

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

# **Rapid contextualization of fragmented scene information in the human visual system**

Daniel Kaiser<sup>1</sup>, Gabriele Inciuraite<sup>2</sup>, Radoslaw M. Cichy<sup>2,3,4</sup>

<sup>1</sup>*Department of Psychology, University of York, York, UK*

<sup>2</sup>*Department of Education and Psychology, Freie Universität Berlin, Berlin, Germany*

<sup>3</sup>*Berlin School of Mind and Brain, Humboldt-Universität Berlin, Berlin, Germany*

<sup>4</sup>*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

21

22

## Abstract

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

Real-world environments are extremely rich in visual information. At any given moment in time, only a fraction of this information is available to the eyes and the brain, rendering naturalistic vision a collection of incomplete snapshots. Previous research suggests that in order to successfully contextualize this fragmented information, the visual system sorts inputs according to spatial schemata, that is knowledge about the typical composition of the visual world. Here, we used a large set of different natural scene fragments to investigate whether this sorting mechanism can operate across the diverse visual environments encountered during real-world vision. We recorded brain activity using electroencephalography (EEG) while participants viewed incomplete scene fragments at fixation. Using representational similarity analysis on the EEG data, we tracked the fragments' cortical representations across time. We found that the fragments' typical vertical location within the environment (top or bottom) predicted their cortical representations, indexing a sorting of information according to spatial schemata. The fragments' cortical representations were most strongly organized by their vertical location at around 200ms after image onset, suggesting rapid perceptual sorting of information according to spatial schemata. Our results identify a cortical sorting process that allows the visual system to efficiently contextualize fragmented inputs. By demonstrating this sorting across a wide range of visually diverse scenes, our study suggests a contextualization mechanism suitable for complex and variable real-world environments.

45

46           *Keywords:* visual perception, scene representation, spatial schema, EEG,

47 representational similarity analysis

48

49

50 **Rapid contextualization of fragmented scene information in the human visual**  
51 **system**

52

53 The visual world around us is structured in predictable ways at multiple levels  
54 (Kaiser et al., 2019a). Natural scenes are characterized by typical distributions of low-  
55 and mid-level visual features (Geisler, 2008; Oliva & Torralba, 2003; Purves et al.,  
56 2011), as well as typical arrangements of high-level contents across the scene (Bar,  
57 2004; Kaiser et al., 2019a; Torralba et al., 2006; Oliva & Torralba, 2007; Vo et al., 2019;  
58 Wolfe et al., 2011). The visual system has adapted to this structure: when multiple  
59 scene elements are arranged in a typical way, cortical processing is more efficient  
60 (Abassi & Papeo, 2019; Baldassano et al., 2016; Bilalic et al., 2019; Gronau et al.,  
61 2008; Kim & Biederman, 2011; Kim et al., 2011; Kaiser et al., 2014, 2019b; Kaiser &  
62 Peelen, 2018; Roberts & Humphreys, 2010). Such results suggest that when multiple  
63 scene elements need to be processed concurrently, cortical processing is strongly  
64 tuned to the typical composition of these elements.

65 In real-life situations, however, we usually do not have access to detailed visual  
66 information about all scene elements at once. Instead, visual inputs are fragmented,  
67 and only incomplete snapshots of the world are available for visual analysis at any  
68 given moment in time. How does the brain assemble a coherent image of the world  
69 from such fragmented inputs? To solve this problem, the visual system may draw  
70 from internal representations of typical scene structure – scene schemata (Mandler,  
71 1984; Minsky, 1975, Rumelhart, 1980) – in order to contextualize the fragmented  
72 inputs with which it is faced. More specifically, schemata may be used to match

73 fragmented visual inputs with their place in the schema: as a result, fragmented visual  
74 information should be sorted according to its typical location within the environment.  
75 This sorting may help to efficiently contextualize visual inputs.

76 A recent study showed that incomplete inputs – fragments of natural scenes –  
77 are indeed sorted according to their typical location in real-world environments  
78 (Kaiser et al., 2019c): In the occipital place area and after 200ms of vision,  
79 representations of scene fragments were organized by their typical vertical location  
80 in the world. For instance, fragments that typically appear in the upper part of a scene  
81 (e.g., a house roof or the ceiling of a room) were represented more similarly to each  
82 other than to fragments that typically appear in the lower part of a scene (e.g., a lawn  
83 or the room's floor). No such organization was found for the fragments' horizontal  
84 location, for which clear schemata are missing (Mandler & Parker, 1976).

85 As a critical limitation, our previous study (Kaiser et al., 2019c) only comprised  
86 six different scenes. However, for this mechanism to be useful in the real world, it has  
87 to operate across huge amounts of vastly different scenes encountered in our  
88 everyday lives. We therefore set out to replicate our findings across a larger and more  
89 diverse set of scene images. Here, we used a set of 210 indoor and outdoor scenes,  
90 which we split into 4 position-specific fragments each, yielding 840 unique scene  
91 fragments (Figure 1a). During an EEG experiment, participants viewed each fragment  
92 centrally and in isolation (Figure 1b). Using representational similarity analysis (RSA;  
93 Kriegeskorte et al., 2018), we then tracked the fragments' cortical representations  
94 across time. As the key result, we found that most prominently after 200ms of visual  
95 processing, the fragments' cortical representations were organized by their vertical  
96 location within the full scene. We conclude that the visual system uses scene

- 97 schemata to sort inputs according to their typical location within the environment,
- 98 supporting the contextualization of fragmented visual information.

