

1 **Semantic scene-object consistency modulates N300/400 EEG**
2 **components, but does not automatically facilitate object**
3 **representations**

4 Lixiang Chen¹, Radoslaw Martin Cichy^{1,*}, Daniel Kaiser^{2,*}

5 ¹ Department of Education and Psychology, Freie Universität Berlin, Berlin, Germany

6 ² Department of Psychology, University of York, York, UK

7 * R.M.C. and D.K. contributed equally to this study

8 **Correspondence:**

9 Lixiang Chen

10 Department of Education and Psychology

11 Freie Universität Berlin

12 Berlin, Germany

13 Email: lixiang.chen@fu-berlin.de

14 **Abstract**

15 During natural vision, objects rarely appear in isolation, but often within a semantically
16 related scene context. Previous studies reported that semantic consistency between
17 objects and scenes facilitates object perception, and that scene-object consistency is
18 reflected in changes in the N300 and N400 components in EEG recordings. Here, we
19 investigate whether these N300/N400 differences are indicative of changes in the
20 cortical representation of objects. In two experiments, we recorded EEG signals while
21 participants viewed semantically consistent or inconsistent objects within a scene; in
22 Experiment 1, these objects were task-irrelevant, while in Experiment 2, they were
23 directly relevant for behavior. In both experiments, we found reliable and comparable
24 N300/400 differences between consistent and inconsistent scene-object combinations.
25 To probe the quality of object representations, we performed multivariate classification
26 analyses, in which we decoded the category of the objects contained in the scene. In
27 Experiment 1, in which the objects were not task-relevant, object category could be
28 decoded from around 100 ms after the object presentation, but no difference in decoding
29 performance was found between consistent and inconsistent objects. By contrast, when
30 the objects were task-relevant in Experiment 2, we found enhanced decoding of
31 semantically consistent, compared to semantically inconsistent, objects. These results
32 show that differences in N300/N400 components related to scene-object consistency do
33 not index changes in cortical object representations, but rather reflect a generic marker
34 of semantic violations. Further, our findings suggest that facilitatory effects between
35 objects and scenes are task-dependent rather than automatic.

36 **Key words**

37 scene-object consistency, N300, object representation, multivariate pattern analysis

38 **Introduction**

39 In the real world, objects rarely appear in isolation, but practically always within a
40 particular scene context (Bar, 2004; Wolfe et al., 2011; Kaiser et al., 2019; Võ et al.,
41 2019). Objects are often semantically related to the scene they appear in: For instance,
42 microwaves usually appear in the kitchen, but practically never in the bathroom.
43 Several behavioral studies have shown that such semantic relations between objects and
44 scenes affect object identification. Early studies using line drawings of scenes and
45 objects found that objects were detected faster and more accurately when they were in
46 a consistent setting than in an inconsistent setting (Palmer, 1975; Biederman et al.,
47 1982; Boyce et al., 1989; Boyce and Pollatsek, 1992). Similar results were recently
48 reported for scene photographs (Davenport and Potter, 2004; Davenport, 2007;
49 Munneke et al., 2013). In line with such findings, eye-tracking studies have shown that
50 inconsistent objects are fixated longer and more often than consistent objects (Võ and
51 Henderson, 2009, 2011; Cornelissen and Võ, 2017), suggesting that objects are
52 perceived more swiftly within a consistent than within an inconsistent scene.
53 Interestingly, such behavioral facilitation effects are also observed when, instead of the
54 object, the scene is task-relevant: Davenport et al. (2004; 2007) reported that scenes
55 were identified more accurately if they contained a consistent foreground object
56 compared to an inconsistent one. These effects suggest that objects and scenes are
57 processed in a highly interactive manner.

58 To characterize the neural basis of these semantic consistency effects, EEG studies have
59 used paradigms in which objects appear within consistent or inconsistent scenes, either
60 simultaneously or sequentially (Ganis and Kutas, 2003; Mudrik et al., 2010; Võ and
61 Wolfe, 2013; Draschkow et al., 2018; Coco et al., 2020). For example, Võ et al. (2013)
62 adopted a sequential design, in which participants first viewed a scene image, followed
63 by a location cue, where then appeared a consistent (e.g., a computer mouse on an office
64 table) or an inconsistent object (e.g., a soap on an office table). They found objects in
65 an inconsistent scene evoked more negative responses than consistent objects in the
66 N300 (around 250-350 ms) and N400 (around 350-600 ms) windows. Several other
67 studies (Mudrik et al., 2010, 2014; Truman and Mudrik, 2018) using a simultaneous
68 design, in which the scene and object were presented simultaneously, reported similar
69 N300 and/or N400 modulations. Critically, the earlier N300 effects are often considered

70 to reflect differences in perceptual processing between typically and atypically
71 positioned objects (Schendan and Maher, 2009; Mudrik et al., 2010; Kumar et al.,
72 2021). On this view, consistency-related differences in EEG waveforms arise as a
73 consequence of differences in the visual analysis of objects and scenes, rather than due
74 to a post-perceptual signaling of (in)consistency.

75 If differences in the N300 waveform indeed index changes in perceptual processing,
76 the N300 ERP effect should be accompanied by differences in the neural representation
77 of the objects. In this study, we put this prediction to the test. Across two experiments,
78 we compared differences in the N300/400 EEG components to multivariate decoding
79 of objects contained in consistent and inconsistent scenes. In both experiments,
80 participants completed a sequential semantic consistency paradigm, in which scenes
81 from 8 different categories were consistently or inconsistently combined with objects
82 from 16 categories. We then examined the influence of scene-object consistency on
83 EEG signals, both when the objects were task-irrelevant (Experiment 1) and when
84 participants performed a recognition task on the objects (Experiment 2). In both
85 experiments, we replicated previously reported ERP effects, with greater N300 and
86 N400 components for inconsistent scene-object combinations, compared to consistent
87 combinations. To probe the quality of object and scene representations, we performed
88 multivariate classification analyses, in which we decoded between the object and scene
89 categories separately for each condition. In Experiment 1, in which the objects were not
90 task-relevant, object category could be decoded from around 100 ms after the object
91 presentation, but no difference in decoding performance was found between consistent
92 and inconsistent objects. In Experiment 2, in which the objects were directly task-
93 relevant, we found enhanced decoding of semantically consistent, compared to
94 semantically inconsistent, objects. In both experiments, we found no differences in
95 scene category decoding between semantically consistent and inconsistent conditions.
96 Together, these results show that differences in N300/N400 components related to
97 scene-object consistency do not necessarily index changes in cortical object
98 representations, but rather reflect a generic marker of semantic violations. Further, they
99 suggest that facilitation effects between objects and scenes are task-dependent rather
100 than automatic.

101 **Materials and Methods**

102 All materials and methods were identical for the two experiments, unless stated
103 otherwise.

104 **Participants** Thirty-two participants (16 males, mean age 26.23 yrs, SD = 2.05 yrs),
105 with normal or corrected-to-normal vision, took part in Experiment 1. Another thirty-
106 two participants (14 males, mean age 26.97 yrs, SD = 1.67 yrs) took part in Experiment
107 2. Participants were paid volunteers or participated for partial course credits. All
108 participants provided written, informed consent prior to participating in the experiment.
109 The experiments were approved by the ethical committee of the Department of
110 Education and Psychology at Freie Universität Berlin and were conducted in
111 accordance with the Declaration of Helsinki.

112 **Stimuli** The stimulus set comprised scene images from 8 categories: beach, bathroom,
113 office, kitchen, gym, street, supermarket, and prairie. The scenes were grouped into 4
114 pairs (beach & bathroom, office & kitchen, gym & street, supermarket & prairie). We
115 chose four objects for each scene pair, two of which were semantically consistent with
116 one scene and two of which were semantically consistent with the other scene. To create
117 semantically inconsistent scene-object combinations, we simply exchanged the objects
118 between the scenes. All combinations of scenes and objects can be found in Table 1.
119 For example, consider the office and kitchen pair: We first chose a computer and a
120 printer as consistent objects for the office, and chose a rice cooker and a microwave as
121 consistent objects for the kitchen. We in turn chose the rice cooker and microwave as
122 inconsistent objects for the office, and the computer and printer as inconsistent objects
123 for the kitchen. We pasted the objects into the scene images using Adobe Photoshop.
124 The object locations were the same across the consistent and inconsistent objects, and
125 they were always in line with the typical position of the consistent object (e.g., a
126 computer was positioned on an office desk in the same way as a rice cooker). We used
127 3 exemplars for each scene category and 3 exemplars for each object, yielding 288
128 unique stimuli. During the experiments, the scenes could also be shown without objects
129 (see below). Fig. 1A shows some examples of the stimuli.

130 **Table 1.** Combinations of scenes and objects for the consistent and inconsistent
131 conditions. Note that scenes were grouped into pairs; for each pair, four objects were
132 used as consistent and inconsistent objects.

