

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Coherent natural scene structure facilitates the extraction of task-relevant object information in visual cortex

Daniel Kaiser¹, Greta Häberle^{2,3,4}, Radoslaw M. Cichy^{2,3,4,5}

¹*Department of Psychology, University of York, York, UK*

²*Department of Education and Psychology, Freie Universität Berlin, Berlin, Germany*

³*Charité – Universitätsmedizin Berlin, Einstein Center for Neurosciences Berlin, Berlin, Germany*

⁴*Humboldt-Universität zu Berlin, Faculty of Philosophy, Berlin School of Mind and Brain, Berlin, Germany*

⁵*Bernstein Center for Computational Neuroscience Berlin, Berlin, Germany*

Correspondence:

Dr. Daniel Kaiser

Department of Psychology

University of York

Heslington, York

YO10 5DD, UK

danielkaiser.net@gmail.com

23

24 **Abstract**

25

26 Looking for objects within complex natural environments is a task everybody
27 performs multiple times each day. In this study, we explore how the brain uses the
28 typical composition of real-world environments to efficiently solve this task. We
29 recorded fMRI activity while participants performed two different categorization tasks
30 on natural scenes. In the object task, they indicated whether the scene contained a
31 person or a car, while in the scene task, they indicated whether the scene depicted
32 an urban or a rural environment. Critically, each scene was presented in an “intact”
33 way, preserving its coherent structure, or in a “jumbled” way, with information
34 swapped across quadrants. In both tasks, participants’ categorization was more
35 accurate and faster for intact scenes. These behavioral benefits were accompanied
36 by stronger responses to intact than to jumbled scenes across high-level visual
37 cortex. To track the amount of object information in visual cortex, we correlated multi-
38 voxel response patterns during the two categorization tasks with response patterns
39 evoked by people and cars in isolation. We found that object information in object-
40 and body-selective cortex was enhanced when the object was embedded in an intact,
41 rather than a jumbled scene. However, this enhancement was only found in the object
42 task: When participants instead categorized the scenes, object information did not
43 differ between intact and jumbled scenes. Together, these results indicate that
44 coherent scene structure facilitates the extraction of object information in a task-
45 dependent way, suggesting that interactions between the object and scene
46 processing pathways adaptively support behavioral goals.

47

48 **Keywords**

49

50 object perception, natural scene categorization, real-world structure, functional
51 magnetic resonance imaging, multivariate pattern analysis

52

53 **1 Introduction**

54

55 Despite the complexity of our everyday environments, perceiving objects embedded
56 in natural scenes is remarkably efficient. This efficiency is illustrated by studies that
57 require participants to categorize objects under conditions of limited visual exposure:
58 For instance, participants can tell whether a scene contains an animal or not from just
59 a single glance (Thorpe et al., 1996; Potter, 1975, 2012), and even when only limited
60 attentional resources are available (Li et al., 2002).

61 The ability to effortlessly make such categorization responses is underpinned
62 by the efficient extraction of object information in visual cortex. Neuroimaging
63 research has shown that the category of task-relevant objects can be reliably
64 decoded from fMRI activity patterns in visual cortex, even when the objects are
65 embedded in complex natural scenes (Peelen et al., 2009; Peelen & Kastner, 2011;
66 Seidl et al., 2012) or movies (Cukur et al., 2013; Nastase et al., 2017; Shahdloo et al.,
67 2020). M/EEG studies demonstrate that object category is represented well within the
68 first 200ms of vision, even when the object is shown under such naturalistic
69 conditions (Cauchoix et al., 2014; Kaiser et al., 2016; VanRullen & Thorpe, 2001;
70 Thorpe et al., 1996). Together, these results highlight that the cortical processing of
71 objects appearing within rich real-world environments is surprisingly efficient.

72 This processing efficiency becomes less surprising if scene context is not just
73 considered as a nuisance that puts additional strain on our visual resources. Indeed,
74 contextual information can facilitate object processing (Bar, 2004): For instance,
75 scene context allows for efficient allocation of attention (Torralba et al., 2006; Wolfe
76 et al., 2011; Vö et al., 2019), or for disambiguating object information under
77 uncertainty (Brandmann & Peelen, 2017; Oliva & Torralba, 2007). Such findings
78 demonstrate that object and scene processing mechanisms interact with each other
79 to enable the efficient processing of object information.

80 Here, we investigated how the coherent spatial structure of the scene context
81 aids the extraction of object information from the scene. To this end, we used a
82 jumbling paradigm, in which we disrupted the scenes' coherent structure by dividing
83 them into multiple rectangular pieces and shuffling those pieces. Classical studies

84 suggest that jumbling drastically impairs participants' ability to categorize both the
85 scene itself (Biederman et al., 1974), and the object embedded within the scene
86 (Biederman et al., 1972, 1973). Such impairments can be linked to changes in cortical
87 scene processing: We have recently shown that scene-selective brain responses are
88 less pronounced and contain less scene category information when the scene is
89 jumbled (Kaiser et al., 2020a, 2020b). However, it is unclear how these changes in
90 scene-selective activations modulate the representation of objects within the scene.

91 In the current study, we thus set out to characterize how the presence of an
92 intact – versus a jumbled – scene context modulates object representations in visual
93 cortex. First, we asked whether cortical object processing is indeed facilitated by the
94 presence of a coherent scene context. Second, we asked whether such facilitation
95 effects depend on the objects being relevant or irrelevant for current behavioral goals.

96 To answer these questions, we recorded fMRI activity while participants
97 categorized objects contained in intact or jumbled scenes. We found that fMRI
98 responses across high-level visual cortex were generally higher for intact scenes than
99 for jumbled scenes, revealing widespread sensitivity to scene structure. When
100 analyzing object category information in multi-voxel response patterns, we found that
101 coherent scene structure enhanced object information in object-selective visual
102 cortex. However, this enhancement was task-specific: When participants categorized
103 the scenes instead of the objects, we found no such enhancement of object
104 information. These results suggest that the visual brain uses coherent real-world
105 structure to more efficiently extract task-relevant object information from complex
106 scenes.

107 **2 Materials and Methods**

108

109 **2.1 Participants**

110 Twenty-five healthy adults (mean age 26.4 years, SD=5.3; 15 female, 10 male)
111 participated. All participants had normal or corrected-to-normal vision. They all
112 provided informed written consent and received either monetary reimbursement or
113 course credits. Procedures were approved by the ethical committee of the
114 Department of Psychology at Freie Universität Berlin and were in accordance with the
115 Declaration of Helsinki.