99

100

## Materials and Methods

101

### 102 **Participants**

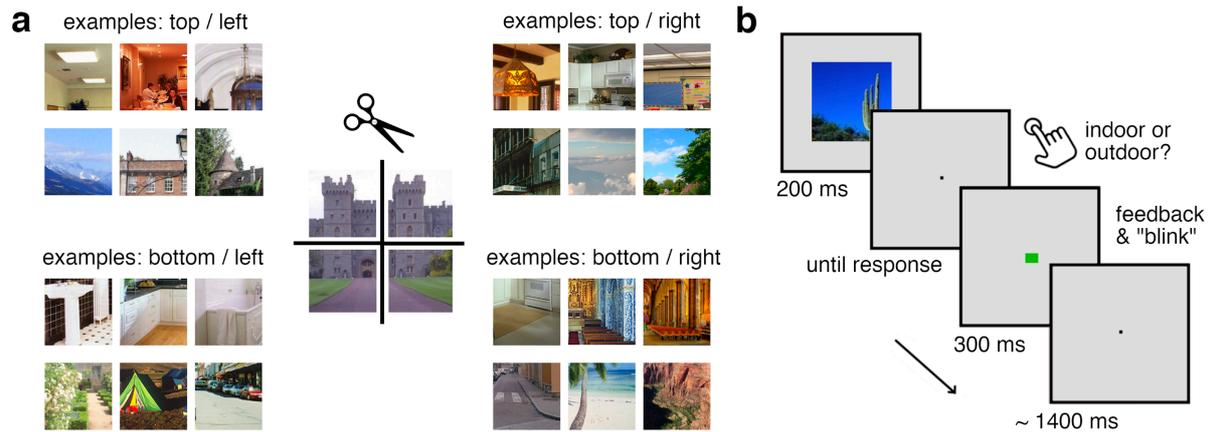
103 Twenty healthy adults (mean age 27.3,  $SD = 4.6$ ; 12 female) participated in the  
104 study. The sample size was identical to the sample size of our previous EEG study  
105 (Kaiser et al., 2019c). All participants had normal or corrected-to-normal vision.  
106 Participants provided informed consent and received monetary reimbursement or  
107 course credits. All procedures were approved by the ethical committee of Freie  
108 Universität Berlin and were in accordance with the Declaration of Helsinki.

109

### 110 **Stimuli**

111 Stimuli were 210 natural scene photographs, taken from an online resource  
112 (Konkle et al., 2010). Half of the stimuli depicted outdoor scenes (bridges, camping  
113 sites, historical buildings, houses, natural environments, streets, and waterfronts) and  
114 half of the stimuli depicted indoor scenes (bathrooms, bedrooms, churches,  
115 classrooms, dining rooms, kitchens, and living rooms). To create position-specific  
116 fragments, each scene was split along the vertical and horizontal axes (Figure 1a),  
117 yielding four fragments of equal size for each scene and 840 fragments in total. The  
118 full stimulus set is available on OSF ([doi.org/10.17605/OSF.IO/D7P8G](https://doi.org/10.17605/OSF.IO/D7P8G)). During the  
119 experiment, these fragments were presented individually and in the center of the  
120 screen ( $5.5^\circ$  by  $5.5^\circ$  visual angle). Participants were not shown the full scene images  
121 prior to the experiment.

122



123

124 *Figure 1.* Stimuli and paradigm. a) To mimic the fragmented nature of natural visual  
125 inputs, we used a set of 210 widely varying indoor and outdoor scene photographs  
126 that each were split into four equally sized fragments (top/left, top/right, bottom/left,  
127 bottom/right). The panel shows representative examples from indoor (upper rows)  
128 and outdoor (lower rows) scenes. b) During the EEG experiment, participants viewed  
129 the individual fragments in the center of the screen while performing an  
130 indoor/outdoor categorization task.

131

## 132 **Experimental paradigm**

133 During the experiment, participants briefly viewed the individual scene  
134 fragments, all presented in the same central location (Figure 1b). Each of the 840  
135 fragments was shown twice during the experiment, yielding 1,680 trials. Trial order  
136 was randomized separately for the first and second half of trials, so that every  
137 fragment appeared once in the first half of the experiment and once in the second  
138 half. On each trial, a single fragment appeared for 200ms and participants were asked  
139 to categorize the fragment as either stemming from an indoor scene or an outdoor  
140 scene using two keyboard buttons. After every response, the fixation cross turned  
141 red or green for 300ms, indicating response correctness. Trials were separated by an  
142 inter-trial interval varying randomly between 1,300ms and 1,500ms. Participants  
143 performed the categorization task well (93% correct responses, SE=1%; 769ms

144 average response time, SE=36ms), with no differences in accuracy or response time  
145 between fragments stemming from the top versus the bottom or from the left versus  
146 the right (all  $t[19] < 1.89$ ,  $p > 0.07$ ). Further, participants were instructed to maintain  
147 central fixation throughout the experiment, and to only blink after they had given a  
148 response. Stimulus presentation was controlled using the Psychtoolbox (Brainard,  
149 1997).

150

### 151 **EEG recording and preprocessing**

152 EEG signals were recorded using an EASYCAP 64-electrode system and a  
153 Brainvision actiCHamp amplifier. Electrodes were arranged in accordance with the  
154 10-10 system. EEG data were recorded at 1000Hz sampling rate and filtered online  
155 between 0.03Hz and 100Hz. All electrodes were referenced online to the Fz electrode.  
156 Offline preprocessing was performed using FieldTrip (Oostenveld et al., 2011). EEG  
157 data were epoched from -200ms to 800ms relative to stimulus onset, and baseline-  
158 corrected by subtracting the mean pre-stimulus signal. Channels and trials containing  
159 excessive noise were removed based on visual inspection. Blinks and eye movement  
160 artifacts were removed using independent component analysis and visual inspection  
161 of the resulting components. The epoched data were downsampled to 200Hz.