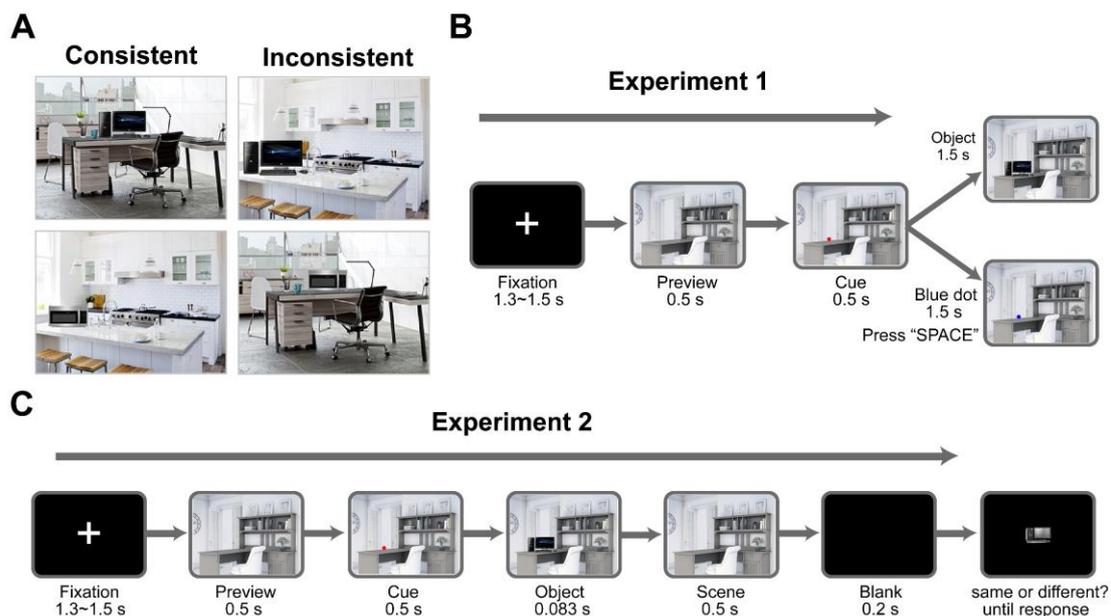
Scenes	Consistent Objects	Inconsistent Objects
bathroom	toilet, washing machine	parasol, deck chair
beach	parasol, deck chair	toilet, washing machine
office	computer, printer	microwave, rice cooker
kitchen	microwave, rice cooker	computer, printer
gym	treadmill, spinning bike	scooter, bus stop sign
road	scooter, bus stop sign	treadmill, spinning bike
supermarket	shopping cart, shop assistant	ostrich, zebra
prairie	ostrich, zebra	shopping cart, shop assistant

133 **Paradigm** Participants were seated in a quiet room. The images were presented on a
134 monitor with a resolution of 1024×768 pixels and a refresh rate of 60 Hz. We adopted
135 a sequential scene-object congruity paradigm similar to Vö et al. (2013). Image
136 presentation and recording of subjects' behavioral responses were controlled using
137 MATLAB and the Psychophysics Toolbox (Brainard, 1997). Each trial began with a
138 fixation cross "+" shown for a random interval between 1.3-1.5 s, after which a scene
139 image (without the critical object) was presented for 500 ms. Next, a red dot cue was
140 presented at a single location in the scene for 500 ms, indicating where the critical object
141 would appear. Participants were instructed to move their eyes to the dot as quickly as
142 possible. To avoid eye movement artifacts during the subsequent object presentation,
143 we told the participants not to move their eyes away from the dot location until the
144 beginning of the next trial. After that, a semantically consistent or inconsistent object
145 appeared at the location of the red dot.

146 In Experiment 1, the objects were not task-relevant. The object simply remained visible
147 together with the scene for 1,500 ms before the next trial started. Participants performed
148 an unrelated attention control task: To ensure that they attended the cued location, we
149 added task trials (10% of trials) during which no object appeared after the cue. Instead,
150 the color of the dot changed from red to blue. Participants were instructed to press the
151 spacebar when they detected the change (detection rate in target trials: 97.8%, SE =
152 0.63%). An example trial for Experiment 1 is shown in Fig. 1B.

153 In Experiment 2, the objects were directly task-relevant. The consistent or inconsistent
154 object appeared only very briefly (83 ms) at the location of the red dot. The scene image
155 remained on the screen for another 500 ms after the object disappeared. After a 200 ms
156 blank interval, participants were asked to perform an object recognition test. During the
157 test, an object was shown on the screen, which was either the same object exemplar
158 they had just seen or a different exemplar of the same category. Test objects were
159 presented in grayscale to increase task difficulty. Participants were asked to determine
160 whether this object exemplar was the one that appeared earlier in the scene. If it was,
161 the participants should press G button, otherwise press H button. The next trial began
162 as soon as participants made a choice. An example trial for Experiment 2 is shown in
163 Fig. 1C.

164 Both experiments included three runs and all 288 unique stimulus images were
165 presented once in random order within each run. Across runs, there were 27 repetitions
166 for each specific scene-object combination.



167

168 **Fig. 1** Experimental design. (A) Examples of consistent and inconsistent scene-object combinations. (B)
169 Trial sequence in Experiment 1. After a fixation interval, a scene without the critical object was
170 presented. Next, a red dot cue was presented, and participants were asked to move their eyes to this
171 location. After that, the critical object appeared at the cued location in the scene. On target trials, the red
172 cue turned blue, and participants were instructed to press spacebar. (C) Trial sequence in Experiment 2.
173 Here, the objects were displayed briefly at the location of the cue. In a subsequent recognition test, an
174 object of the same category appeared on the screen. Participants were asked to determine whether this
175 object exemplar was the one that appeared earlier in the scene.

176 **EEG Recording and Preprocessing**

177 EEG signals were recorded using an EASYCAP 64-electrode system and a Brainvision
178 actiCHamp amplifier in both Experiments. Electrodes were arranged according to the
179 10-10 system. EEG data were recorded with a sample rate of 1000 Hz and filtered online
180 between 0.03 Hz and 100 Hz. All electrodes were referenced online to the Fz electrode
181 and re-referenced offline to the average of data from all channels. Offline data
182 preprocessing was performed using FieldTrip (Oostenveld et al., 2011). EEG data were
183 segmented into epochs from -100 ms to 800 ms relative to the onset of the critical object
184 and baseline corrected by subtracting the mean signal prior to the object onset. To track
185 the temporal representations of scenes, EEG data were segmented into epochs from -
186 1100 ms to 800 ms relative to the onset of the object and baseline corrected by
187 subtracting the mean signal prior to the scene presentation (-100-0 ms relative to the
188 scene onset). Channels and trials containing excessive noise were removed by visual
189 inspection. Blinks and eye movement artifacts were removed using independent
190 component analysis and visual inspection of the resulting components. The epoched
191 data were downsampled to 200 Hz.

192 **ERP Analyses**

193 To replicate semantic consistency ERP effect reported in previous scene studies (e.g.,
194 Mudrik et al., 2014; V̇o et al., 2013), we performed ERP analyses using FieldTrip. For
195 these analyses, the preprocessed data were additionally band-pass filtered at 0.1-30 Hz.
196 In accordance with V̇o et al. (2013), we chose 9 electrodes (FC1, FCz, FC2, C1, Cz,
197 C2, CP1, CPz, and CP2) located in the mid-central region for further ERP analysis. We
198 first averaged the evoked responses across these electrodes and then averaged these
199 mean responses separately for the consistent and inconsistent conditions and each
200 participant.

201 **Decoding Analyses**

202 We performed two complementary multivariate decoding analyses to track temporal
203 representations of objects and scenes across time. First, to track representations of
204 objects and investigate how consistent or inconsistent scene contexts affect objects
205 processing, we performed decoding analyses between two consistent and inconsistent

206 objects separately within each scene at each time point from -100 ms to 800 ms relative
207 to the onset of the object. For example, we performed classification analyses to either
208 differentiate printers (consistent) from computers (consistent) in office scenes, or to
209 differentiate printers (inconsistent) from computers (inconsistent) in kitchen scenes, at
210 each time point. Second, to track the impact of consistent or inconsistent objects on
211 scenes representations, we performed decoding analyses to discriminate between the
212 eight scene categories separately for consistent and inconsistent conditions at each time
213 point from -100 ms to 1800 ms relative to the onset of the scene (-1100 ms to 800 ms
214 relative to the onset of the object). For all decoding analyses, we adopted two
215 approaches: standard timeseries decoding (Boring et al., 2020; Kaiser and Nyga, 2020),
216 using data from a sliding time window, and cumulative decoding (Ramkumar et al.,
217 2013; Kaiser et al., 2020a), using aggregated data from all elapsed time points. The two
218 approaches are detailed in the following paragraphs.

