiological activity was used to estimate the perceived valence of websites' images.

269 Features extracted from the epochs in which participants were engaged in viewing images from the IAPS
270 dataset were used as training data of an MLP Regressor Solver = "sgd", $\alpha = 0.0001$, number of hidden
271 layers = 100), with the standardized valence values of the images, provided within the IAPS dataset, used
272 as training labels. To reduce the number of input features, the best six physiological components were
273 identified through Principal Component Analysis, standard scaled and fed to the model.

274 Once trained, the average accuracy of the model was tested, using bootstrapping (N=100) against real
275 IAPS' valence value.

276 Finally, the model was fed with features extracted from the portions of signals where participants were
277 rating website pictures, in order to obtain an estimated implicit valence value (implicit rating) for each
278 website.

279

280 2.6.4 Preliminary analysis

281 **2.6.4.1 Expertise and exposure**

282 To compare differences between high and low expertise, and between high- and low-exposure users, each
283 participant was assigned to one of the two groups for each classification. Assignment was done by defining
284 a threshold that allowed the authors to obtain groups of similar sizes.

285

286 **2.6.4.2 Gender differences**

287 Previous studies highlighted the fact that males and females rate some websites with significantly different
288 scores. To omit gender of participants from the model, we conducted a preliminary analysis to identify if,
289 within our dataset, some of the websites received significantly different behavioral ratings by males and
290 females and subsequently removed those websites from our analysis.

291

292 **2.7 Predicting perceived visual aesthetic**

293 Prediction of the perceived visual aesthetic from estimated website features has been performed using
294 two different machine learning models, a generalized linear model, as implemented in *statsmodel* (Seabold

⁵ Pypi: <https://pypi.org/project/prettywebsite/>

Github: <https://github.com/Gabrock94/PrettyWebsite>

295 and Perktold, 2010) (GLM) and a recursive partitioning (Decision Tree, min samples per leaf=100, max
296 depth=5, max features=5), as implemented in Scikit-learn (Pedregosa et al., 2011).

297 In total, 24 different machine learning models were trained and tested using bootstrapping ($N = 1000$). For
298 each rating and physiological measure of aesthetic judgments, the two classifiers were used (GLM and
299 Decision Tree). Each classifier was tested 6 times, 3 times predicting values on a 5-point scale, and 3 times
300 predicting values on a binomial scale (1-3, 4-5). Finally, each of these 3 models was tested three times: one
301 with no reference to participants' expertise or exposure, one with expertise (high/low) as a factor of the
302 model and one with exposure (high/low) as a factor of the model.

303

3 RESULTS

304 Out of the data recorded from 59 participants, data of one participant ($N = 1$) was removed because of
305 technical issues in collected physiological samples. Therefore all the data described below are based on 58
306 participants ($N = 58$, $F = 33$, $M = 25$, *Mean age*: 21.4 ± 2.2).

307

3.1 Expertise and Exposure

308 With regard to expertise, ten participants ($N=10$) reported to have developed, and twelve ($N = 12$) to have
309 managed at least one website. Of the above, five ($N = 5$) reported having both developed and owned at
310 least one website. More than half of the participants ($N = 35$) reported having at least basic knowledge
311 of one or more programming languages. Participants who had at least basic knowledge of two or more
312 programming languages and have developed or managed a website were assigned to the "expert group" (N
313 = 29). Three websites (AVI 78, AVI 42 and AVI 128) received significantly different ratings by experts
314 and non-experts (Table 1). Thumbnails of those websites are reported in Supplementary Material (Figure
315 S1). Results from the t-tests showed that the power of the test was medium and only ratings given to AVI
316 128 were statistically significant using the KS-test. Of the three pages, the first two were rated on average
317 higher by expert users while the last was rated higher by the non-expert users.

318

Table 1. p-value (and power) of Student's t, F, Kolmogorov-Smirnov tests and means of the ratings for the Image that have been rated significantly different by the expert and non expert groups.