162

### 163 **Measuring representational similarity**

164 To track the representations of individual fragments across time, we used  
165 representational similarity analysis (RSA; Kriegeskorte et al., 2008). First, we created  
166 neural representational dissimilarity matrices (RDMs) for each time point in the EEG  
167 epochs (5ms resolution), reflecting the pairwise dissimilarity of the fragments' brain

168 representations. Second, we modeled the organization of the neural RDMs in a  
169 regression approach (Proklova et al., 2016, 2019), which allowed us to track when  
170 representations are explained by the fragments' vertical and horizontal location within  
171 the full scene as well as the scene's category.

172 To construct neural RDMs, we computed the pairwise dissimilarity of all  
173 fragments at each time point using the CoSMoMvPA toolbox (Oosterhof et al., 2016).  
174 This analysis was done separately for each participant. For this, we used response  
175 patterns across 17 posterior electrodes (Kaiser et al., 2019c) in our EEG montage (O1,  
176 O2, Oz, PO3, PO4, PO7, PO8, POz, P1, P2, P3, P4, P5, P6, P7, P8, Pz). For each of  
177 the 840 fragments, we computed the response pattern to the fragment by averaging  
178 across the two repetitions of the fragment. If after trial rejection during preprocessing  
179 only one trial was left for a fragment, we used the response pattern from this one trial.  
180 If after preprocessing no trial was left for a fragment, this fragment was excluded from  
181 the analysis (i.e., removed from all RDMs of the respective participant). Neural  
182 dissimilarity was computed by correlating the response patterns to each individual  
183 fragment in a pairwise fashion and subtracting the resulting correlations from 1,  
184 yielding an index of neural dissimilarity (0: minimum dissimilarity, 2: maximum  
185 dissimilarity). Computing this index for each pairwise comparison of fragments, we  
186 obtained an 840-by-840 neural RDM for each time point.

187

### 188 **Modelling representational similarity**

189 To quantify how well the neural organization is explained by the fragments'  
190 vertical and horizontal location within the full scene and by the original scene's  
191 category (indoor vs. outdoor), we modeled the neural RDMs in a general linear model

192 (GLM) with three predictors: (1) a vertical location RDM, in which each pair of  
193 conditions is assigned either a value of 0, if the fragments stem from the same vertical  
194 location (e.g., both from the top), or the value 1, if they stem from different vertical  
195 locations (e.g., one from the top and one from the bottom), (2) a horizontal location  
196 RDM, in which each pair of conditions is assigned either a value of 0, if the fragments  
197 stem from the same horizontal location (e.g., both from the left), or the value 1, if they  
198 stem from different horizontal locations (e.g., one from the left and one from the right),  
199 and (3) a category RDM, in which each pair of conditions is assigned either a value of  
200 0, if the fragments stem from the same scene category (e.g., both from indoor  
201 scenes), or the value 1, if they stem from different categories (e.g., one from an indoor  
202 scene and one from an outdoor scene).

203 GLMs were constructed with the neural RDMs as the regression criterion and  
204 the vertical and horizontal location RDMs as well as the category RDM as predictors.  
205 For these GLMs, the neural RDMs and predictor RDMs were vectorized by selecting  
206 all lower off-diagonal elements – the rest of the entries, including the diagonal, were  
207 discarded. Values for the neural RDMs were z-scored. Estimating this GLM yielded  
208 three beta weights for each time point and participant. We subsequently tested these  
209 beta weights across participants against zero, which reveal whether the fragments’  
210 vertical location, horizontal location, and their category significantly explained the  
211 neural organization at each time point.

212 We additionally performed two control analyses. In the first control analysis,  
213 we assessed if the location-based organization in the fragments’ neural  
214 representations can be explained by simple visual features. To do so, we used the  
215 fragments’ similarity in image space to predict their neural organization. We

216 computed the fragments' image similarity by correlating their pixel values. We then  
217 constructed a pixel RDM that contained the fragments' pairwise dissimilarity (1-  
218 correlation) in pixel values. To control for image similarity, we ran a GLM in which we  
219 predicted neural RDMs as a function of the pixel RDM. In another, second GLM, we  
220 then modelled the residuals of the first model using the vertical location, horizontal  
221 location, and category predictors (see above). This two-stage approach allowed us  
222 to quantify how much vertical and horizontal location information as well as category  
223 information remained unaccounted for by the fragments' pixel similarity.

224 In the second control analysis, we aimed at eliminating visual and conceptual  
225 features that are common to either indoor or outdoor scenes (e.g., the top fragments  
226 from outdoor scenes often show blue skies). We thus constructed RDMs for  
227 horizontal and vertical location information which only contained comparisons across  
228 indoor and outdoor scenes. These RDMs were constructed in the same way as  
229 explained above, but now all comparisons within the same scene type (e.g.,  
230 comparisons of different indoor scene fragments) were removed. We then repeated  
231 the original GLM analysis (see above) with these restricted RDMs, allowing us to see  
232 if the organization according to vertical location persists when only between-category  
233 comparisons are considered.

