

1           **Fixation-related brain potentials during semantic integration**  
2                           **of object-scene information**

3  
4                           Moreno I. Coco<sup>1,2</sup>, Antje Nuthmann<sup>3</sup> and Olaf Dimigen<sup>4</sup>

5  
6           1. School of Psychology, The University of East London, London, UK

7           2. Faculdade de Psicologia, Universidade de Lisboa, Lisbon, Portugal

8           3. Institute of Psychology, Christian-Albrechts-Universität zu Kiel, Kiel, Germany

9           4. Department of Psychology, Humboldt-Universität zu Berlin, Berlin, Germany

10  
11  
12  
13  
14  
15   This research was supported by the Leverhulme Trust (grant ECF-014-205) and Fundação  
16   para a Ciência e Tecnologia (grant PTDC/PSI-ESP/30958/2017) to M. I. C, while he was a  
17   Research Fellow at the University of Edinburgh. The authors thank Benedikt Ehinger for  
18   helpful discussions on EEG deconvolution techniques. Please address correspondence to  
19   moreno.cocoi@gmail.com or olaf.dimigen@hu-berlin.de.

20

21 **Abstract**

22 In vision science, a particularly controversial topic is whether and how quickly the  
23 semantic information about objects is available outside foveal vision. Here, we aimed at  
24 contributing to this debate by co-registering eye-movements and EEG while participants  
25 viewed photographs of indoor scenes that contained a semantically consistent or inconsistent  
26 target object. Linear deconvolution modelling was used to analyse the event-related  
27 potentials (ERP) evoked by scene onset as well as the fixation-related potentials (FRPs)  
28 elicited by the fixation on the target object ( $t$ ) and by the preceding fixation ( $t-1$ ). Object-  
29 scene consistency did not influence the probability of immediate target fixation or the ERP  
30 evoked by scene onset, which suggests that object-scene semantics was not accessed  
31 immediately. However, during the subsequent scene exploration, inconsistent objects were  
32 prioritized over consistent objects in extrafoveal vision (i.e., looked at earlier) and were more  
33 effortful to process in foveal vision (i.e., looked at longer). In FRPs, we demonstrate a  
34 fixation-related N300/N400 effect, whereby inconsistent objects elicit a larger frontocentral  
35 negativity than consistent objects. In line with the behavioural findings, this effect was  
36 already seen in FRPs aligned to the pre-target fixation  $t-1$  and persisted throughout fixation  $t$ ,  
37 indicating that the extraction of object semantics can already begin in extrafoveal vision.  
38 Taken together, the results emphasize the usefulness of combined EEG/eye-movement  
39 recordings for understanding the mechanisms of object-scene integration during natural  
40 viewing.

41

42 *Keywords:* object-scene integration; foveal and peripheral vision; semantic processing;  
43 fixation-related potentials, eye tracking, N300/N400, regression-ERPs

44

## Introduction

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

In our daily activities, for example when we search for something in a room, our attention is mostly oriented to objects. The time course of object recognition and the role of overt attention in this process are therefore topics of considerable interest in the visual sciences. In the context of real-world scene perception, the question of what constitutes an object is a more complex question than intuition would suggest (e.g., Wolfe, Alvarez, Rosenholtz, Kuzmova, & Sherman, 2011). An object is likely a hierarchical construct (e.g., Feldman, 2003), with both low-level features (e.g., visual saliency) and high-level properties (e.g., semantics) contributing to its identity. Accordingly, when a natural scene is inspected with eye-movements, the observer's attentional selection is thought to be based either on objects (e.g., Nuthmann & Henderson, 2010), image features (saliency; Itti, Koch, & Niebur, 1998) or some combination of the two (e.g., Stoll, Thrun, Nuthmann, & Einhäuser, 2015).

An early and uncontroversial finding is that the recognition of objects is mediated by their semantic consistency. For example, an object that the observer would not expect to occur in a particular scene (e.g., a toothbrush in a kitchen) is recognized less accurately (e.g., Biederman, 1972; Davenport & Potter, 2004; Fenske, Aminoff, Gronau, & Bar, 2006) and looked at for longer than an expected object (e.g., Cornelissen & Vö, 2017; De Graef, Christiaens, & D'Ydewalle, 1990; Henderson, Weeks Jr, & Hollingworth, 1999).

What is more controversial, however, is the exact time course along which the meaning of an object is processed and how this semantic processing then influences the overt allocation of visual attention (see Wu, Wick, & Pomplun, 2014, for a review). Two interrelated questions are at the core of this debate: (1) How much time is needed to access the meaning of objects after a scene is displayed and (2) can object semantics be extracted before the object is overtly attended, that is, while the object is still outside of high-acuity foveal vision ( $> 1^\circ$  eccentricity) or even in the periphery ( $> 5^\circ$  eccentricity).

69 Evidence that the meaning of not-yet-fixated objects can capture overt attention comes  
70 from experiments that have used sparse displays of several standalone objects (e.g., Belke,  
71 Humphreys, Watson, Derrick, Meyer, & Telling, 2008; Cimminella, Della Sala & Coco, in  
72 press; Moores, Laiti, & Chelazzi, 2003; Nuthmann, de Groot, Huettig, & Olivers, 2019). For  
73 example, across three different experiments Nuthmann et al. found that the very first saccade  
74 in the display was directed more frequently to objects that were semantically related to a  
75 target object rather than to unrelated objects.

76 Whether such findings generalize to objects embedded in real-world scenes is  
77 currently an open research question. The size of the visual span – that is the area of the visual  
78 field from which observers can take in useful information (see Rayner, 2014 for review) – is  
79 large in scene viewing. For object-in-scene search, it corresponded to approximately 8° in  
80 each direction from fixation (Nuthmann, 2013). This opens up the possibility that both low-  
81 level and high-level object properties can be processed outside the fovea. This is the case for  
82 low-level visual features: objects that are highly salient (i.e., visually distinct) are  
83 preferentially selected for fixation (e.g., Stoll et al., 2015). But what about high-level  
84 semantic information? If extra-foveal semantic processing takes place, then objects that are  
85 inconsistent with the scene context (which are also thought to be more informative, Antes,  
86 1974) should be fixated earlier in time than consistent ones (Loftus & Mackworth, 1978;  
87 Mackworth & Morandi, 1967).

88 However, results from eye-movement studies on this issue have been mixed. A  
89 number of studies have indeed reported evidence for an inconsistent object advantage (e.g.,  
90 Bonitz & Gordon, 2008; Borges, Fernandes, & Coco, 2019; LaPointe & Milliken, 2016;  
91 Loftus & Mackworth, 1978; Underwood, Templeman, Lamming, & Foulsham, 2008).  
92 Among these studies, only Loftus & Mackworth (1978) have reported evidence for *immediate*  
93 extrafoveal attentional capture (i.e. within the first fixation) by object-scene semantics. In this

94 study, which used relatively sparse line drawings of scenes, the mean amplitude of the  
95 saccade into the critical object was more than  $7^\circ$ , suggesting that viewers could process  
96 semantic information based on peripheral information obtained in a single fixation. Several  
97 other studies, however, have failed to find any advantage for inconsistent objects in attracting  
98 overt attention (e.g., De Graef, Christiaens, & D'Ydewalle, 1990; Henderson, Weeks, &  
99 Hollingworth, 1999; Vö & Henderson, 2009, 2011). In these experiments, only measures of  
100 foveal processing – such as gaze duration – were influenced by object-scene consistency,  
101 with longer fixations times on inconsistent than on consistent objects.

102         Interestingly, a similar controversy exists in the literature on eye guidance in sentence  
103 reading. Although some degree of parafoveal processing during reading is uncontroversial, it  
104 is less clear whether semantic information is acquired from the parafovea (Andrews &  
105 Veldre, 2019, for review). Most evidence from studies involving readers of English has been  
106 negative (e.g., Rayner, Balota, & Pollatsek, 1986), whereas results from reading German  
107 (e.g., Hohenstein & Kliegl, 2014) and Chinese (e.g., Yan, Richter, Shu, & Kliegl, 2009)  
108 suggest that parafoveal processing can advance up to the level of semantic processing.

109         The processing of object-scene inconsistencies and its time course have also been  
110 investigated in electrophysiological studies (e.g., Ganis & Kutas, 2003; Mudrik, Lamy, &  
111 Deouell, 2010). In event-related potentials (ERPs), it is commonly found that scene-  
112 inconsistent objects elicit a larger negative brain response compared to consistent ones. This  
113 long-lasting negative shift typically starts as early as 200-250 ms after stimulus onset (e.g.,  
114 Mudrik, Shalgi, Lamy, & Deouell, 2014; Draschkow, Heikel, Vö, Fiebach, & Sassenhagen,  
115 2018) and has its maximum at frontocentral scalp sites, in contrast to the centroparietal N400  
116 effect for words (e.g., Kutas & Federmeier, 2011). The effect was found for objects that  
117 appeared at a cued location after the scene background was already shown (Ganis & Kutas,  
118 2003), for objects that were photoshopped into the scene (Mudrik, Lamy, & Deouell, 2010;

119 Mudrik, et al., 2014; Coco, Araujo, & Petersson, 2017), and for objects that were part of  
120 realistic photographs (Võ & Wolfe, 2013). These ERP effects of object-scene consistency  
121 have typically been subdivided into two distinct components: N300 and N400. The earlier  
122 part of the negative response, usually referred to as N300, has been taken to reflect the  
123 context-dependent difficulty of object identification, whereas the later N400 has been linked  
124 to semantic integration processes after the object is identified (e.g., Dyck & Brodeur, 2015).  
125 The present study was not designed to differentiate between these two subcomponents,  
126 especially considering that their scalp distribution is strongly overlapping or even  
127 topographically indistinguishable (Draschkow et al., 2018). Thus, for reasons of simplicity,  
128 we will in most cases simply refer to all frontocentral negativities as “N400”.

129         One limiting factor of existing ERP studies is that the data were gathered using steady-  
130 fixation paradigms in which the free exploration of the scene through eye-movements was not  
131 permitted. Instead, the critical object was typically large and/or located relatively close to the  
132 centre of the screen, and ERPs were time-locked to the onset of the image (e.g., Mudrik et al.,  
133 2010). Due to these limitations, it remains unclear whether foveation of the object is a  
134 necessary condition for processing object-scene consistencies, or whether such processing  
135 can at least begin in extrafoveal vision.

136         In the current study, we used fixation-related potentials (FRPs), that is EEG  
137 waveforms aligned to fixation onset, to shed new light on the controversial findings of the  
138 role of foveal versus extrafoveal vision in extracting object semantics, while providing  
139 insights into the patterns of brain activity that underlie them (for reviews about FRPs see  
140 Dimigen, Sommer, Hohlfeld, Jacobs, & Kliegl, 2011; Nikolaev, Meghanathan, & van  
141 Leeuwen, 2016).

142         FRPs have previously been used to investigate the brain-electric correlates of natural  
143 reading, as opposed to serial word presentation, helping researchers to provide finer details

144 about the online processing of linguistic features (such as word predictability, Kretzschmar,  
145 Bornkessel-Schlesewsky, & Schlesewsky, 2009; Kliegl, Dambacher, Dimigen, Jacobs, &  
146 Sommer, 2012) or the dynamics of the perceptual span during reading (e.g., parafovea-on-  
147 fovea effects, Niefind & Dimigen, 2016). More recently, the co-registration method has also  
148 been applied to investigate active visual search (e.g., Devillez, Guyader, & Guerin-Dugue,  
149 2015; Kamienkowski, Ison, Quiroga, & Sigman, 2012; Kaunitz et al., 2014), object  
150 identification (Rämä & Baccino, 2010), and affective processing in natural scene viewing  
151 (Simola, Le Fevre, Torniaainen, & Baccino, 2015).

152         In the present study, we simultaneously recorded eye-movements and FRPs during the  
153 viewing of real-world scenes to distinguish between three alternative hypotheses on object-  
154 scene integration that can be derived from the literature: (A) one glance of the scene is  
155 sufficient to extract object semantics from extrafoveal vision (e.g., Loftus & Mackworth,  
156 1978), (B) extrafoveal processing of object-scene semantics is possible but takes some time to  
157 unfold (e.g., Bonitz & Gordon, 2008; Underwood et al., 2008), and (C) the processing of  
158 object semantics requires foveal vision, that is, a direct fixation of the critical object (e.g., De  
159 Graef et al., 1990; Henderson et al., 1999; Vö & Henderson, 2009). We note that these  
160 possibilities are not mutually exclusive, an issue we elaborate on in the *General Discussion*.

161         For the behavioural data, these hypotheses translate as follows: under (A), the  
162 probability of immediate target fixation should reveal that already the first saccade on the  
163 scene goes more often towards inconsistent than consistent objects. Under (B), there should  
164 be no effect on the first eye-movement, but the latency to first fixation on the critical object  
165 should be shorter for inconsistent than consistent objects. Under (C), only fixation times on  
166 the critical object itself should differ as a function of object-scene consistency, with longer  
167 gaze durations on inconsistent objects.

168 For the electrophysiological data analysis, we used a novel regression-based analysis  
169 approach (linear deconvolution modelling; Cornelissen, Sassenhagen, & Vö, 2019, Dandekar,  
170 Privitera, Carney, & Klein, 2011; Dimigen & Ehinger, 2019; Ehinger & Dimigen, 2018;  
171 Smith & Kutas, 2015b), which allowed us to control for the confounding influences of  
172 overlapping potentials and oculomotor covariates during natural viewing, which can  
173 otherwise distort the neural responses. In the EEG, hypothesis (A) can be tested by computing  
174 the ERP time-locked to the onset of the scene on the display, following the traditional  
175 approach. Given that the critical objects in our study were not placed directly in the centre of  
176 the screen from which observers started their exploration of the scene, any effect of object-  
177 scene congruency in this ERP would suggest that object semantics is rapidly processed in  
178 extrafoveal vision, even before the first eye-movement is generated, in line with Loftus &  
179 Mackworth, 1978. Under hypothesis (B) we would not expect to see an effect in the scene-  
180 onset ERP. Instead, we should find a negative brain potential (N400) for inconsistent as  
181 compared to consistent objects in the FRP aligned to the fixation that *precedes* the one that  
182 first lands on the critical object. Finally, if (C) is correct, an N400 for inconsistent objects  
183 should only arise once the critical object is foveated, i.e., in the FRP aligned to the target  
184 fixation (fixation  $t$ ). In contrast, no consistency effects should appear in the scene-onset ERP  
185 or in the FRP aligned to the pre-target fixation (fixation  $t-1$ ). To preview the results, both the  
186 eye-movement as well as the EEG data lend support for hypothesis (B).

