

1 **Title:** Automatic encoding of a view-centered background image in the macaque temporal lobe

2 **Author Names and affiliations:** He Chen<sup>1</sup>, Yuji Naya<sup>1,2,3</sup>

3 <sup>1</sup> School of Psychological and Cognitive Sciences, Peking University, No. 52, Haidian Road,  
4 Haidian District, Beijing 100805, China

5 <sup>2</sup> IDG/McGovern Institute for Brain Research at Peking University, No. 52, Haidian Road,  
6 Haidian District, Beijing 100805, China

7 <sup>3</sup> Beijing Key Laboratory of Behavior and Mental Health, Peking University, No. 52, Haidian  
8 Road, Haidian District, Beijing 100805, China

9

10 **Corresponding Author:** Yuji Naya

11 Address: School of Psychological and Cognitive Sciences, Peking University, No. 52, Haidian  
12 Road, Wang Kezhen Building, Room 1707, Haidian District, Beijing 100805, China

13 E-mail: yujin@pku.edu.cn

14

15

16

17

18

19

20

21 **Abstract**

22 Perceptual processing along the ventral visual pathway to the hippocampus is hypothesized to be  
23 substantiated by signal transformation from retinotopic space to relational space, which  
24 represents interrelations among constituent visual elements. However, our visual perception  
25 necessarily reflects the first person's perspective based on the retinotopic space. To investigate  
26 this two-facedness of visual perception, we compared neural activities in the temporal lobe  
27 (anterior inferotemporal cortex, perirhinal and parahippocampal cortices, and hippocampus)  
28 between when monkeys gazed on an object and when they fixated on the screen center with an  
29 object in their peripheral vision. We found that in addition to the spatially invariant object signal,  
30 the temporal lobe areas automatically represent a large-scale background image, which specify  
31 the subject's viewing location. These results suggest that a combination of two distinct visual  
32 signals on relational space and retinotopic space may provide the first person's perspective  
33 serving for perception and presumably subsequent episodic memory.

34

35

36

## 37 **Introduction**

38           Visual information of our external world could be once decomposed into “what” and  
39 “where” before we attained its mental representation as the first person’s perspective  
40 (Eichenbaum, Yonelinas, & Ranganath, 2007; Palombo et al., 2015; Tulving, 2002). For several  
41 decades, it has been considered that the perception of these two visual features proceeds  
42 exclusively through the ventral and dorsal pathways (Goodale & Milner, 1992; Haxby et al.,  
43 1991; Mishkin & Ungerleider, 1982). Instead of this widespread dichotomy, contemporary visual  
44 neuroscience research suggests a presence of spatial information in the ventral pathway for  
45 perception (Chen & Naya, 2019; Connor & Knierim, 2017; Russell A. Epstein & Julian, 2013;  
46 Freud, Plaut, & Behrmann, 2016; Hong, Yamins, Majaj, & DiCarlo, 2016; Kornblith, Cheng,  
47 Ohayon, & Tsao, 2013; Mormann et al., 2017; Schenk, 2010). For instance, neurons in the  
48 inferotemporal (IT) cortex (TEO and TEd) of non-human primates exhibited preferential  
49 responses to scene-like stimuli rather than object-like stimuli (Kornblith et al., 2013; Vaziri,  
50 Carlson, Wang, & Connor, 2014). The response pattern of scene-selective IT neurons may be  
51 comparable to an activation pattern in the parahippocampal place area detected in human  
52 functional imaging studies (R. Epstein & Kanwisher, 1998; Julian & Epstein, 2013). The  
53 parahippocampal place area is located within the parahippocampal cortex (PHC) of the medial  
54 temporal lobe (MTL), which receives inputs from the early stages of the ventral pathway  
55 including the TEO and posterior TEd in addition to inputs from the dorsal pathway, and provides  
56 spatial information to the hippocampus (HPC) – a candidate of the final brain region for scene  
57 perception (Burgess, 2008) - via the medial/posterior entorhinal cortex (ERC) (Rolls, 2018). On  
58 the other hand, neurons in the IT cortex also represent location information of an object within a  
59 scene either at population-coding level (Hong et al., 2016) or at single-neuron level (Chen &

60 Naya, 2019). It is worth noting that while most neurophysiological studies had shown a spatial  
61 invariance of object at single-neuron level during monkeys' fixating on the center of a display  
62 under either the passive-viewing task (Hong et al., 2016; Kobatake & Tanaka, 1994) or delayed  
63 matching-to-sample (object) task (Miyashita & Chang, 1988; Nakamura, Matsumoto, Mikami, &  
64 Kubota, 1994), our recent study demonstrated equivalent or even more neurons exhibiting  
65 location signal compared with object signal in the ventral part of the anterior inferotemporal  
66 cortex (TEv) and its downstream MTL area (e.g., perirhinal cortex, PRC) during an item-location  
67 retention (ILR) task requiring monkeys to encode both identity and location of a sample object  
68 using a foveal vision (Fig. 1). Importantly, the location-selective activity during the ILR task  
69 could not be explained by the animals' eye-positions themselves (Chen & Naya, 2019).

70         Considering that different gaze positions cause a substantial difference in the large-scale  
71 visual input in the ILR task using the foveal-view (F-V) condition (Fig. 1D), the most  
72 straightforward explanation for the robust location signal might be that a substantial number of  
73 neurons in the IT cortex and MTL areas are driven by the retinotopic signal including parafoveal  
74 vision, which would not only serve for recognizing a scene (Connor & Knierim, 2017; Dilks,  
75 Julian, Kubilius, Spelke, & Kanwisher, 2011; Kornblith et al., 2013; Vaziri et al., 2014) but also  
76 signal a particular location in the scene (Chen & Naya, 2019; Hong et al., 2016). An alternative  
77 explanation would be the location information of an object is coded into internal spatial  
78 relationships within a large complex stimulus including an object and its background regardless  
79 of their absolute retinotopic positions. In other words, the IT cortex and MTL areas would  
80 represent object location by transforming representations of the object and its background on the  
81 retinotopic space (Zhaoping, 2019) into those on the "relational space" (Connor & Knierim,  
82 2017). In this case, the location signal in the ILR task would be sensitive to the task demand

83 requiring the animals to retain the object location for a following action rather than a retinotopic  
84 image depending on the animals' gaze position.

85         To address this question and investigate characteristics of spatial information in the  
86 ventral pathway and its downstream (i.e., MTL areas), we examined single-unit activities and  
87 local-field potentials (LFPs) from the TEv and MTL subregions during an object stimulus  
88 presented randomly at one of the quadrants on the display in a peripheral-view (P-V) as well as  
89 in an F-V condition (Fig. 1). In the P-V condition, animals were required to fixate on a central  
90 dot and obtain the location and item-identity information of the sample object using their  
91 peripheral vision (Fig. 1A). We compared the location effects between the two-view conditions  
92 by testing two rhesus macaques, and found that regardless of the task demands for encoding of  
93 an object and its location, there were much more abundant location signal in the F-V condition  
94 compared with the P-V condition on all the recording regions of the two monkeys.

## 95 **Results**

96           We collected data in both F-V and P-V conditions from two rhesus macaques (Fig. 1B).  
97   During the recording, Monkey A was required to encode an identity of a sample stimulus and its  
98   location actively for a subsequent response (i.e., ILR task). We reported the single-unit data in  
99   the F-V condition of the ILR task in the previous study (Chen & Naya, 2019); here, we refer to  
100   the ILR task as an “active-encoding task.” On the other hand, monkey F was only required to  
101   fixate on a small white dot, viewing a sample stimulus passively (“passive-encoding task”) in  
102   both view conditions (Figs. 1A&B). We did not record from the single monkeys in both  
103   encoding tasks because it would be difficult to exclude both explicit and implicit influences of  
104   learning the active-encoding task on the cognitive process in the passive-encoding task. We used  
105   the same six visual objects (yellow Chinese characters, radius = 3°) as sample stimuli for both  
106   monkeys through all the recording sessions (Fig. 1C). It should be noted that the retinotopic  
107   images differed entirely between F-V and P-V conditions although a position of a sample  
108   stimulus was identical relative to the external world including a large square background (48°  
109   each side) on the display between the two view conditions (Figs. 1C&D). This two-by-two  
110   experimental design (“F-V vs. P-V” × “active-encoding vs. passive-encoding”) allowed us to  
111   compare the neural signals in the F-V condition with those in the P-V condition in the animals  
112   with different task demands. Monkey A performed the active-encoding task at high  
113   performances in both F-V ( $96.2 \pm 3.7$  %, 454 sessions) and P-V ( $92.2 \pm 7.1$  %, 477 sessions)  
114   conditions.