234

### 235 **Statistical testing**

236 To test whether GLM beta weights were significantly greater than zero, we used a  
237 threshold-free cluster enhancement procedure (Smith and Nichols, 2009) and  
238 multiple-comparison correction based on a sign-permutation test (with null  
239 distributions created from 10,000 bootstrapping iterations), as implemented in

240 CoSMoMVPA (Oosterhof et al., 2016). The resulting statistical maps were thresholded  
241 at  $z > 1.64$  (i.e.,  $p_{corr} < .05$ ). For all peaks in the time series, we additionally report results  
242 of conventional one-sided t-tests against zero. To estimate the robustness of peak  
243 latencies we performed a bootstrapping analysis. In this analysis, we created 1000  
244 samples of 20 randomly chosen datasets each (with possible repetitions). For each  
245 random sample, we computed the peak latency (i.e., the highest beta estimate in the  
246 average time course). We then computed a confidence interval ( $ci$ ) by selecting the  
247 central 95% of the distribution across the 1000 random samples. Given the clear two-  
248 peak structure in vertical location information, we performed the bootstrapping  
249 analysis separately for the early and late peaks, by splitting the data for each random  
250 sample along the minimum beta value between 100ms and 200ms.

251

## 252 **Data Availability**

253 Data are publicly available on OSF ([doi.org/10.17605/OSF.IO/D7P8G](https://doi.org/10.17605/OSF.IO/D7P8G)).

254

255

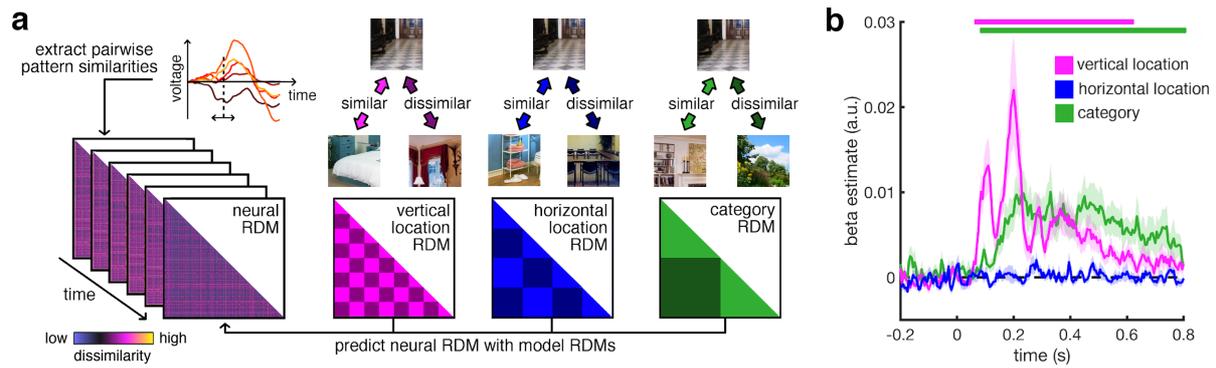
## Results

256

257 To model the fragments' cortical representations across time, we ran a GLM  
258 analysis with three predictors capturing the fragments' vertical and horizontal  
259 locations within the full scene and the full scene's category (Figure 2a). This analysis  
260 revealed three key insights. First, the fragments' cortical organization was explained  
261 by their vertical location within the scene (Figure 2b), from 70ms to 615ms (peaking  
262 at 110ms, peak  $t[19]=4.07$ ,  $p<0.001$ ,  $p_{corr}<0.05$ ,  $ci=[85ms, 115ms]$ , and at 200ms,  
263 peak  $t[19]=3.47$ ,  $p=0.001$ ,  $p_{corr}<0.05$ ,  $ci=[190ms, 210ms]$ ). This suggests that  
264 fragmented scene information is sorted by its typical origin within the visual world.  
265 Second, the fragments' horizontal location did not significantly predict their neural  
266 organization, suggesting that the more rigid real-world location along the vertical axis  
267 is more strongly reflected in cortical signals. Third, the fragments' category (i.e.,  
268 whether a fragment stems from an indoor or an outdoor scene) was also reflected in  
269 their neural organization, from 90ms to 800ms (peaking at 330ms, peak  $t[19]=3.06$ ,  
270  $p=0.007$ ,  $p_{corr}<0.05$ ,  $ci=[210ms, 635ms]$ ). This finding supports previous studies  
271 showing that scene category can be rapidly decoded from EEG signals (Dima et al.,  
272 2018; Kaiser et al., 2019c; Lowe et al., 2019).

273 Together, these results suggest that fragmented visual information is  
274 organized with respect to its typical vertical location in the world. To test how flexible  
275 this organization is with respect to visual and conceptual properties of individual  
276 scenes, we additionally performed two control analyses.