## 187 **Methods**

### 188 ***Design and task overview***

189 We designed a short-term visual working memory change detection task, illustrated in  
190 Figure 1 and 2. During the *study phase*, participants were exposed to photographs of indoor  
191 scenes (e.g., a bathroom), each of which contained a target object that was either semantically



192 *consistent* (e.g., toothpaste) or *inconsistent* (e.g., a flashlight) with the scene context. In the  
193 following *recognition phase*, after a short retention interval of 900 ms, the same scene was  
194 shown again, but in half of the trials either the identity, the location, or both the identity and  
195 location of the target object had changed relative to the study phase.

196 Insert Figure 1 and 2 about here

197 The participants' task was to indicate with a keyboard press whether or not a change had  
198 happened to the scene (see also LaPointe & Milliken, 2016). All eye-movement and EEG  
199 analyses in the present article focus on the semantic consistency manipulation of the target  
200 object during the *study phase*.

### 201 ***Participants***

202 Twenty-four participants (9 male) between the ages of 18 and 33 (mean: 25.0 years)  
203 took part in the experiment after providing written informed consent. They were compensated  
204 with £7 per hour. All participants had normal or corrected-to-normal vision. Data from an  
205 additional two participants was recorded but removed from the analysis due to excessive  
206 scalp muscle (EMG) activity or skin potentials in the raw EEG. Ethics approval was obtained  
207 from the Psychology Research Ethics Committee of the University of Edinburgh.

### 208 ***Apparatus and Recording***

209 Scenes were presented on a 19" CRT monitor (Iiyama Vision Master Pro 454) at a  
210 vertical refresh rate of 75 Hz. At the viewing distance of 60 cm, each scene subtended  $35.8^\circ \times$   
211  $26.9^\circ$  (width  $\times$  height). Eye-movements were recorded monocularly from the dominant eye  
212 using an SR Research EyeLink 1000 desktop-mounted system at a sampling rate of 1000 Hz.  
213 Eye dominance for each participant was determined with a parallax test. A chin and forehead  
214 rest was used to stabilize the participant's head. Nine-point calibrations were run at the

215 beginning of each session and whenever the participant's fixation deviated by  $> 0.5^\circ$   
216 horizontally or  $> 1^\circ$  vertically from a drift correction point presented at trial onset.

217 The EEG was recorded from 64 active electrodes at a sampling rate of 512 Hz using  
218 BioSemi ActiveTwo amplifiers. Four electrodes, located near the left and right canthus and  
219 above and below the right eye, recorded the electro-oculogram (EOG). All channels were  
220 referenced against the BioSemi common mode sense (CMS; active electrode) and grounded  
221 to a passive electrode. The BioSemi hardware is DC coupled and applies digital low-pass  
222 filtering through the A/D-converter's decimation filter, which has a 5th order sinc response  
223 with a -3 dB point at 1/5th of the sample rate (corresponding approximately to a 100 Hz low-  
224 pass filter).

225 Offline, the EEG was re-referenced to the average of all scalp electrodes and filtered  
226 using EEGLAB's (Delorme & Makeig, 2004) Hamming-windowed sinc FIR filter  
227 (`pop_eegfiltnew.m`) with default settings. The lower edge of the filter's passband was set to  
228 0.2 Hz (with -6 dB attenuation at 0.1 Hz) and the upper edge to 30 Hz (with -6 dB attenuation  
229 at 33.75 Hz). Eye tracking and EEG were synchronized using shared triggers sent via the  
230 parallel port of the stimulus presentation PC to the two recording computers. Synchronization  
231 was performed offline using the EYE-EEG extension (v0.8) for EEGLAB (Dimigen et al.,  
232 2011). All datasets were aligned with a mean synchronization error  $\leq 2$  ms as computed based  
233 on trigger alignment after synchronization.

### 234 ***Materials & Rating***

235 Stimuli consisted of 192 colour photographs of indoor scenes (e.g., bedrooms,  
236 bathrooms, offices). Real target objects were placed in the physical scene, before each picture  
237 was taken with a tripod under controlled lighting conditions and with a fixed aperture (i.e.,  
238 there was no photo-editing). One scene is shown in Figure 1; miniature version of all stimuli  
239 used in the present study are found as part of the *Supplementary Materials*. Of the 192 scenes,

240 96 were conceived as change items and 96 as no-change items. Each one of the 96 change  
241 scenes was created in four versions. In particular, each scene (e.g., a bathroom) was  
242 photographed with two alternative target objects in it, one that was consistent with the scene  
243 context (e.g., a toothbrush) and one that was not (e.g., a flashlight). Moreover, each of these  
244 two objects was placed at two alternative locations (left or right side) within the scene (e.g.,  
245 either on the sink or on the bathtub). Accordingly, three types of change were implemented  
246 during the recognition phase (*Congruency, Location, and Both*); see *Procedure* section below.

247 Each of the 96 no-change scenes was also a real photograph with either a consistent or  
248 an inconsistent object in it, which was again located in either the left or right half of the  
249 scene. Across the 96 no-change scenes, factors consistency (consistent vs. inconsistent  
250 objects) and location (left and right) were also balanced. However, each no-change scene was  
251 unique, that is, we did not create four different versions of each no-change scene. The data of  
252 the 96 no-change scenes, which were originally conceived to be filler trials, was included to  
253 improve the signal-to-noise ratio of the EEG analyses, as these scenes also had a balanced  
254 distribution of inconsistent and consistent objects.

255 As explained above, each scene contained a critical object that was either consistent or  
256 inconsistent with the scene context. Object consistency was assessed in a pre-test rating study  
257 by eight naïve participants who were not involved in any other aspect of the study. Each  
258 participant rated all of the no-changes scenes as well as one of the four versions of each  
259 change-scene (counterbalanced across raters). Together with each scene, raters saw a box  
260 with a cropped image of the critical object. They were asked (a) to write down the name for  
261 the displayed object, and (b) to respond to the question “*How likely is it that this object would*  
262 *be found in this room?*” using a six-point Likert scale (1-6). For the object naming, a mean  
263 naming agreement of 96.35% was obtained. Furthermore, consistent objects were judged as  
264 significantly more likely (mean = 5.78, SD = ± 0.57) to appear in the scene than inconsistent

265 objects ( $1.88 \pm 1.11$ ), as confirmed by an independent-samples Kruskal-Wallis  $H$ -test ( $\chi^2(1) =$   
266  $616.09, p < .001$ ).

267 In addition, we ensured that there was no difference between consistent and  
268 inconsistent objects on three important low-level variables: object size (pixels square),  
269 distance from the centre of the scene (degrees of visual angle) and mean visual saliency of the  
270 object as computed using the Adaptive Whitening Saliency model (Garcia-Diaz, Fdez-Vidal,  
271 Pardo, & Dosil, 2012). Table 1 provides additional information about the target object. Paired  
272  $t$ -tests showed no significant difference between consistency conditions in object size,  $t(476)$   
273  $= -1.2, p = 0.2$ ; visual saliency,  $t(476) = 0.1, p = 0.9$ ; and distance from the centre,  $t(476) =$   
274  $0.48, p = 0.6$ .

275 The position of each target object was marked with an invisible rectangular bounding  
276 box, which was used to implement the gaze contingency mechanism (described in the  
277 *Procedure* section below) and to determine whether a fixation was inside the target object.  
278 The average width of the bounding box was  $6.1^\circ \pm 2.0$  for consistent objects and  $6.1^\circ \pm 2.1$   
279 for inconsistent objects (see Table 1); the average height was  $5.1^\circ \pm 1.8$  and/or  $5.4^\circ \pm 2.2$ ,  
280 respectively. The average distance of the object centroid from the centre of the scene was  
281  $12.1^\circ (\pm 2.8)$  for consistent and  $11.7^\circ (\pm 3.0)$  for inconsistent objects.

## 282 ***Procedure***

283 A schematic representation of the task is shown in Figure 2. Each trial started with a  
284 drift correction of the eye-tracker. Afterwards, the study scene was presented (e.g., a  
285 bathroom). The display duration of the study scene was controlled by a gaze-contingent  
286 mechanism that ensured that participants fixated the target object (e.g., toothbrush or  
287 flashlight) at least once during the trial. Specifically, the study scene disappeared on average  
288 2000 ms (with a random jitter of  $\pm 200$  ms, drawn from a uniform distribution) after the  
289 participant's eyes left the invisible bounding box of the target object (and provided that the

290 target had been fixated for at least 150 ms). The jittered delay of about 2000 ms was  
291 implemented to prevent participants from learning to associate the last fixated object during  
292 the study phase with the changed object during the recognition phase. If the participant did  
293 not fixate the target object within 10 s, the study scene disappeared from the screen and the  
294 retention interval was triggered, which lasted for 900 ms.

295 In the following recognition phase (data not analysed here), the scene was presented  
296 again, either with (50% of trials) or without (50% of trials) a change to an object in the scene.  
297 Three types of object changes occurred with equal probability: Location, Consistency, or  
298 Both. In the (a) *Location* condition, the target object changed its position and moved either  
299 from left to right or from right to left to another plausible location within the scene (e.g., a  
300 toothbrush was placed elsewhere within the scene). In the (b) *Consistency* condition, the  
301 object remained in the same location, but was replaced with another object of opposite  
302 semantic consistency (e.g., the toothbrush was replaced by a flashlight or vice-versa). Finally,  
303 in the (c) *Both* condition, the object was both replaced and moved within the scene (e.g., a  
304 toothbrush was replaced by a flashlight at a different location).

305 During the recognition phase, participants had to indicate whether they noticed any  
306 kind of change within the scene by pressing the arrow keys on the keyboard. Afterwards, the  
307 scene disappeared and the next trial began. If participants did not respond within 10 s, a  
308 missing response was recorded.

309 The type of change between trials was fully counterbalanced using a Latin Square  
310 rotation. Specifically, the 96 change trials were distributed across 12 different lists,  
311 implementing the different types of change. This implies that each participant was exposed to  
312 an equal number of consistent and inconsistent change trials. The 96 no-change trials also  
313 comprised an equal number of consistent and inconsistent scenes and were the same for each  
314 participant. All 192 trials were presented in a randomized order. These trials were preceded

315 by four practice trials at the start of the session. Written instructions were given to explain the  
316 task which took 20-40 minutes to complete. The experiment was implemented using the SR  
317 Research Experiment Builder software.

### 318 *Data preprocessing*

#### 319 *Eye-movement events and data exclusion*

320 Fixations and saccades events were extracted from the raw gaze data using the SR Research  
321 Data Viewer software, which performs saccade detection based on velocity and acceleration  
322 thresholds of  $30^\circ \text{ s}^{-1}$  and  $9,500^\circ \text{ s}^{-2}$ , respectively. To provide directly comparable results for  
323 eye-movement behaviour and FRP analyses, we discarded all trials on which we did not have  
324 clean data from both recordings. Specifically, from a total of 4,608 trials (24 participants  $\times$   
325 192 trials), we excluded 10 trials (0.2%) because of machine error (i.e., no data was recorded  
326 for those trials), 689 trials (15.0%) because the participant responded incorrectly after the  
327 recognition phase and 494 trials (10.7%) because the target object was not fixated during the  
328 study phase. Finally, we removed an additional 97 trials (2.1%) for which the target fixation  
329 overlapped with intervals of the EEG that contained non-ocular artefacts (see below). The  
330 final dataset therefore comprised 3,318 unique trials: 1,567 for the consistent condition and  
331 1,751 for the inconsistent condition. Per participant, this corresponded to an average of 65.3  
332 trials ( $\pm 6.9$ , range = 48-78) for consistent and 73.0 trials ( $\pm 6.9$ , range = 59-82) for  
333 inconsistent items. Due to the fixation check, participants were always fixating at the screen  
334 centre when the scene appears on the display. This on-going central fixation was removed  
335 from all analyses.

#### 336 *EEG ocular artefact correction*

337 EEG recordings during free viewing are contaminated by three types of ocular  
338 artefacts (e.g., Dimigen, 2018; Plöchl, Ossandón, & König, 2012) which need to be removed  
339 to get at the genuine brain activity. Here we applied an optimized variant (Dimigen, 2018) of

340 Independent Component Analysis (ICA, Jung et al., 1998), which uses the information  
341 provided by the eye-tracker to objectively identify ocular ICA components (Plöchl et al.,  
342 2012).

343 In a first step, we created optimized ICA training data by high-pass filtering a copy of  
344 the EEG at 2 Hz (Winkler, Debener, Müller, & Tangermann, 2015; Dimigen, 2018) and  
345 segmenting it into epochs lasting from scene onset until 3 s thereafter. This high pass-filtered  
346 training data was entered into an extended Infomax ICA using EEGLAB, and the resulting  
347 unmixing weights were then transferred to the original (i.e. less strictly filtered) recording  
348 (Viola, Debener, Thorne, & Schneider, 2010). From this original EEG dataset, we then  
349 removed all independent components whose time course varied more strongly during saccade  
350 intervals (defined as lasting from -20 ms before saccade onset until 20 ms after saccade  
351 offset) than during fixations with the threshold for the variance ratio (saccade/fixation, see  
352 Plöchl et al., 2012) set to 1.3. The artefact-corrected continuous EEG was then back-projected  
353 to the sensor space. For a validation of the ICA procedure, please refer to Supplementary  
354 Figure S1.