219 ***Timeseries decoding*** Timeseries decoding analyses were performed using Matlab and
220 CoSMoMVPA (Oosterhof et al., 2016). To increase the power of our timeseries
221 decoding, the analysis was performed on a sliding time window (50 ms width), with a
222 5 ms resolution. This approach thus not only utilizes data from current time point, but
223 the data from 5 time points before and after the current time point.

224 Considering excessive data dimensionality may harm classification, we adopted
225 principal component analysis (PCA) to reduce the dimensionality of the data
226 (Grootswagers et al., 2017; Kaiser et al., 2020a; Kaiser and Nyga, 2020). For each
227 classification, a PCA was performed on all data from the training set, and the PCA
228 solution was projected onto data from the testing set. For each PCA, we retained the set
229 of components explaining 99% of the variance in the training set data.

230 The classification was performed separately for each time point from -100 to 800 ms
231 (from -1100 to 800 ms for scene decoding), using LDA classifiers with 10-fold cross-
232 validation. Specifically, the EEG data from all epochs were first allocated to 10 folds
233 randomly. LDA classifiers were then trained on data from 9 folds and then tested on
234 data from the left-out fold. The amount of data in the training set was always balanced
235 across conditions. For each object decoding analysis, the training set included ~48
236 trials, and the testing set included ~6 trials; for each scene decoding analysis, the
237 training set included ~384 trials and the testing set included ~32 trials. The

238 classification was done repeatedly until every fold was left out once. For each time
239 point, the accuracies were averaged across the 10 repetitions.

240 **Cumulative decoding** We also performed cumulative decoding analyses, which takes
241 into account the data of all time points before the current time point in the epoch for
242 classifications (Ramkumar et al., 2013; Kaiser et al., 2020a). For example, for the first
243 time point in the epoch, the classifier was trained and tested on response patterns at this
244 time point in the epoch; at the second time point in the epoch, the classifier was trained
245 and tested on response patterns at the first and second time point in the epoch; and at
246 the last time point in the epoch, the classifier was trained and tested on response patterns
247 at all time points in the epoch. This decoding approach uses larger amounts of data that
248 span multiple time points, so that it may provide additional sensitivity for detecting
249 effects that are transported by variations in the time domain. As for the timeseries
250 decoding, LDA classifiers with 10-fold cross-validation were used for classifications
251 and PCA was adopted to reduce the dimensionality of the data for each classification
252 step across time.

253 **Statistics**

254 For Experiment 2, we used paired *t*-tests to compare participants' accuracy and
255 response times when they were asked to recognize consistent and inconsistent objects.

256 For ERP analyses, we used paired *t*-tests to compare the averaged EEG responses
257 evoked by consistent and inconsistent scene-object combinations, at each time point.

258 For decoding analyses, we used one-sample *t*-tests to compare decoding accuracies
259 against chance level and paired *t*-tests to compare decoding accuracies between the
260 consistent and inconsistent conditions, at each time point.

261 Multiple-comparison corrections were performed using FDR ($p < 0.05$), and only
262 clusters of at least 5 consecutive significant time points (i.e., 25 ms) were considered.

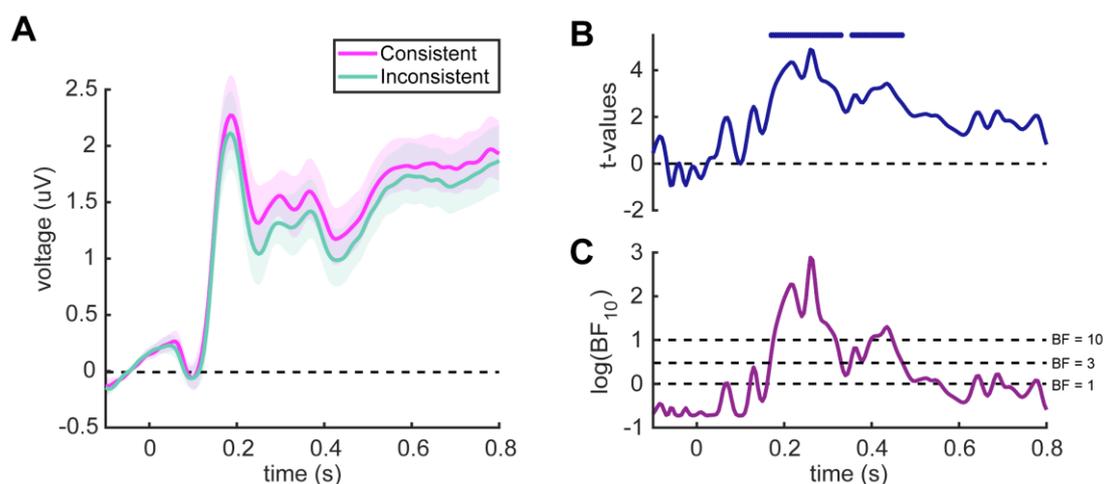
263 We also calculated Bayes factors (Rouder et al., 2009) for all analyses.

264 Results

265 Experiment 1

266 ERP signals indexing scene-object consistency

267 To track the influence of scene-object consistency on EEG responses, we first analyzed
268 EEG waveforms in mid-central electrodes. In this analysis, we found more negative
269 responses evoked by inconsistent scene-object combinations than consistent
270 combinations, which emerged at 170-330 ms (peak: $t = 4.884$, $BF_{10} = 765.26$) and 355-
271 470 ms (peak: $t = 3.429$, $BF_{10} = 20.06$) (Fig. 2). These results demonstrate larger N300
272 and N400 components evoked by inconsistent scenes, which is in line with previous
273 findings (Mudrik et al., 2010, 2014; Vö and Wolfe, 2013).



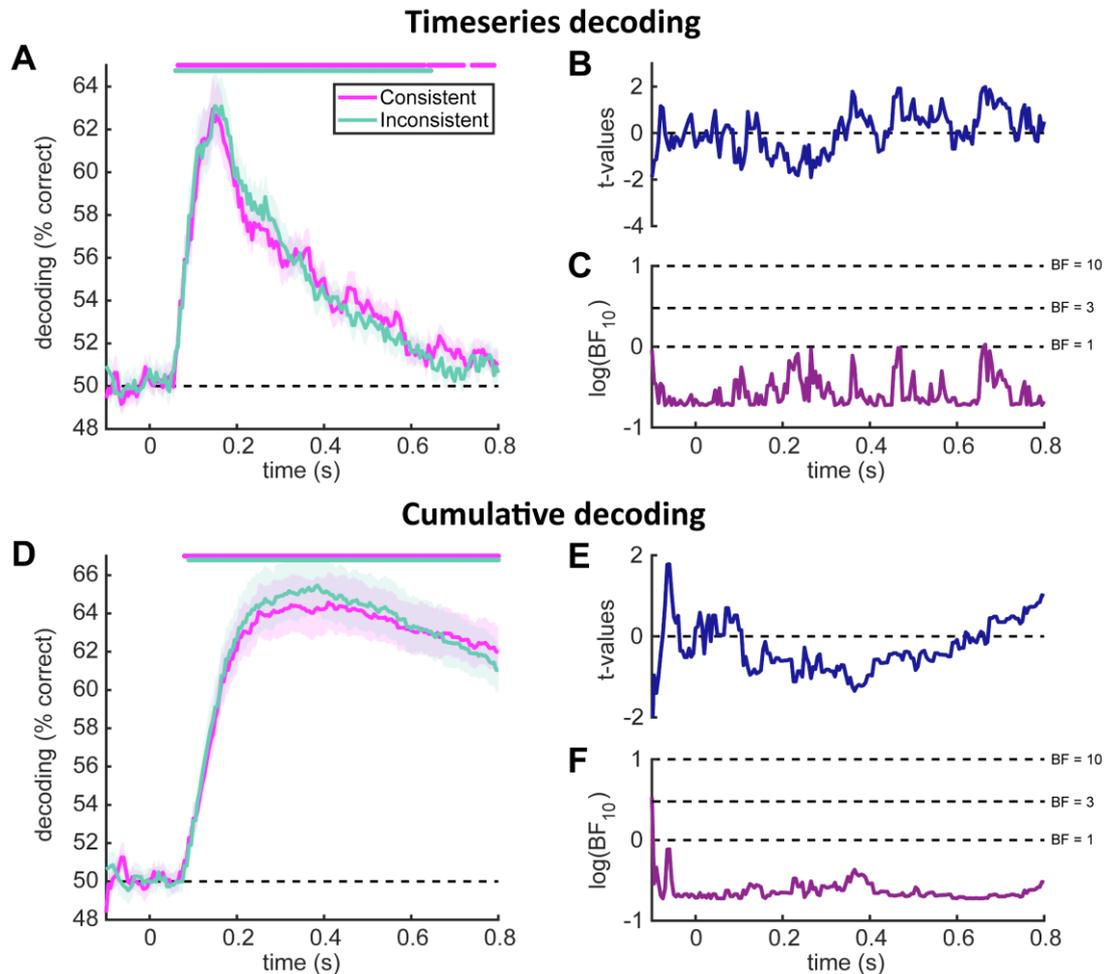
274

275 **Fig. 2** Event-related potentials (ERPs) in Experiment 1. (A) ERPs recorded from the mid-central region
276 for consistent and inconsistent scene-object combinations. Error margins represent standard errors. (B)
277 t -values for the comparisons between consistent and inconsistent conditions. Line markers denote
278 significant differences between conditions ($p < 0.05$, FDR-corrected). (C) Bayes factors (BF_{10}) for the
279 comparisons between consistent and inconsistent conditions. For display purposes, the BF_{10} values were
280 log-transformed. Dotted lines show low ($BF_{10} = 1$), moderate ($BF_{10} = 3$), and high ($BF_{10} = 10$) evidence
281 for a difference between conditions. In line with previous reports, these results show that scene-object
282 consistency is represented in evoked responses 170-330 ms and 355-470 ms after the object onset.