Image	p-val. (t)	Power (t)	p-val. (F)	p-val. (K-S)	Avg. E.	Avg. Non-E.
78.png	0.0118	0.682	0.281	0.0706	3.8	3.4
42.png	0.0366	0.536	0.713	0.1951	2.6	2.3
128.png	0.0136	0.689	0.489	0.074	2.6	3.0

319 With respect to exposure, 4 participants browsed the web using only either a laptop or desktop, and 11
320 browsed the web using only mobile devices. Almost half of the participants reported browsing up to 5
321 different websites per day, and 30 reported browsing more than 5 websites per day. More than half of the
322 participants spent less than 3 hours browsing websites ($N = 39$). Half of the participants ($N = 28$) reported
323 spending the majority of their time on a single website, such as Facebook or Twitter. Participants who
324 reported browsing 10 or more websites per day and who indicated browsing the web for more than 2 hours
325 per day were assigned to the "high exposure" group ($N = 27$). Two websites received significantly different
326 ratings by the "high exposure" and "low exposure" groups (Table 2). A previes of those websites is reported
327 in Supplementary Material (Figure S2).

Table 2. p-value (and power) of Student's t, F, Kolmogorov-Smirnov's tests and means of the ratings for the Image that have been rated significantly different by the high exposure and low exposure groups.

Image	p-val. (t)	Power (t)	p-val. (F)	p-val. (K-S)	Avg. E.	Avg. Non-E.
98.png	0.0439	0.621	0.239	0.338	2.7	3.1
20.png	0.0077	0.826	0.194	0.0453	3.7	3.3

328 3.2 Gender

329 Five websites received significantly different ratings from male and female participants. Results are
 330 reported in Table 3, while thumbnails of the images are reported in Supplementary Material (Figure S3).
 331 These websites were therefore removed from subsequent analysis.

Table 3. Results of t-test, Kolmogorov-Smirnov test and Fisher's test of websites with statistically significant differences between males and females.

Image	p-value (t-test)	Power (t-test)	p-value (F-test)	p-value (KS-test)
36.png	0.0017	0.764	0.339	0.041
101.png	0.0074	0.618	0.198	0.087
132.png	0.0236	0.506	0.423	0.224
66.png	0.0133	0.576	0.381	0.256
76.png	0.004	0.695	0.261	0.194

332 3.3 Website features

333 A visual feature — e.g. Visual complexity — can be estimated using different methods. In this work,
 334 where more than one algorithm was available, we adopted the most prominent. Therefore, the index of
 335 visual complexity used was based on the QDT (as opposed to the images' weight, $R^2 = 0.5$ between the
 336 two indexes), brightness was estimated from the *BT.709*⁶ index (as opposed to the *BT.601*⁷, ($R^2 = 1.0$)),
 337 and colorfulness was extracted from the HSV colorscheme (as opposed to the RGB colorscheme, ($R^2 =$
 338 0.58)), as done by Yendrikhovskij et al. (1998).

339 For our predicted models, we used *Symmetry*, *Colorfulness* (HSV), *Visual Complexity* (Quadratic Tree
 340 decomposition), *brightness* (BT709) and *number of Images* - automatically evaluated applying the Space-
 341 based decomposition and OCR- as independent variables.

342 3.4 Estimation of image valence from viewers' physiological activity

343 To obtain an estimation of an image's valence from a viewer's physiological activity, we used an MLP
 344 Regressor. First, extracted physiological features were used to estimate the valence of IAPS images.
 345 Average accuracy of MLP Regressor, tested against real IAPS' valence value, is 97.9% ($\sigma = 0.004$).
 346 Implicit ratings of website stimuli were estimated for 2792 epochs from 44 different participants (Mean
 347 number of stimuli per participant = 63.5 ± 14.5).