### 115 *Gaze-related location signal*

116           We first investigated single-unit activities signaling location information. Figure 2A  
117   shows an example of TEv neurons that were recorded in the active-encoding task. The neuron

118 showed the largest responses when the animal fixated on the *position I* (*top right* on the large  
119 square background). Although the responses once decayed, the neuron responded strongly when  
120 an item stimulus was presented as a sample stimulus at the same *position I* in the F-V condition.  
121 We examined the neuronal responses during 80-1000 ms after the onset of sample presentation  
122 (sample period) using a two-way ANOVA with item identities (six items) and locations (four  
123 locations) as main effects. The neuron showed a significant location effect [ $P < 0.0001$ ,  $F(3,156)$   
124 = 18.98] but not for item identities of sample stimuli [ $P = 0.309$ ,  $F(5,156) = 1.21$ ]. In contrast to  
125 the strong location-selectivity in the F-V condition, the same TEv neuron did not show location-  
126 selective activities in the P-V condition during the sample period [ $P = 0.183$ ,  $F(3,157) = 1.63$ ].

127 Figure 2B shows an example of PRC neurons that also exhibited location-selective activities  
128 only in the F-V conditions. This neuron signaled location information only after sample  
129 presentation in the F-V condition, suggesting that the presence of location signal in the F-V  
130 condition cannot be necessarily explained by preceding location-selective activity before sample  
131 presentation. We examined the prevalence of location signal in the two view conditions among  
132 the recording regions by calculating proportions of neurons with significant ( $P < 0.01$ , two-way  
133 ANOVA) location-selective activities during the sample period in each area. All recording  
134 regions contained significantly ( $P < 0.0016$  in each region,  $\chi^2$  test) larger proportions of location-  
135 selective cells in the F-V condition (24%, TE; 27%, PRC; 21%, HPC; 20%, PHC) than the P-V  
136 condition (7%, TE; 10%, PRC; 7%, HPC; 4%, PHC) (Fig. 3A). These results indicated that the  
137 location information in the ventral pathway and MTL areas were sensitive to the view  
138 conditions, although the same task-relevant information was required for a following action in  
139 the active-encoding task. The robust location signal only in the F-V condition implicates that the

140 temporal lobe areas represent a visual image, which subjects view rather than the goal-directed  
141 spatial information related with an action plan.

142         The different sensitivity to the two view conditions was also observed for the location  
143 signal in the passive-encoding task (Figs. 2C&D). Similar to the active-encoding task, we found  
144 a substantial number of neurons exhibiting location effect (29%, TE; 15%, PRC; 21%, HPC;  
145 34%, PHC) under the F-V condition (Fig. 3B). This result indicates that the location-selective  
146 response in the active-encoding task did not result from the task requirement, in which the  
147 animal was required to maintain actively a location of a sample stimulus. Compared with the F-V  
148 condition, the number of location-selective cells decreased dramatically under the P-V condition  
149 in all areas (8%, TE; 4%, PRC; 0%, HPC; 10%, PHC) (Fig. 3B). These results are also consistent  
150 with the single-unit results in the active-encoding task, and suggest that the gaze-sensitive  
151 location signal is automatically encoded by neurons in the TEv and MTL. The marked reduction  
152 of location signal in the P-V condition during either active or passive-encoding task argued  
153 against the possibility that the location-selective cells distinguish the structural organization of  
154 large objects with internal structures (e.g., a large grey square with a small letter at its top-left vs.  
155 at its bottom-right) which would be represented by the relational rather than the retinotopic space  
156 (Connor & Knierim, 2017).

157         The most straight-forward interpretation of fewer active location-selective cells under the  
158 P-V condition may be that fixating on the center of the display reduces attention to a sample  
159 stimulus and attenuates the response of location-selective cells, which showed robust location  
160 signals in the F-V condition. If this situation applies, we would then expect that neurons with  
161 stronger location selectivity in the F-V condition would show relatively stronger location  
162 selectivity in the P-V condition (i.e., a positive correlation). To test this possibility, we estimated



163 strengths of location signals for neurons with location-selective activity in either F-V or P-V  
164 condition using  $F$  values indicating a location effect in the two-way ANOVA. Notably, we  
165 observed a negative correlation in amplitudes of the  $F$  values between the conditions in all areas  
166 during either active-encoding (Spearman rank correlation = -0.24 among 229 neurons across  
167 areas,  $P = 0.0003$ , two-tailed) (Fig. 4A) or passive-encoding task (Spearman rank correlation = -  
168 0.20 among 71 neurons,  $P = 0.090$ , two-tailed) (Fig. 4B). These results suggest that the weak  
169 location signal in the P-V condition was not due to the attenuated attention to a sample item. A  
170 reasonable interpretation of the negatively correlated location signal might be that separate visual  
171 inputs on the retinae drive different ensembles of neurons between the two view conditions (Fig.  
172 1D). This interpretation is consistent with the significant reduction in the proportion of location-  
173 selective cells from the F-V to the P-V condition (Fig. 3) because a retinotopic shift of a large  
174 background square ( $48^\circ$ , each side, Fig. 1C) in the F-V condition (Fig. 1D, left) would drive  
175 more neurons than that of a small sample stimulus ( $3^\circ$ , radius) in the P-V condition (Fig. 1D,  
176 right). Collectively, the TEv and MTL areas may automatically signal large-scale background  
177 information represented on the retinotopic space, which necessarily reflects a perspective that a  
178 subject is viewing.

### 179 *Task-dependent item signal*

180 In contrast to the dramatic difference in the location-selective activity between the F-V  
181 and P-V conditions, neurons in the temporal lobe showed consistent item-selective responses  
182 between the two view conditions during the active-encoding task (Fig. S1). In all recording  
183 regions except for the PHC, we found a substantial number of item-selective cells under the P-V  
184 condition (TE 23%, PRC 22%, HPC 27%, and PHC 2%) as well as F-V condition (TE 14%,  
185 PRC 22%, HPC 32%, and PHC 3%) (Fig. 3A, bottom). These results are consistent with

186 previous studies indicating the spatial invariance of object representation (Kobatake & Tanaka,  
187 1994; Miyashita & Chang, 1988; Nakamura et al., 1994), which would be obtained by  
188 transforming it from the retinotopic space into the relational space along the ventral pathway  
189 (Connor & Knierim, 2017). In contrast to the location signal, the signal strengths of item  
190 information positively correlated between the F-V and P-V conditions (Figs. 4&D). These results  
191 indicate distinct processing between the item and its background (i.e., location signal) regarding  
192 their sensitivity to the view conditions. Interestingly, the number of item-selective cells was  
193 negligible in all areas under both view conditions in the passive-encoding task (F-V condition:  
194 TE 6%, PRC 0%, HPC 2%, PHC 3%; P-V condition: TE 2%, PRC 2%, HPC 3% PHC 1%; Fig.  
195 3B, bottom), which contrasts to the substantial number of item-selective cells in the active-  
196 encoding task. The inconsistency in the item signal between the two tasks suggests that the  
197 object representation depends on the task demand, which required the subject to maintain an item  
198 identity of a sample stimulus for the following action.