116

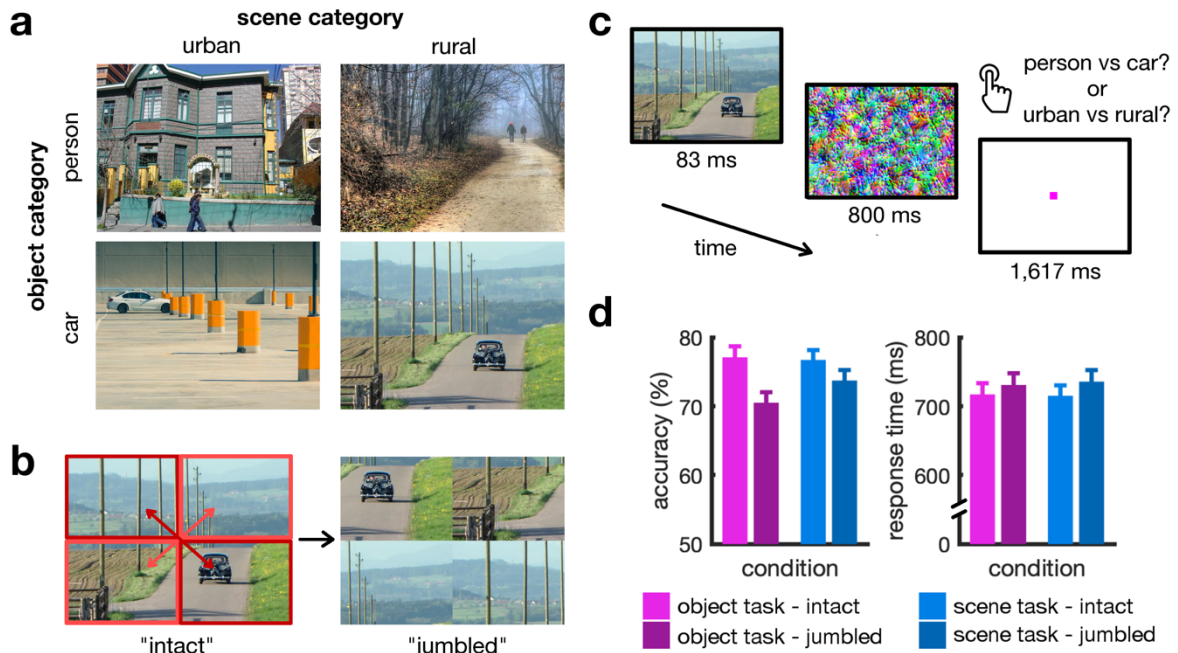
117 **2.2 Stimuli**

118 The stimulus set consisted of colored natural scene photographs (640×480 pixels
119 resolution). Scenes were selected to cover three independent manipulations. First,
120 each scene contained one of two object categories: half of the scenes contained a
121 person (or multiple people), whereas the other half contained a car (or multiple cars).
122 Second, the person or car appeared equally often in each of the quadrants of the
123 scene. Third, each scene belonged to one of two scene categories: half of the scenes
124 depicted urban environments, the other half depicted rural environments. For each
125 possible combination of these factors (e.g., a person appearing in the bottom left
126 quadrant of a rural scene), 10 unique scene exemplars were available, yielding 160
127 scenes in total (2 object categories × 4 object locations × 2 scene categories × 10
128 exemplars). During the experiment, the scenes could be presented in their original
129 orientation or mirrored along their vertical axis (as in Kaiser et al., 2016), yielding a
130 total of 320 different scene stimuli. Example scenes are shown in Figure 1a.

131 To manipulate scene structure, we either presented the scenes in a coherent,
132 “intact” condition or in an incoherent, “jumbled” condition. Jumbled scenes were
133 generated by shuffling the four quadrants of the image in a crisscrossed way (Figure
134 1b). This manipulation solely affected the scene’s structure, but not the people or
135 cars contained in the scene: First, as the objects never straddled the boundary
136 between quadrants, the objects themselves always remained unaltered. Second, as
137 the objects appeared equally often in each quadrant before jumbling the scenes, they
138 also appeared equally often in each quadrant after jumbling them.

139 In total, 640 scene images were used, which covered 320 intact scenes and
140 320 jumbled scenes. Additionally, 200 colored texture masks (Kaiser et al., 2016) were
141 used to visually mask the scenes during the experiment (see below).

142



143

144 **Figure 1. Stimuli, paradigm, and behavioral results.** **a)** Stimuli consisted of natural
145 scene images from two categories: urban or rural environments. Each of the scenes
146 contained one of two object categories: people or cars. **b)** During the experiment,
147 these scenes were shown in an unaltered way ("intact" condition) or with their
148 quadrants intermixed ("jumbled" condition). The jumbled scenes were created by
149 shuffling the quadrants in a crisscrossed way, as illustrated. **c)** Participants viewed
150 each scene briefly, followed by a visual mask. In separate runs, they performed two
151 different tasks: They were either asked to indicate whether the scene contained a
152 person or a car ("object task") or whether the scene depicted an urban or a rural
153 environment ("scene task"). **d)** In both tasks, scene structure impacted behavioral
154 performance: Participants were significantly more accurate and faster for the intact
155 scenes than for the jumbled scenes. Error bars represent standard errors of the mean.

156

157 **2.3 Experimental Paradigm**

158 Each participant completed four experimental runs of 17 minutes each. Each run
159 contained 320 experimental trials, corresponding to 320 unique scene stimuli. For

160 half of the participants, the even runs only contained the original scenes, while the
161 odd runs only contained the horizontally mirrored scenes; for the other half of the
162 participants, the odd runs only contained the original scenes, while the even runs only
163 contained the horizontally mirrored scenes. Each of the scenes was presented once
164 during the run, with order fully randomized.

165 On each trial, the scene was presented for 83ms, immediately followed by a
166 visual mask (chosen randomly from the 200 available masks) for 800ms. Images were
167 shown within a black rectangle (10deg X 7.5deg visual angle). After an inter-trial
168 interval of 1,617ms, during which a pink fixation dot was shown, the next trial started.
169 An example trial is illustrated in Figure 1c. In addition to the experimental trials, each
170 run contained 80 fixation-only trials, during which only the fixation dot was displayed.
171 Runs started and ended with a brief fixation period.