⁶ ITU-R Recommendation BT.709

⁷ ITU-R Recommendation BT.601

348 3.4.1 Predicting perceived visual aesthetic ratings

349 Average predictive accuracy of GLMs and Decision Trees are reported in Table 4. For the 5-point scale,
 350 no differences in the average prediction accuracy were reported between the explicit ratings and implicit
 351 appraisals when using GLM. However, for the same scale, when using recursive partitioning, tree-based
 352 models showed that implicit appraisals predicted better (73%) as compared to explicit ratings (60%). On
 353 the other hand, for the binary rating estimation, using GLM, the prediction of explicit ratings outperformed
 354 that of implicit appraisals by almost 10 percentage points. When using a decision tree, no differences were
 355 found in the performance when applied to the two different types of ratings.

Table 4. Comparison between average accuracy of implicit appraisals and explicit ratings from website features by model, presence of expertise/exposure factors and type of prediction (1-5 points or binary).

Prediction	Model	Expertise / Exposure	Average accuracy	
			Explicit	Implicit
5-Points	GLM	None	67.7%	67.7%
		Exposure	62.1%	64.0%
		Expertise	62.8%	69.2%
	Decision Tree	None	60.1%	72.9%
		Exposure	60.1%	68.2%
		Expertise	61.3%	66.0%
Binary	GLM	None	97.1%	89.9%
		Exposure	97.7%	86.6%
		Expertise	96.9%	87.7%
	Decision Tree	None	87.6%	87.1%
		Exposure	88.5%	86.5%
		Expertise	88.0%	85.0%

4 DISCUSSION

356 4.1 Prediction of perceived visual aesthetic

357 Our results showed that by using automatic estimable features from still images of web-pages, regressive
 358 models can be used to predict with reasonably high accuracy if a page will be explicitly and/or implicitly
 359 perceived as aesthetically pleasant, thereby supporting our first hypothesis.

360 Our finding provides further support to the existing literature whereby visual properties have been found
 361 to play a role in aesthetic judgment. More importantly, this finding provides insight into the predictive
 362 capabilities of these visual properties on both explicit and implicit aesthetic judgments and how they can
 363 be utilized effectively depending on the type of scale the researcher prefers.

364 More specifically, depending on the type of desired data, different models can be selected to predict the
 365 different ratings. When a 5-point scale is preferred, either GLM or Decision Tree model can be used to
 366 both predict explicit ratings or implicit appraisals. However, when a binary scale is preferred, using GLM

367 will provide a higher prediction accuracy for explicit ratings than for implicit ratings while using Decision
368 Tree will provide similar accuracy for both implicit appraisals and explicit ratings.

369 **4.2 Does the level of exposure to different websites influences perceived visual** 370 **appeal?**

371 With regards to our second hypothesis, we predicted that the level of exposure to different web pages
372 moderates users' aesthetic judgments.

373 Despite the fact that two websites received significantly different ratings by participants of the high and
374 low-exposure groups, the addition of the level of expertise as a factor of our regressive model resulted in
375 no significant increase in their accuracy. We can, therefore, conclude that our second hypothesis, within
376 given limits of the number of websites and participants, is not confirmed.

377

378 **4.3 Does the level of expertise play a role in perceived visual appeal?**

379 For our third hypothesis, we predicted that participants' expertise in the design and development of web
380 pages affects their aesthetic judgment. Similar to the comparison between highly-exposed and low-exposed
381 participants, three websites received significantly different ratings by participants of the two groups, but
382 the addition of the users' expertise as a factor of the models led to no significant improvement of their
383 accuracy, hence not supporting our third hypothesis.