277



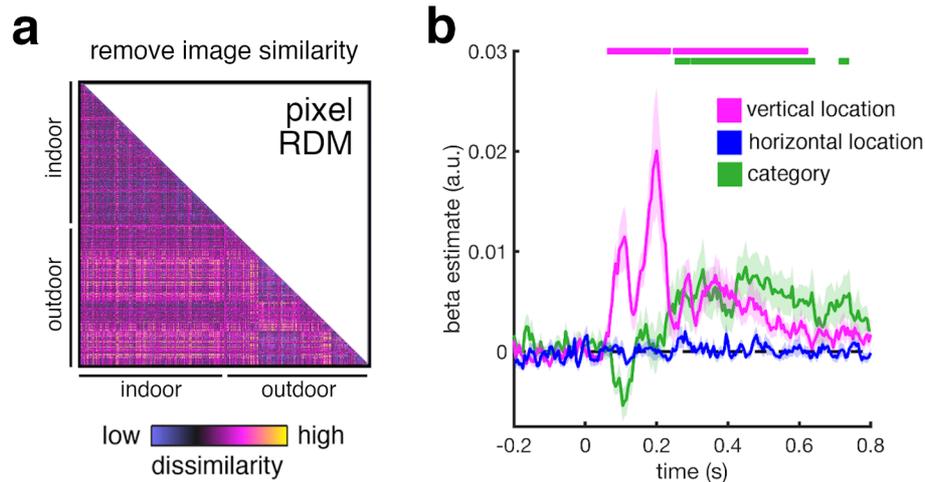
278

279 *Figure 2.* Analysis approach and main result. a) We first extracted neural RDMs from  
280 EEG signals in a time-resolved manner. That is, for each time point in an epoch we  
281 correlated the response patterns evoked by each one fragment with the response  
282 pattern evoked by each other fragment, yielding an 840-by-840 matrix of pairwise  
283 neural dissimilarities. The neural RDMs were then modeled as a combination of three  
284 predictor RDMs that captured the fragments' dissimilarity in vertical location (e.g.  
285 different fragments from the same location were considered similar), horizontal  
286 location (e.g. different fragments from the same location were considered similar),  
287 and category (e.g. different fragments from the same category were considered  
288 similar). Estimating this model for each time point yielded three time courses of beta  
289 estimates, indicating how well the neural organization matched each of the predicted  
290 organizations. b) The fragments' vertical location (but not their horizontal location)  
291 predicted neural organization between 70ms and 615ms, suggesting a sorting of  
292 information according to typical real-world locations. Additionally, the fragments'  
293 category predicted their neural organization between 90ms and 800ms. Significance  
294 markers denote  $p_{corr} < 0.05$ . Shaded margins represent standard errors of the mean.

295

296 In the first control analysis, we tested whether simple image features can  
297 explain the vertical location organization in the neural data. To quantify such simple  
298 visual features, we constructed RDMs from the fragments' pixel values (Figure 3a)  
299 and removed these features from the neural data before performing the GLM analysis  
300 (see Materials and Methods).

301



302

303 *Figure 4. Controlling for image similarity.* a) We computed the fragments' low-level  
304 visual similarity by computing pairwise dissimilarities (1-correlation) between their  
305 pixel values, yielding an 840-by-840 element pixel RDM. We then removed the pixel  
306 RDM from the neural organization (see Materials and Methods) before performing the  
307 GLM analysis as above, i.e., modelling the fragments' neural organization as a  
308 function of their vertical and horizontal locations and their category. d) Removing the  
309 fragments' low-level visual similarity did not abolish vertical location information,  
310 which remained significant between 70ms and 615ms, suggesting that the sorting of  
311 fragments according to their vertical location in the world is flexible with regards to  
312 the low-level features. Notably, after removing the pixel RDM, category information  
313 emerged only after 260ms, suggesting that the more rapidly emerging vertical  
314 location information does not depend on image properties that distinguish the indoor  
315 and outdoor scenes. Significance markers denote  $p_{corr} < 0.05$ . Shaded margins  
316 represent standard errors of the mean.

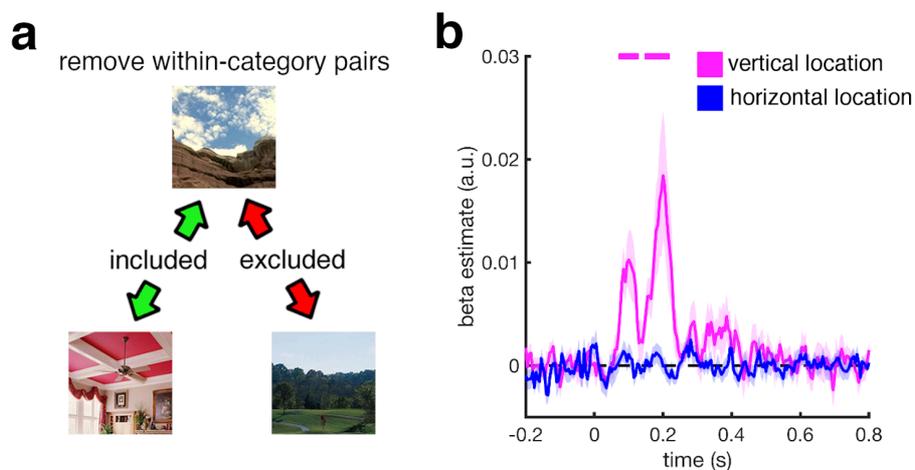
317

318 After removing the fragments' pixel similarity from the neural RDM, we still  
319 found significant vertical location information (Figure 3b), between 70ms and 615ms  
320 (peaking at 110ms, peak  $t[19]=3.74$ ,  $p < 0.001$ ,  $p_{corr} < 0.05$ ,  $ci=[85ms, 115ms]$ , and at  
321 200ms, peak  $t[19]=3.21$ ,  $p=0.002$ ,  $p_{corr} < 0.05$ ,  $ci=[190ms, 210ms]$ ), but no significant  
322 horizontal location information. This shows that controlling for simple visual features  
323 cannot explain the fragments' neural organization by vertical location. Category

324 information was still present in the neural data, but now emerged only later, between  
325 260ms and 640ms as well as between 720ms and 730ms (peaking at 450ms, peak  
326  $t[19]=3.23$ ,  $p=0.002$ ,  $p_{corr}<0.05$ ,  $ci=[260ms, 720ms]$ ). This result suggests that during  
327 early visual processing vertical location information remains prominent, even if during  
328 this time no category information can be observed in the neural data.