283 Tracking object representations in consistent and inconsistent scenes

284 Having established reliable ERP differences between consistent and inconsistent scene-
285 object combinations, we were next interested if these differences were accompanied by
286 differences in how well the consistent and inconsistent objects were represented. We
287 performed timeseries and cumulative decoding analyses between two consistent or

288 inconsistent objects separately within each scene at each time point from -100 to 800
289 ms relative to the onset of the object. In both analyses, we found highly similar decoding
290 performances for both consistent and inconsistent objects. Specifically, there was
291 significant decoding between consistent objects, which emerged at 65-790 ms in the
292 timeseries decoding (Fig. 3A), and between 80 and 800 ms in the cumulative decoding
293 (Fig. 3D), and there was significant decoding between inconsistent objects in both the
294 timeseries decoding (60-645 ms; Fig. 3A) and cumulative decoding (90-800 ms; Fig.
295 3D). No significant differences in decoding accuracy were found between consistent
296 and inconsistent objects. Hence, despite the reliable ERP differences between
297 consistent and inconsistent scene-object stimuli, there was no evidence for an automatic
298 facilitation from the scene to the semantically consistent object.



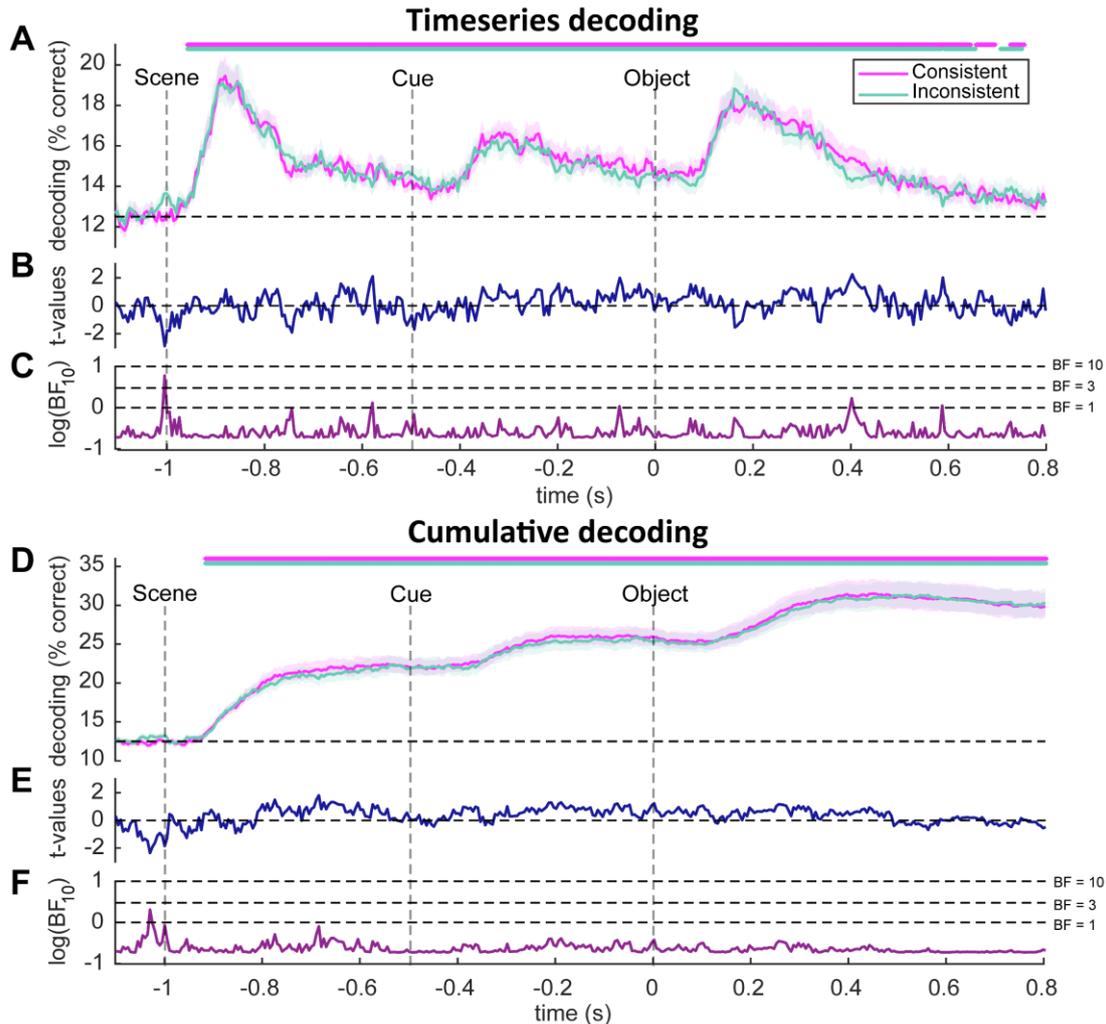
299

300 **Fig. 3** Decoding for the consistent or inconsistent objects within each scene in Experiment 1. (A)
301 Timeseries decoding results, separately for consistent and inconsistent objects. Line markers denote
302 significant above-chance decoding ($p < 0.05$, FDR-corrected). (B) t -values for the comparisons between
303 consistent and inconsistent conditions. (C) Bayes factors (BF₁₀) for the comparisons between consistent
304 and inconsistent conditions. For display purposes, the BF₁₀ values were log-transformed. Dotted lines
305 show low (BF₁₀ = 1), moderate (BF₁₀ = 3), and high (BF₁₀ = 10) evidence for a difference between
306 conditions. (D) Cumulative decoding results, separately for consistent and inconsistent objects. Line

307 markers denote significant above-chance decoding ($p < 0.05$, FDR-corrected). (**E, F**) t -values and Bayes
308 factors (BF_{10}) for the comparisons between consistent and inconsistent conditions, as in (B, C). These
309 results show robust decoding for consistently and inconsistently placed objects. However, despite the
310 reliable differences between consistent and inconsistent scene-object combinations in ERP signals, object
311 decoding was highly similar between the consistent and inconsistent conditions.

312 *Tracking the representation of scenes with consistent and inconsistent objects*

313 Although scene-object consistency does not automatically facilitate the representation
314 of the objects, there may still be an opposite cross-facilitation effect where the
315 consistent object enhances scene representations. To test this possibility, we performed
316 decoding analyses to discriminate between the eight scene categories separately for the
317 consistent and inconsistent conditions from -100 ms to 1800 ms relative to the onset of
318 the scene. We found significant decoding between scenes with a consistent object in
319 both the timeseries decoding (50-1645 ms, 660-695 ms, and 1765-1800 ms) and
320 cumulative decoding (85-1800 ms) analyses. Significant decoding between scenes that
321 contained inconsistent objects was also found in both the timeseries decoding (50-1585
322 ms, 1595-1655 ms and 1710-1750 ms) and cumulative decoding (85-1800 ms). These
323 results are consistent with previous findings (Lowe et al., 2018; Kaiser et al., 2020b),
324 which suggest scene category can be decoded within 100 ms (Fig. 4). However, no
325 significant differences were found between these scenes with consistent and
326 inconsistent objects. These results suggest that scene category can be decoded in a
327 temporally sustained way, but semantically consistent objects have no facilitatory effect
328 on scene representations.



329
330 **Fig. 4** Decoding between scenes with consistent or inconsistent objects in Experiment 1. (A) Timeseries
331 decoding results, separately for scene with consistent and inconsistent objects. Line markers denote
332 significant above-chance decoding ($p < 0.05$, FDR-corrected). (B) t -values for the comparisons between
333 consistent and inconsistent conditions. (C) Bayes factors (BF_{10}) for the comparisons between consistent
334 and inconsistent conditions. For display purposes, the BF_{10} values were log-transformed. Dotted lines
335 show low ($BF_{10} = 1$), moderate ($BF_{10} = 3$), and high ($BF_{10} = 10$) evidence for a difference between
336 conditions. (D) Cumulative decoding results, separately for scene with consistent and inconsistent
337 objects. Line markers denote significant above-chance decoding ($p < 0.05$, FDR-corrected). (E, F) t -
338 values and Bayes factors (BF_{10}) for the comparisons between consistent and inconsistent conditions, as
339 in (B, C). These results show that scene category can be well decoded across time, but the consistency
340 of embedded objects has no facilitatory effects on scene representations.