172 In two of the four runs, participants were asked to categorize the object
173 contained in each scene as either a person or a car (“object task”). In the other two
174 runs, participants were asked to categorize the scene as either a rural or an urban
175 environment (“scene task”). Participants were instructed to respond as accurately
176 and quickly as possible, with an emphasis on accuracy. Button-press responses were
177 recorded during the whole inter-trial interval (i.e., until 2,500s after stimulus onset).
178 The four runs were alternating between the object and scene tasks. The task in the
179 first run was counter-balanced between participants. Notably, physical stimulation
180 was completely identical across the object and scene tasks.

181 All stimuli were back-projected onto a translucent screen mounted to the head
182 end of the scanner bore. Participants viewed the stimulation through a mirror
183 attached to the head coil. Stimulus presentation was controlled using the
184 Psychtoolbox (Brainard, 1997).

185

186 **2.4 Benchmark Localizer Paradigm**

187 In addition to the experimental runs, each participant completed a benchmark
188 localizer run, which was designed to obtain “benchmark” patterns in response to
189 people and cars in isolation (Peelen et al., 2009; Peelen & Kastner, 2011). During this
190 run, participants viewed images of bodies, cars, and scrambled images of bodies and
191 cars. For each of the three categories, 40 images were used. All images were different

192 than the ones used in the main experiment. These images were presented in a block
193 design. Each block lasted 20 seconds and contained 20 images of one of the three
194 categories, or only a fixation cross. Images were presented for 500ms (5deg × 5deg
195 visual angle), separated by a 500ms inter-stimulus interval. The benchmark localizer
196 run consisted of a total of 24 blocks (6 blocks for each of the three stimulus
197 categories, and 6 fixation-only blocks). Four consecutive blocks always contained the
198 four different conditions in random order. Participants were instructed to respond to
199 one-back image repetitions (i.e., two identical images back-to-back), which
200 happened once during each non-fixation block. The benchmark localizer run lasted
201 8:30 minutes and was completed halfway through the experiment, after two of the
202 four experimental runs.

203

204 **2.5 fMRI recording and preprocessing**

205 MRI data was acquired using a 3T Siemens Tim Trio Scanner equipped with a 12-
206 channel head coil. T2*-weighted gradient-echo echo-planar images were collected
207 as functional volumes (TR=2s, TE=30ms, 70° flip angle, 3mm³ voxel size, 37 slices,
208 20% gap, 192mm FOV, 64×64 matrix size, interleaved acquisition). Additionally, a T1-
209 weighted anatomical image (MPRAGE; 1mm³ voxel size) was obtained.

210 Preprocessing and hemodynamic response modelling was performed using
211 SPM12 (www.fil.ion.ucl.ac.uk/spm/). Functional volumes were realigned and
212 coregistered to the anatomical image. Further, the T1 image was segmented to obtain
213 transformation parameters to standard MNI-305 space.

214 Functional data from each experimental run were modelled in a general linear
215 model (GLM) with 16 experimental predictors (2 object categories × 4 object locations
216 × 2 scene categories). Additionally, we included the six movement regressors
217 obtained during realignment. Data from the benchmark localizer run were modelled
218 in a GLM with three experimental predictors (person, car, scrambled) and six
219 movement regressors.

220

221 **2.6 Region of interest definition**

222 We restricted fMRI analyses to five regions of interest (ROIs): early visual cortex (V1),
223 object-selective lateral occipital cortex (LO), body-selective extrastriate body area

224 (EBA), scene-selective occipital place area (OPA), and scene-selective
225 parahippocampal place area (PPA). ROIs masks were defined using group-level
226 activation masks from functional brain atlases (for V1: Wang et al., 2015; for LO, EBA,
227 OPA, and PPA: Julian et al., 2012). ROIs were defined separately for each
228 hemisphere. All ROI masks were inverse-normalized into individual-participant space
229 using the parameters obtained during T1 segmentation. Voxel counts in individual-
230 participant space amounted to 248/271 (V1; left/right), 929/947 (LO), 402/443 (EBA),
231 26/47 (OPA), and 140/105 (PPA). Notably, the LO and EBA ROIs overlapped to some
232 extent (300/406 voxels overlap, left/right); the inclusion of the EBA allowed us to see
233 whether the results hold in a smaller cortical region with a narrower category
234 preference for bodies. As we did not have any hypothesis related to hemispheric
235 differences, all results for the left- and right-hemispheric ROIs were averaged before
236 statistical analysis.

237

238 **2.7 Univariate analysis**

239 Response magnitudes during the experimental runs were analyzed separately for
240 each ROI. We first averaged beta values across the two object-task and scene-task
241 runs, respectively. We then averaged beta values across object categories, object
242 locations, and scene categories. This way, we obtained response magnitudes for four
243 conditions: (1) responses to intact scenes in the object task, (2) responses to jumbled
244 scenes in the object task, (3) responses to intact scenes in the scene task, and (4)
245 responses to jumbled scenes in the scene task. These four conditions allowed us to
246 separately estimate the effects of task (object task versus scene task) and scene
247 structure (intact versus jumbled) on neural responses across the five ROIs.

248

249 **2.8 Multivariate pattern analysis**

250 Multivariate pattern analysis (MVPA) was carried out in CoSMoMVPA (Oosterhof et
251 al., 2016). Our MVPA approach closely followed similar fMRI studies that investigated
252 the representation of objects in natural scenes (Peelen et al., 2009; Peelen & Kastner,
253 2011). We first computed a one-sample t-contrasts for every condition against
254 baseline. In the benchmark localizer run, there were 2 such t-contrasts (one for people
255 versus baseline, and one for cars versus baseline). In the object task and scene task

256 runs, there were 16 t-contrasts each (one contrast for each experimental condition
257 against baseline, reflecting 2 object categories \times 4 object locations \times 2 scene
258 categories). For each of the three tasks (benchmark localizer, object task, and scene
259 task), the resulting t-values were normalized for each voxel by subtracting the average
260 t-value across conditions. For each ROI, multi-voxel response patterns were
261 constructed by concatenating the t-values across all voxels belonging to the ROI.