329 In the second control analysis, we more explicitly tested whether the sorting  
330 of fragmented information is independent of visual category information by testing if  
331 the sorting generalized across indoor and outdoor scenes. In this analysis, we  
332 restricted our models to comparisons between indoor and outdoor scenes (i.e., all  
333 comparisons within the same scene category were removed from all RDMs). The  
334 comparisons between indoor and outdoor scenes share fewer visual and conceptual  
335 properties than the comparisons of scenes from the same category (Figure 4a): For  
336 example, two fragments from the upper part of outdoor scenes often share both  
337 visual properties (e.g., they tend to be blue-colored) and conceptual content (e.g.,  
338 they tend to contain the same objects, such as clouds of tree tops).

339



340

341 *Figure 4.* Controlling for similarity within scene category. a) In this analysis, we  
342 removed all comparisons between the fragments of the same category from the

343 neural and predictor RDMs. This allowed us to control for visual and conceptual  
344 features shared by fragments stemming from the same location (e.g., fragments from  
345 the upper part of outdoor scenes often contain skies and clouds). d) Removing the  
346 same-category comparisons did not abolish vertical location information, which  
347 remained significant between 80ms and 210ms. This indicates that the sorting of  
348 fragments according to their vertical location in the world is flexible with regards to  
349 visual and conceptual attributes shared by the scenes of the same category.  
350 Significance markers denote  $p_{corr} < 0.05$ . Shaded margins represent standard errors of  
351 the mean.

352

353 This analysis revealed significant vertical location information (Figure 4b), from  
354 80ms to 120ms (peaking at 100ms, peak  $t[19]=3.63$ ,  $p < 0.001$ ,  $p_{corr} < 0.05$ ,  $ci=[80ms,$   
355  $120ms]$ ) and from 155ms to 210ms (peaking at 200ms, peak  $t[19]=2.98$ ,  $p=0.004$ ,  
356  $p_{corr} < 0.05$ ,  $ci=[190ms, 215ms]$ ), but no horizontal location information. This suggests  
357 that the cortical sorting of information according to the fragments' vertical location  
358 occurs similarly for visually and conceptually diverse scenes. Note that – as within-  
359 category comparisons were removed – no category information could be computed  
360 in this analysis.

361 Together, our results suggest that fragmented information is sorted according  
362 to its typical vertical location in the world, providing a mechanism for the  
363 contextualization of incomplete visual information. Even when controlling for visual  
364 and conceptual scene attributes, this mechanism can operate rapidly and most  
365 strongly determines the cortical organization after 200ms of visual processing.

366

367

368

## Discussion

369

370           During natural vision the brain is constantly faced with incomplete snapshots  
371 of the world from which it needs to infer the structure of the whole scene. Here, we  
372 show that in order to meet this challenge, the visual system rapidly contextualizes  
373 incoming information according to its typical place in the world: after 200ms of  
374 processing, fragmented scene information is sorted according to its real-world  
375 location. By using a large stimulus set (comprising 840 unique fragments) we provide  
376 compelling evidence that this mechanism supports spatial contextualization across  
377 diverse visual environments.

378           Which features allow the visual system to make such inferences about a  
379 fragment's typical position within the environment? Solidifying our previous results  
380 (Kaiser et al., 2019c), the strongest organization according to vertical location became  
381 apparent at around 200ms after onset. At this time higher-level scene attributes –  
382 such as the scene's clutter or openness (Cichy et al., 2017; Harel et al., 2016) – are  
383 analyzed, suggesting that the sorting of information according to real-world locations  
384 is determined by more complex scene properties, rather than low-level visual  
385 features. This is consistent with previous fMRI results which demonstrated a vertical  
386 location organization in the occipital place area, but not in early visual cortex (Kaiser  
387 et al., 2019c). However, we additionally found a very rapid onset of vertical location  
388 information with the first peak shortly after 100ms, which suggests that the visual  
389 features extracted during early visual analysis are also diagnostic of a fragment's  
390 typical location in the world. Such features could comprise particular distributions of

391 spatial frequency content or texture information (Dima et al., 2018; Groen et al., 2013,  
392 2017). Alternatively, these early effects could reflect the rapid analysis of scene  
393 geometry (Henriksson et al., 2019). Future studies need to isolate the contribution of  
394 different visual features to the sorting of visual information by real-world location at  
395 different processing times.

396 Previous research has demonstrated that the representation of individual  
397 naturalistic stimuli depends on whether their current position in the visual field  
398 matches their typical position in the world (Chan et al., 2010; de Haas et al., 2016;  
399 Mannion, 2015; Kaiser & Cichy, 2018; Kaiser et al., 2018). For instance, when face  
400 parts (e.g., an eye) or objects (e.g., a lamp) are presented in their typically experienced  
401 visual-field position (e.g., the upper visual field), they evoke more efficient cortical  
402 representations (de Haas et al., 2016; Kaiser & Cichy, 2018). This suggests that  
403 across diverse visual contents cortical representations of a stimulus are entwined with  
404 preferences for its typical location (Kaiser & Haselhuhn, 2017; Kaiser et al., 2019a).  
405 Here we show that the pairing of visual representations and location information is  
406 apparent even when objects do not appear in their expected locations: although in  
407 the current study all fragments were presented in the same central location, their  
408 representations were still organized by their typical location. This shows that even in  
409 the absence of location information the brain can use the characteristic spatial  
410 distribution of visual contents to organize their representation in an efficient way. A  
411 related effect of real-world structure on individual object representations was  
412 previously reported in face processing, where representations of individual face  
413 fragments are grouped according to their position within a face (Henriksson et al.,  
414 2015). Together, these results suggest that information across different types of

415 fragmented visual contents can be contextualized on the basis of real-world  
416 structure.