### 199 *Population-coding analysis*

200 The analyses based on the spike-firing data of individual neurons indicated substantially  
201 stronger location signal in the F-V condition compared with the P-V condition regardless of the  
202 task demands. One remaining question might be whether the location signal could be represented  
203 equivalently between the two view conditions by population coding. To test this possibility, we  
204 conducted the “representational similarity analyses” (RSA) (Kriegeskorte, Mur, & Bandettini,  
205 2008); we first constructed a population vector consisting of firing rates of all recorded neurons  
206 in each area as its elements. In each combination of view condition and encoding-type, there  
207 were twenty-four (six items  $\times$  four locations) of  $n$ -dimensional population vectors. “ $n$ ” indicates  
208 a number of the recorded neurons in each area. We then calculated correlation coefficients

209 between the population vectors, indicating the similarity level of neural representations between  
210 trial-types with different item-location combinations. Figures 5A and 5B displayed the similarity  
211 level of neural representations in the HPC during the sample presentation period in the active-  
212 encoding and passive-encoding tasks, respectively. In both tasks, the representational similarities  
213 between trial-types with same locations (e.g., location 1 item 1 & location 1 item 2) were  
214 substantially larger than the similarities between trial-types with different locations (e.g., location  
215 1 item 1 & location 2 item 2) in the F-V condition ( $P < 0.001$  in both tasks, one-side, simulation  
216 test), suggesting that the HPC represents the item location that the animals were viewing,  
217 regardless of the task demands. In contrast to the F-V condition, the HPC's discriminability in  
218 the location of a sample stimulus was considerably diminished in the P-V condition (Figs.  
219 5A&B). In the RSA, other recorded regions also showed the marked reduction of the location  
220 signal in the P-V condition compared with the F-V condition in both tasks (Figs. 5C&D).  
221 Together, consistent with the analyses based on the single neurons, the analyses examining the  
222 population coding suggest that the temporal lobe areas represent the location information more  
223 robustly in the F-V condition than the P-V condition. As to the item signal, the RSA also  
224 provided the results which were consistent with the results of the single-neuron-based analyses  
225 (Fig. S2).

#### 226 *LFP activity depending on both view-condition and task-demand*

227 In addition to spiking data, we investigated the LFP activity during the sample period.  
228 Figure 6A shows the differential spectrums between the viewing conditions (F-V condition  
229 minus P-V condition) in each recording region under the active-encoding task (*left column*) and  
230 passive-encoding task (*right column*). During the early sample presentation period (0-300 msec  
231 after sample onset), there is an enhanced beta-band activity (1- 25 Hz) expressed non-selectively

232 across the brain regions and tasks (Fig. 6B). This higher beta-band activity in the F-V condition  
233 is consistent with preceding literature indicating that larger beta-band activity is observed when  
234 the current cognitive or perceptual status should be actively maintained (i.e. the sample stimulus  
235 appears at the same position as with the fixation period in the F-V condition) than when the  
236 current state is disrupted by an unexpected event (i.e. the sample stimulus appears randomly at  
237 one out of the four positions in the P-V condition) (Engel & Fries, 2010). A view-condition  
238 dependent LFP activity was also observed in a gamma-band (30-80 Hz) during the late sample  
239 presentation period (350-800 msec after sample onset) (Fig. 6A). In contrast to the widely  
240 distributed beta-band, the gamma-band activity was selectively expressed only in the PRC and  
241 HPC when a sample item and its location were encoded actively by the foveal vision (Fig. 6B),  
242 in which situation both the item and location signals appeared robustly in these brain regions  
243 (Figs. 3&5&S2). These results may implicate that the increased gamma-band activity is related  
244 with the interaction between the item and location signals, which reportedly occurs in the PRC  
245 and HPC but not in TEv nor PHC (Chen & Naya, 2019).

## 246 **Discussion**

247 The present study provides single-unit data showing robust spatial information in the TEv  
248 and MTL areas, which signaled a particular location where the animals were viewing (F-V  
249 condition) rather than an object position presented in the peripheral view (P-V condition). These  
250 results were shown for each of the recording regions by the independent analyses for each of the  
251 two monkeys, indicating the very robust animal consistency. In addition, this animal consistency  
252 was confirmed even though the two animals were tested in different task demands (i.e., active-  
253 encoding and passive-encoding of an object and its location), which manifests the robustness of  
254 the present findings showing an existence of the location signal characterized by the clear

255 difference in its sensitivity to the two view conditions. These new findings suggest that the  
256 location signal in the primate temporal lobe areas may represent a view-centered background  
257 image, which could specify the current gaze position within a scene (Fig. 7). This view-centered  
258 background may be automatically represented in the temporal lobe areas because it was observed  
259 in the passive-encoding task as well as the active-encoding task. The TEv and MTL areas except  
260 for the PHC also signaled object information. However, in contrast to the background  
261 information, the object information was represented regardless of the view conditions when it  
262 was actively encoded. These results from the single-neuron-based analyses were confirmed by  
263 population-coding analyses. Taken together, the present study suggests that the ventral pathway  
264 and its downstream in the MTL signal not only spatially-invariant object information but also  
265 view-centered background information, which may automatically locate the object in a scene  
266 when it is viewed by the foveal vision.

267         One naïve question on the gaze-related location signal might be whether the location  
268 signal could be explained by non-visual sensory/motor information, which reflects the animals'  
269 eye positions relative to their heads. Our previous study indicated that neurons in the TEv and  
270 MTL areas responded differently to the same gaze positions depending on the position of the  
271 large background square within the display (leftward or rightward) (Chen & Naya, 2019),  
272 suggesting that the gaze-related location signal reflects visual inputs rather than  
273 somatosensory/motor-related information of the gaze itself. In the present study, we  
274 characterized the location signal, which were widely distributed over the temporal lobe areas, by  
275 revealing the underlying visual inputs not to be represented on the relational space, but instead  
276 on the retinotopic space (i.e., view-centered background). An important question about the view-  
277 centered background information on the retinotopic space might be whether it only reflects the

278 parafoveal vision or not. In the present study, the location-selective activity depends on the  
279 parafoveal vision of the background, which shows an edge of the large grey square or the display  
280 frame. However, some neurons exhibited location-selective activities only after sample  
281 presentation in the F-V condition (Figs. 2B-D) (10.8% and 8.5% across areas in the active and  
282 passive-encoding tasks), which suggest an existence of neuronal population that represent the  
283 view-centered background including foveal vision as well as parafoveal vision. The view-  
284 centered background signal in the present study may explain response patterns of “spatial view  
285 cells” in the HPC (and posterior PRC) reported by Rolls (Rolls, Robertson, & Georges-François,  
286 1997). The spatial view cells show selective responses to a particular location where an animal  
287 views regardless of its standing position. This allocentric coding property of the spatial view  
288 cells could be due to similar visual inputs when an animal views the same location from different  
289 positions.

290 In spite of the location signal which may reflect the background information on the  
291 retinotopic space, the object signal was detected regardless of its retinotopic position in the  
292 active encoding task (Fig. 7), which confirmed the preceding literature showing the spatial  
293 invariant of object representation in the IT cortex (Miyashita & Chang, 1988; Nakamura et al.,  
294 1994). The representation of an object may be explained by a spatial relationship among the  
295 internal elements of it, which necessarily accompany its transformation from the retinotopic  
296 space into the relational space (Connor & Knierim, 2017). The present study suggests that  
297 neurons in the temporal lobe signal the location information of an object as its background image  
298 represented on the retinotopic space (Fig. 7) rather than an interrelation between the object and  
299 any other spatial structure such as a large gray square behind it. Based on the present  
300 experimental set up, the background image encoded by neurons in the TEv and MTL areas

301 should cover larger than 30 degrees in the visual angle (diameter) to include the edge of the large  
302 gray square background, which may cause different responses according to the gaze positions.  
303 As well as the object signal, the large-scale background image is reportedly processed along the  
304 ventral pathway (Kornblith et al., 2013; Vaziri et al., 2014). One remaining question is whether  
305 the processing of the background image in the ventral pathway imparts more generalized spatial  
306 features (e.g., field, valley, forest), which may be represented on the relational space and serve  
307 for recognizing an entire scene (e.g., suburb rather than modern city) regardless of the gaze  
308 positions (e.g., an eagle over the valley).