355 In a next step, intervals with residual non-ocular artefacts (e.g., EMG bursts) were  
356 detected by shifting a 2000 ms moving window in steps of 100 ms across the continuous  
357 recording. Whenever the voltages within the window exceeded a peak-to-peak threshold of  
358 100  $\mu$ V in at least one of the channels, all data within the window was marked as “bad” and  
359 subsequently excluded from analysis. Within the linear deconvolution framework (see  
360 below), this can easily be done by setting all predictors to zero during these bad EEG  
361 intervals (Smith & Kutas, 2015b), meaning that the data in these intervals will not affect the  
362 computation.

363 *Analysis*364 *Eye-movement data*

365 *Dependent measures:* Behavioural analyses focused on four eye-movement measures  
 366 commonly reported in the semantic consistency literature: (a) cumulative probability of  
 367 having fixated the target object as a function of the ordinal fixation number, (b) the  
 368 probability of immediate object fixation, (c) the latency to first fixation on the target object,  
 369 and (d) the gaze duration on the target object (cf. Võ & Henderson, 2009).

370 *Linear-mixed effect modelling:* Eye-movement data were analysed using linear mixed-  
 371 effects models (LMM) and generalized linear mixed-effects models (GLMM) as implemented  
 372 in the `lme4` package in R (Bates et al., 2015). The only exception was the cumulative  
 373 probability of first-fixations on the target for which a generalized linear model (GLM) was  
 374 used. One advantage of (G)LMM modelling is that it allows one to simultaneously model the  
 375 intrinsic variability of both participants and scenes (e.g., Nuthmann & Einhäuser, 2015).

376 In all analyses, the main predictor was the *Consistency* of the critical object (contrast  
 377 coding: Consistent = -0.5, Inconsistent = 0.5) in the study scene. In the (G)LMMs, *Participant*  
 378 (24) and *Scene* (192) were included as random intercepts<sup>1</sup>. The cumulative probability of object  
 379 fixation was analysed using a GLM with a binomial (probit) link. This model included the  
 380 *Ordinal Number of Fixation* on the scene as a predictor; it was entered as a continuous variable  
 381 ranging from 1 to a maximum of 28 (the 99<sup>th</sup> quantile).

382 In the tables of results, we report the beta coefficients, *t*-values (LMM), *z*-values  
 383 (GLMM), and *p*-values for each model. The level of significance was calculated from an *F*-test

---

<sup>1</sup> We did not include random slopes for two reasons: For *Participant*, the inclusion of a random slope led to a small variance and perfect correlation between intercept and slope. For the random effect *Scene*, only the change trials were fully counterbalanced in terms of location and consistency, meaning that the slope for *Consistency* could not be estimated for the no-change trials.



384 based on the Satterthwaite approximation to the effective degrees of freedom (Satterthwaite,  
385 1946), where  $p$ -values in GLMMs are based on asymptotic Wald tests.

### 386 *Electrophysiological data*

387 *Linear Deconvolution Modelling (first level of analysis)*: EEG measurements during active  
388 vision are associated with two major methodological problems: overlapping potentials and  
389 low-level signal variability (Dimigen & Ehinger, 2019). Overlapping potentials arise from the  
390 rapid pace of active information sampling through eye-movements, which causes the neural  
391 responses that are evoked by subsequent fixations on the stimulus to overlap with each other.  
392 Because the average fixation duration usually varies between conditions, this changing  
393 overlap can easily confound the measured waveforms. A related issue is the mutual overlap  
394 between the ERP elicited by the initial presentation of the stimulus and the FRPs evoked by  
395 the subsequent fixations on it. This second type of overlap is especially important in  
396 experiments like ours, in which the critical fixations occurred at different latencies after scene  
397 onset in the two experimental conditions.

398 The problem of signal variability refers to the fact that low-level visual and  
399 oculomotor variables can also influence the morphology of the predominantly visually-  
400 evoked fixation-related neural responses (e.g., Dimigen et al., 2011; Kristensen, Rivet, &  
401 Guerin-Dugué, 2017; Nikolaev et al., 2016). The most relevant of these variables, which is  
402 known to modulate the entire FRP waveform, is the amplitude of the saccade that precedes  
403 fixation onset (e.g., Dandekar et al., 2011; Thickbroom, Knezevič, Carroll, & Mastaglia,  
404 1991). One option for controlling the effect of saccade amplitude is to include it as a  
405 continuous covariate in a massive univariate regression-model (Smith & Kutas, 2015a,  
406 2015b), in which a separate regression model is computed for each EEG time point and  
407 channel (Weiss, Knakker & Vidnyánszky, 2016). However, this method does not account for  
408 overlapping potentials.

409           An approach that allows one to simultaneously control for overlapping potentials and  
410 low-level covariates is deconvolution within the linear model (for tutorial reviews see  
411 Dimigen & Ehinger, 2019; Smith & Kutas, 2015a, 2015b), sometimes also called continuous-  
412 time regression (Smith & Kutas, 2015b). Initially developed to separate overlapping BOLD  
413 responses (e.g., Serences, 2004), linear deconvolution has also been applied to separate  
414 overlapping potentials in ERP (Smith & Kutas, 2015b) and FRP paradigms (Cornelissen et  
415 al., 2019; Dandekar et al., 2011; Ehinger & Dimigen, 2018; Kristensen, et al. 2017). Another  
416 elegant property of this approach is that the ERPs elicited by scene onset and the FRPs  
417 elicited by fixations on the scene can be disentangled and simultaneously estimated in the  
418 same regression model. The benefits of deconvolution are illustrated in more detail in  
419 Supplementary Figures S2 and S3.

420           Here, we applied this technique by using the new *unfold* toolbox (Ehinger & Dimigen,  
421 2018), which represents the first-level analysis and provides us with the partial effects (i.e.  
422 the beta coefficients or “regression-ERPs”, Smith & Kutas, 2015a, 2015b) for each predictor  
423 of interest. In a first step, both stimulus onset events and fixation onset events were included  
424 as stick functions (also called finite impulse responses, FIR) in the design matrix of the  
425 regression model. To account for overlapping activity from adjacent experimental events, the  
426 design matrix was then time-expanded in a time window between -300 and +800 ms around  
427 each stimulus and fixation onset event. Time-expansion means that the time points within this  
428 window are added as predictors to the regression model. Because the temporal distance  
429 between subsequent events in the experiment is variable, it is possible to disentangle their  
430 overlapping responses. Time-expansion with stick functions is explained in Serences (2004)  
431 and Ehinger & Dimigen (2018, see their Figure 2). The model was run on EEG data sampled  
432 at the original 512 Hz, that is, no down-sampling was performed.

433           Using Wilkinson notation, the model formula for scene onset events was defined as:

434 
$$\text{ERP} \sim 1 + \text{Consistency}$$

435 In this formula, the beta coefficients for the intercept (1) capture the shape of the overall  
 436 waveform of the stimulus-ERP in the consistent condition, which was used as the reference  
 437 level, whereas those for *Consistency* capture the differential effect of presenting an  
 438 inconsistent object in the scene (relative to a consistent object) on the ERP. The coefficients  
 439 for the predictor *Consistency* are therefore analogous to a difference waveform in a traditional  
 440 ERP analysis (Smith & Kutas, 2015a, 2015b) and would reveal if semantic processing  
 441 already occurs immediately after the initial presentation of the scene.

442 In the same regression model, we also included the onsets of all fixations made on the  
 443 scene. Fixation onsets were modelled with the formula

444 
$$\text{FRP} \sim 1 + \text{Consistency} * \text{Type} + \text{Sacc\_Amplitude}$$

445 Thus, we predicted the FRP for each time-point as a function of the semantic *Consistency* of  
 446 the target object (Consistent vs. Inconsistent; Consistent as reference level) in interaction with  
 447 the *Type* of fixation (Critical fixation vs. Non-target fixation; Non-target fixation as reference  
 448 level). In this model, any FRP consistency effects elicited by the pre-target or target fixation  
 449 would appear as an interaction between *Consistency* and fixation *Type*. In addition, we  
 450 included the incoming *Saccade Amplitude* (in degrees of visual angle) as a continuous linear  
 451 covariate to control for the effect of saccade size on the FRP waveform<sup>2</sup>. Thus, the full model  
 452 was as follows:

---

<sup>2</sup> Other low-level variables, such as local image features in the currently foveated image region (e.g., luminance, spatial frequency) are also known to modulate the FRP waveform. In the model presented here, we did not include these other covariates because (1) their influence on the FRP waveform is small compared to that of saccade amplitude and (2) the properties of the target object (such as its visual saliency) did not differ between the two levels of object consistency (see *Materials and Rating*). For reasons of simplicity, saccade amplitude was included as a linear predictor in the current model, although its influence on the FRP becomes non-linear for large saccades (e.g., Dandekar et al., 2011)..



467 Both runs of the model (the one for  $t-1$  and  $t$ ) also yield an estimate for the scene  
468 onset-ERP, but because the results for the scene-onset ERP were virtually identical, we  
469 present the betas from the first run of the model.

470 The average number of events entering the model per participant was 65.7 and 73.6 for  
471 scene onsets (consistent and inconsistent condition, respectively), 864.2 and 887.9 for non-  
472 target fixations ( $nt$ ), 58.3 and 59.8 for pre-target fixations ( $t-1$ ), and 63.8 and 70.6 for target  
473 fixations ( $t$ ).

474 *Baseline placement for FRPs:* Another challenging issue for free-viewing EEG  
475 experiments is the choice of an appropriate neutral baseline interval for the FRP waveforms  
476 (Dimigen et al., 2011; Nikolaev et al., 2016). Baseline placement is particularly relevant for  
477 experiments on extrafoveal processing where we do not know in advance when EEG  
478 differences will arise, and whether they may already develop prior to fixation onset.

479 For the pre-target fixation  $t-1$  and non-target fixations  $nt$ , we used a standard baseline  
480 interval by subtracting the mean channel voltages between -200 and 0 ms before the event  
481 (note that the saccadic spike potential ramping up at the end of this interval was almost  
482 completely removed by our ICA procedure; see Figure A1). For fixation  $t$ , we cannot use  
483 such a baseline because semantic processing may already be ongoing by the time the target  
484 object is fixated. Thus, to apply a neutral baseline to fixation  $t$ , we subtracted the mean  
485 channel voltages in the 200 ms interval before the *preceding* fixation  $t-1$  also from the FRP  
486 aligned to the target fixations  $t$  (see Nikolaev et al., 2016 for similar procedures). The scene-  
487 onset ERP was corrected with a standard pre-stimulus baseline (-200 to 0 ms).

488 *Group statistics for EEG (second level of analysis):* To perform second-level group  
489 statistics, averaged EEG waveforms at the single-subject level (“regression-ERPs”) were  
490 reconstructed from the beta coefficients of the linear deconvolution model. These regression-  
491 based ERPs are directly analogous to subject-level averages in a traditional ERP analysis

492 (Smith & Kutas, 2015a). We then used two complementary statistical approaches to examine  
493 consistency effect in the EEG: linear mixed models and a cluster-based permutation test.

494 *LMM in a-priori defined time windows.* LMM were used to provide hypothesis-based  
495 testing motivated by existing literature. Specifically, we adopted the spatiotemporal  
496 definitions by Vö & Wolfe (2013) and compared the consistent and inconsistent condition in  
497 the time windows from 250 – 350 ms (early effect) and 350 – 600 ms (late effect) at a mid-  
498 central region-of-interest (ROI) of nine electrodes (comprising FC1, FCz, FC2, C1, Cz, C2,  
499 CP1, CPz, and CP2). Because the outputs provided by the linear deconvolution model (the  
500 first-level analysis) are already aggregated at the level of subject-averages, the only predictor  
501 included in these LMMs was the *Consistency* of the object. Furthermore, to minimize the risk  
502 of Type I error (Barr, Levy, Scheepers, & Tily, 2013) we started with a random effect  
503 structure with *Participant* as random intercept and slope for the *Consistency* predictor. This  
504 random effect structure was then evaluated and backwards-reduced using the `step` function  
505 of the `lmerTest` package (Kuznetsova, Brockhoff, & Christensen, 2017) to retain the model  
506 that was justified by the data, i.e., it converged, and it was parsimonious in the number of  
507 parameters (Matuschek, Kliegl, Vasishth, Baayen, & Bates, 2017).

508 *Cluster permutation tests.* It is still largely unknown to what extent the topography of  
509 traditional ERP effects translates to natural viewing. Therefore, in order to test for  
510 consistency effects across all channels and time points, we additionally applied the  
511 Threshold-Free Cluster Enhancement (TFCE) procedure developed by Smith & Nichols  
512 (2009) and adapted for EEG data by Mensen & Khatami (2013,  
513 [http://github.com/Mensen/ept\\_TFCE-matlab](http://github.com/Mensen/ept_TFCE-matlab)). In a nutshell, TFCE is a non-parametric  
514 permutation test that controls for multiple comparisons across time and space, while  
515 maintaining relatively high sensitivity (e.g. compared to a Bonferroni correction). Its  
516 advantage over previous cluster permutation tests (e.g., Maris & Oostenveld, 2007) is that it

517 does not require the experimenter to set an arbitrary cluster-forming threshold. In the first  
518 stage of the TFCE procedure, a raw statistical measure (here:  $t$ -values) is weighted according  
519 to the support provided by clusters of similar values at surrounding electrodes and time  
520 points. In the second stage, these cluster-enhanced  $t$ -values are then compared to the  
521 maximum cluster-enhanced values observed under the null hypotheses (based on  $n=2000$   
522 random permutations of the data). In the present manuscript (Figures 4 and 5), we not only  
523 report the global result of the test, but also plot the spatiotemporal extent of the first-stage  
524 clusters, since they provide some indication about which time points and electrodes likely  
525 contributed to the overall significant effect established by the test. Please note, however, that  
526 unlike the global test result, these first-stage values are not stringently controlled for false  
527 positives and do not establish precise effect onset or offsets (Sassenhagen & Draschkow,  
528 2019). We report them here as a descriptive statistic.