262 To obtain an index of object discriminability (i.e., how discriminable people and
263 cars in scenes are based on multi-voxel response patterns), we performed a cross-
264 correlation MVPA. The goal of this analysis was to quantify how “person-like” or “car-
265 like” the cortical representation of each of the scenes was, thereby isolating the
266 amount of object category information in visual cortex. To this end, we correlated
267 multi-voxel response patterns evoked by people and cars in isolation (from the
268 benchmark localizer) with response patterns evoked by people and cars contained in
269 a scene (from one of the experimental tasks). These correlations were Fisher-
270 transformed. To quantify object discriminability, we then subtracted the correlations
271 between different categories (e.g., person in isolation and car within a scene) from
272 correlations between the same categories (e.g., person in isolation and person within
273 a scene). This yielded an index of category-discriminability, with values greater than
274 zero indicating that the two categories are represented differently (Haxby et al., 2001).

275 Before performing this analysis, response patterns in the main experiment
276 were averaged across object locations and scene categories. This way, we obtained
277 an index of object category-discriminability for four separate conditions: (1) category-
278 discriminability for intact scenes in the object task, (2) category-discriminability for
279 jumbled scenes in the object task, (3) category-discriminability for intact scenes in
280 the scene task, and (4) category-discriminability for jumbled scenes in the scene task.
281 The resulting four conditions allowed us to estimate the effects of scene structure on
282 the quality of object representations in visual cortex, both when the objects were
283 task-relevant and task-irrelevant.

284

285 **2.9 Statistical testing**

286 To compare behavioral performance, univariate responses, and multi-voxel pattern
287 information across conditions, we used repeated-measures ANOVAs and paired-
288 sample t-tests.

289

290 ***2.10 Data availability***

291 Data are publicly available on OSF (doi.org/10.17605/osf.io/gs2t5). Other materials
292 are available from the corresponding author upon request.

293

294

295 **3 Results**

296

297 ***3.1 Coherent scene structure facilitates the perception of objects within scenes***

298 We first analyzed participants' behavioral performance in the object and scene tasks,
299 separately for the intact and jumbled scenes (Figure 1d). In the object task,
300 participants' categorization (person versus car) of objects within the intact scenes
301 was more accurate, $t(24)=8.28$, $p<.001$, and faster, $t(24)=3.26$, $p=.0033$, compared to
302 the jumbled scenes. In the scene task, participants' categorization (rural versus urban)
303 of the intact scenes was more accurate, $t(24)=4.77$, $p<.001$, and faster, $t(24)=3.26$,
304 $p=.0033$, compared to the jumbled scenes. These results are in line with classical
305 findings on object and scene categorization in jumbling paradigms (Biederman, 1972;
306 Biederman et al., 1973, 1974), showcasing that scene jumbling has a profound impact
307 on perception.

308 Further, when directly comparing the two tasks, we did not find differences in
309 accuracy, $F(1,24)=3.13$, $p=.090$, or response times, $F(1,24)=0.04$, $p=.84$. Any
310 differences in neural responses are therefore unlikely to reflect differences in task
311 difficulty, and therefore attentional engagement, between the two tasks.

312 Together, these results demonstrate that jumbling similarly impairs the
313 perception of the scene and the objects contained in it, demonstrating a cross-
314 facilitation between scene and object vision that can be observed on the behavioral
315 level.

316

317 ***3.2 Scene structure impacts univariate responses across object- and scene- 318 selective cortex***

319 To quantify the effects of scene jumbling on the neural level, we first ran univariate
320 analyses. In these analyses, we compared fMRI response magnitudes across the
321 intact and jumbled scenes and across the two tasks (Figure 2). To do so, we
322 performed a 2×2 repeated measures ANOVA with the factors scene structure (intact
323 versus jumbled) and task (object task versus scene task). The analysis was performed
324 separately and in turn for each of the five ROIs: V1, LO, EBA, OPA, and PPA. Detailed
325 results for these analyses can be found in Table 1.

326 In V1, responses were comparable across all conditions, all $F < 1.25$, $p > .27$,
 327 suggesting that V1 is not sensitive to typical scene composition.

328 In all extrastriate ROIs, we found a main effect of scene structure, which
 329 indicated stronger responses to intact than to jumbled scenes, all $F(1,24) > 7.95$,
 330 $p < .010$. Comparing this effect across regions, we found that it was more pronounced
 331 in the scene-selective regions, OFA versus LO/EBA, both $F(1,24) > 31.17$, $p < .001$, and
 332 PPA versus LO/EBA, both $F(1,24) > 35.55$, $p < .001$. This finding confirms our previous
 333 fMRI results, which revealed particularly strong effects of scene jumbling in scene-
 334 selective areas of visual cortex (Kaiser et al., 2020a).

335

336 **Table 1:** *Univariate responses, analyzed in a 2x2 repeated measures ANOVA with the*
 337 *factors scene structure (intact versus jumbled) and task (object task versus scene*
 338 *task). Significant effects are highlighted in bold.*

ROI	Main Effect Scene Structure	Main Effect Task	Interaction Effect Structure × Task
V1	$F(1,24)=1.25$, $p=.28$	$F(1,24)<0.01$, $p=.98$	$F(1,24)=0.09$, $p=.76$
LO	$F(1,24)=9.74$, $p=.005$	$F(1,24)=0.04$, $p=.85$	$F(1,24)=0.97$, $p=.33$
EBA	$F(1,24)=7.95$, $p=.009$	$F(1,24)=0.21$, $p=.65$	$F(1,24)=2.46$, $p=.13$
OPA	$F(1,24)=27.18$, $p<.001$	$F(1,24)=0.09$, $p=.77$	$F(1,24)=0.97$, $p=.34$
PPA	$F(1,24)=48.02$, $p<.001$	$F(1,24)=6.51$, $p=.017$	$F(1,24)=0.51$, $p=.48$