384 **4.4 Limitations**

385 As is common to all experimental studies, limitations are inevitable and should be mentioned. With
386 respect to the physiological measurements utilized to assess participants' implicit appraisals, we are
387 unable to control for the participants' physiological state at the beginning of each session. Despite the fact
388 that a baseline correction is applied during feature extraction, possible differences in pre-experimental
389 physiological arousal and valence may still be present and should be taken into consideration. Next, it
390 should also be noted that participants' explicit ratings and expertise/exposure responses are all self-reported
391 measurements and the social desirability factor could affect the reliability of the reporting. Participants
392 may feel the social pressure in not stating the truth to questions about the amount of time they spend on the
393 Internet and about the average number of pages browsed per day. Another important consideration is that
394 the number of existing websites are almost impossible to be determined and, as such, the usage of a limited
395 number of websites may not be suitable if used pages are not representative of the whole dataset. Thus, it is
396 ideal for further studies to be conducted to determine how representative the used pages are of the entire
397 dataset.

398 Across different age groups and different cultures, the daily usage of websites can vary greatly and as such,
399 future studies should also take these factors into consideration and select appropriate thresholds for their
400 sample. More specifically, future studies can consider including participants with a broader range in design
401 and development knowledge, time spent browsing pages and number of different pages browsed per day.
402 Different indicators of expertise and exposure can then be considered. In addition, our sample is not a
403 perfect representation of the actual age range of Internet users. Therefore, future studies should also involve
404 younger and older participants in order to test the reliability of our models on a more varied sample.

5 CONCLUSION

405 In this work, we investigated the possibility of predicting both implicit and explicit user aesthetic judgment
406 of websites from visual properties while considering expertise and exposure as possible predictive factors.
407 Results showed that by investigating the visual properties of web pages, it is possible to predict, with a
408 good degree of accuracy, if a website will be perceived -explicitly or implicitly- as aesthetically pleasing by
409 possible users. Although differences in ratings given by experts and non-experts as well as high-exposure

410 and low-exposure users have been found, the accuracy of predictive models was not enhanced by the
411 addition of expertise and exposure as factors.
412 Findings from this study will help designers uncover the most critical aspects that they should consider in
413 sketching the layout of digital interfaces.

CONFLICT OF INTEREST STATEMENT

414 The authors declare that the research was conducted in the absence of any commercial or financial
415 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

416 G.G. and G.E. conceived and planned the experiments; G.G. conducted the study and analyzed the data;
417 G.G., M.H.B. and G.E. discussed and interpreted the results; G.G. wrote the original draft; all the authors
418 reviewed and edited the submitted version.

FUNDING

419 This study was supported by the NAP-SUG program of the Nanyang Technological University and by the
420 Singapore Ministry of Education Academic Research Grants - Tier1.

ACKNOWLEDGMENTS

421 Thank to Chiara Iannaccone, Giulia Garbin, and Mengyu Lim for their help.

DATA AVAILABILITY STATEMENT

422 The raw and processed data, a copy of the python packages and of the scripts used in this work can be found
423 in the data repository of the Nanyang Technological University [https://doi.org/10.21979/N9/
424 YCDXNE](https://doi.org/10.21979/N9/YCDXNE).

REFERENCES

- 425 Batista, D., Silva, H., and Fred, A. (2017). Experimental characterization and analysis of the bitalino
426 platforms against a reference device. In *2017 39th Annual International Conference of the IEEE
427 Engineering in Medicine and Biology Society (EMBC)* (IEEE), 2418–2421
- 428 Birtolo, C., Pagano, P., and Troiano, L. (2009). Evolving colors in user interfaces by interactive genetic
429 algorithm. In *Nature & Biologically Inspired Computing, 2009. NaBIC 2009. World Congress on* (IEEE),
430 349–355
- 431 Blackburn, K. and Schirillo, J. (2012). Emotive hemispheric differences measured in real-life portraits
432 using pupil diameter and subjective aesthetic preferences. *Experimental Brain Research*
- 433 Bölte, J., Hösker, T. M., Hirschfeld, G., and Thielsch, M. T. (2017). Electrophysiological correlates of
434 aesthetic processing of webpages: a comparison of experts and laypersons. *PeerJ* 5, e3440
- 435 Bornstein, R. F. and D'agostino, P. R. (1992). Stimulus recognition and the mere exposure effect. *Journal
436 of personality and social psychology* 63, 545
- 437 Bucy, E. P., Lang, A., Potter, R. F., and Grabe, M. E. (1999). Formal features of cyberspace: Relationships
438 between Web page complexity and site traffic. *Journal of the American Society for Information Science*
439 50, 1246–1256
- 440 Cela-Conde, C. J., Agnati, L., Huston, J. P., Mora, F., and Nadal, M. (2011). The neural foundations of
441 aesthetic appreciation. *Prog. Neurobiol.* 94, 39–48