309 In addition to the view conditions testing the representation spaces (i.e., relational vs.  
310 retinotopic), the object and the background signals showed differential sensitivity patterns to the  
311 task demands in the present study. The background signal was encoded irrespective of the task  
312 demand while the object signal was encoded only in the active-encoding task. The automatic  
313 encoding of the background signal suggests that when we direct our gaze toward an object to  
314 obtain its high-resolution image, we would spontaneously receive the spatial information, which  
315 would be assigned to the object (Chen & Naya, 2019). One remaining problem about the object  
316 signal might be whether the lack of item-selective activity in the passive-encoding task is due to  
317 the present stimulus set (i.e., Chinese character) because the IT neurons reportedly respond to  
318 object stimuli such as face stimuli in a passive-viewing task (Kiani, Esteky, Mirpour, & Tanaka,  
319 2007; Tsao, Freiwald, Knutsen, Mandeville, & Tootell, 2003). Compared with a natural object  
320 such as a face stimulus, a fabricated two-dimensional stimuli used in the present study may not  
321 bring about a bottom-up attention to be perceived as an object. In the active-encoding task, the  
322 monkey learned the Chinese characters to discriminate one from another. The repetitive training  
323 in the active-encoding task might form a long-term learning effect on the stimulus to induce the

324 bottom-up attention, which may lead a transformation of representations of Chinese-characters  
325 from the retinotopic space into the relational space. Although we cannot address if the attention  
326 was derived from the bottom-up or the top-down, the attention-dependent object signal and the  
327 attention-independent background signal may derive from a figure-background segmentation,  
328 which reportedly occurred at the V4, a start point of the ventral pathway (Roe et al., 2012).  
329 Previous studies have focused on the object information which is filtered, and implicated that the  
330 object representation is transformed from the retinotopic space into the relational space with the  
331 increase of neurons' receptive fields along the ventral pathway (Connor & Knierim, 2017). We  
332 hypothesize that the background information, which is filtered-out at the figure-ground  
333 segmentation, spreads into the ventral pathway with its representation remaining on the  
334 retinotopic space rather than the relational space. Our previous report has demonstrated that the  
335 two distinct signals, which are segmented from the same retinal image, are integrated step-by-  
336 step from the TEv, PRC to HPC (Chen & Naya, 2019). From the ventral stream to the MTL  
337 areas, the strongest integration effect was found in the PRC at the single neurons level. This  
338 integration process may be related with the largest gamma-band LFP activity in the PRC, which  
339 was observed when the monkey gazed at an object to encode its identity and location information  
340 actively (Fig. 6).

341 In the present study, the PHC represents the view-centered background signal whose  
342 property is similar to that in the TEv and other MTL areas including the PRC. Considering the  
343 heavier projections from the posterior parietal cortex to the PHC compared with the AIT cortex  
344 including the PRC (Kravitz, Saleem, Baker, & Mishkin, 2011), the PHC may also process the  
345 spatial information related with the eye/self-movement. Contributions of the PHC to scene  
346 construction process may become apparent when a subject perceives the environment by moving



347 their gazes (Zhang & Naya, 2019) in which multiple views should be coordinated according to  
348 the eye/self-movements, beyond encoding a single snapshot focusing on one object which was  
349 investigated in the present study. We propose a future study to investigate how the past multiple  
350 views influence on the present view to build the current first person's perspective (Eichenbaum  
351 et al., 2007; Palombo et al., 2015; Tulving, 2002), which may be related with an encoding of  
352 episodic memory.

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

## 370 **Materials and Methods**

### 371 *Subjects*

372 Two male monkeys (*Macaca mulatta*) (9.3 kg, monkey A; 10.1 kg, monkey F) were  
373 used for the experiments. All procedures and treatments were performed in accordance with the  
374 NIH Guide for the Care and Use of Laboratory Animals and were approved by the Institutional  
375 Animal Care and Use Committee (IACUC) of Peking University.

### 376 *Behavioral task*

377 We trained monkey A on a foveal-view/F-V condition of an active-encoding task with six  
378 visual items (Fig. 1). During both training and recording sessions, monkeys performed the task  
379 under dim light in an electromagnetic shielded room (length \* width \* height = 160 cm \* 120 cm  
380 \* 222 cm). The task began with an encoding phase, which was initiated by the animal pulling a  
381 lever and fixating on a white square (0.6 ° of visual angle) presented within one of the four  
382 quadrants (12.5 ° from the center) of a touch screen (3MTM MicroTouch™ Display M1700SS,  
383 17 inch, horizontal viewing angle: ~59 °, vertical viewing angle: ~49 °) with a custom-made  
384 metal frame (diagonal size: 22 inch, horizontal viewing angle: ~72 °, vertical viewing angle: ~71  
385 °) situated ~28 cm from the subjects. Eye position was monitored using an infrared digital  
386 camera with a 120 Hz sampling frequency (ETL-200, ISCAN) placed next to the left edge of the  
387 touch screen. The eye position calibration was conducted before starting each recording session  
388 (Monkey logic). After a 0.6 s fixation, one of the six items (3.0 °, radius) was presented in the  
389 same quadrant as a sample stimulus for 0.3 s, followed by another 0.7 s fixation on the white  
390 square. An additional 0.017 s, reflecting the design of software and hardware controlling the  
391 behavioral task was added to each trial event. If the fixation was successfully maintained  
392 (typically, < 2.5 °), the encoding phase ended with the presentation of a single drop of water.

393           The encoding phase was followed by a blank interphase delay interval of 0.7-1.4 s during  
394    which no fixation was required. Then, the response phase was initiated with a fixation dot  
395    presented at the center of the screen. One of the six items was then presented at the center for 0.3  
396    s as a cue stimulus. After another 0.5 s delay period, five discs were presented as choices,  
397    including a blue disc in each quadrant and a green disc at the center. When the cue stimulus was  
398    the same as the sample stimulus, the subject was required to choose by touching the blue disc in  
399    the same quadrant as the sample (i.e., match condition). Otherwise, the subject was required to  
400    choose the green disc (i.e., nonmatch condition). If the animal made the correct choice, four to  
401    eight drops of water were given as a reward; otherwise, an additional 4 s was added to the  
402    standard intertrial interval (1.5-3 s). During the trial, a large gray square (48 ° on each side, RGB  
403    value: 50, 50, 50, luminance: 3.36 cd/m<sup>2</sup>) was presented at the center of the display (backlight  
404    luminance: 0.22 cd/m<sup>2</sup>) as a background. After the end of a trial, all stimuli disappeared and the  
405    entire screen displayed light-red color during the inter-trial interval. The start of a new trial was  
406    indicated by the re-appearance of the large gray square on the display, upon which the monkey  
407    could start to pull the lever triggering an appearance of a white fixation dot. In the match  
408    condition, sample stimuli were pseudorandomly chosen from six well-learned visual items, and  
409    each item was presented pseudorandomly within the four quadrants, resulting in 24 (6 × 4)  
410    different configuration patterns. In the nonmatch condition, the position of the sample stimulus  
411    was randomly chosen from the four quadrants, and the cue stimulus was randomly chosen from  
412    the five items that differed from the sample stimulus. The match and nonmatch conditions were  
413    randomly presented at a ratio of 4:1, resulting in 30 (24+6) different configuration patterns. The  
414    same six stimuli were used during all recording sessions.

415 In addition to the F-V condition, we tested the neuronal responses of monkey A in the  
416 peripheral-view/P-V condition of the active-encoding task. In this view condition, fixation on the  
417 center of the display was required during the encoding phase (Fig. 1). Other parameters were the  
418 same as those in the F-V condition of the active-encoding task. Correct performance under F-V  
419 condition:  $97.5 \pm 2.6\%$  in the match trials and  $90.8 \pm 8.1\%$  in the nonmatch trials ( $n = 454$   
420 sessions); P-V condition:  $94.3 \pm 6.2\%$  in the match trials and  $84.1 \pm 10.8\%$  in the nonmatch trials  
421 ( $n = 478$  sessions).

422 We tested the neuronal responses of monkey F in both F-V and P-V condition of a  
423 passive-encoding task, in which the task sequence and requirement were same as the encoding  
424 phase of the active-encoding task but without a lever-pulling requirement (no interphase delay  
425 interval and response phase). The configuration of visual stimuli (such as visual angles,  
426 configuration patterns, and others) was same as that for monkey A. We tested the neuronal  
427 response of both monkey A and monkey F in the F-V and P-V conditions in a block manner.