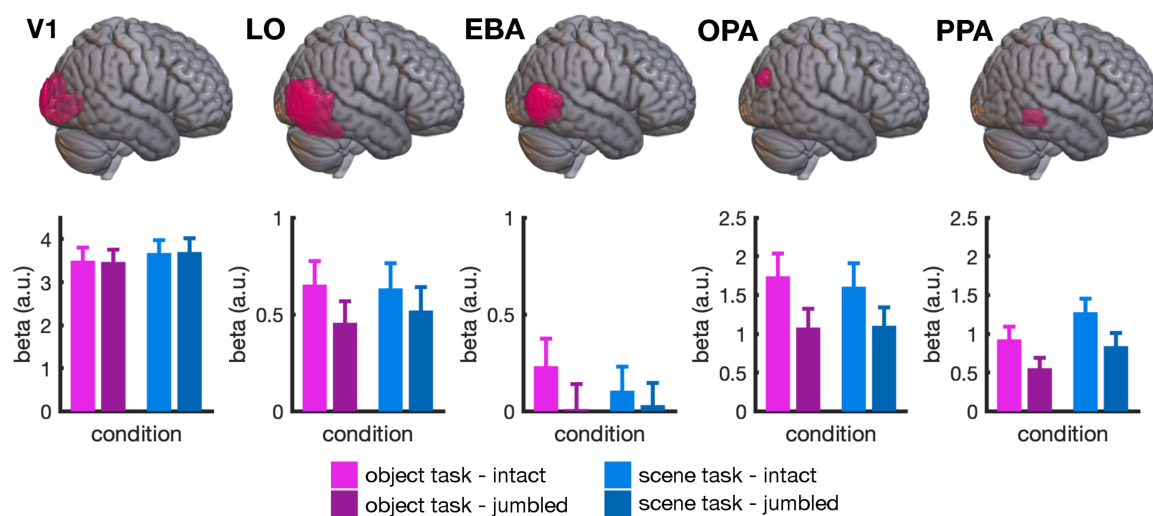
339

340 In all ROIs, scene structure affected univariate responses similarly across the
 341 two tasks, as indexed by no significant interaction effects, all $F < 2.46$, $p > 0.12$. This
 342 pattern of results mirrors the pattern observed in behavior, where scene jumbling
 343 produced comparable effects in the object and scene tasks.

344 PPA was the only region that additionally showed an effect of task,
 345 $F(1,24)=6.51$, $p=.017$, with stronger responses in the scene task compared to the
 346 object task. This suggests an increased importance of computations in higher-level
 347 scene-selective cortex when scene attributes were behaviorally relevant.

348 Having established that scene structure enhanced cortical responses across
 349 object- and scene-selective cortex, and similarly for both tasks, we next asked how
 350 scene structure contributes to the extraction of object information – both when the
 351 objects are behaviorally relevant and when they are not.

352



353

354 **Figure 2. Univariate results.** In all extrastriate regions, but not in V1, we found a
355 significant main effect of scene structure: Intact scenes led to significantly stronger
356 responses than jumbled scenes. This effect was comparable across the two tasks and
357 most pronounced in scene-selective ROIs. PPA was the only region that additionally
358 showed a modulation by task, with significantly stronger responses when participants
359 were categorizing the scenes, compared to when they were categorizing the objects
360 within them. For illustration purposes, ROI masks are shown on the right hemisphere
361 of a standard-space template using MRICroGL (Li et al., 2016); the displayed results
362 are averaged across ROIs in both hemispheres. Error bars represent standard errors
363 of the mean.

364

365 **3.2 Coherent scene structure enhances task-relevant object information in** 366 **multi-voxel response patterns**

367 To understand how the coherent spatial structure of the scene impacts cortical object
368 processing, we performed a cross-correlation multivariate pattern analysis (MVPA).
369 In this analysis, we correlated the multi-voxel response patterns evoked by objects
370 embedded in scenes (from the object and scene tasks) with the patterns evoked by
371 the objects in isolation (from the benchmark localizer) (Figure 3a). This approach
372 allowed us to quantify how “person-like” or “car-like” the cortical representation of
373 each of the scenes was, thereby isolating the amount of object information present
374 in visual cortex. When object information is operationalized in this way, it can be
375 separated from differences in the scene context (as in the benchmark localizer no

376 scene context is presented) and task-related differences (as in the benchmark
377 localizer participants perform a different task).

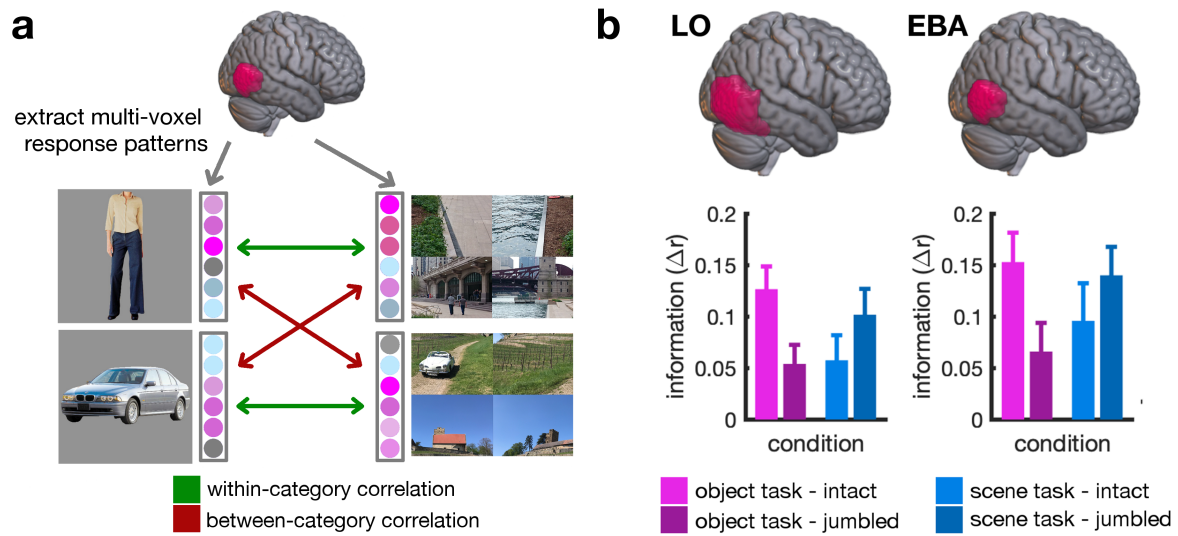
378 To quantify object information, we computed a cross-correlation measure by
379 subtracting correlations between different categories (e.g., person in isolation and car
380 within a scene) from correlations between the same categories (e.g., person in
381 isolation and person within a scene) (Figure 3a). This measure was computed
382 separately for each the object and scene tasks, the intact and jumbled scenes, and
383 all ROIs.

384 To test whether multi-voxel response patterns contained any information at all
385 about the object contained in the scenes, we first averaged the cross-correlation
386 measure across all conditions. We then tested whether the average category
387 information was significantly different from zero, separately for each ROI. As
388 expected, people and cars could be reliably discriminated from response patterns in
389 the object-selective LO, $t(24)=7.56$, $p<.001$, and body-selective EBA, $t(24)=8.00$,
390 $p<.001$, but not from response patterns in V1, $t(24)=0.80$, $p=.43$, or scene-selective
391 OPA, $t(24)=0.49$, $p=.63$, and PPA, $t(24)=0.70$, $p=.49$.