529 Insert Figure 3 and Tables 1, 2, 3 about here

530 Finally, for purely descriptive purposes and to provide a-priori information for future  
531 studies, we also plot the 95% between-subject confidence interval for the consistency effects  
532 at the central ROI (corresponding to sample-by-sample paired  $t$ -testing without correction for  
533 multiple comparisons; see also Mudrik et al., 2014) in Figures 4 and 5.

## 534 **Results**

### 535 ***Task performance (change detection task)***

536 Following the recognition phase, participants pressed a button to indicate whether or  
537 not a change had taken place within the scene. Response accuracy in this task was high (mean  
538 = 85.0%± 35.7%) and did not differ as a function of whether the study scene contained a  
539 consistent (84.6%± 36.1%) or an inconsistent (85.3%± 35.3%) target object.

540 *Eye-movement behaviour*

541 Figure 3A shows the cumulative probability of having fixated the target object as a  
542 function of the ordinal number of fixation and semantic consistency, and Table 2 reports the  
543 corresponding GLM model coefficients. We found a significant main effect of *Consistency*;  
544 overall, inconsistent objects were looked at with a higher probability than consistent objects.  
545 As expected, the cumulative probability of looking at the critical object increased as a  
546 function of the *Ordinal Number of Fixation*. There was also a significant interaction between  
547 the two variables.

548 Complementing this global analysis, we analysed the very first eye-movement during  
549 scene exploration to assess whether observers had immediate extrafoveal access to object-  
550 scene semantics (Loftus & Mackworth, 1978). The mean probability of immediate object  
551 fixation was 12.77%; we observed a numeric advantage of inconsistent objects over  
552 consistent objects (Figure 3B) but this difference was not significant (Table 3). The latency to  
553 first fixation on the target object is another measure to capture the potency of an object in  
554 attracting early attention in extrafoveal vision (e.g., Underwood & Foulsham, 2006; Vö &  
555 Henderson, 2009). This measure is defined as the time elapsed between the onset of the scene  
556 image and the first fixation on the critical object. Importantly, this latency was significantly  
557 shorter for inconsistent as compared to consistent objects (Figure 3C, Table 3).

558 Moreover, we analysed gaze duration as a measure of foveal object processing time  
559 (e.g., Henderson et al., 1999). First-pass gaze duration for a critical object is defined as the  
560 sum of all fixation durations from first entry to first exit. On average, participants looked  
561 longer at inconsistent (520 ms) than consistent objects (409 ms) before leaving the target  
562 object for the first time, and this difference was significant (Table 3). Table 1 summarizes  
563 additional oculomotor characteristics in the two conditions of object consistency.

564





590 5B; see also Figure A6). The scalp distribution of this difference, shown in Figure 6, closely  
591 resembled the frontocentral N400 (and N300) previously reported in ERP studies on object-  
592 scene consistency (e.g., Mudrik et al., 2014; Vö & Wolfe, 2013). In the LMM analyses  
593 conducted on the mid-central ROI, this effect was marginally significant ( $p < 0.1$ ) for the  
594 early time window (250 to 350 ms), but became highly significant between 350 and 600 ms  
595 ( $p < 0.001$ , Table 4). The TFCE test across all channels and time points revealed a significant  
596 effect of consistency on the pre-target FRP ( $p < 0.05$ ). Figure 5C also shows the extents of the  
597 underlying spatiotemporal clusters, computed in the first stage of the TFCE procedure.  
598 Between 372 ms and 721 ms after fixation onset, we observed a cluster of 14 frontocentral  
599 electrodes that was shifted slightly to the left hemisphere. This N400 modulation on the pre-  
600 target fixation could be seen even in traditionally-averaged FRP waveforms without any  
601 control of overlapping potentials (see Supplementary Figure S3). In summary, we were able  
602 to measure a significant frontocentral N400 modulation during natural scene viewing that  
603 already emerged in FRPs aligned to the pre-target fixation.

604 On average, the target fixation  $t$  occurred at a median latency of 240 ms ( $\pm 18$  ms)  
605 after fixation  $t-1$ , as marked by the vertical dashed line in Figure 5B. If we take the extent of  
606 the cluster from the TFCE tests as a rough approximation for the likely onset of the effect in  
607 the FRP, this means that, on average, at the time when the ERP consistency effect started  
608 (372 ms) the eyes had been looking at the target object for only 132 ms (372 minus 240 ms).

609 **Target fixation,  $t$ .** An anterior N400 effect was also clearly visible in the FRP aligned  
610 to fixation  $t$ . In the LMM analysis at the central ROI, the effect was significant in both the  
611 early (250-350 ms,  $p < 0.01$ ) and late window (350-600 ms,  $p < 0.05$ ; see Table 4). However,  
612 compared to the effect aligned to the pre-target fixation, this N400 was significant at only a  
613 few electrodes in the TFCE statistic (Cz, FCz, and FC1; see Figure 6). Aligned to the target  
614 fixation  $t$ , the N400 also peaked extremely early, with the maximum of the difference curve

615 already observed at 200 ms after fixation onset (Figure 5G). Qualitatively, a frontocentral  
616 negativity was already visible much earlier than that, within the first 100 ms after fixation  
617 onset (Figure 5D). The TFCE permutation test confirmed an overall effect of consistency ( $p <$   
618 0.05) on the target-locked FRP. Figure 5G also shows the extents of the underlying first-stage  
619 clusters. For the target fixation, clusters only extended across a brief interval between 151 and  
620 263 ms after fixation onset, an interval during which the N400 effect also reached its peak.

621 Figure 5E shows that, numerically, voltages at the central ROI were more negative in  
622 the inconsistent condition during the baseline interval already, that is, before the critical  
623 object was fixated. To understand the role of activity already present before fixation onset, we  
624 repeated the FRP analyses for fixation  $t$  after applying a standard baseline correction, with the  
625 baseline placed immediately before the target fixation itself (-200 to 0 ms). This way, we  
626 eliminate any weak N400-like effects that may have already been on-going before target  
627 fixation onset. Interestingly, in the resulting FRP waveforms, the target-locked N400 effects  
628 were weakened: The N400 effect now failed to reach significance in the TFCE statistic and in  
629 the LMM analysis for the second window (350 to 600 ms; see last row of Table 4) and only  
630 remained significant for the early window (250 to 350 ms). This indicates that some N400-  
631 like negativity was already ongoing before target fixation onset. To summarize, we found no  
632 immediate influences of object-scene consistency in ERPs time-locked to scene onset.  
633 However, N400 consistency effects were found in FRPs aligned to the target fixation ( $t$ ) and  
634 in those aligned to the pre-target fixation ( $t-1$ ).

## 635 Discussion

636 Substantial research in vision science has been devoted to understanding the  
637 behavioural and neural mechanisms underlying object recognition (e.g., Biederman, 1972;  
638 Loftus & Mackworth, 1978). At the core of this debate are the type of object features that are  
639 accessed (e.g., low-level vs. high-level), the time-course of their processing (e.g., pre-

640 attentive vs. attentive), and the region of the visual field in which these features can be  
641 acquired (e.g., foveal vs. extrafoveal). A particularly controversial topic is whether and how  
642 quickly the semantic properties of objects are available outside foveal vision.

643 In the current study, we approached these questions from a new perspective by co-  
644 registering eye-movements and EEG while participants freely inspected images of real-world  
645 scenes in which a critical object was either consistent or inconsistent with the scene context.  
646 As a novel finding, we demonstrate a fixation-related N400 effect during natural scene  
647 viewing. Moreover, behavioural and electrophysiological measures converge to suggest that  
648 the extraction of object-scene semantics can already begin in extrafoveal vision, before the  
649 critical object is fixated.

650 It is a rather undisputed finding that inconsistent objects, such as a flashlight in a  
651 bathroom, require increased processing when selected as targets of overt attention.  
652 Accordingly, several eye-movement studies have reported longer gaze durations on  
653 inconsistent than consistent objects, probably reflecting the greater effort required to resolve  
654 the conflict between object meaning and scene context (e.g., Cornelissen & Võ, 2017; De  
655 Graef et al., 1990; Henderson et al., 1999). In addition, a number of traditional ERP studies  
656 using steady-fixation paradigms have found that inconsistent objects elicit a larger negative  
657 brain response at frontocentral channels (an N300/N400 complex) as compared to consistent  
658 objects (e.g., Coco et al., 2017; Ganis & Kutas, 2003; Mudrik et al., 2010).

659 However, previous research with eye-movements remained inconclusive on whether  
660 semantic processing can take place prior to foveal inspection of the object. Evidence in  
661 favour of extrafoveal processing of object-scene semantics comes from studies in which  
662 inconsistent objects were selected for fixation earlier than consistent ones (e.g., Borges et al.,  
663 2019; LaPointe & Milliken, 2016; Underwood et al., 2008). However, other studies have not  
664 found evidence for earlier selection of inconsistent objects (e.g., De Graef et al., 1990;

665 Henderson et al., 1999; Võ & Henderson, 2009, 2011). Extrafoveal and peripheral vision are  
666 known to be crucial for saccadic programming (e.g., Nuthmann, 2014). Therefore, any  
667 demonstration that semantic information can act as a source of guidance for fixation selection  
668 in scenes implies that some semantic processing must have occurred prior to its fixation, that  
669 is, in extrafoveal vision.

670 ERPs are highly sensitive to semantic processing (Kutas & Federmeier, 2011) and  
671 provide excellent temporal resolution to investigate the time course of object processing.  
672 However, an obvious limitation of existing ERP studies is that observers were not allowed to  
673 explore the scene with saccadic eye-movements, thereby constraining their normal attentional  
674 dynamics. Instead, the critical object was usually large and/or placed near the point of  
675 fixation. Hence, these studies were unable to establish whether semantic processing can take  
676 place prior to its foveal inspection.

677 In the current study, we addressed this problem by simultaneously recording  
678 behavioural and brain-electric correlates of object processing. Specifically, we analysed  
679 different eye-movement responses that tap into extrafoveal and foveal processing along with  
680 FRPs time-locked to the first fixation on the critical object ( $t$ ) and the fixation preceding it  
681 ( $t-1$ ). We also analysed the scene-onset ERP evoked by the trial-initial presentation of the  
682 image. Recent advances in linear deconvolution methods for EEG (e.g., Ehinger & Dimigen,  
683 2018) allowed us to disentangle the overlapping brain potentials produced by the scene onset  
684 and the subsequent fixations, and to control for the modulating influence of saccade  
685 amplitude on the FRP.

686 The eye-movement behaviour showed no evidence for hypothesis (A) as outlined in  
687 the *Introduction*, according to which semantic information can exert an immediate effect on  
688 eye-movement control (Loftus & Mackworth, 1978). Specifically, the mean probability of  
689 immediate object fixation was fairly low (12.8%) and not modulated by *Consistency*. Instead,

690 the data lends support to hypothesis (B) according to which extrafoveal processing of object-  
691 scene semantics is possible but takes some time to unfold. In particular, the results for the  
692 latency to first fixation of the critical object show that inconsistent objects were, on average,  
693 looked at sooner than consistent objects (cf. Bonitz & Gordon, 2008; Underwood et al.,  
694 2008). At the same time, we observed longer gaze durations on inconsistent objects,  
695 replicating previous findings (e.g., De Graef et al., 1990; Henderson et al., 1999; Vö &  
696 Henderson, 2009). Thus, we found behavioural evidence for the extrafoveal processing of  
697 object-scene (in)consistencies, but also differences in the subsequent foveal processing.

698         The question then remains why existing eye-movement studies have provided very  
699 different results, ranging from rapid processing of semantic information in peripheral vision  
700 to a complete lack of evidence for extrafoveal semantic processing. Researchers have  
701 suggested that the outcome may depend on factors related to the critical object or the scene in  
702 which it is located. Variables that may (or may not) facilitate the appearance of the  
703 incongruency effect include visual saliency (e.g., Henderson et al., 1999; Underwood &  
704 Foulsham, 2006), image clutter (Henderson & Ferreira, 2004), and the critical object's size  
705 and eccentricity (Gareze & Findlay, 2007). Therefore, an important question for future  
706 research is to identify the specific conditions under which extrafoveal semantic information  
707 can be extracted, or when the three outlined hypotheses and/or outcomes would prevail.

708         Returning to the present data, the FRP waveforms showed a negative shift over frontal  
709 and central scalp sites when participants fixated a scene-inconsistent object. This result is in  
710 agreement with traditional ERP studies that have shown an frontocentral N300/N400  
711 complex after passive foveal stimulation (e.g., Coco et al., 2017; Ganis & Kutas, 2003;  
712 Mudrik et al., 2014; Vö & Wolfe, 2013) and extends this finding for the first time to a natural  
713 viewing situation with eye-movements. Regarding the time course, the present data suggest

714 that the effect was already initiated during the preceding fixation ( $t-1$ ), but then carried on  
715 through fixation ( $t$ ) on the target object.