- 442 Cox, D. and Cox, A. D. (2002). Beyond first impressions: The effects of repeated exposure on consumer
443 liking of visually complex and simple product designs. *Journal of the Academy of Marketing Sciences*
444 Cyr, D. (2008). Modeling web site design across cultures: relationships to trust, satisfaction, and e-loyalty.
445 *Journal of Management Information Systems* 24, 47–72
- 446 Cyr, D., Head, M., and Larios, H. (2010). Colour appeal in website design within and across cultures: A
447 multi-method evaluation. *Int. J. Hum. Comput. Stud.* 68, 1–21
- 448 de Jong, M. A. (1972). A physiological approach to aesthetic preference. I. Paintings. *Psychother.*
449 *Psychosom.* 20, 360–365
- 450 de Jong, M. A., van Mourik, K. R., and Schellekens, H. M. C. (1973). A Physiological Approach to
451 Aesthetic Preference. *Psychother. Psychosom.* 22, 46–51
- 452 [Dataset] Gabrieli, G. (2019a). Gabrock94/prettywebsite: Version 0.0.3 - build 1. doi:10.5281/zenodo.
453 3187316
- 454 [Dataset] Gabrieli, G. (2019b). Gabrock94/physiology: Version 0.0.9 - build 3. doi:10.5281/zenodo.2622204
- 455 Gabrieli, G., Azhari, A., and Esposito, G. (2020). Physiology: A python package for physiological feature
456 extraction. In *Neural Approaches to Dynamics of Signal Exchanges* (Springer). 395–402
- 457 Guerreiro, J., Martins, R., Silva, H., Lourenço, A., and Fred, A. L. (2013). Bitalino—a multimodal platform
458 for physiological computing. In *ICINCO (1)*. 500–506
- 459 Hall, R. H. and Hanna, P. (2004). The impact of web page text-background colour combinations on
460 readability, retention, aesthetics and behavioural intention. *Behav. Inf. Technol.* 23, 183–195
- 461 Harmon-Jones, E. and Allen, J. J. B. (2001). The role of affect in the mere exposure effect: Evidence from
462 psychophysiological and individual differences approaches. *Personality and social psychology bulletin*
- 463 Hasler, D. and Suesstrunk, S. E. (2003). Measuring colorfulness in natural images. In *Human vision and*
464 *electronic imaging VIII* (International Society for Optics and Photonics), vol. 5007, 87–96
- 465 Ivory, M. Y., Sinha, R. R., and Hearst, M. A. (2001). Empirically validated web page design metrics. In
466 *Proceedings of the SIGCHI conference on Human factors in computing systems* (ACM), 53–60
- 467 Jacobs, K. W. and Hustmyer Jr, F. E. (1974). *Effects of four psychological primary colors on GSR, heart*
468 *rate and respiration rate*, vol. 38 (SAGE Publications Sage CA: Los Angeles, CA)
- 469 Karvonen, K. (2000). The beauty of simplicity. In *Proceedings on the 2000 conference on Universal*
470 *Usability* (ACM), 85–90
- 471 Kim, J., Lee, J., and Choi, D. (2003). Designing emotionally evocative homepages: an empirical study of
472 the quantitative relations between design factors and emotional dimensions. *Int. J. Hum. Comput. Stud.*
473 59, 899–940
- 474 Kim, J. and Moon, J. Y. (1998). Designing towards emotional usability in customer interfaces—
475 trustworthiness of cyber-banking system interfaces. *Interact. Comput.* 10, 1–29
- 476 Lang, P. and Bradley, M. M. (2007). The international affective picture system (iaps) in the study of
477 emotion and attention. *Handbook of emotion elicitation and assessment* 29
- 478 Lang, P. J., Greenwald, M. K., Bradley, M. M., and Hamm, A. O. (1993). Looking at pictures: Affective,
479 facial, visceral, and behavioral reactions. *Psychophysiology* 30, 261–273
- 480 Lin, Y.-C., Yeh, C.-H., and Wei, C.-C. (2013). How will the use of graphics affect visual aesthetics? a
481 user-centered approach for web page design. *International Journal of Human-Computer Studies* 71,
482 217–227
- 483 Lindgaard, G., Fernandes, G., Dudek, C., and Brown, J. (2006). Attention web designers: You have 50
484 milliseconds to make a good first impression! *Behav. Inf. Technol.* 25, 115–126
- 485 Maughan, L., Gutnikov, S., and Stevens, R. (2007). Like more, look more. look more, like more: The
486 evidence from eye-tracking. *Journal of Brand management* 14, 335–342
-