392 Given that we only found robust object information in LO and EBA, we only
393 performed further analyses for these two regions (Figure 3b). Data were again
394 analyzed in a 2x2 ANOVA with factors scene structure (intact vs jumbled) and task
395 (object task vs scene task), separately for LO and EBA.

396 When analyzing the amount of object information contained in LO response
397 patterns, we found a significant interaction between task and scene structure,
398 $F(1,24)=5.63$, $p=.026$: When participants performed the object task, object
399 information in LO was more pronounced for objects embedded in intact compared to
400 jumbled scenes, $t(24)=2.65$, $p=.014$. This effect was absent when participants
401 performed the scene task, $t(24)=1.22$, $p=.24$. A similar interaction effect was found in
402 the EBA, $F(1,24)=5.19$, $p=.032$: Object information was again enhanced for intact
403 scenes during the object task, $t(24)=2.30$, $p=.030$, but not during the scene task,
404 $t(24)=0.92$, $p=.37$. These results demonstrate that coherent scene structure indeed
405 enhances object representations in visual cortex. However, this enhancement
406 depends on the behavioral relevance of the object: When scene category, rather than
407 object category, was task-relevant, no such enhancement was observed.

408



409

410 **Figure 3. Cross-correlation MVPA logic and results.** **a)** To measure object
411 discriminability, we extracted multi-voxel response patterns for each ROI, separately
412 for objects in isolation (from the benchmark localizer) and objects appearing within
413 the scenes (from the main experiment). We then computed within- and between-
414 category correlations. By subtracting the between-category from the within-category
415 correlations, we obtained an index of category information (Δr). **b)** In both LO and
416 EBA, category information was significantly higher for objects that were embedded in
417 intact scenes than for objects embedded in jumbled scenes. However, this was only
418 true when participants performed the object task; when they performed the scene
419 task, no significant difference in object category information was observed when
420 comparing intact and jumbled scenes. For illustration purposes, ROI masks are shown
421 on the right hemisphere of a standard-space template using MRIcroGL (Li et al., 2016);
422 the displayed results are averaged across ROIs in both hemispheres. Error bars
423 represent standard errors of the mean.

424

425

426 **4 Discussion**

427

428 ***4.1 Coherent scene structure facilitates task-relevant object processing***

429 In this study, we shed light on neural object processing in situations where the object
430 is either embedded within a coherent, intact scene or an incoherent, jumbled scene.
431 Consistent with classical studies (Biederman, 1972; Biederman et al., 1973, 1974),
432 our participants were more accurate and faster in perceiving intact, compared to
433 jumbled scenes, both when performing an object categorization task and a scene
434 categorization task. Consistent with our own recent fMRI work (Kaiser et al., 2020a),
435 intact scenes yielded stronger neural responses than jumbled scenes across high-
436 level visual cortex. Importantly, our current results show that scene structure also
437 matters when it comes to the neural representation of objects within the scene: When
438 analyzing the amount of object information contained in multi-voxel response
439 patterns in object and body-selective visual cortex, we found an enhancement of
440 object information when the objects were embedded within intact scenes, compared
441 to jumbled scenes. Critically, this enhancement only emerged in the object
442 categorization task, suggesting that coherent scene structure facilitates the
443 extraction of object information only when the objects are relevant for current
444 behavioral goals.

445

446 ***4.2 Interactions between object and scene processing are mediated by scene*** 447 ***structure***

448 Our findings support the view that the scene and object processing pathways are not
449 functionally separate, but that scene information can aid the extraction of object
450 information (Brandmann & Peelen, 2017). Theories of contextual facilitation propose
451 that scene structure is analyzed rapidly, potentially based on coarse low-spatial
452 frequency information (Bar, 2004; Bar et al., 2006). This idea is consistent with the
453 observation that an initial representation of scene meaning – the scene’s “gist” – can
454 be extracted from just a single glance (Greene & Oliva, 2009; Oliva & Torralba, 2006,
455 2007). Contextual facilitation theories argue that detailed object analysis is facilitated
456 by this more readily available information about scene gist (Bar, 2004; Hochstein &

457 Ahissar, 2002). Informing object analysis through the analysis of coarse scene
458 properties may be particularly useful when perception is challenged by the presence
459 of many distracter items and limited visual exposure. Probing perception with such a
460 challenging task, our study shows that the cross-facilitation between object and
461 scene processing is mediated by the scene's structural coherence: When the analysis
462 of scene gist is disrupted by jumbling the scene, contextual information cannot
463 amplify object processing in the same way as it can for intact scenes.

464 The enhanced extraction of object information from the intact scenes suggests
465 that useful information about scene gist is extracted less efficiently from the jumbled
466 scenes. Indeed, the rapid analysis of scene gist depends on our priors about typical
467 scene composition (Csathó et al., 2015; Greene et al., 2015). Neuroimaging studies
468 suggest that the cortical scene processing network is tuned to these priors (Kaiser et
469 al., 2020a; Torralbo et al., 2013), and that the early extraction of properties like the
470 scene's basic-level category depends on the structural coherence of the scene
471 (Kaiser et al., 2020b). Jumbling is a strong manipulation in the sense that it disrupts
472 multiple aspects of the scene's spatial coherence at the same time: it disrupts the
473 spatial positioning of individual pieces of information in visual space (Kaiser et al.,
474 2018; Mannion, 2015), the positioning of objects relative to each other (Kaiser et al.,
475 2019; Kaiser & Peelen, 2018), as well as the typical geometry of the scene (Dillon et
476 al., 2018; Spelke & Lee, 2012). Future research is needed to disentangle these
477 different factors, and how much they each contribute to the facilitation of object
478 representation.

479 Although jumbling is a strong manipulation that conflates multiple factors of
480 scene structure, it preserves critical characteristics of the objects: First, the objects
481 remain completely unaltered across the intact and jumbled scenes. Second, the
482 objects' absolute positions in visual space were matched across the intact and
483 jumbled scenes. Finally, each object's local visual context remains constant across
484 the intact and jumbled scenes. These properties allow us to attribute differences in
485 object representations to facilitates effects from cortical scene analysis: If the visual
486 brain would not take global scene context into account and would only analyze the
487 objects in their local visual surroundings, our paradigm should yield comparable
488 results for structurally coherent, intact scenes and incoherent, jumbled scenes.