716 As a cautionary note, we emphasize that it is not trivial to unambiguously ascribe  
717 typical N400 (and N300) effects in the EEG to either extrafoveal or foveal processing. The  
718 reason is that these canonical congruency effects only begin 200-250 ms after stimulus onset  
719 (Draschkow et al., 2018; Mudrik et al., 2010). This means that even a purely extrafoveal  
720 effect would be almost impossible to measure during the pre-target fixation ( $t-1$ ) itself, since  
721 it would only emerge at a time when the eyes are already moving to the target object. That  
722 being said, three properties of the observed FRP consistency effect suggest that it was already  
723 initiated during the pre-target fixation:

724 First, due to the temporal jitter introduced by variable fixation durations, an effect that  
725 only arises in foveal vision should be the most robust in the FRP averages aligned to fixation  
726  $t$ , but latency-jittered and attenuated in those aligned to fixation  $t-1$ . However, the opposite  
727 was the case: At least qualitatively, a frontocentral N400 effect was seen at more electrodes  
728 (Figure 6) and for longer time intervals (Figure 5) in the FRP aligned to the pre-target fixation  
729 as compared to the actual target fixation. The second argument for extrafoveal contributions  
730 to the effect is the forward-shift in its time course. Relative to fixation  $t$ , the observed N400  
731 occurred almost instantly: As the effect topographies in Figure 5H show, the frontocentral  
732 negativity for inconsistent objects was qualitatively visible within the first 100 ms after  
733 fixation onset and the effect reached its peak after just 200 ms. Clusters underlying the TFCE  
734 test were also restricted to an early time range between 151 and 263 ms after fixation onset  
735 and therefore to a much earlier interval to what we would expect from the canonical N300 or  
736 N400 effect elicited by foveal stimulation.

737 Of course, it is possible that even purely foveal N400 effects may emerge earlier  
738 during active scene exploration with eye movements as compared to the latencies established

739 in traditional ERP research. For example, it is reasonable to assume that during natural vision,  
740 observers pre-process some low-level (non-semantic) features of the soon-to-be fixated object  
741 in extrafoveal vision (cf. Nuthmann, 2017). This non-semantic preview benefit might then  
742 speed up the timeline of foveal processing (including the latency of semantic access) once the  
743 object is fixated (cf. Dimigen, et al., 2012, for reading). Moreover, if eye movements are  
744 permitted, observers have more time to build a representation of the scene before they foveate  
745 the target, and this increased contextual constraint may also affect the N400 timing (but see  
746 Kutas & Hillyard, 1984). Importantly, however, neither of these two accounts could explain  
747 why the N400 effect is stronger – rather than much weaker – in the waveforms aligned to  
748 fixation  $t-1$  as compared to fixation  $t$ . The fact that the eye movement data also provided  
749 clear evidences in favour of extra-foveal processing further strengthens our interpretation of  
750 the N400 timing.

751 Finally, we found that the N400 consistency effect aligned to the target fixation ( $t$ )  
752 became weaker (and non-significant in two out of the three statistical measures considered) if  
753 the baseline interval for the FRP analysis was placed directly before this target fixation.  
754 Again, this indicates that at least a weak frontocentral negativity in the inconsistent condition  
755 was already present during the baseline period before the target was fixated. Together, these  
756 results are difficult to reconcile with a pure foveal processing account and are more consistent  
757 with the notion that semantic processing of the object was at least initiated in extrafoveal  
758 vision (and then continued after it was foveated).

759 Crucially, we did not find any effect of target consistency in the traditional ERP  
760 aligned to scene onset. In line with the behavioural results, this goes against the most extreme  
761 hypothesis A postulating that object semantics can be extracted from peripheral vision  
762 already at the first glance of a scene (Loftus & Mackworth, 1978). Similarly, there was no  
763 effect of consistency on the FRPs evoked by the non-target fixations on the scene (Figure 4);



764 this was also the case in a control analysis that only included non-target fixations that  
765 occurred earlier than  $t-1$  and at an extrafoveal distance between  $3^\circ$  and  $7^\circ$  from the target  
766 object (see Supplementary Figure S10). All these analyses suggest that the semantic  
767 information of the critical object started during fixation  $t-1$ . However, from any given fixation  
768 there are many candidate locations that could potentially be chosen for the next saccade (cf.  
769 Tatler, Brockmole, & Carpenter, 2017). Thus, it is conceivable that observers may have  
770 partially acquired semantic information of the critical object outside foveal vision prior to  
771 fixation  $t-1$ , but without selecting it as a saccade target. Such reasoning leaves open the  
772 possibility that observers may have already picked up some information about the target  
773 object's semantics during these occasions.

774         Taken together, our behaviour and electrophysiological findings are consistent with  
775 the claim formulated in hypothesis B that objects can be recognized outside of the fovea or  
776 even in the visual periphery, at least to some degree. Indirectly, our results also speak to the  
777 debate about the unit of saccade targeting and, by inference, attentional selection during scene  
778 viewing. Finding effects of object-scene semantics on eye guidance is evidence in favour of  
779 object- and meaning-based, rather than image-based guidance of attention in scenes (e.g.,  
780 Hwang, Wang, & Pomplun, 2011; Henderson, Hayes, Peacock, & Rehrig, 2019).

781         In sum, our findings converge to suggest that the visual system is capable of accessing  
782 semantic features of objects in extrafoveal vision to guide attention towards objects that do  
783 not fit to the scene's overall meaning. They also highlight the utility of investigating  
784 attentional and neural mechanisms in parallel to uncover the mechanisms underlying object  
785 recognition during the unconstrained exploration of naturalistic scenes.

786

787

788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818

### References

- Andrews, S., & Veldre, A. (2019). What is the most plausible account of the role of parafoveal processing in reading? *Language and Linguistics Compass*, *13*(7), e12344.
- Antes, J. R. (1974). The time course of picture viewing. *Journal of Experimental Psychology*, *103*(1), 62.  
<https://doi.org/10.1037/h0036799>
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*(3), 255–278.  
<https://doi.org/10.1016/j.jml.2012.11.001>
- Bates, D., Mächler, M., Bolker, B. M., Walker, S. C., Machler, M., Bolker, B. M., ... Walker, S. C. (2015). Fitting Linear Mixed-Effects Models using lme4. *Journal of Statistical Software*, *67*(1), 1–48.  
<https://doi.org/10.18637/jss.v067.i01>
- Belke, E., Humphreys, G. W., Watson, Derrick, G., Meyer, A. S., & Telling, A. L. (2008). Top-down effects of semantic knowledge in visual search are modulated by cognitive but not perceptual load. *Perception & Psychophysics*, *70*(8), 1444–1458. <https://doi.org/10.3758/PP>
- Biederman, I. (1972). Perceiving real-world scenes. *Science*, *177*, 77–80.  
<https://doi.org/10.1126/science.177.4043.77>
- Bonitz, V. S., & Gordon, R. D. (2008). Attention to smoking-related and incongruous objects during scene viewing. *Acta Psychologica*, *129*(2), 255–263. <https://doi.org/10.1016/j.actpsy.2008.08.006>
- Borges, M. T., Fernandes, E. G., & Coco, M. I. (2019). Age-related differences during visual search: the role of contextual expectations and cognitive control mechanisms. *Aging, Neuropsychology, and Cognition*.  
<https://doi.org/10.1080/13825585.2019.1632256>
- Brouwer, A.-M., Reuderink, B., Vincent, J., van Gerven, M. A. J., & van Erp, J. B. F. (2013). Distinguishing between target and nontarget fixations in a visual search task using fixation-related potentials. *Journal of Vision*, *13*(3), :17, 1-17. <https://doi.org/10.1167/13.3.17>
- Cimminella, F., Della Sala, S., & Coco, M. I. (in press). Extra-foveal Processing of Object Semantics Guides Early Overt Attention During Visual Search. *Attention, Perception, & Psychophysics*.  
<https://doi.org/10.3758/s13414-019-01906-1>
- Coco, M. I., Araujo, S., & Petersson, K. M. (2017). Disentangling stimulus plausibility and contextual congruency: Electro-physiological evidence for differential cognitive dynamics. *Neuropsychologia*, *96*, 150–163. <https://doi.org/10.1016/j.neuropsychologia.2016.12.008>

- 819 Cornelissen, T. H. W., Sassenhagen, J., & Vö, M. L. H. (2019). Improving free-viewing fixation-related EEG  
 820 potentials with continuous-time regression. *Journal of Neuroscience Methods*, *313*, 77–94.  
 821 <https://doi.org/10.1016/j.jneumeth.2018.12.010>
- 822 Cornelissen, T. H. W., & Vö, M. L.-H. (2017). Stuck on semantics: Processing of irrelevant object-scene  
 823 inconsistencies modulates ongoing gaze behavior. *Attention, Perception, & Psychophysics*, *79*(1), 154–  
 824 168. <https://doi.org/10.3758/s13414-016-1203-7>
- 825 Dandekar, S., Privitera, C., Carney, T., & Klein, S. A. (2011). Neural saccadic response estimation during  
 826 natural viewing. *Journal of Neurophysiology*, *107*(6), 1776–1790. <https://doi.org/10.1152/jn.00237.2011>
- 827 Davenport, J. L., & Potter, M. C. (2004). Scene consistency in object and background perception. *Psychological*  
 828 *Science*, *15*(8), 559–564. <https://doi.org/10.1111/j.0956-7976.2004.00719.x>
- 829 De Graef, P., Christiaens, D., & D’Ydewalle, G. (1990). Perceptual effects of scene context on object  
 830 identification. *Psychological Research*, *52*(4), 317–329. <https://doi.org/10.1007/BF00868064>
- 831 Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics  
 832 including independent component analysis. *Journal of Neuroscience Methods*, *134*, 9–21.  
 833 <https://doi.org/10.1016/j.jneumeth.2003.10.009>
- 834 Devillez, H., Guyader, N., & Guerin-Dugue, A. (2015). An eye fixation-related potentials analysis of the P300  
 835 potential for fixations onto a target object when exploring natural scenes. *Journal of Vision*, *15*((13); 20),  
 836 1–31. <https://doi.org/10.1167/15.13.20>
- 837 Dimigen, O. (2018). Optimized ICA-based removal of ocular EEG artifacts from free viewing experiments.  
 838 *BioRxiv*, (0), 446955. <https://doi.org/10.1101/446955>
- 839 Dimigen, O., & Ehinger, B. (2019). Analyzing combined eye-tracking/EEG experiments with (non)linear  
 840 deconvolution models. *BioRxiv*. <https://doi.org/10.1101/735530>
- 841 Dimigen, O., Kliegl, R., & Sommer, W. (2012). Trans-saccadic parafoveal preview benefits in fluent reading: A  
 842 study with fixation-related brain potentials. *NeuroImage*, *62*(1), 381–393.  
 843 <https://doi.org/10.1016/j.neuroimage.2012.04.006>
- 844 Dimigen, O., Sommer, W., Hohlfeld, A., Jacobs, A. M., & Kliegl, R. (2011). Coregistration of Eye Movements  
 845 and EEG in Natural Reading: Analyses and Review. *Journal of Experimental Psychology: General*,  
 846 *140*(4), 552–572. <https://doi.org/10.1037/a0023885>
- 847 Draschkow, D., Heikel, E., Vö, M. L.-H., Fiebach, C. J., & Sassenhagen, J. (2018). No evidence from MVPA  
 848 for different processes underlying the N300 and N400 incongruity effects in object-scene processing.

- 849 *Neuropsychologia*, 120, 9–17. <https://doi.org/10.1016/J.NEUROPSYCHOLOGIA.2018.09.016>
- 850 Dyck, M., & Brodeur, M. B. (2015). ERP evidence for the influence of scene context on the recognition of  
851 ambiguous and unambiguous objects. *Neuropsychologia*, 72, 43–51.  
852 <https://doi.org/10.1016/j.neuropsychologia.2015.04.023>
- 853 Ehinger, B., & Dimigen, O. (2018). Unfold: An integrated toolbox for overlap correction, non-linear modeling,  
854 and regression-based EEG analysis. *BioRxiv*, 360156. <https://doi.org/10.1101/360156>
- 855 Feldman, J. (2003). What is a visual object? *Trends in Cognitive Sciences*, 7(6), 252–256.  
856 [https://doi.org/10.1016/S1364-6613\(03\)00111-6](https://doi.org/10.1016/S1364-6613(03)00111-6)
- 857 Fenske, M. J., Aminoff, E., Gronau, N., & Bar, M. (2006). Chapter 1 Top-down facilitation of visual object  
858 recognition: object-based and context-based contributions. *Progress in Brain Research*, 155 B, 3–21.  
859 [https://doi.org/10.1016/S0079-6123\(06\)55001-0](https://doi.org/10.1016/S0079-6123(06)55001-0)
- 860 Ganis, G., & Kutas, M. (2003). An electrophysiological study of scene effects on object identification. *Cognitive*  
861 *Brain Research*, 16(2), 123–144. [https://doi.org/10.1016/S0926-6410\(02\)00244-6](https://doi.org/10.1016/S0926-6410(02)00244-6)
- 862 Garcia-Diaz, A., Fdez-Vidal, X. R., Pardo, X. M., & Dosil, R. (2012). Saliency from hierarchical adaptation  
863 through decorrelation and variance normalization. *Image and Vision Computing*, 30(1), 51–64.  
864 <https://doi.org/10.1016/j.imavis.2011.11.007>
- 865 Gareze, L., & Findlay, J. M. (2007). Absence of scene context effects in object detection and eye gaze capture.  
866 In R. P. G. van Gompel & R. L. Hill (Eds.), *Eye movements: A window on mind and brain* (pp. 537–562).  
867 Oxford: Elsevier.
- 868 Hauk, O., Davis, M. H., Ford, M., Pulvermüller, F., & Marslen-Wilson, W. D. (2006). The time course of visual  
869 word recognition as revealed by linear regression analysis of ERP data. *NeuroImage*, 30(4), 1383–1400.  
870 <https://doi.org/10.1016/j.neuroimage.2005.11.048>
- 871 Henderson, J. M., & Ferreira, F. (2004). Scene perception for psycholinguists. In J. M. Henderson & F. Ferreira  
872 (Eds.), *The interface of language, vision, and action: Eye movements and the visual world* (pp. 1–58).  
873 New York: Psychology Press.
- 874 Henderson, J. M., Hayes, T. R., Peacock, C. E., & Rehrig, G. (2019). Meaning and attentional guidance in  
875 scenes: A review of the meaning map approach. *Vision*, 3(2), 19. <https://doi.org/10.3390/vision3020019>
- 876 Henderson, J. M., Weeks, Jr, P. A., & Hollingworth, A. (1999). The effects of semantic consistency on eye  
877 movements during complex scene viewing. *Journal of Experimental Psychology: Human Perception and*  
878 *Performance*, 25(1), 210–228. <https://doi.org/10.1037//0096-1523.25.1.210>