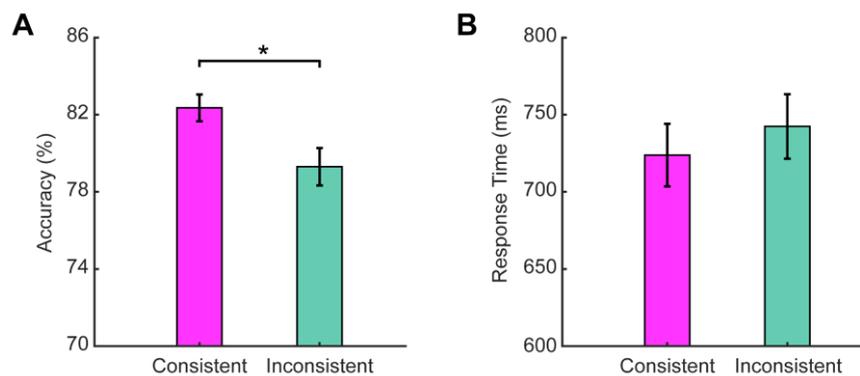
341 Experiment 2

342 In Experiment 1, we did not find differences in object and scene representations
343 between the consistent and inconsistent object-scene combinations, despite robust ERP
344 differences between the two conditions. However, the objects and scenes were both not
345 task-relevant in Experiment 1 – although participants spatially attended the object
346 location, the objects' features were not important for solving the task. Under such
347 conditions, object representations may not benefit from semantically consistent context

348 to the same extent as when object features are critical for solving the task. In Experiment
349 2, we therefore made the objects task-relevant.

350 *Behavioral object recognition in semantically consistent and inconsistent scenes*

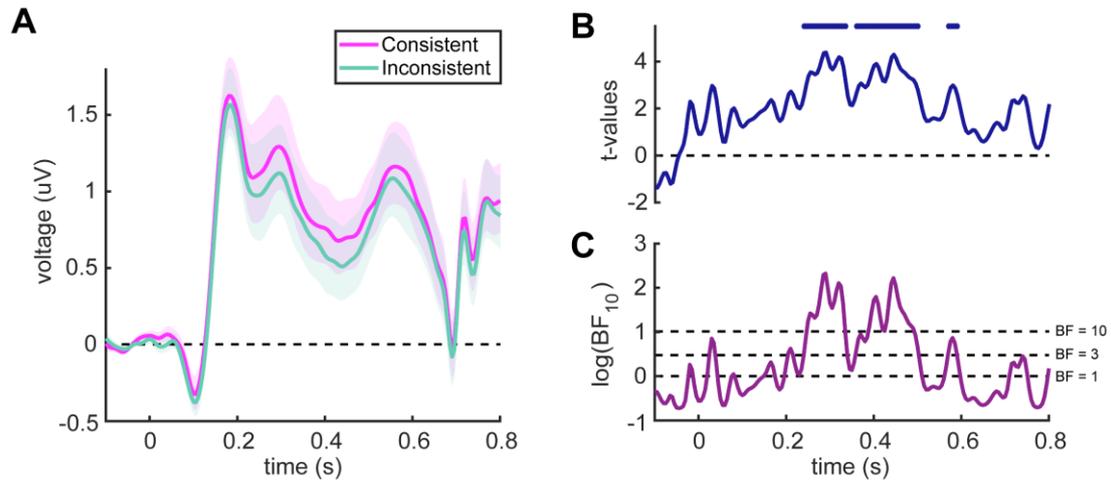
351 In Experiment 2, participants performed a recognition task, in which they were asked
352 to report whether a test object was identical to the one they had previously seen in the
353 scene (Fig. 1C). In line with previous findings (Davenport and Potter, 2004; Davenport,
354 2007; Munneke et al., 2013), we found that objects were recognized more accurately
355 when they were embedded in consistent scenes than in inconsistent scenes (mean
356 accuracy: consistent = 82.36%, inconsistent = 79.30%; $t = 2.598$, $p = 0.011$; Fig. 5A).
357 These results suggest that semantically consistent scenes can enhance the recognition
358 of objects. There was no difference in response times between two conditions (mean
359 response time: consistent = 723.8 ms, inconsistent = 742.4 ms; $t = -0.648$, $p = 0.519$;
360 Fig. 5B).



361
362 **Fig. 5** Accuracy (A) and response time (B) in consistent and inconsistent conditions in Experiment 2.
363 Consistent objects were recognized more accurately than inconsistent object, but there was no significant
364 difference in response time between two conditions. Error bars represent standard error of the mean. *
365 represents $p < 0.05$.

366 *ERP signals indexing scene-object consistency*

367 Inconsistent scene-object combinations evoked more negative responses in mid-central
368 electrodes than consistent combinations at 240-335 ms (peak: $t = 4.385$, $BF_{10} = 210.85$),
369 360-500 ms (peak: $t = 4.291$, $BF_{10} = 165.94$), and 570-590 ms (peak: $t = 2.986$, $BF_{10} =$
370 7.36) (Fig. 6). The results suggest larger N300 and N400 components evoked by
371 semantically inconsistent scene-object combinations, replicating the findings from
372 Experiment 1.

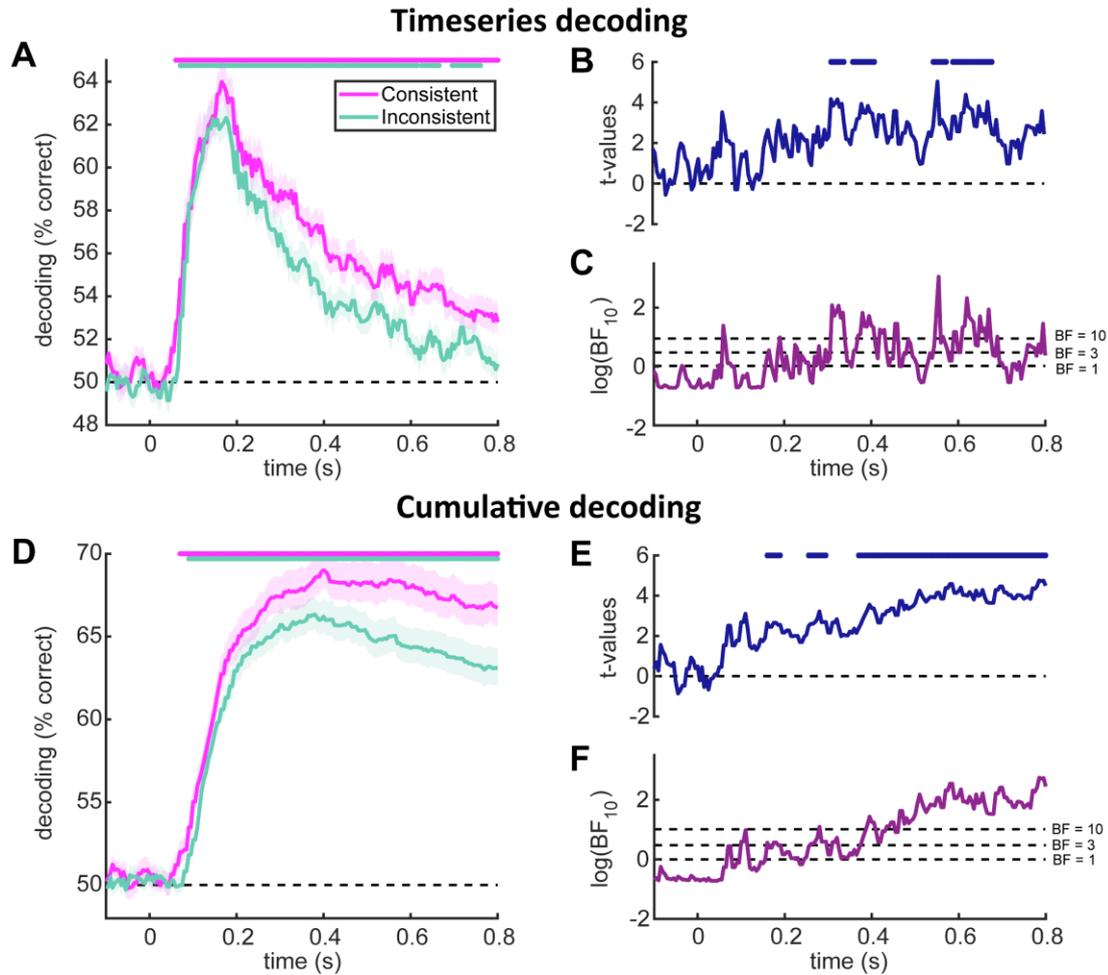


373

374 **Fig. 6** Event-related potentials (ERPs) in Experiment 2. **(A)** ERPs recorded from the mid-central region
375 for consistent and inconsistent scene-object combinations. Error margins represent standard errors. **(B)**
376 *t*-values for the comparisons between consistent and inconsistent conditions. Line markers denote
377 significant differences between conditions ($p < 0.05$, FDR-corrected). **(C)** Bayes factors (BF₁₀) for the
378 comparisons between consistent and inconsistent conditions. For display purposes, the BF₁₀ values were
379 log-transformed. Dotted lines show low (BF₁₀ = 1), moderate (BF₁₀ = 3), and high (BF₁₀ = 10) evidence
380 for a difference between conditions. Similar to Experiment 1, inconsistent scene-object combinations
381 evoked more negative responses at 240-335 ms, 360-500 ms, and 570-590 ms after the object onset
382 relative to consistent combinations.