#### 428 *Electrophysiological recording*

429 Following initial behavioral training, animals were implanted with a head post and  
430 recording chamber under aseptic conditions using isoflurane anesthesia. To record single-unit  
431 activity, we used a 16-channel vector array micrILRobe (V1 X 16-Edge, NeuroNexus), 16-  
432 channel U-Probe (Plexon), tungsten tetrode probe (Thomas RECORDING), or a single-wire  
433 tungsten microelectrode (Alpha Omega), which was advanced into the brain using a hydraulic  
434 Microdrive (MO-97A, Narishige) (Naya & Suzuki, 2011). The microelectrode was inserted  
435 through a stainless steel guide tube positioned in a customized grid system on the recording  
436 chamber. Neural signals for single units were collected (low-pass, 6 kHz; high-pass, 200 Hz) and  
437 digitized (40 kHz) (OmniPlex Neural Data Acquisition System, Plexon). These signals were then

438 sorted using an offline sorter provided by the OmniPlex system. We did not attempt to prescreen  
439 isolated neurons. Instead, once we isolated any neuron, we started to record its activity. The  
440 location of microelectrodes in target areas was guided by individual brain atlases from MRI  
441 scans (3T, Siemens). We also constructed individual brain atlases based on the  
442 electrophysiological properties around the tip of the electrode (e.g., gray matter, white matter,  
443 sulcus, lateral ventricle, and bottom of the brain). The recording sites were estimated by  
444 combining the individual MRI atlases and physiological atlases (Naya, Chen, Yang, & Suzuki,  
445 2017). To record LFPs, we used neural signals from the same electrodes as we used for the  
446 recording of spikes. However, the signals were collected using different filters (low-pass, 200  
447 Hz; high-pass, 0.05 Hz), and digitized at 1 kHz.

448         The recording sites in monkey A covered an area between 5 and 24 mm anterior to the  
449 interaural line (right hemisphere). The recording sites in monkey F covered an area between 6.6  
450 and 23.4 mm anterior to the interaural line (right hemisphere). The recording sites in HPC  
451 appeared to cover all its subdivisions (i.e., dentate gyrus, CA3, CA1, and subicular complex).  
452 The recording sites in PHC focused on approximately the lateral 2/3. The recording sites in PRC  
453 appeared to cover areas 35 and 36 from the fundus of the rhinal sulcus to the medial lip of the  
454 anterior middle temporal sulcus (amts). The border of PRC's caudal limit (PHC's rostral limit)  
455 was determined according to the rostral limit of the occipital temporal sulcus and the caudal limit  
456 of the rhinal sulcus (Suzuki & Amaral, 2003). In monkey A, the caudal limit of the recording  
457 sites in PRC is 2 mm posterior to the caudal limit of its rhinal sulcus and 1 mm anterior to the  
458 rostral limit of the occipital temporal sulcus. In monkey F, the caudal limit of the recording sites  
459 in PRC is 0 mm posterior to the caudal limit of its rhinal sulcus and 0 mm anterior to the rostral

460 limit of the occipital temporal sulcus. The recording sites in TE were limited to its ventral area,  
461 including both banks of the amts.

#### 462 *Data analysis*

463 All neuronal data were analyzed using MATLAB (MathWorks) with custom written  
464 programs, including the statistics toolbox. For responses before sample presentation, we tested  
465 each neuron's firing rate during the 700 ms period before the sample stimulus onset, including  
466 the 100 ms before the fixation start, as the monkeys typically started fixation 160-170 ms after  
467 fixation dot presentation. For responses during/after sample presentation, the firing rate during  
468 the period extending from 80 to 1000 ms after sample onset was tested. For responses before  
469 sample presentation, we evaluated the effects of "location" for each neuron using one-way  
470 ANOVA ( $P < 0.01$ ). For sample responses, we evaluated the effects of "location" and "item" for  
471 each neuron using two-way ANOVA with interactions ( $P < 0.01$  for each). We analyzed neurons  
472 that we tested in at least 60 trials (10 trials for each stimulus, 15 trials for each location).

473

474 **Acknowledgments:** We thank E.T. Rolls, W.A. Suzuki, M. Zhang, K.W. Koyano, C. Yang for  
475 helpful comments and S. Xue for expert animal care. Funding: The present study was funded by  
476 National Natural Science Foundation of China Grant 31421003 & 31871139 (to Y.N.).

477 **Author Contributions:** Y.N. designed the experiments. H.C. performed the experiments. H.C.  
478 and Y.N. analyzed data and wrote the manuscript.

479 **Declaration of Interests:** The authors declare no competing financial interests.

480

481

482 **References**

- 483 Burgess, Neil. (2008). Spatial Cognition and the Brain. *Ann N Y Acad Sci*, 1124(1), 77-97.  
484 doi:10.1196/annals.1440.002
- 485 Chen, He, & Naya, Yuji. (2019). Forward Processing of Object-Location Association from the  
486 Ventral Stream to Medial Temporal Lobe in Nonhuman Primates. *Cerebral cortex (New*  
487 *York, N.Y. : 1991)*. doi:10.1093/cercor/bhz164
- 488 Connor, C. E., & Knierim, J. J. (2017). Integration of objects and space in perception and  
489 memory. *Nat Neurosci*, 20(11), 1493-1503. doi:10.1038/nn.4657
- 490 Dilks, Daniel D., Julian, Joshua B., Kubiilius, Jonas, Spelke, Elizabeth S., & Kanwisher, Nancy.  
491 (2011). Mirror-image sensitivity and invariance in object and scene processing pathways.  
492 *Journal of Neuroscience*, 31(31), 11305-11312. doi:10.1523/JNEUROSCI.1935-11.2011
- 493 Eichenbaum, H., Yonelinas, A. P., & Ranganath, C. (2007). The medial temporal lobe and  
494 recognition memory. *Annu Rev Neurosci*, 30, 123-152.  
495 doi:10.1146/annurev.neuro.30.051606.094328
- 496 Engel, A. K., & Fries, P. (2010). Beta-band oscillations--signalling the status quo? *Curr Opin*  
497 *Neurobiol*, 20(2), 156-165. doi:10.1016/j.conb.2010.02.015
- 498 Epstein, Russell, & Kanwisher, Nancy. (1998). The Parahippocampal Place Area: A Cortical  
499 Representation of the Local Visual Environment. *Neuroimage*, 7(4), S341-S341.  
500 doi:10.1016/S1053-8119(18)31174-1
- 501 Epstein, Russell A, & Julian, Joshua B. (2013). Scene Areas in Humans and Macaques. *Neuron*,  
502 79(4), 615-617. doi:10.1016/j.neuron.2013.08.001

- 503 Freud, Erez, Plaut, David C., & Behrmann, Marlene. (2016). ‘What’ Is Happening in the Dorsal  
504 Visual Pathway. *Trends in Cognitive Sciences*, 20(10), 773-784.  
505 doi:10.1016/j.tics.2016.08.003
- 506 Goodale, M. A., & Milner, A. D. (1992). Separate Visual Pathways for Perception and Action.  
507 *Trends Neurosci*, 15(1), 20-25. doi:Doi 10.1016/0166-2236(92)90344-8
- 508 Haxby, J. V., Grady, C. L., Horwitz, B., Ungerleider, L. G., Mishkin, M., Carson, R. E., . . .  
509 Rapoport, S. I. (1991). Dissociation of object and spatial visual processing pathways in  
510 human extrastriate cortex. *Proc Natl Acad Sci U S A*, 88(5), 1621-1625.  
511 doi:10.1073/pnas.88.5.1621
- 512 Hong, H., Yamins, D. L., Majaj, N. J., & DiCarlo, J. J. (2016). Explicit information for category-  
513 orthogonal object properties increases along the ventral stream. *Nat Neurosci*, 19(4), 613-  
514 622. doi:10.1038/nn.4247
- 515 Julian, J., & Epstein, R. (2013). The Landmark Expansion Effect: Navigational Relevance  
516 Influences Memory of Object Size. *J Vis*, 13(9), 49-49. doi:10.1167/13.9.49
- 517 Kiani, Roozbeh, Esteky, Hossein, Mirpour, Koorosh, & Tanaka, Keiji. (2007). Object Category  
518 Structure in Response Patterns of Neuronal Population in Monkey Inferior Temporal  
519 Cortex. *J Neurophysiol*, 97(6), 4296-4309. doi:10.1152/jn.00024.2007
- 520 Kobatake, E., & Tanaka, K. (1994). Neuronal selectivities to complex object features in the  
521 ventral visual pathway of the macaque cerebral cortex. *J Neurophysiol*, 71(3), 856-867.  
522 doi:10.1152/jn.1994.71.3.856
- 523 Kornblith, Simon, Cheng, Xueqi, Ohayon, Shay, & Tsao, Doris Y. (2013). A Network for Scene  
524 Processing in the Macaque Temporal Lobe. *Neuron*, 79(4), 766-781.  
525 doi:10.1016/j.neuron.2013.06.015