417         How does this contextualization mechanism aid perception under naturalistic  
418 conditions? The mechanism may be particularly beneficial across a variety of  
419 situations where visual inputs are incomplete. Such situations include partially  
420 occluded objects, fast-changing and dynamic environments, and fragmented  
421 information arising from eye movements across a scene. In each of these situations,  
422 matching the input with its typical position in the context of the current environment  
423 can facilitate the understanding of the incomplete information available at every point  
424 in time. Future studies need to connect the rapid sorting process described here and  
425 behavioral benefits in the aforementioned situations.

426

427

428

### **Acknowledgements**

429

430 D.K. and R.M.C. are supported by Deutsche Forschungsgemeinschaft (DFG) grants  
431 (KA4683/2-1, CI241/1-1, CI241/3-1). R.M.C. is supported by a European Research  
432 Council Starting Grant (ERC-2018-StG 803370). The authors declare no conflict of  
433 interest.

434

435

### **Author Contributions**

436

437 D. K. and R.M.C. designed research, D.K. and G.I. acquired data, D.K. and G.I.  
438 analyzed data, D.K., G.I., and R.M.C. interpreted results, D.K. prepared figures, D.K.  
439 drafted manuscript, D.K., G.I., and R.M.C. edited and revised manuscript. All authors  
440 approved the final version of the manuscript.

441

442

## References

443

444 Abassi, E., & Papeo, L. (2019). The representation of two-body shapes in the human  
445 visual cortex. *Journal of Neuroscience*, doi.org/10.1523/JNEUROSCI.1378-  
446 19.2019

447 Baldassano, C., Beck, D. M., & Fei-Fei, L. (2017). Human-object interactions are more  
448 than the sum of their parts. *Cerebral Cortex*, *27*, 2276-2288.

449 Bar, M. (2004). Visual objects in context. *Nature Reviews Neuroscience*, *5*, 617-629.

450 Bilalic, M., Lindig, T., & Turella, L. (2019). Parsing rooms: the role of the PPA and RSC  
451 in perceiving object relations and spatial layout. *Brain Structure and Function*,  
452 *224*, 2505-2524.

453 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, *10*, 433-436.

454 Chan, A. W., Kravitz, D. J., Truong, S., Arizpe, J., & Baker, C.I. (2010). Cortical  
455 representations of bodies and faces are strongest in commonly experienced  
456 configurations. *Nature Neuroscience*, *13*, 417-418.

457 Cichy, R. M., Khosla, A., Pantazis, D., & Oliva, A. (2017). Dynamics of scene  
458 representations in the human brain revealed by magnetoencephalography and  
459 deep neural networks. *NeuroImage*, *153*, 346-358.

460 de Haas, B., Schwarzkopf, D. S., Alvarez, I., Lawson, R. P., Henriksson, L.,  
461 Kriegeskorte, N., & Rees, G. (2016). Perception and processing of faces in the  
462 human brain is tuned to typical facial feature locations. *Journal of*  
463 *Neuroscience*, *36*, 9289-9302.

- 464 Dima, D. C., Perry, G., & Singh, K. D. (2018). Spatial frequency supports the  
465 emergence of categorical representations in visual cortex during natural scene  
466 perception. *Neuroimage*, 179, 102-116.
- 467 Geisler, W. S. (2008). Visual perception and the statistical properties of natural  
468 scenes. *Annual Review of Psychology*, 59, 167-192.
- 469 Groen, I. I., Ghebreab, S., Prins, H., Lamme, V. A., & Scholte, H.S. (2013). From image  
470 statistics to scene gist: evoked neural activity reveals transition from low-level  
471 natural image structure to scene category. *Journal of Neuroscience*, 33, 18814-  
472 18824.
- 473 Groen, I. I., Silson, E. H., & Baker, C. I. (2017). Contributions of low- and high-level  
474 properties to neural processing of visual scenes in the human brain.  
475 *Philosophical Transactions of the Royal Society B – Biological Sciences*, 372.
- 476 Gronau, N., Neta, M., & Bar, M. (2008). Integrated contextual representation for  
477 objects' identities and their locations. *Journal of Cognitive Neuroscience*, 20,  
478 371-388
- 479 Harel, A., Groen, I. I. A., Kravitz, D. J., Deouell, L. Y., & Baker, C. I. (2016). The  
480 temporal dynamics of scene processing: A multifaceted EEG investigation.  
481 *eNeuro*, 3, ENEURO.0139-16.2016.
- 482 Henriksson, L., Mur, M., & Kriegeskorte, N. (2015). Faciotopy – A face-feature map  
483 with face-like topography in the human occipital face area. *Cortex*, 72, 156-  
484 167.
- 485 Henriksson, L., Mur, M., & Kriegeskorte, N. (2019). Rapid invariant encoding of scene  
486 layout in human OPA. *Neuron*, 103, 161-171.