383 *Tracking object representations in consistent and inconsistent scenes*

384 To test whether semantically consistent scenes facilitate object representations
385 differently from semantically inconsistent scenes when the objects are task-relevant,
386 we performed both timeseries and cumulative decoding analyses, where we classified
387 two consistent or inconsistent objects within each scene at each time point from -100
388 to 800 ms relative to the onset of the object. We found significant decoding for both
389 consistent objects (timeseries decoding: 60-800 ms; cumulative decoding: 70-800 ms)
390 and inconsistent objects (timeseries decoding: 70-760 ms; cumulative decoding: 90-
391 800 ms). Critically, we found the consistent objects were decoded more accurately than
392 inconsistent objects in both the timeseries decoding (310-410 ms and 545-680 ms) and
393 cumulative decoding analyses (160-190 ms, 255-295 ms, and 370-800 ms) (Fig. 7).
394 These results suggest that scene-object consistency can facilitate cortical object
395 representations – but only when the objects are task-relevant. Our data show that such
396 effects arise at least from around 300ms, although the more sensitive cumulative
397 decoding suggests that such effects may be seen much earlier, even within the first
398 200ms of processing. As the current evidence for such early effects is only moderately
399 strong, the exact timing of such effects needs to be confirmed in future studies.



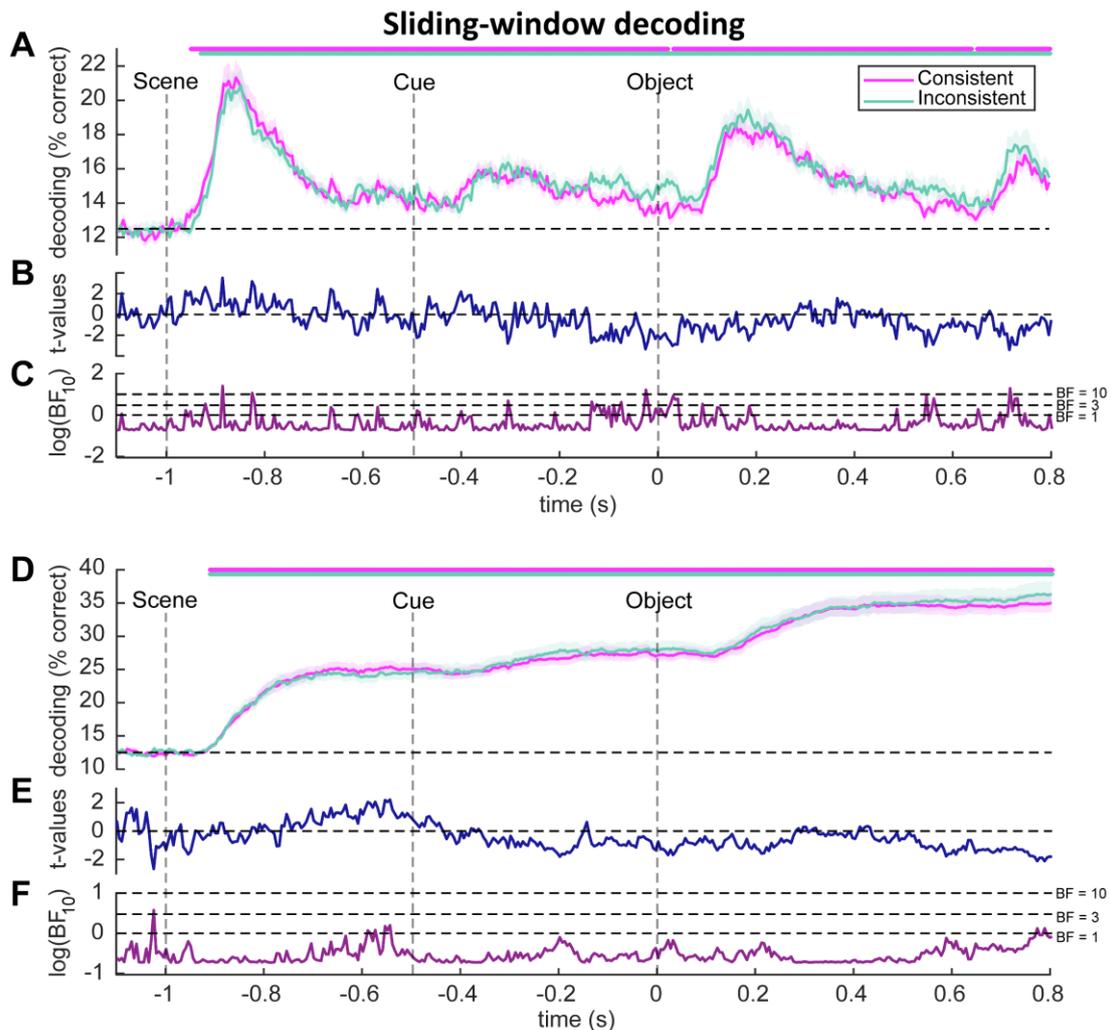
400

401 **Fig. 7** Decoding for the consistent or inconsistent objects within each scene in Experiment 2. **(A)**
 402 Timeseries decoding results, separately for consistent and inconsistent objects. Line markers denote
 403 significant above-chance decoding ($p < 0.05$, FDR-corrected). **(B)** t -values for the comparisons between
 404 consistent and inconsistent conditions. Line markers denote significant differences between the
 405 consistent and inconsistent conditions ($p < 0.05$, FDR-corrected). **(C)** Bayes factors (BF₁₀) for the
 406 comparisons between consistent and inconsistent conditions. For display purposes, the BF₁₀ values were
 407 log-transformed. Dotted lines show low (BF₁₀ = 1), moderate (BF₁₀ = 3), and high (BF₁₀ = 10)
 408 evidence for a difference between conditions. **(D)** Cumulative decoding results, separately for consistent and
 409 inconsistent objects. Line markers denote significant above-chance decoding ($p < 0.05$, FDR-corrected).
 410 **(E, F)** t -values and Bayes factors (BF₁₀) for the comparisons between consistent and inconsistent
 411 conditions, as in (B, C). These results are markedly different from Experiment 1: Scenes can indeed
 412 facilitate the cortical representations of consistent objects when the objects are task-relevant.

413 *Tracking the representation of scenes with consistent and inconsistent objects*

414 As in Experiment 1, we also tested whether semantically consistent objects can
 415 facilitate scene representations. We performed timeseries and cumulative decoding
 416 analyses to discriminate between the eight scene categories separately for the consistent
 417 and inconsistent conditions from -100 ms to 1800 ms relative to the onset of the scene.
 418 We found very similar results as the Experiment 1, with significant decoding for both
 419 consistent scenes (timeseries decoding: 50-1800 ms; cumulative decoding: 90-1800 ms)

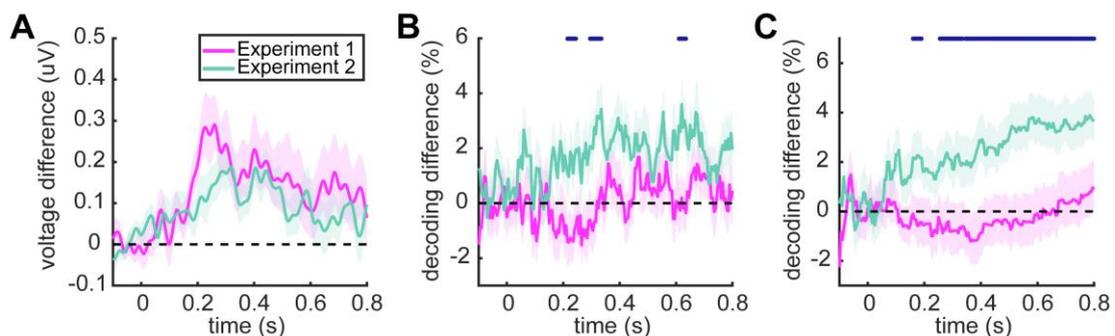
420 and inconsistent scenes (timeseries decoding: 70-1800 ms; cumulative decoding: 90-
421 1800 ms), but no difference in decoding performance between consistent and
422 inconsistent conditions (Fig. 8). These results suggest that facilitation effects between
423 scenes and objects are not mutual, but that they likely depend on behavioral goals: once
424 the objects were task relevant, we found a facilitation effect originating from
425 semantically consistent scenes.