489

490 **4.2 Attention mediates contextual facilitation effects**

491 Unlike task-relevant objects, task-irrelevant objects were not processed differently as
492 a function of scene coherence. This finding shows that contextual facilitation of object
493 processing is not an automatic process. On the contrary, interactions between the
494 object and scene processing systems seem to be mediated by attention. This
495 observation fits well with previous results from studies on object detection in natural
496 scenes. Compared to task-relevant objects, multi-voxel response patterns in visual
497 cortex contain far less information about unattended objects (Peelen et al., 2009;
498 Peelen & Kastner, 2011). Further, MEG decoding results suggest strong differences
499 in the representation of attended and unattended object categories (Kaiser et al.,
500 2016): Particularly at early stages of processing, within the first 200ms after stimulus
501 onset, the category of unattended objects is represented less accurately. Beyond the
502 visual brain, differences in task demands also affect more widespread activations
503 across the cortex (Cukur et al., 2013; Harel et al., 2014; Hebart et al., 2018; Nastase
504 et al., 2017), potentially causing substantial task-related changes in processing
505 dynamics. One such change may be an alteration of the crosstalk between
506 representations in different visual domains. Our data indeed suggests that the
507 exchange of information between the object and scene processing pathways is not
508 mandatory, but rather constitutes an adaptive mechanism for improving task
509 performance. Under this view, interactions between the scene and object processing
510 pathways may be specifically “switched on” when objects are part of current
511 attentional templates (Battistoni et al., 2017; Peelen & Kastner, 2011). The specific
512 mechanism underlying this adaptive control of the crosstalk between scene and
513 object processing needs further investigation.

514 How does the apparent importance of attention tie in with previous studies that
515 reported a cross-facilitation between the object and scene-processing systems
516 (Brandmann & Peelen, 2017, 2019)? While these studies did not use object
517 categorization tasks, they still explicitly asked participants to attend to the objects
518 appearing within the scene (either by asking them to memorize them or through one-
519 back tasks). In our scene categorization task, the situation was entirely different, as
520 the objects were completely irrelevant for solving the task. In fact, this orthogonality

521 of object and scene category in our design may have introduced an active
522 suppression of object information when participants performed the scene
523 categorization task. Previous studies suggest that task-irrelevant distracter objects
524 can be suppressed effectively and quickly (Seidl et al., 2012; Hickey et al., 2019).
525 During the scene task, we indeed found numerically better object representations for
526 jumbled scenes. This tentative reversal of the facilitation effect could hint at a more
527 efficient suppression of object information when the object is embedded in a
528 structurally coherent scene. Although largely speculative at this point, this assertion
529 could be tested in future experiments that include additional conditions in which the
530 scene and the objects are similarly task-relevant.

531

532 **4.4 Conclusion**

533 In conclusion, our results show that the object and scene processing pathways can
534 interact to facilitate the processing of task-relevant object information embedded in
535 coherent scenes. However, such interactions are not mandatory. They rather seem
536 to be guided by current behavioral goals. Our findings therefore suggest that the
537 visual brain adaptively exploits coherent scene context to resolve object perception
538 in challenging real-world situations.

539

540 **Acknowledgements**

541

542 We thank Sina Schwarze for helping with the stimulus collection. D.K. and R.M.C. are
543 supported by Deutsche Forschungsgemeinschaft (DFG) grants (KA4683/2-1,
544 CI241/1-1, CI241/3-1). G.H. is supported by a PhD fellowship of the Einstein Center
545 for Neurosciences Berlin. R.M.C. is supported by a European Research Council
546 Starting Grant (ERC-2018-StG 803370).

547

548 **References**

549

550 Bar M. (2004). Visual objects in context. *Nat Neurosci*, 5, 617-629.

551 Bar M, Kassam KS, Ghuman AS, Boshyan J, Schmid AM, Dale AM, Hämäläinen MS,
552 Marinkovic K, Schacter DL, Rosen BR, Halgren E. (2006). Top-down facilitation
553 of visual recognition. *Proc Natl Acad Sci USA*, 103, 449-454.

554 Battistoni E, Stein T, Peelen MV. (2017). Preparatory attention in visual cortex. *Ann N*
555 *Y Acad Sci*, 1396, 92-107.

556 Biederman I. (1972). Perceiving real-world scenes. *Science*, 177, 77-80.

557 Biederman I, Glass AL, Stacy EW. (1973). Searching for objects in real-world scenes.
558 *J Exp Psychol*, 97, 22-27.

559 Biederman I, Rabinowitz JC, Glass AL, Stacy EW. (1974). On the information
560 extracted from a glance at a scene. *J Exp Psychol*, 103, 597-600.

561 Brandman T, Peelen MV. (2017). Interaction between scene and object processing
562 revealed by human fMRI and MEG decoding. *J Neurosci*, 37, 7700-7710.

563 Brandman T, Peelen MV. (2019). Signposts in the fog: objects facilitate scene
564 representations in left scene-selective cortex. *J Cogn Neurosci*, 31, 390-400.

565 Brainard DH. (1997). The psychophysics toolbox. *Spat Vis*, 10, 433-436.

566 Cauchoix M, Barragan-Jason G, Serre T, Barbeau EJ. (2014). The neural dynamics of
567 face detection in the wild revealed by MVPA. *J Neurosci*, 34, 846-854.

568 Csathó Á, van der Linden D, Gács B. (2015). Natural scene recognition with increasing
569 time-on-task: the role of typicality and global image properties. *Q J Exp*
570 *Psychol*, 68, 814-828.

571 Cukur T, Nishimoto S, Huth AG, Gallant JL. (2013). Attention during natural vision
572 warps semantic representation across the human brain. *Nat Neurosci* 16, 763-
573 770.

574 Dillon MR, Persichetti AS, Spelke ES, Dilks DD. (2018). Places in the brain: bridging
575 layout and object geometry in scene-selective cortex. *Cereb Cortex*, 28, 2365-
576 2374.