- 487 Kaiser, D., & Cichy, R. M. (2018). Typical visual-field locations enhance processing in  
488 object-selective channels of human occipital cortex. *Journal of*  
489 *Neurophysiology*, *120*, 848-853.
- 490 Kaiser, D., Häberle, G., & Cichy, R. M. (2019b) Cortical sensitivity to natural scene  
491 structure. *Human Brain Mapping*, doi.org/10.1002/hbm.24875
- 492 Kaiser, D., & Haselhuhn, T. (2017). Facing a regular world: How spatial object structure  
493 shapes visual processing. *Journal of Neuroscience*, *37*, 1965-1967.
- 494 Kaiser, D., Moeskops, M. M., & Cichy, R. M. (2018). Typical retinotopic locations  
495 impact the time course of object coding. *NeuroImage*, *176*, 372-379.
- 496 Kaiser, D., & Peelen, M. V. (2018). Transformation from independent to integrative  
497 coding of multi-object arrangements in human visual cortex. *NeuroImage* *169*,  
498 334-341.
- 499 Kaiser, D., Quek, G. L., Cichy, R. M., & Peelen, M. V. (2019a). Object vision in a  
500 structured world. *Trends in Cognitive Sciences*, *23*, 672-685.
- 501 Kaiser, D., Stein, T., & Peelen, M. V. (2014). Object grouping based on real-world  
502 regularities facilitates perception by reducing competitive interactions in visual  
503 cortex. *Proceedings of the National Academy of Sciences USA*, *111*, 11217-  
504 11222.
- 505 Kaiser, D., Turini, J., & Cichy, R. M. (2019c). A neural mechanism for contextualizing  
506 fragmented inputs during naturalistic vision. *eLife*, *8*, e48182.
- 507 Kim, J. G., & Biederman, I. (2011). Where do objects become scenes? *Cerebral*  
508 *Cortex*, *21*, 1738-1746.
- 509 Kim, J. G., Biederman, I., & Juan, C. H. (2011). The benefit of object interactions arises  
510 in the lateral occipital cortex independently of attentional modulation from the

- 511 intraparietal sulcus: a transcranial magnetic stimulation study. *Journal of*  
512 *Neuroscience*, 31, 8320-8324.
- 513 Konkle, T., Brady, T. F., Alvarez, G. A., & Oliva, A. (2010). Scene memory is more  
514 detailed than you think: the role of categories in visual long-term memory.  
515 *Psychological Science*, 21, 1551-1556.
- 516 Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity  
517 analysis – connecting the branches of systems neuroscience. *Frontiers in*  
518 *Systems Neuroscience*, 2, 4.
- 519 Lowe, M. X., Rajsic, J., Ferber, S., & Walther, D. B. (2018). Discriminating scene  
520 categories from brain activity within 100 milliseconds. *Cortex*, 106, 275-287.
- 521 Mandler, J. M. (1984). Stories, scripts and scenes: aspects of schema theory. L.  
522 Erlbaum.
- 523 Mandler, J. M., & Parker, R. E. (1976). Memory for descriptive and spatial information  
524 in complex pictures. *Journal of Experimental Psychology: Human Learning,*  
525 *Memory, & Cognition*, 2, 38-48.
- 526 Mannion, D. J. (2015). Sensitivity to the visual field origin of natural image patches in  
527 human low-level visual cortex. *PeerJ*, 3, e1038.
- 528 Minsky, M. (1975). A framework for representing knowledge. In: The psychology of  
529 computer vision. Winston, P. (ed), McGraw-Hill.
- 530 Oliva, A., & Torralba, A. (2003). Statistics of natural image categories. *Network*, 14,  
531 391-412.
- 532 Oliva, A., & Torralba, A. (2007). The role of context in object recognition. *Trends in*  
533 *Cognitive Sciences*, 11, 520-527.

- 534 Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J. M. (2011). FieldTrip: Open source  
535 software for advanced analysis of MEG, EEG, and invasive  
536 electrophysiological data. *Computational Intelligence and Neuroscience*, 2011,  
537 156869.
- 538 Oosterhof, N. N., Connolly, A. C., & Haxby, J. V. (2016). CoSMoMVPA: Multi-modal  
539 multivariate pattern analysis of neuroimaging data in Matlab/GNU Octave.  
540 *Frontiers in Neuroinformatics*, 10, 20.
- 541 Proklova, D., Kaiser, D., & Peelen, M. V. (2016). Disentangling representations of  
542 object shape and object category in human visual cortex: the animate-  
543 inanimate distinction. *Journal of Cognitive Neuroscience*, 28, 680-692.
- 544 Proklova, D., Kaiser, D., & Peelen, M. V. (2019). MEG sensor patterns reflect  
545 perceptual but not categorical similarity of animate and inanimate objects.  
546 *Neuroimage*, 192, 167-177.
- 547 Purves, D., Wojtach, W. T., & Lotto, R. B. (2011). Understanding vision in wholly  
548 empirical terms. *Proceedings of the National Academy of Sciences USA*, 108,  
549 15588-15595.
- 550 Roberts, K. L., & Humphreys, G. W. (2010). Action relationships concatenate  
551 representations of separate objects in the ventral visual cortex. *NeuroImage*  
552 52, 1541-1548.
- 553 Rumelhart, D. E. (1980). Schemata: the building blocks of cognition. In: Theoretical  
554 issues in reading comprehension. Spiro R. J., et al. (eds.), L. Erlbaum.
- 555 Smith, S. M., & Nichols, T. E. (2009). Threshold-free cluster enhancement: addressing  
556 problems of smoothing, threshold dependence and localisation in cluster  
557 inference. *NeuroImage*, 44, 83-98.

- 558 Torralba, A., Oliva, A., Castelhana, M. S., & Henderson, J. M. (2006). Contextual  
559 guidance of eye movements and attention in real-world scenes: the role of  
560 global features in objects search. *Psychological Review*, 113, 766-786.
- 561 Võ, M. L.-H., Boettcher, S. E. P., & Draschkow, D. (2019). Reading scenes: How scene  
562 grammar guides attention and aids perception in real-world environments.  
563 *Current Opinion in Psychology*, 29, 205-210.
- 564 Wolfe, J. M., Võ, M. L.-H., Evans, K. K., & Greene, M. R. (2011). Visual search in  
565 scenes involves selective and nonselective pathways. *Trends in Cognitive  
566 Sciences*, 15, 77-84.