- 487 Miniukovich, A. and De Angeli, A. (2014). Quantification of interface visual complexity. In *Proceedings*
488 *of the 2014 International Working Conference on Advanced Visual Interfaces - AVI '14* (New York, New
489 York, USA: ACM Press), 153–160. doi:10.1145/2598153.2598173
- 490 Moshagen, M. and Thielsch, M. T. (2010). Facets of visual aesthetics. *International journal of human-*
491 *computer studies* 68, 689–709
- 492 Müller, M., Höfel, L., Brattico, E., and Jacobsen, T. (2010). Aesthetic judgments of music in experts and
493 laypersons—an erp study. *International Journal of Psychophysiology* 76, 40–51
- 494 Oliphant, T. E. (2006). *A guide to NumPy*, vol. 1 (Trelgol Publishing USA)
- 495 Oliphant, T. E. (2007). Python for scientific computing. *Computing in Science & Engineering* 9, 10–20
- 496 Orr, M. G. and Ohlsson, S. (2005). Relationship between complexity and liking as a function of expertise.
497 *Music Perception: An Interdisciplinary Journal*
- 498 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn:
499 Machine learning in Python. *Journal of Machine Learning Research* 12, 2825–2830
- 500 Pihko, E., Virtanen, A., Saarinen, V.-M., Pannasch, S., Hirvenkari, L., Tossavainen, T., et al. (2011).
501 Experiencing art: the influence of expertise and painting abstraction level. *Frontiers in human*
502 *neuroscience* 5, 94
- 503 Quispel, A., Maes, A., and Schilperoord, J. (2016). Graph and chart aesthetics for experts and laymen in
504 design: The role of familiarity and perceived ease of use. *Information Visualization* 15, 238–252
- 505 Ray, G., Kaplan, A. Y., and Jovanov, E. (1997). Morphological variations in eeg during music-induced
506 change in consciousness. In *Engineering in Medicine and Biology Society, 1997. Proceedings of the*
507 *19th Annual International Conference of the IEEE (IEEE)*, vol. 1, 227–230
- 508 Reinecke, K., Yeh, T., Miratrix, L., Mardiko, R., Zhao, Y., Liu, J., et al. (2013). Predicting users' first
509 impressions of website aesthetics with a quantification of perceived visual complexity and colorfulness.
510 In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (ACM), 2049–2058
- 511 Robbins, J. N. (2012). *Learning Web Design: A Beginner's Guide to {HTML}, {CSS}, {JavaScript}, and*
512 *Web Graphics* (“O'Reilly Media, Inc.”)
- 513 Seabold, S. and Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with python. In
514 *Proceedings of the 9th Python in Science Conference* (Scipy), vol. 57, 61
- 515 Seckler, M., Opwis, K., and Tuch, A. N. (2015a). Linking objective design factors with subjective
516 aesthetics: An experimental study on how structure and color of websites affect the facets of users' visual
517 aesthetic perception. *Computers in Human Behavior* 49, 375–389. doi:10.1016/j.chb.2015.02.056
- 518 Seckler, M., Opwis, K., and Tuch, A. N. (2015b). Linking objective design factors with subjective
519 aesthetics: An experimental study on how structure and color of websites affect the facets of users' visual
520 aesthetic perception. *Computers in Human Behavior* 49, 375–389
- 521 Tuch, A. N., Bargas-Avila, J. A., Opwis, K., and Wilhelm, F. H. (2009). Visual complexity of websites:
522 Effects on users' experience, physiology, performance, and memory. *Int. J. Hum. Comput. Stud.* 67,
523 703–715
- 524 Tuch, A. N., Presslauer, E. E., Stöcklin, M., Opwis, K., and Bargas-Avila, J. A. (2012). The role
525 of visual complexity and prototypicality regarding first impression of websites: Working towards
526 understanding aesthetic judgments. *International Journal of Human Computer Studies* 70, 794–811.
527 doi:10.1016/j.ijhcs.2012.06.003
- 528 Ulrich Kirk, M. S. C. N. N., Martin Skov (2009). Brain correlates of aesthetic expertise: A parametric fMRI
529 study. *Brain and Cognition*
- 530 Van Rossum, G. and Drake, F. L. (2011). *The python language reference manual* (Network Theory Ltd.)