- 577 Greene MR, Botros AP, Beck DM, Fei-Fei L. (2015). What you see is what you expect:
578 rapid scene understanding benefits from prior experience. *Atten Percept*
579 *Psychophys*, 77, 1239-1251.
- 580 Greene MR, Oliva A. (2009). Recognition of natural scenes from global properties:
581 seeing the forest without representing the trees. *Cogn Psychol*, 58, 137-176.
- 582 Harel A, Kravitz DJ, Baker CI. (2014). Task context impacts visual object processin
583 differentially across the cortex. *Proc Natl Acad Sci USA*, 111, 962-971.
- 584 Haxby JV, Gobbini IM, Furey ML, Ishai A, Schouten JL, Pietrini P. (2001). Distributed
585 and overlapping representations of faces and objects in ventral temporal
586 cortex. *Science*, 293, 2425-2430.
- 587 Hebart MN, Bankson BB, Harel A, Baker CI, Cichy RM. (2018). *eLife*, 7, e32816.
- 588 Hickey C, Pollicino D, Bertazzoli G, Barbaro L. (2019). Ultrafast object detection in
589 naturalistic vision relies on ultrafast distractor suppression. *J Cogn Neurosci*,
590 31, 1563-1572.
- 591 Hochstein S, Ahissar M. (2002). View from the top: hierarchies and reverse hierarchies
592 in the visual system. *Neuron*, 36, 791-804.
- 593 Julian JB, Fedorenko E, Webster J, Kanwisher N. (2012). An algorithmic method for
594 functionally defining regions of interest in the ventral visual pathway.
595 *Neuroimage*, 60, 2357-2364.
- 596 Kaiser D, Cichy RM. (2018). Typical visual-field locations enhance processing in
597 object-selective channels of human occipital cortex. *J Neurophysiol*, 120, 848-
598 853.
- 599 Kaiser D, Häberle G, Cichy RM. (2020a). Cortical sensitivity to natural scene structure.
600 *Hum Brain Mapp*, 41, 1286-1295.
- 601 Kaiser D, Häberle G, Cichy RM. (2020b). Real-world structure facilitates the rapid
602 emergence of scene category information in visual brain signals. *J*
603 *Neurophysiol*, 124, 145-151.
- 604 Kaiser D, Oosterhof NN, Peelen MV. (2016). The neural dynamics of attentional
605 selection in natural scenes. *J Neurosci*, 36, 10522-10528.
- 606 Kaiser D, Quek GL, Cichy RM, Peelen MV. (2019). Object vision in a structured world.
607 *Trends Cogn Sci*, 23, 672-685.

- 608 Kaiser D, Peelen MV. (2018). Transformation from independent to integrative coding
609 of multi-object arrangements in human visual cortex. *Neuroimage*, 169, 334-
610 341.
- 611 Li FF, VanRullen R, Koch C, Perona P. (2002). Rapid natural scene categorization in
612 the near absence of attention. *Proc Natl Acad Sci USA*, 99, 9596-9601.
- 613 Li X, Morgan PS, Ashburner J, Smith J, Rorden C. (2016). The first step for
614 neuroimaging data analysis: DICOM to NIfTI conversion. *J Neurosci Methods*,
615 264, 47-56.
- 616 Mannion DJ. (2015). Sensitivity to the visual field origin of natural image patches in
617 human low-level visual cortex. *PeerJ*, 3, e1038.
- 618 Nastase SA, Connolly AC, Oosterhof NN, Halchenko YO, Guntupalli JS, Visconti di
619 Oleggio Castello M, Gobbini I, Haxby JV. (2017). Attention selectively reshapes
620 the geometry of distributed semantic representation. *Cereb Cortex*, 27, 4277-
621 4291.
- 622 Oliva A, Torralba A. (2006). Building the gist of a scene: the role of global image
623 features in recognition. *Prog Brain Res*, 155, 23-36.
- 624 Oliva A, Torralba A. (2007). The role of context in object recognition. *Trends Cogn Sci*,
625 11, 520-527.
- 626 Oosterhof NN, Connolly AC, Haxby JV. (2016). CoSMoMMPA: Multi-modal
627 multivariate pattern analysis of neuroimaging data in Matlab/GNU Octave.
628 *Front Neuroinform*, 10, 20.
- 629 Peelen MV, Fei-Fei L, Kastner S. (2009). Neural mechanisms of rapid natural scene
630 categorization in human visual cortex. *Nature*, 460, 94-97.
- 631 Peelen MV, Kastner S. (2011). A neural basis for real-world visual search in human
632 occipitotemporal cortex. *Proc Natl Acad Sci USA*, 108, 12125-12130.
- 633 Potter MC. (1975). Meaning in visual search. *Science*, 187, 965-966.
- 634 Potter MC. (2012). Recognition and memory for briefly presented scenes. *Front*
635 *Psychol*, 3, 32.
- 636 Seidl KN, Peelen MV, Kastner S. (2012). Neural evidence for distracter suppression
637 during visual search in real-world scenes. *J Neurosci*, 32, 11812-11819.
- 638 Shahdloo M, Celik E, Cukur T. (2020). Biased competition in semantic representation
639 during natural visual search. *Neuroimage*, 216, 116383.

- 640 Spelke ES, Lee SA. (2012). Core systems of geometry in animal minds. *Philos Trans*
641 *R Soc Lond B Biol Sci*, 367, 2784-2793.
- 642 Thorpe S, Fize D, Marlot C. (1996). Speed of processing in the human visual system.
643 *Nature*, 381, 520-522.
- 644 Torralba A, Oliva A, Castelhana MS, Henderson JM. (2006). Contextual guidance of
645 eye movements and attention in real-world scenes: the role of global features
646 in object search. *Psychol Rev*, 113, 766-786.
- 647 Torralbo A, Walther DB, Chai B, Caddigan E, Fei-Fei L, Beck DM. (2013). Good
648 exemplars of natural scene categories elicit clearer patterns than bad
649 exemplars but not greater BOLD activity. *Plos One*, 8, e58594.
- 650 VanRullen R, Thorpe SJ. (2001). The time course of visual processing: from early
651 perception to decision-making. *J Cogn Neurosci*, 13, 454-461.
- 652 Võ MLH, Boettcher SEP, Draschkow D. (2019). Reading scenes: How scene grammar
653 guides attention and aids perception in real-world environments. *Curr Opin*
654 *Psychol*, 29, 205-210.
- 655 Wang L, Mruczek RE, Arcaro MJ, Kastner S. (2015). Probabilistic maps of visual
656 topography in human cortex. *Cereb Cortex*, 25, 3911-3931.
- 657 Wolfe JM, Võ ML-H, Evans KK, Greene MR. (2011). Visual search in scenes involves
658 selective and nonselective pathways. *Trends Cogn Sci*, 15, 77-84.