426
427 **Fig. 8** Decoding between scenes with consistent or inconsistent objects in Experiment 2. (A) Timeseries
428 decoding results, separately for scene with consistent and inconsistent objects. Line markers denote
429 significant above-chance decoding ($p < 0.05$, FDR-corrected). (B) t -values for the comparisons between
430 consistent and inconsistent conditions. (C) Bayes factors (BF_{10}) for the comparisons between consistent
431 and inconsistent conditions. For display purposes, the BF_{10} values were log-transformed. Dotted lines
432 show low ($BF_{10} = 1$), moderate ($BF_{10} = 3$), and high ($BF_{10} = 10$) evidence for a difference between
433 conditions. (D) Cumulative decoding results, separately for scene with consistent and inconsistent
434 objects. Line markers denote significant above-chance decoding ($p < 0.05$, FDR-corrected). (E, F) t -
435 values and Bayes factors (BF_{10}) for the comparisons between consistent and inconsistent conditions, as
436 in (B, C). The results show that consistent embedded objects do not automatically facilitate the
437 representation of scenes.

438 *Comparison across experiments*

439 The pattern of results across our experiments revealed reliable ERP effects that are
440 independent of task-relevance, but multivariate decoding demonstrated that
441 representational facilitation effects can only be observed when the objects are task
442 relevant. To statistically quantify this pattern, we directly compared the ERP and
443 decoding results between two experiments. For each participant, we computed the
444 difference between the consistent and inconsistent conditions, and then compared these
445 differences across experiments using independent-samples *t*-tests. For the ERP results,
446 we found no statistical differences across the experiments (all $p > 0.05$, FDR-corrected;
447 Fig. 9A), suggesting that N300/400 effects emerge independently of the task-relevance
448 of the objects. On the flipside, object representations benefitted more strongly from
449 semantically consistent context when the objects were directly task-relevant: In the
450 timeseries decoding, differences between the two experiments emerged at 215-245 ms,
451 295-335 ms, and 610-635 ms, suggesting that during these time points task-relevance
452 enhances the effect of semantically consistent scene context (Fig. 9B). Clear evidence
453 for this effect was also found in the cumulative decoding, where the effect of
454 semantically consistent scene context was stronger in Experiment 2 between 160-190
455 ms and 255-800 ms (Fig. 9C). There were no significant differences in scene decoding
456 between the two experiments. Together this shows that N300/400 ERP differences
457 emerge independently of task relevance, suggesting that they do not index changes in
458 perceptual object representation. By contrast, multivariate decoding reveals that the
459 changes in object representations are modulated by task-relevance: when the objects
460 are critical for behavior, semantically consistent scenes more strongly enhance their
461 cortical representation.



462

463 **Fig. 9** Comparisons of ERP and decoding effects between experiments. (A) Differences in ERP effects
464 (consistent – inconsistent) between experiments. (B) Differences in decoding effects (consistent –
465 inconsistent) for object in timeseries decoding analyses between experiments. Line markers denote
466 significant differences between two experiments ($p < 0.05$). (C) Differences in decoding effects for object
467 in cumulative decoding analyses between experiments. Line markers denote significant differences

468 between two experiments ($p < 0.05$). The results shows that the N300/400 effects emerge independently
469 of the task relevance of the objects, but facilitations from scenes to the representations of objects are
470 task-dependent.

471 **Discussion**

472 In this study, we used EEG to investigate how scene-object consistency affects the
473 quality of object and scene representations. In two experiments, we replicated previous
474 scene-object consistency ERP effects (Mudrik et al., 2010, 2014; Vö and Wolfe, 2013),
475 showing that inconsistent scene-object combinations evoked more negative responses
476 in the N300 and N400 windows than consistent combinations. Critically, multivariate
477 decoding analyses revealed whether these scene-object consistency effects in ERPs
478 were accompanied by changes in the quality of cortical object and scene
479 representations. We found that task-irrelevant consistent and inconsistent objects were
480 decoded equally well in Experiment 1, despite pronounced ERP differences in the
481 N300/400 range. When the objects were task-relevant in Experiment 2, we observed a
482 comparable N300/400 ERP effect, which now was accompanied by enhanced object
483 decoding. Across both experiments, we found no significant differences in scene
484 category decoding between consistent and inconsistent conditions. These results
485 suggest that the N300/N400 ERP effects are not necessarily indicative of enhanced
486 perceptual representations. Further, they suggest that facilitations between objects and
487 scenes are task-dependent rather than automatic.

488 **N300 effects do not index changes in perceptual processing**

489 The N300 effects found in the study replicated previous findings in studies of scene-
490 object consistency (Vö and Wolfe, 2013; Mudrik et al., 2014; Draschkow et al., 2018;
491 Truman and Mudrik, 2018; Coco et al., 2020). Particularly the early N300 effects were
492 often interpreted as reflecting differences in perceptual processing (Schendan and
493 Maher, 2009; Mudrik et al., 2010; Dyck and Brodeur, 2015; Sauv e et al., 2017; Kumar
494 et al., 2021). Such findings are often explained through models of contextual facilitation
495 (Bar, 2004; Bar et al., 1996), which propose that object representations are refined by
496 more readily available information about the consistent context. Specifically, when a
497 scene is presented, gist-consistent schemas are rapidly activated through non-selective
498 processing channels (Wolfe et al., 2011). By comparing this rapidly available scene gist
499 to incoming visual information, perceptual uncertainty in object recognition is reduced.

500 However, if the object does not match the scene gist, its identification should be
501 impeded. It was argued that this mismatch between inconsistent objects and the pre-
502 activated schemas elicits a larger N300 amplitude, signifying a prediction error that
503 occurs during perceptual object analysis (Kumar et al., 2021).

504 Our data challenge this interpretation. We show that enhanced N300 amplitudes are
505 observed independently of changes in object representation. We found a reliable N300
506 differences between consistent and inconsistent objects, which were highly similar for
507 task-relevant and task-irrelevant objects (if anything, the effect for task-relevant objects
508 was smaller). By contrast, object information was similar for consistent and
509 inconsistent objects when they were not task-relevant; only when they were task-
510 relevant, we found that scene-object consistency facilitated object representations. In is
511 worth noting that both task-relevant and task-irrelevant objects within the scenes could
512 be decoded reliably and with high accuracy in both experiments, which is in line with
513 previous reports (Kaiser et al., 2016); our results therefore cannot be attributed to a
514 failure to decode the objects in the first place.

515 The pattern of results obtained in our study is therefore inconsistent with the N300
516 indexing a change in perceptual representations. Our results are rather consistent with
517 an interpretation that views the N300 as a general marker of inconsistency or a purely
518 attentional response. On this view, N300 differences are post-perceptual in nature,
519 possibly reflecting differences in attention. Contrary to the N300, consistency-related
520 differences in the N400 time window are commonly interpreted as a marker of
521 differences in post-perceptual semantic processing (Võ and Wolfe, 2013; Truman and
522 Mudrik, 2018). In fact, a recent study has shown that N400 effects are qualitatively
523 similar to N300 effects (Draschkow et al., 2018), further supporting the view that N300
524 differences are not directly indicative of changes in perceptual encoding.

525 **Semantic consistency only facilitates task-relevant representations**

526 Our results suggest that cross-facilitation effects between objects and scenes are not
527 automatic but task-dependent. Consistent objects were only decoded better than
528 inconsistent objects in Experiment 2 where they were directly task-relevant, suggesting
529 that semantically consistent scenes only facilitate object processing when the objects
530 are critical for behavior. Further, decoding between the different scenes was similar for

531 scenes that contained consistent and inconsistent objects. As the scenes were never task-
532 relevant, this supports the view that that mutual influences between scene and object
533 representations are only observed when they support ongoing behavior.

534 Several previous neuroimaging studies reported a cross-facilitation between scene and
535 object processing (Brandman and Peelen, 2017, 2019; Kaiser et al., 2021), reporting
536 that scenes enhance the cortical representation of objects (Brandman and Peelen, 2017;
537 Kaiser et al., 2021), and objects facilitate the representation of scenes (Brandman and
538 Peelen, 2019). In these studies, participants were asked to attend the objects or scenes
539 by memorizing them, completing repetition detection tasks, or categorization tasks.
540 One recent study directly compared cross-facilitation effects between objects and
541 scenes under different task demands (Kaiser et al., 2021). In this study, spatially
542 consistent scene context facilitated object representation more than spatially
543 inconsistent scene context when the objects were task-relevant. When participants
544 instead performed a task on the scene, object representations were comparable for the
545 spatially consistent and inconsistent scene contexts.