- 526 Kravitz, D. J., Saleem, K. S., Baker, C. I., & Mishkin, M. (2011). A new neural framework for  
527 visuospatial processing. *Nat Rev Neurosci*, *12*(4), 217-230. doi:10.1038/nrn3008
- 528 Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity analysis -  
529 connecting the branches of systems neuroscience. *Front Syst Neurosci*, *2*, 4.  
530 doi:10.3389/neuro.06.004.2008
- 531 Mishkin, M., & Ungerleider, L. G. (1982). Contribution of striate inputs to the visuospatial  
532 functions of parieto-preoccipital cortex in monkeys. *Behav Brain Res*, *6*(1), 57-77.  
533 doi:10.1016/0166-4328(82)90081-x
- 534 Miyashita, Y., & Chang, H. S. (1988). Neuronal correlate of pictorial short-term memory in the  
535 primate temporal cortex. *Nature*, *331*(6151), 68-70. doi:10.1038/331068a0
- 536 Mormann, Florian, Kornblith, Simon, Cerf, Moran, Ison, Matias J., Kraskov, Alexander, Tran,  
537 Michelle, . . . Fried, Itzhak. (2017). Scene-selective coding by single neurons in the  
538 human parahippocampal cortex. *Proc Natl Acad Sci U S A*, *114*(5), 1153-1158.  
539 doi:10.1073/pnas.1608159113
- 540 Nakamura, K., Matsumoto, K., Mikami, A., & Kubota, K. (1994). Visual response properties of  
541 single neurons in the temporal pole of behaving monkeys. *J Neurophysiol*, *71*(3), 1206-  
542 1221. doi:10.1152/jn.1994.71.3.1206
- 543 Naya, Y., Chen, H., Yang, C., & Suzuki, W. A. (2017). Contributions of primate prefrontal  
544 cortex and medial temporal lobe to temporal-order memory. *Proc Natl Acad Sci U S A*,  
545 *114*(51), 13555-13560. doi:10.1073/pnas.1712711114
- 546 Naya, Y., & Suzuki, W. A. (2011). Integrating what and when across the primate medial  
547 temporal lobe. *Science*, *333*(6043), 773-776. doi:10.1126/science.1206773

- 548 Palombo, Daniela J., Alain, Claude, Söderlund, Hedvig, Khuu, Wayne, Levine, Brian,  
549 Institutionen för, psykologi, . . . Samhällsvetenskapliga, fakulteten. (2015). Severely  
550 deficient autobiographical memory (SDAM) in healthy adults: A new mnemonic  
551 syndrome. *Neuropsychologia*, 72, 105-118. doi:10.1016/j.neuropsychologia.2015.04.012
- 552 Roe, Anna W, Chelazzi, Leonardo, Connor, Charles E, Conway, Bevil R, Fujita, Ichiro, Gallant,  
553 Jack L, . . . Vanduffel, Wim. (2012). Toward a Unified Theory of Visual Area V4.  
554 *Neuron*, 74(1), 12-29. doi:10.1016/j.neuron.2012.03.011
- 555 Rolls, Edmund T. (2018). The storage and recall of memories in the hippocampo-cortical system.  
556 *Cell Tissue Res*, 373(3), 577-604. doi:10.1007/s00441-017-2744-3
- 557 Rolls, Edmund T., Robertson, Robert G., & Georges - François, Pierre. (1997). Spatial View  
558 Cells in the Primate Hippocampus. *European Journal of Neuroscience*, 9(8), 1789-1794.  
559 doi:10.1111/j.1460-9568.1997.tb01538.x
- 560 Schenk, Thomas. (2010). Visuomotor robustness is based on integration not segregation. *Vision*  
561 *Res*, 50(24), 2627-2632. doi:10.1016/j.visres.2010.08.013
- 562 Suzuki, W. A., & Amaral, D. G. (2003). Perirhinal and parahippocampal cortices of the macaque  
563 monkey: cytoarchitectonic and chemoarchitectonic organization. *J Comp Neurol*, 463(1),  
564 67-91. doi:10.1002/cne.10744
- 565 Tsao, D. Y., Freiwald, W. A., Knutsen, T. A., Mandeville, J. B., & Tootell, R. B. (2003). Faces  
566 and objects in macaque cerebral cortex. *Nat Neurosci*, 6(9), 989-995. doi:10.1038/nn1111
- 567 Tulving, E. (2002). Episodic memory: from mind to brain. *Annu Rev Psychol*, 53, 1-25.  
568 doi:10.1146/annurev.psych.53.100901.135114

- 569 Vaziri, Siavash, Carlson, Eric T, Wang, Zhihong, & Connor, Charles E. (2014). A Channel for  
570 3D Environmental Shape in Anterior Inferotemporal Cortex. *Neuron*, 84(1), 55-62.  
571 doi:10.1016/j.neuron.2014.08.043
- 572 Zhang, Bo, & Naya, Yuji. (2019). *Object-Based Cognitive Map in the Human Hippocampus and*  
573 *Medial Prefrontal Cortex*: bioRxiv 680199; doi: <https://doi.org/10.1101/680199>
- 574 Zhaoping, Li. (2019). A new framework for understanding vision from the perspective of the  
575 primary visual cortex. *Curr Opin Neurobiol*, 58, 1-10. doi:10.1016/j.conb.2019.06.001