- 879 Hohenstein, S., & Kliegl, R. (2014). Semantic preview benefit during reading. *Journal of Experimental*  
 880 *Psychology: Learning Memory and Cognition*, *40*(1), 166–190. <https://doi.org/10.1037/a0033670>
- 881 Hwang, A. D., Wang, H. C., & Pomplun, M. (2011). Semantic guidance of eye movements in real-world scenes.  
 882 *Vision Research*, *51*(10), 1192–1205. <https://doi.org/10.1016/j.visres.2011.03.010>
- 883 Itti, L., Koch, C., & Niebur, E. (1998). A model of saliency-based visual attention for rapid scene analysis. *IEEE*  
 884 *Transaction on Pattern Analysis and Machine Learning*, *20*(11), 1254–1259.  
 885 <https://doi.org/10.1109/TPAMI.2012.125>
- 886 Jung, T.-P., Humphries, C., Lee, T.-W., Makeig, S., McKeown, M. J., Iragui, V., & Sejnowski, T. J. (1998).  
 887 Extended ICA removes artifacts from electroencephalographic recordings. *Advances in Neural*  
 888 *Information Processing Systems*, *10*, 894–900.
- 889 Kamienkowski, J. E., Ison, M. J., Quiroga, R. Q., & Sigman, M. (2012). Fixation-related potentials in visual  
 890 search: A combined EEG and eye tracking study. *Journal of Vision*, *12*(7), :4, 1-20.  
 891 <https://doi.org/10.1167/12.7.4>
- 892 Kaunitz, L. N., Kamienkowski, J. E., Varatharajah, A., Sigman, M., Quiroga, R. Q., & Ison, M. J. (2014).  
 893 Looking for a face in the crowd: Fixation-related potentials in an eye-movement visual search task.  
 894 *NeuroImage*, *89*, 297–305. <https://doi.org/10.1016/j.neuroimage.2013.12.006>
- 895 Kliegl, R., Dambacher, M., Dimigen, O., Jacobs, A. M., & Sommer, W. (2012). Eye movements and brain  
 896 electric potentials during reading. *Psychological Research*. <https://doi.org/10.1007/s00426-011-0376-x>
- 897 Kretzschmar, F., Bornkessel-Schlesewsky, I., & Schlewsky, M. (2009). Parafoveal versus foveal N400s  
 898 dissociate spreading activation from contextual fit. *NeuroReport*, *20*(18), 1613–1618.  
 899 <https://doi.org/10.1097/WNR.0b013e328332c4f4>
- 900 Kristensen, E., Guerin-Dugué, A., & Rivet, B. (2017). Regularization and a general linear model for event-  
 901 related potential estimation. *Behavior Research Methods*, *49*(6), 2255–2274.  
 902 <https://doi.org/10.3758/s13428-017-0856-z>
- 903 Kristensen, E., Rivet, B., & Guerin-Dugué, A. (2017). Estimation of overlapped Eye Fixation Related Potentials:  
 904 The General Linear Model, a more flexible framework than the ADJAR algorithm. *Journal of Eye*  
 905 *Movement Research*, *10*(1), 1–27. <https://doi.org/http://dx.doi.org/10.16910/jemr.10.1.7>
- 906 Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: finding meaning in the N400 component of  
 907 the event-related brain potential (ERP). *Annual Review of Psychology*, *62*(1), 621–647.  
 908 <https://doi.org/10.1146/annurev.psych.093008.131123>

- 909 Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest Package: Tests in Linear Mixed  
 910 Effects Models. *Journal of Statistical Software*, 82(13). <https://doi.org/10.18637/jss.v082.i13>
- 911 LaPointe, M. R. P., & Milliken, B. (2016). Semantically incongruent objects attract eye gaze when viewing  
 912 scenes for change. *Visual Cognition*, 24(1), 63–77. <https://doi.org/10.1080/13506285.2016.1185070>
- 913 Loftus, G. R., & Mackworth, N. H. (1978). Cognitive determinants of fixation location during picture viewing.  
 914 *Journal of Experimental Psychology. Human Perception and Performance*, 4(4), 565–572.  
 915 <https://doi.org/10.1037/0096-1523.4.4.565>
- 916 Mackworth, N. H., & Morandi, A. J. (1967). The gaze selects information details within pictures. *Perception*  
 917 *and Psychophysics*, 2(11), 547–552. <https://doi.org/10.3758/BF03210264>
- 918 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data. *Journal of*  
 919 *Neuroscience Methods*, 164(1), 177–190. <https://doi.org/10.1016/j.jneumeth.2007.03.024>
- 920 Matuschek, H., Kliegl, R., Vasishth, S., Baayen, H., & Bates, D. (2017). Balancing Type I error and power in  
 921 linear mixed models. *Journal of Memory and Language*, 94, 305–315.  
 922 <https://doi.org/10.1016/j.jml.2017.01.001>
- 923 Mensen, A., & Khatami, R. (2013). Advanced EEG analysis using threshold-free cluster-enhancement and non-  
 924 parametric statistics. *NeuroImage*, 67, 111–118. <https://doi.org/10.1016/j.neuroimage.2012.10.027>
- 925 Moores, E., Laiti, L., & Chelazzi, L. (2003). Associative knowledge controls deployment of visual selective  
 926 attention. *Nature Neuroscience*, 6(2), 182–189. <https://doi.org/10.1038/mn996>
- 927 Mudrik, L., Lamy, D., & Deouell, L. Y. (2010). ERP evidence for context congruity effects during simultaneous  
 928 object-scene processing. *Neuropsychologia*, 48(2), 507–517.  
 929 <https://doi.org/10.1016/j.neuropsychologia.2009.10.011>
- 930 Mudrik, L., Shalgi, S., Lamy, D., & Deouell, L. Y. (2014). Synchronous contextual irregularities affect early  
 931 scene processing: Replication and extension. *Neuropsychologia*, 56, 447–458.  
 932 <https://doi.org/10.1016/j.neuropsychologia.2014.02.020>
- 933 Niefind, F., & Dimigen, O. (2016). Dissociating parafoveal preview benefit and parafovea-on-fovea effects  
 934 during reading: A combined eye tracking and EEG study. *Psychophysiology*, 53(12), 1784–1798.  
 935 <https://doi.org/10.1111/psyp.12765>
- 936 Nikolaev, A. R., Meghanathan, R. N., & van Leeuwen, C. (2016). Combining EEG and eye movement recording  
 937 in free viewing: Pitfalls and possibilities. *Brain and Cognition*, 107, 55–83.  
 938 <https://doi.org/10.1016/j.bandc.2016.06.004>

- 939 Nuthmann, A. (2013). On the visual span during object search in real-world scenes. *Visual Cognition*, *21*(7),  
 940 803–837. <https://doi.org/10.1080/13506285.2013.832449>
- 941 Nuthmann, A. (2014). How do the regions of the visual field contribute to object search in real-world scenes?  
 942 Evidence from eye movements. *Journal of Experimental Psychology: Human Perception and*  
 943 *Performance*, *40*(1), 342–360. <https://doi.org/10.1037/a0033854>
- 944 Nuthmann, A., de Groot, F., Huettig, F., & Olivers, C. N. L. (2019). Extrafoveal attentional capture by object  
 945 semantics. *Plos One*, *14*(5), e0217051. <https://doi.org/10.1371/journal.pone.0217051>
- 946 Nuthmann, A., & Einhäuser, W. (2015). A new approach to modeling the influence of image features on fixation  
 947 selection in scenes. *Annals of the New York Academy of Sciences*, *1339*(1), 82–96.  
 948 <https://doi.org/10.1111/nyas.12705>
- 949 Nuthmann, A., & Henderson, J. M. (2010). Object-based attentional selection in scene viewing. *Journal of*  
 950 *Vision*, *10*(8), :20, 1-19. <https://doi.org/10.1167/10.8.20>
- 951 Plöchl, M., Ossandón, J. P., & König, P. (2012). Combining EEG and eye tracking: identification,  
 952 characterization, and correction of eye movement artifacts in electroencephalographic data. *Frontiers in*  
 953 *Human Neuroscience*, *6*(October), 1–23. <https://doi.org/10.3389/fnhum.2012.00278>
- 954 Rämä, P., & Baccino, T. (2010). Eye fixation-related potentials (EFRPs) during object identification. *Visual*  
 955 *Neuroscience*, *27*(5–6), 187–192. <https://doi.org/10.1017/S0952523810000283>
- 956 Rayner, K. (2014). The gaze-contingent moving window in reading: Development and review. *Visual Cognition*,  
 957 *22*(3), 242–258. <https://doi.org/10.1080/13506285.2013.879084>
- 958 Rayner, K., Balota, D. A., & Pollatsek, A. (1986). Against parafoveal semantic preprocessing during eye  
 959 fixations in reading. *Canadian Journal of Psychology*, *40*(4), 473–483. <https://doi.org/10.1037/h0080111>
- 960 Sassenhagen, J., & Draschkow, D. (2019). Cluster-based permutation tests of MEG/EEG data do not establish  
 961 significance of effect latency or location. *Psychophysiology*, *56*(6), 1–8.  
 962 <https://doi.org/10.1111/psyp.13335>
- 963 Satterthwaite, F. E. (1946). An approximate distribution of estimates of variance components. *Biometrics*  
 964 *Bulletin*, *2*(6), 110–114. <https://doi.org/10.2307/3002019>
- 965 Serences, J. T. (2004). A comparison of methods for characterizing the event-related BOLD timeseries in rapid  
 966 fMRI. *NeuroImage*, *21*(4), 1690–1700. <https://doi.org/10.1016/j.neuroimage.2003.12.021>
- 967 Simola, J., Le Fevre, K., Torniaainen, J., & Baccino, T. (2015). Affective processing in natural scene viewing:  
 968 Valence and arousal interactions in eye-fixation-related potentials. *NeuroImage*, *106*, 21–33.

- 969 <https://doi.org/10.1016/j.neuroimage.2014.11.030>
- 970 Smith, N. J., & Kutas, M. (2015a). Regression-based estimation of ERP waveforms: I. The rERP framework.  
 971 *Psychophysiology*, 52(2), 157–168. <https://doi.org/10.1111/psyp.12317>
- 972 Smith, N. J., & Kutas, M. (2015b). Regression-based estimation of ERP waveforms: II. Non-linear effects,  
 973 overlap correction, and practical considerations. *Neuropsychology*, 52(2), 169–181.  
 974 <https://doi.org/10.1111/psyp.12320>
- 975 Smith, S. M., & Nichols, T. E. (2009). Threshold-free cluster enhancement: Addressing problems of smoothing,  
 976 threshold dependence and localisation in cluster inference. *NeuroImage*, 44(1), 83–98.  
 977 <https://doi.org/10.1016/j.neuroimage.2008.03.061>
- 978 Stoll, J., Thrun, M., Nuthmann, A., & Einhäuser, W. (2015). Overt attention in natural scenes: Objects dominate  
 979 features. *Vision Research*, 107, 36–48. <https://doi.org/10.1016/j.visres.2014.11.006>
- 980 Thickbroom, G. W., Knezević, W., Carroll, W. M., & Mastaglia, F. L. (1991). Saccade onset and offset lambda  
 981 waves: relation to pattern movement visually evoked potentials. *Brain Research*, 551(1–2), 150–156.  
 982 [https://doi.org/10.1016/0006-8993\(91\)90927-N](https://doi.org/10.1016/0006-8993(91)90927-N)
- 983 Underwood, G., & Foulsham, T. (2006). Visual saliency and semantic incongruency influence eye movements  
 984 when inspecting pictures. *Quarterly Journal of Experimental Psychology (2006)*, 59(11), 1931–1949.  
 985 <https://doi.org/10.1080/17470210500416342>
- 986 Underwood, G., Templeman, E., Lamming, L., & Foulsham, T. (2008). Is attention necessary for object  
 987 identification? Evidence from eye movements during the inspection of real-world scenes. *Consciousness*  
 988 *and Cognition*, 17(1), 159–170. <https://doi.org/10.1016/j.concog.2006.11.008>
- 989 Ušćumlić, M., & Blankertz, B. (2016). Active visual search in non-stationary scenes: coping with temporal  
 990 variability and uncertainty. *Journal of Neural Engineering*, 13(1), 1–12. [https://doi.org/10.1088/1741-](https://doi.org/10.1088/1741-2560/13/1/016015)  
 991 [2560/13/1/016015](https://doi.org/10.1088/1741-2560/13/1/016015)
- 992 Viola, F. C., Debener, S., Thorne, J., & Schneider, T. R. (2010). Using ICA for the analysis of multi-channel  
 993 EEG data. In *Simultaneous EEG and fMRI: Recording, Analysis, and Application: Recording, Analysis,*  
 994 *and Application* (pp. 121–133).
- 995 Vö, M. L.-H., & Henderson, J. M. (2009). Does gravity matter? Effects of semantic and syntactic  
 996 inconsistencies on the allocation of attention during scene perception. *Journal of Vision*, 9(3), :24, 1-15.  
 997 <https://doi.org/10.1167/9.3.24>
- 998 Vö, M. L.-H., & Henderson, J. M. (2011). Object-scene inconsistencies do not capture gaze: evidence from the