546 These results are in line with the current study, in which semantically consistent scene
547 context only facilitated perceptual object processing when it was beneficial for the task
548 at hand. Our findings therefore support a view where the visual system uses contextual
549 information in a flexible and strategic way: When scene context is beneficial for the
550 current task demands, the visual system harnesses contextual information to enhance
551 object representations. Conversely, if the current task does not benefit from contextual
552 information, no cross-facilitation between object and scene processing is found.

553 **Conclusions**

554 In the study, we investigated how scene-object consistency affects scene and object
555 representations. Our results suggest that differences in the N300/N400 components
556 related to scene-object consistency do not directly index differences in perceptual
557 representations, but rather reflect a generic marker of semantic violations. Further, they
558 suggest that facilitation effects between objects and scenes are task-dependent rather
559 than automatic. Our findings highlight that there are multiple markers of semantic
560 consistency that reflect different underlying brain mechanisms. How these mechanisms
561 interact to support efficient real-world vision needs to be explored in future studies.

562 **Acknowledgments**

563 D.K. and R.M.C. are supported by the Deutsche Forschungsgemeinschaft (DFG) grants
564 (CI241/1-1, CI241/3-1, CI241/7-1, KA4683/2-1). R.M.C. is supported by the European
565 Research Council (ERC) grant (803370). L.C. is supported by the Chinese Scholarship
566 Council (CSC).

567 **References**

- 568 Bar M (2004) Visual objects in context. *Nat Rev Neurosci* 5:617–629.
- 569 Bar M, Ullman S (1996) Spatial Context in Recognition. *Perception* 25:343–352.
- 570 Biederman I, Mezzanotte RJ, Rabinowitz JC (1982) Scene perception: Detecting and judging
571 objects undergoing relational violations. *Cognit Psychol* 14:143–177.
- 572 Boring MJ, Ridgeway K, Shvartsman M, Jonker TR (2020) Continuous decoding of cognitive
573 load from electroencephalography reveals task-general and task-specific correlates. *J*
574 *Neural Eng* 17:056016.
- 575 Boyce SJ, Pollatsek A (1992) Identification of objects in scenes: the role of scene background
576 in object naming. *J Exp Psychol Learn Mem Cogn* 18:531.
- 577 Boyce SJ, Pollatsek A, Rayner K (1989) Effect of background information on object
578 identification. *J Exp Psychol Hum Percept Perform* 15:556.
- 579 Brainard DH (1997) The psychophysics toolbox. *Spat Vis* 10:433–436.
- 580 Brandman T, Peelen MV (2017) Interaction between scene and object processing revealed by
581 human fMRI and MEG decoding. *J Neurosci* 37:7700–7710.
- 582 Brandman T, Peelen MV (2019) Signposts in the fog: Objects facilitate scene representations
583 in left scene-selective cortex. *J Cogn Neurosci* 31:390–400.
- 584 Coco MI, Nuthmann A, Dimigen O (2020) Fixation-related brain potentials during semantic
585 integration of object–scene information. *J Cogn Neurosci* 32:571–589.
- 586 Cornelissen TH, Vö ML (2017) Stuck on semantics: Processing of irrelevant object-scene
587 inconsistencies modulates ongoing gaze behavior. *Atten Percept Psychophys* 79:154–
588 168.
- 589 Davenport JL (2007) Consistency effects between objects in scenes. *Mem Cognit* 35:393–401.
- 590 Davenport JL, Potter MC (2004) Scene consistency in object and background perception.
591 *Psychol Sci* 15:559–564.
- 592 Draschkow D, Heikel E, Vö ML, Fiebach CJ, Sassenhagen J (2018) No evidence from MVPA
593 for different processes underlying the N300 and N400 incongruity effects in object-
594 scene processing. *Neuropsychologia* 120:9–17.
- 595 Dyck M, Brodeur MB (2015) ERP evidence for the influence of scene context on the
596 recognition of ambiguous and unambiguous objects. *Neuropsychologia* 72:43–51.
- 597 Ganis G, Kutas M (2003) An electrophysiological study of scene effects on object
598 identification. *Cogn Brain Res* 16:123–144.
- 599 Grootswagers T, Wardle SG, Carlson TA (2017) Decoding dynamic brain patterns from evoked
600 responses: A tutorial on multivariate pattern analysis applied to time series

- 601 neuroimaging data. *J Cogn Neurosci* 29:677–697.
- 602 Kaiser D, Häberle G, Cichy RM (2020a) Real-world structure facilitates the rapid emergence
603 of scene category information in visual brain signals. *J Neurophysiol* 124:145–151.
- 604 Kaiser D, Häberle G, Cichy RM (2021) Coherent natural scene structure facilitates the
605 extraction of task-relevant object information in visual cortex. *NeuroImage*:118365.
- 606 Kaiser D, Inciuraitė G, Cichy RM (2020b) Rapid contextualization of fragmented scene
607 information in the human visual system. *NeuroImage* 219:117045.
- 608 Kaiser D, Nyga K (2020) Tracking cortical representations of facial attractiveness using time-
609 resolved representational similarity analysis. *Sci Rep* 10:1–10.
- 610 Kaiser D, Oosterhof NN, Peelen MV (2016) The neural dynamics of attentional selection in
611 natural scenes. *J Neurosci* 36:10522–10528.
- 612 Kaiser D, Quek GL, Cichy RM, Peelen MV (2019) Object vision in a structured world. *Trends*
613 *Cogn Sci* 23:672–685.
- 614 Kumar M, Federmeier KD, Beck DM (2021) The N300: An Index for Predictive Coding of
615 Complex Visual Objects and Scenes. *Cereb Cortex Commun* 2:tgab030.
- 616 Lowe MX, Rajsic J, Ferber S, Walther DB (2018) Discriminating scene categories from brain
617 activity within 100 milliseconds. *Cortex* 106:275–287.
- 618 Mudrik L, Lamy D, Deouell LY (2010) ERP evidence for context congruity effects during
619 simultaneous object–scene processing. *Neuropsychologia* 48:507–517.
- 620 Mudrik L, Shalgi S, Lamy D, Deouell LY (2014) Synchronous contextual irregularities affect
621 early scene processing: Replication and extension. *Neuropsychologia* 56:447–458.
- 622 Munneke J, Brentari V, Peelen M (2013) The influence of scene context on object recognition
623 is independent of attentional focus. *Front Psychol* 4:552.
- 624 Oostenveld R, Fries P, Maris E, Schoffelen J-M (2011) FieldTrip: open source software for
625 advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput*
626 *Intell Neurosci* 2011.
- 627 Oosterhof NN, Connolly AC, Haxby JV (2016) CoSMoMVPA: multi-modal multivariate
628 pattern analysis of neuroimaging data in Matlab/GNU Octave. *Front Neuroinformatics*
629 10:27.
- 630 Palmer SE (1975) The effects of contextual scenes on the identification of objects. *Mem Cognit*
631 3:519–526.
- 632 Ramkumar P, Jas M, Pannasch S, Hari R, Parkkonen L (2013) Feature-specific information
633 processing precedes concerted activation in human visual cortex. *J Neurosci* 33:7691–
634 7699.
- 635 Sauvė G, Harmand M, Vanni L, Brodeur MB (2017) The probability of object–scene co-
636 occurrence influences object identification processes. *Exp Brain Res* 235:2167–2179.
- 637 Schendan HE, Maher SM (2009) Object knowledge during entry-level categorization is
638 activated and modified by implicit memory after 200 ms. *Neuroimage* 44:1423–1438.
- 639 Truman A, Mudrik L (2018) Are incongruent objects harder to identify? The functional
640 significance of the N300 component. *Neuropsychologia* 117:222–232.
- 641 Vő ML, Boettcher SE, Draschkow D (2019) Reading scenes: how scene grammar guides
642 attention and aids perception in real-world environments. *Curr Opin Psychol* 29:205–
643 210.
- 644 Vő ML, Henderson JM (2009) Does gravity matter? Effects of semantic and syntactic
645 inconsistencies on the allocation of attention during scene perception. *J Vis* 9:24.1-15.

- 646 V̄o ML, Henderson JM (2011) Object–scene inconsistencies do not capture gaze: evidence
647 from the flash-preview moving-window paradigm. *Atten Percept Psychophys*
648 73:1742–1753.
- 649 V̄o ML, Wolfe JM (2013) Differential electrophysiological signatures of semantic and syntactic
650 scene processing. *Psychol Sci* 24:1816–1823.
- 651 Wolfe JM, V̄o ML, Evans KK, Greene MR (2011) Visual search in scenes involves selective
652 and nonselective pathways. *Trends Cogn Sci* 15:77–84.
- 653