576

577

578

579

580

581

582

583

584

585

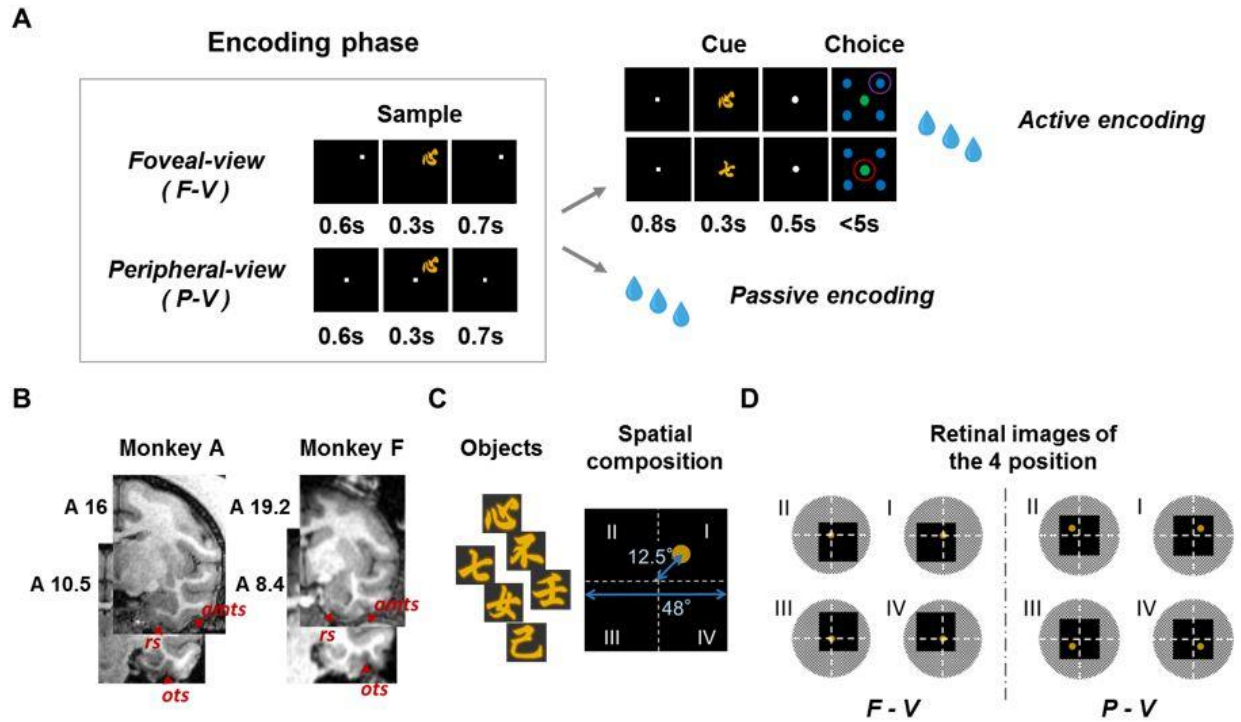
586

587

588

589 **Figures**

590



591

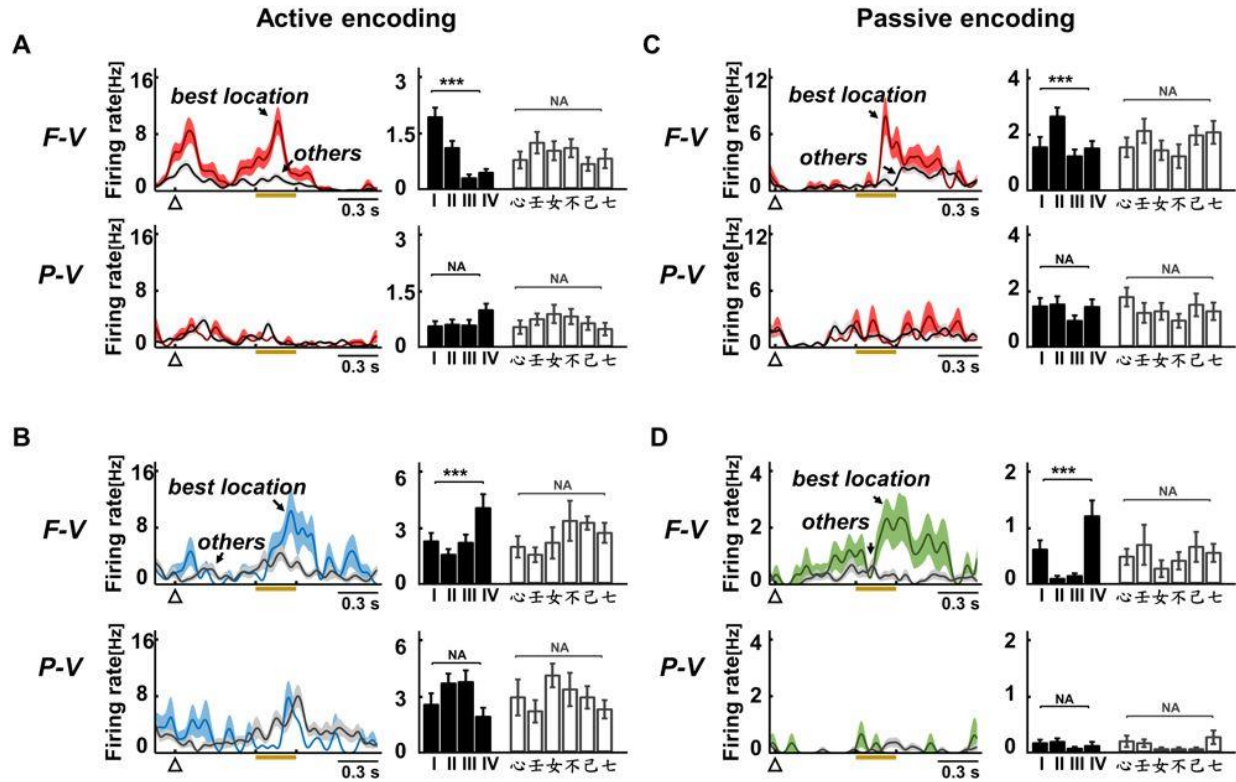
592 **Fig. 1. Encoding of location and item in two view conditions**

593 (A) Schematic diagram of location and item encoding in the F-V and P-V conditions of the  
 594 active-encoding and passive-encoding tasks. In the active-encoding task, the cue stimulus was  
 595 the same as the sample stimulus during the encoding phase in the match trial (Top), while the  
 596 two stimuli differed in the nonmatch trial (Bottom). Red circles indicate correct answers.

597 Passive-encoding task consisted of only the encoding phase of the active-encoding task. (B)

598 Example of coronal sections from monkey A and monkey F. The sections from Monkey A are 16  
 599 mm and 10.5 mm anterior to the interaural line and include the hippocampus (HPC),  
 600 parahippocampal cortex (PHC), perirhinal cortex (PRC), and area TE (TE). amts, anterior middle  
 601 temporal sulcus; ots, occipital temporal sulcus; rs, rhinal sulcus. Coronal sections from monkey F

602 are 19.2 mm and 8.4 mm anterior to the interaural line. (C) Six object stimuli were used in the  
603 task, and an example of spatial composition during the sample period is shown. A yellow disk  
604 indicates an object position. (D) Schematic diagram of visual inputs to the retinae during the  
605 sample period; white dashed lines indicate the horizontal and vertical meridians of the visual  
606 field.  
607

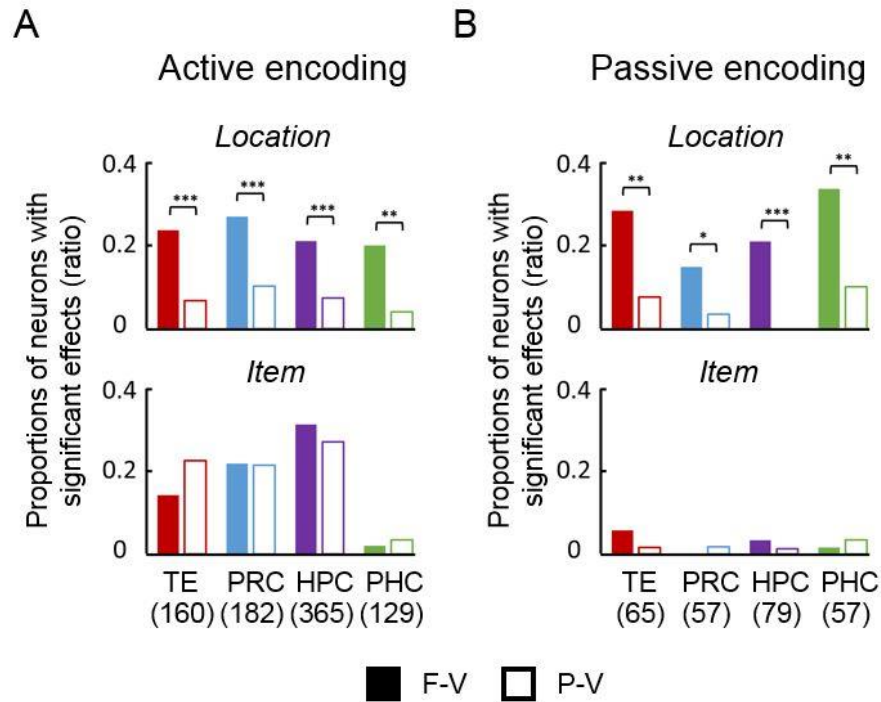


608

609 **Fig. 2. Responses of the location-selective cells in the active-encoding and passive-encoding**  
 610 **task**

611 (A) Example of the location-selective cells from TE in the F-V and P-V condition of the active-  
 612 encoding task. (Left) Spike-density functions (SDFs) ( $\sigma = 20$  ms) indicating the firing rates  
 613 under two conditions (best location and the average of other three locations). (Right) Bar graph  
 614 indicating the mean firing rate during sample period (80-1000 ms after sample on) under each  
 615 location and each item. (B) Example of the location-selective cells from PRC in the F-V and P-V  
 616 condition of the active-encoding task. (C-D) Examples of the location-selective cells in TE (C)  
 617 and PHC (D) in the F-V and P-V conditions of the passive-encoding task.