- 999 flash-preview moving-window paradigm. *Attention, Perception & Psychophysics*, 73(6), 1742–1753.  
1000 <https://doi.org/10.3758/s13414-011-0150-6>
- 1001 Vő, M. L.-H., & Wolfe, J. M. (2013). Differential electrophysiological signatures of semantic and syntactic  
1002 scene processing. *Psychological Science*, 24(9), 1816–1823. <https://doi.org/10.1177/0956797613476955>
- 1003 Weiss, B., Knakker, B., & Vidnyánszky, Z. (2016). Visual processing during natural reading. *Scientific Reports*,  
1004 6, 26902. <https://doi.org/10.1038/srep26902>
- 1005 Winkler, I., Debener, S., Müller, K.-R., & Tangermann, M. (2015). On the influence of high-pass filtering on  
1006 ICA-based artifact reduction in EEG-ERP. *Proceedings of the Annual International Conference of the*  
1007 *IEEE Engineering in Medicine and Biology Society, EMBS, 2015-Novem*, 4101–4105.  
1008 <https://doi.org/10.1109/EMBC.2015.7319296>
- 1009 Wolfe, J. M., Alvarez, G. A., Rosenholtz, R., Kuzmova, Y. I., & Sherman, A. M. (2011). Visual search for  
1010 arbitrary objects in real scenes. *Attention, Perception, and Psychophysics*, 73(6), 1650–1671.  
1011 <https://doi.org/10.3758/s13414-011-0153-3>
- 1012 Wu, C. C., Wick, F. A., & Pomplun, M. (2014). Guidance of visual attention by semantic information in real-  
1013 world scenes. *Frontiers in Psychology*. <https://doi.org/10.3389/fpsyg.2014.00054>
- 1014 Yan, M., Richter, E. M., Shu, H., & Kliegl, R. (2009). Readers of Chinese extract semantic information from  
1015 parafoveal words. *Psychonomic Bulletin and Review*, 16(3), 561–566.  
1016 <https://doi.org/10.3758/PBR.16.3.561>
- 1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028

1029

## Tables

1030 Table 1

		Consistent	Inconsistent
		Mean $\pm$ SD	Mean $\pm$ SD
Eye movement behaviour	Ordinal fixation number of first target fixation	7.7 $\pm$ 6.0	6.1 $\pm$ 5.4
	Fixation duration ( $t-2$ ), in ms	220.7 $\pm$ 105	212.9 $\pm$ 95
	Fixation duration ( $t-1$ ), in ms	207.6 $\pm$ 96	197 $\pm$ 91
	Fixation duration ( $t$ ), in ms	261.6 $\pm$ 146	263.3 $\pm$ 136
	Gaze duration on target, in ms	408.5 $\pm$ 367.1	519.1 $\pm$ 373.6
	Number of re-fixations on target	1.7 $\pm$ 2	2.2 $\pm$ 2.1
	Duration of re-fixations on target, in ms	238.9 $\pm$ 121.8	250.2 $\pm$ 135.7
	Fixation duration ( $t+1$ ), in ms	245.3 $\pm$ 148	243.7 $\pm$ 146
	Incoming saccade amplitude to $t-1$ ( $^{\circ}$ )	6.1 $\pm$ 5.2	5.9 $\pm$ 4.8
	Incoming saccade amplitude to $t$ ( $^{\circ}$ )	8.4 $\pm$ 5.2	8.2 $\pm$ 4.8
	Incoming saccade amplitude to $t+1$ ( $^{\circ}$ )	9.5 $\pm$ 5.9	10.1 $\pm$ 5.8
	Distance of fixation $t-1$ from closest edge of target ( $^{\circ}$ )	6.7 $\pm$ 9.92	6.3 $\pm$ 9.77
	Number of fixations after first encountering target object until end of trial	7.3 $\pm$ 2.1	7.2 $\pm$ 1.7
	Duration of fixations after first encountering target object (until end of trial)	254.6 $\pm$ 120.4	251.6 $\pm$ 118.8
Target object properties	Distance of target object centre from screen centre ( $^{\circ}$ )	12.1 $\pm$ 2.8	11.7 $\pm$ 3
	Mean visual saliency (AWS model)	0.35 $\pm$ 0.16	0.37 $\pm$ 0.16
	Width ( $^{\circ}$ )	6.1 $\pm$ 2	6.1 $\pm$ 2.1
	Height ( $^{\circ}$ )	5.1 $\pm$ 1.8	5.4 $\pm$ 2.2
	Area (degrees of visual angle squared)	16.1 $\pm$ 8.7	17.3 $\pm$ 11.4

*Note.* Target object size and distance to target are based on the bounding box around the object. The fixation  $t+1$  is the first fixation after leaving the bounding box of the target object.

*Table 1.* Eye movement behaviour in the task and properties of the target object.

1031  
 1032  
 1033 Table 2  
 1034

Predictor	Cumulative probability of First Fixation <sup>1035</sup>			
	$\beta$	<i>SE</i>	<i>z</i> -value	Pr(>  <i>z</i>  ) <sup>1036</sup>
Intercept	-1.04	0.02	-49.91	0.00001 <sup>1037</sup>
Nr. Fixation	-2.00	0.05	-35.9	0.00001 <sup>1038</sup>
Consistency	0.17	0.03	5.8	0.00001 <sup>1039</sup>
Consistency × Nr. Fixation	-0.71	0.09	-7.9	0.00001 <sup>1040</sup>

1043 *Table 2.* Cumulative probability of having fixated the critical object as a function of the  
 1044 ordinal number of fixations on the scene (binomial probit). The centred predictors are  
 1045 *Consistency* (Consistent: -0.5, Inconsistent: 0.5) and *Number of Fixation*  
 1046

1047 Table 3

Predictor	Probability of immediate fixation			Latency to first fixation			Gaze Duration		
	$\beta$	$SE$	$z$	$\beta$	$SE$	$t$	$\beta$	$SE$	$t$
Intercept	-2.93	0.19	-14.73***	1,904.4	83.8	22.7***	400.1	20.97	19.08***
Consistency	0.21	0.15	1.38	-246.4	64.0	-3.85***	105.0	20.77	7.08***

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ 

1048 *Table 3.* Probability of immediate fixation and latency to first fixation. The simple coded  
 1049 predictor is *Consistency* (Consistent = -0.5, Inconsistent = 0.5). We report the  $\beta$ , the standard  
 1050 error, the  $z$ - (for binomial link) and  $t$ -value. Asterisks indicate significant predictors.

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068 Table 4

1069

Type of Event	Analysis window	$\beta$	$SE$	$t$ -value
<i>Scene onset</i>	Early (250-350 ms)	0.28	0.39	0.72
	Late (350-600 ms)	0.34	0.39	0.87
<i>nt</i>	Early (250-350 ms)	-0.06	0.07	-0.75
	Late (350-600 ms)	-0.09	0.08	-1.10
<i>t -1</i>	Early (250-350 ms)	-0.28	0.15	-1.77*
	Late (350-600 ms)	-0.46	0.12	-3.76***
<i>t</i>	Early (250-350 ms)	-0.52	0.17	-3.03**
	Late (350-600 ms)	-0.38	0.15	-2.43*
<i>t</i> (control analysis with baseline before fixation <i>t</i> )	Early (250-350 ms)	-0.34	0.16	-2.20*
	Late (350-600 ms)	-0.20	0.17	-1.10

(\*)  $p < 0.1$ , \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

1090 *Table 4.* Mixed-effects models with maximal random structure for the FRP at the mid-central  
1091 region-of-interest (comprising electrodes FC1, FCz, FC2, C1, Cz, C2, CP1, CPz, and CP2)  
1092 for two temporal windows of analysis (Early, 250-350 ms and Late, 350-600 ms) as predicted  
1093 by *Consistency* (Consistent = -0.5, Inconsistent = 0.5) of which we report the  $\beta$  the standard  
1094 error, and the  $t$ -value. Asterisks indicate the level of significance.

1095

1096

1097

1098

1099

1100

1101

1102

1103

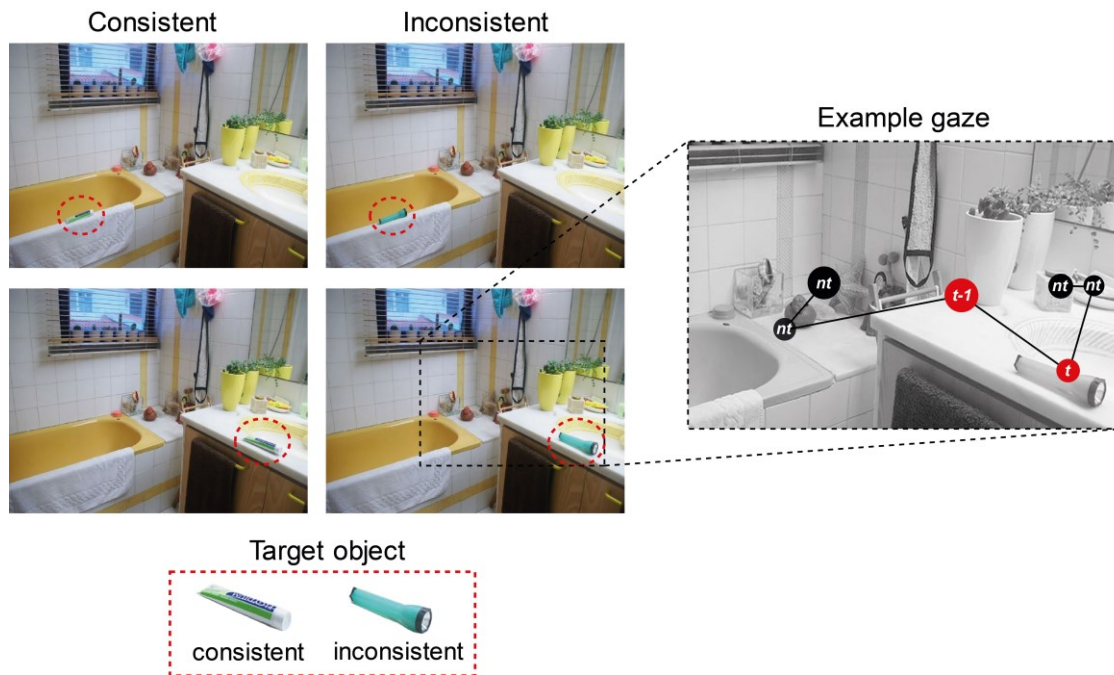
1104

1105

## Figures

1106

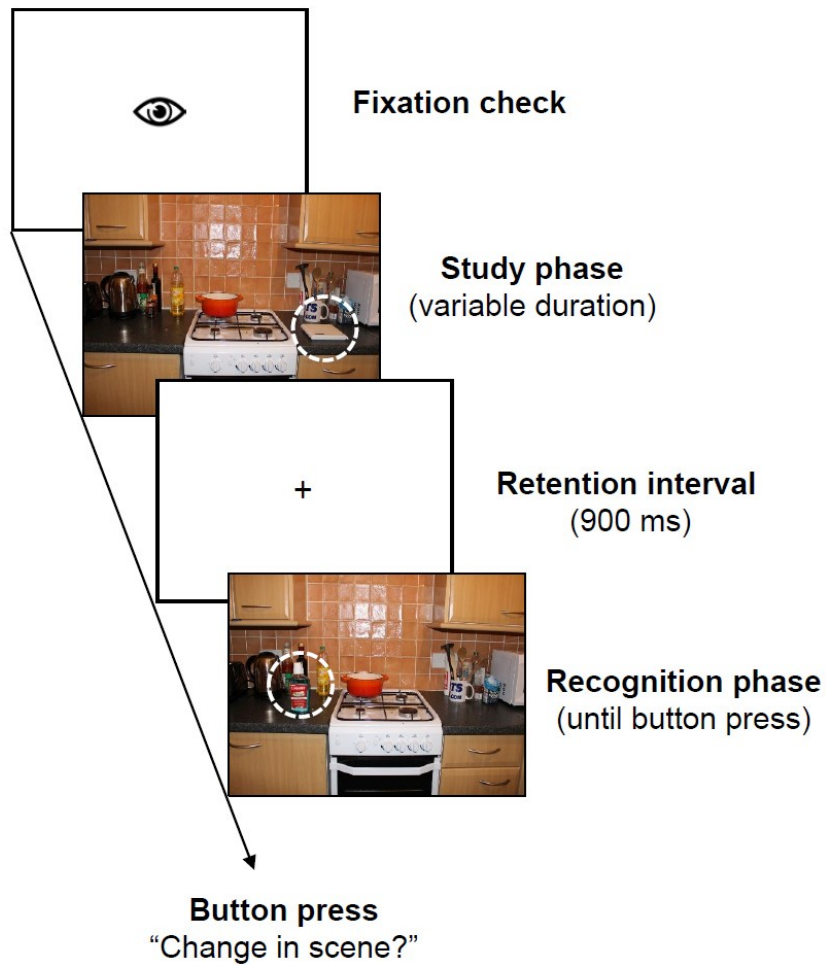
Figure 1



*Figure 1.* Example stimuli and conditions in the study. Participants viewed photographs of indoor scenes that contained a target object (highlighted with a red circle) that was either semantically consistent (here: toothpaste) or semantically inconsistent (here: flashlight) with the context of the scene. The target object could be placed at different locations within the scene, either on the left or the right side. The example gaze path plotted on the right illustrates the three type of fixations analysed in the study: (a) *t-1*; the fixation preceding the first fixation to the target object, (b) *t*; the first fixation to the target and (c) *nt*; all other (non-target) fixations. Fixation duration is proportional to the diameter of the circle, which is red for the critical fixations, and black for the non-target fixations.