- 531 Wang, M. and Li, X. (2016). Effects of the aesthetic design of icons on app downloads: evidence from an
532 android market. *Electr. Commerce Res.* 17, 83–102
- 533 Wang, S. and Ding, R. (2012). A qualitative and quantitative study of color emotion using valence-arousal.
534 *Frontiers of Computer Science*
- 535 Winkielman, P. and Cacioppo, J. T. (2001). Mind at ease puts a smile on the face: Psychophysiological
536 evidence that processing facilitation elicits positive affect. *J. Pers. Soc. Psychol.* 81, 989–1000
- 537 Yanulevskaya, V., Uijlings, J., Bruni, E., Sartori, A., Zamboni, E., Bacci, F., et al. (2012). In the eye of the
538 beholder: employing statistical analysis and eye tracking for analyzing abstract paintings. In *Proceedings*
539 *of the 20th ACM international conference on Multimedia* (ACM), 349–358
- 540 Yendrikhovskij, S., Blommaert, F. J., and de Ridder, H. (1998). Optimizing color reproduction of natural
541 images. *Color and Imaging Conference* 1998, 140–145
- 542 Zajonc, R. B. (1968). Attitudinal effects of mere exposure. *Journal of personality and social psychology*
- 543 Zhang, P. and Von Dran, G. M. (2000). Satisfiers and dissatisfiers: a two-factor model for website design
544 and evaluation. *Journal of the American Society for Information Science and Technology* 51, 1253–1268.
545 doi:10.1002/1097-4571(2000)9999:9999<::AID-ASII1039>3.0.CO;2-O
- 546 Zhang, P., Von Dran, G. M., Blake, P., and Pipithsuksunt, V. (2001). Important Design Features in Different
547 Web Site Domains: An Empirical Study of User Perceptions. *e-Service Journal* 1, 77–91
- 548 Zheng, X. S., Chakraborty, I., Lin, J. J.-W., and Rauschenberger, R. (2009). Correlating low-level image
549 statistics with users - rapid aesthetic and affective judgments of web pages. In *Proceedings of the 27th*
550 *international conference on Human factors in computing systems - CHI 09* (New York, New York, USA:
551 ACM Press), 1. doi:10.1145/1518701.1518703
-