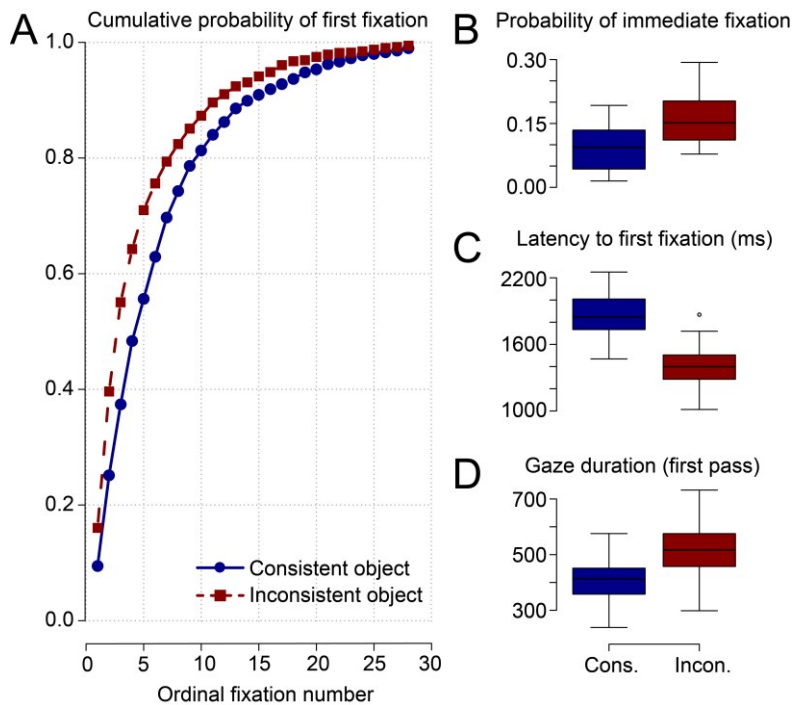
1107

Figure 2



*Figure 2.* Trial scheme. Following a drift correction, the study scene appeared. The display duration of the scene was controlled by a gaze-contingent mechanism and it disappeared on average 2000 ms after the target object was fixated. In the following retention interval, only a fixation cross was presented. During the recognition phase, the scene was presented again until participants pressed a button to indicate whether or not a change had occurred within the scene. All analyses in the present paper focus on eye-movement and EEG data collected during the study phase.

Figure 3



*Figure 3.* Eye-movement correlates of early overt attention towards consistent and inconsistent critical objects. **A.** Cumulative probability of fixating the critical object as a function of the ordinal fixation number on the scene. Blue solid line = consistent object; red dashed line = inconsistent object. **B.** Probability of fixating the critical object immediately, that is with the first fixation after scene onset. **C.** Latency until fixating the critical object for the first time. **D.** First-pass gaze duration for the critical object, i.e. the sum of all fixation durations from first entry to first exit. Whiskers of the boxplots (B, C, D) represent the 25<sup>th</sup> and 75<sup>th</sup> percentile of the measure (lower and upper quartiles). Dots indicates observations lying beyond the extremes of the whiskers.

1109

1110



Figure 4

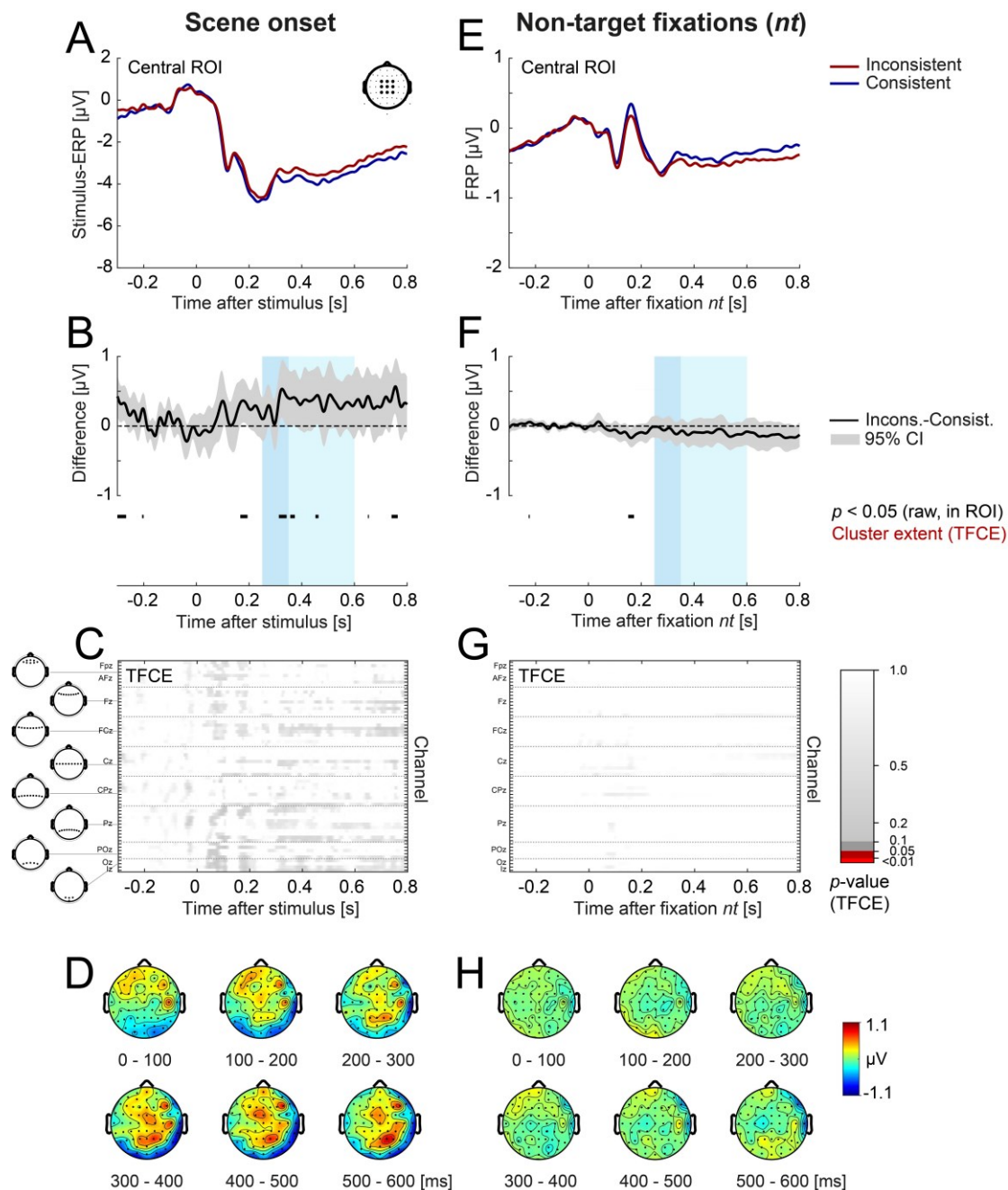
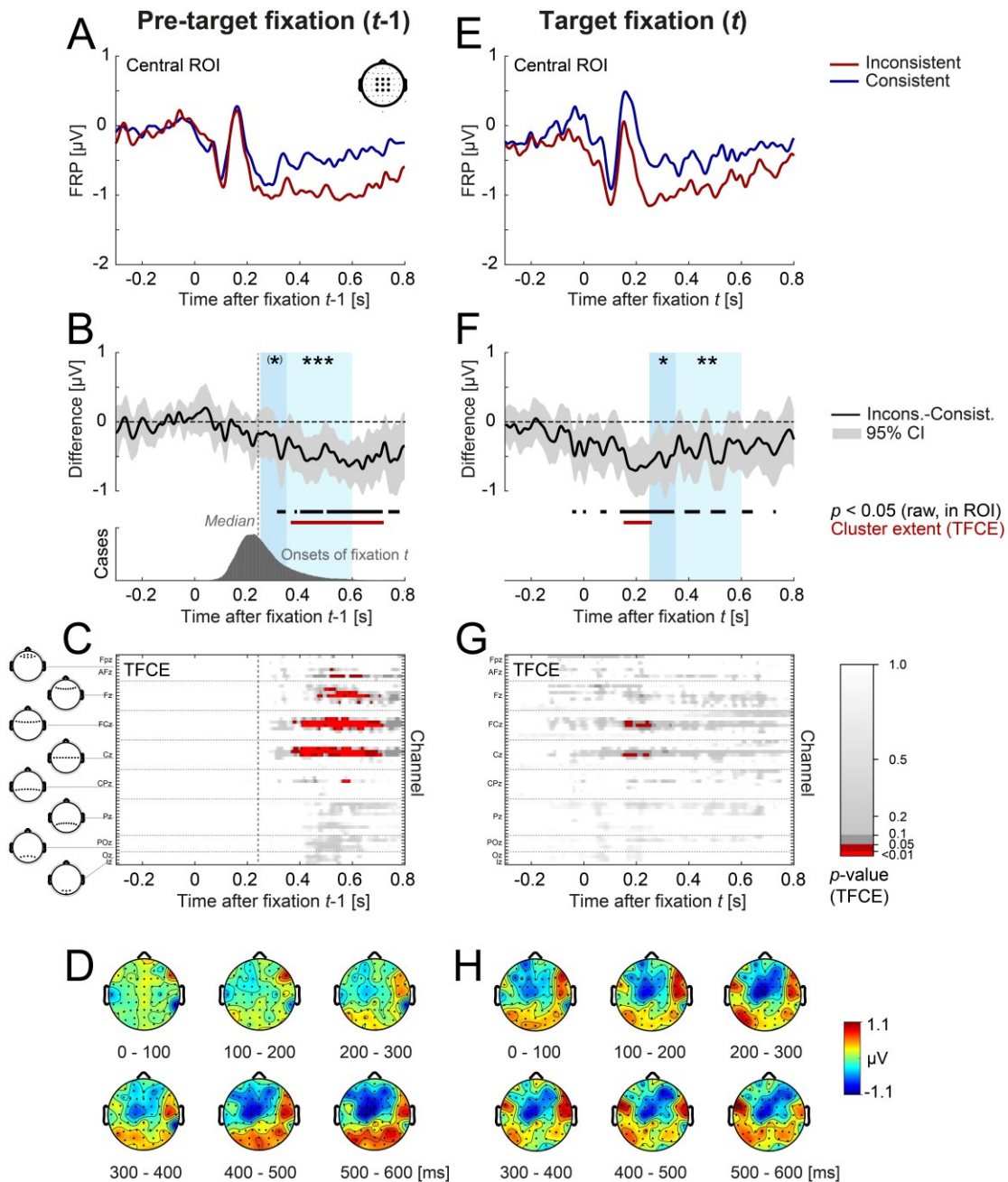


Figure 4. Stimulus-ERP aligned to scene onset (left panels) and FRP aligned to non-target fixations (right panels) as a function of object-scene consistency. **A and E.** Grand-average ERP/FRP at the central region-of-interest (comprising electrodes FC1, FCz, FC2, C1, Cz, C2, CP1, CPz, CP2). Red lines represent the *Inconsistent* condition, blue lines represent the *Consistent* condition. **B and F.** Corresponding difference waves (inconsistent minus consistent) at the central ROI. Grey shading illustrates the 95% confidence interval (without correction for multiple comparisons) of the difference wave with values outside the CI also marked in black below the curve. The two windows used for LMM statistics (250-350 and 350-600 ms) are indicated in light

blue. **C and G.** Extent of the spatiotemporal clusters underlying the cluster-based permutation statistic (TFCE) computed across all electrodes/time points. There were no significant ( $p < 0.05$ ) effects. **D and H.** Scalp topographies of the consistency effect (inconsistent minus consistent) averaged across successive 100 ms time windows. Object-scene consistency had no significant effects on the stimulus-ERP or on the FRP elicited by non-target fixations, neither in the LMM statistic, nor in the cluster permutation test.

1111

Figure 5



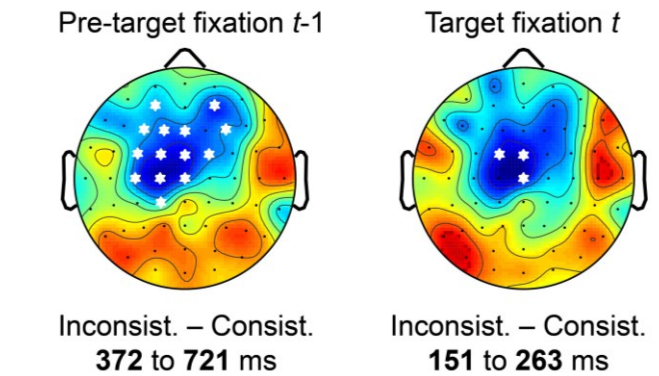
*Figure 5.* Grand-average FRP elicited by pre-target fixation (left panels) and target fixation (right panels) as a function of object-scene consistency. **(A, E)** Grand-average FRPs at the central ROI. **(B, F)** Difference waves at the central ROI. In panel B, the grey distribution shows the onset of fixation  $t$  relative to the onset of the pre-target fixation  $t-1$ , with the vertical dotted line indicating the mean latency (260 ms). **(C, G)** Results of cluster-based permutation testing (TFCE). The extent of the clusters from the first stage of the permutation test (marked in red) provides some indication which spatiotemporal features of the waveforms likely contributed to the overall significant effect of consistency. The temporal extent of the clusters is also illustrated by the red

bars in panels B and F. **(D, H)** Scalp topographies of the consistency effect (inconsistent minus consistent) across successive 100 ms time windows. A frontocentral N400 effect emerged in the FRP time-locked to fixation  $t-1$  and reached significance shortly after the eyes had moved on to fixation  $t$ . This effect then continued during fixation  $t$  reaching a maximum 200 ms after the start of the target fixation.

1112

1113

Figure 6



*Figure 6.* Scalp distribution of frontocentral N400 effects in the time windows significant in the TFCE statistic (see also Figure 5). White asterisks highlight the spatial extent of the clusters observed in the first stage of the TFCE permutation test for both intervals. In the FRP aligned to the pre-target fixation (left), clusters extended from 372 to 721 ms and across 13 frontocentral channels. In the FRP aligned to the target fixation (right), clusters extended from 151 and 263 ms at three frontocentral channels.

1114  
 1115