618



619

620 **Fig. 3. Proportions of location-selective and item-selective cells**

621 (A) Proportions of location-selective cells (Top) and item-selective cells (Bottom) during the

622 sample period (80-1000 ms after sample on) in the F-V (filled bars) and P-V conditions (open

623 bars) in the active-encoding task. Numbers of recorded neurons (tested in both view conditions)

624 are indicated in parentheses. \*\* $P < 0.0016$ ,  $\chi^2 = 10.0$  for PHC, d.f. = 1. \*\*\* $P < 0.0001$ .  $\chi^2 =$

625 19.5, 20.0, and 28.3 for TE, PRC, and HPC, respectively. (B) Proportions of location-selective

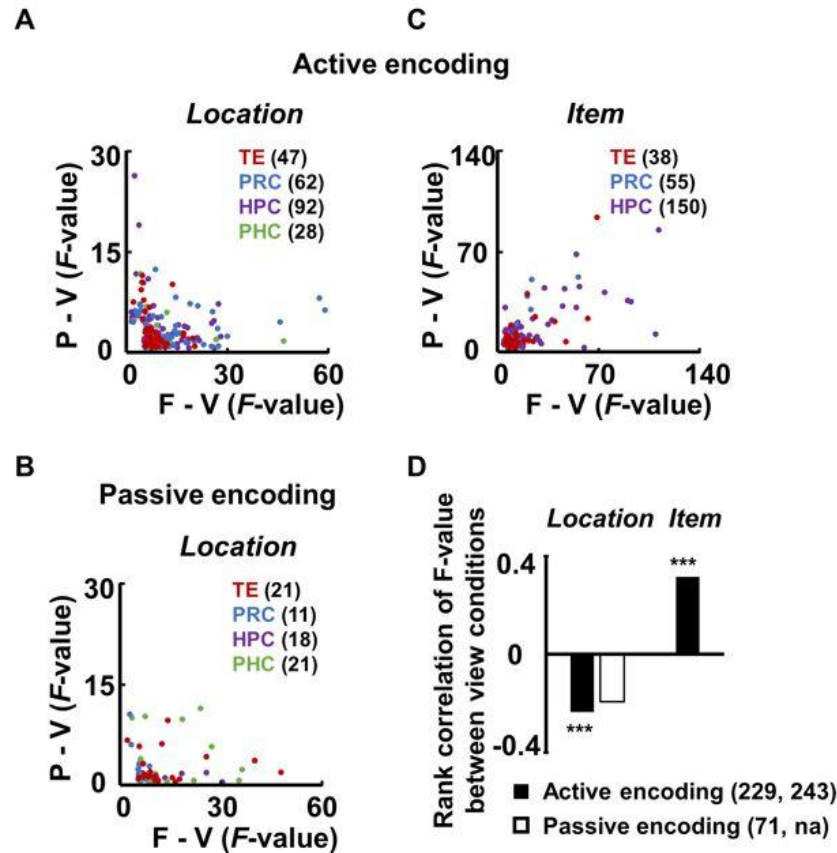
626 cells (Top) and item-selective cells (Bottom) during the sample period in the F-V (filled bars)

627 and P-V conditions (open bars) in the passive-encoding task. \* $P < 0.026$ ,  $\chi^2 = 4.9$  for PRC, d.f. =

628 1. \*\* $P < 0.005$ .  $\chi^2 = 11.1$  and 8.7 for TE and PHC, respectively. \*\*\* $P < 0.0001$ .  $\chi^2 = 20.3$  for

629 HPC.

630



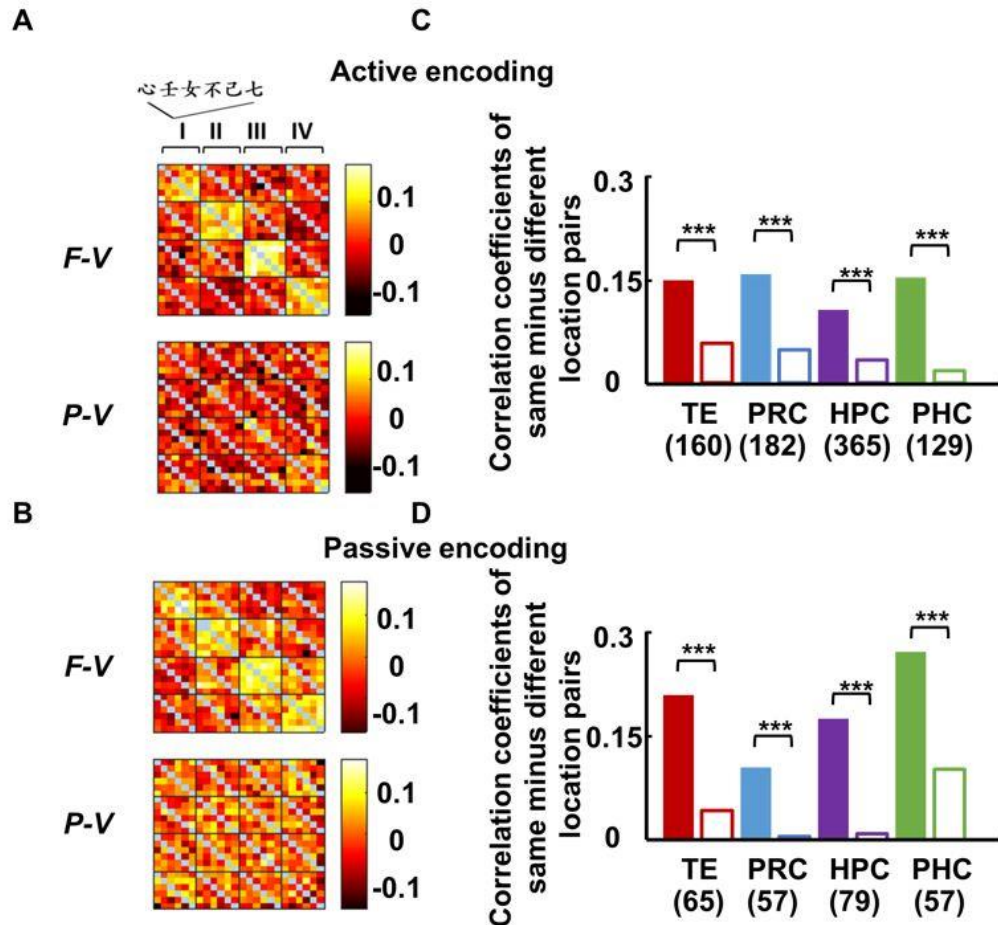
631

632 **Fig. 4. Location and item signal intensity between the two view conditions**

633 (A) Location effect of the location-selective cells in the F-V and P-V conditions of the active-  
 634 encoding task. F values in the P-V condition are plotted against those in the F-V condition for  
 635 location-selective cells in either of the two view conditions. Neurons showing significant effects  
 636 in either of the two conditions were used for the calculation of the F values. Numbers of the  
 637 location-selective cells used for final calculation in each region are indicated in parentheses. (B)  
 638 Location effect in the F-V and P-V conditions of the passive-encoding task. (C) Item effect of  
 639 the item-selective cells in the two view conditions of the active-encoding task. The axis ranges in  
 640 A-C were adjusted for display purpose, which included majorities of the data sets (A: 97.8%, B:  
 641 98.6%, C: 99.6%). (D) Correlation of the signal intensity between the two view conditions. Data  
 642 from MTL and TEv were merged in the active-encoding and passive-encoding tasks,



643 respectively. The total numbers of location-selective and item-selective cells used for final  
644 calculation are indicated in parentheses (left and right, respectively). na, not accountable. P =  
645 0.0003, 0.0000(3.4E-07) and 0.09;  $\rho = -0.24, 0.32, \text{ and } -0.20$ ; d.f. = 227, 241 and 69 for the  
646 active-encoding (location), active-encoding (item), and passive-encoding (location), respectively.  
647 Spearman's rank correlation, two-tailed.  
648



649

650 **Fig. 5. Location effects at population level**

651 (A-B) Correlation coefficients of each pair out of the full 24 (four locations × six items)\* 24

652 (four locations × six items) population vectors in the HPC under the F-V and P-V conditions of

653 the (A) active-encoding and (B) passive-encoding tasks. Correlation coefficients of dummy data

654 sets with location labels randomly shuffled (n=1000) were subtracted from the raw correlation

655 coefficients. All recorded neurons from HPC were used in this analysis. Pearson's linear

656 correlation coefficient. (C-D) Difference value between the mean correlation coefficient under

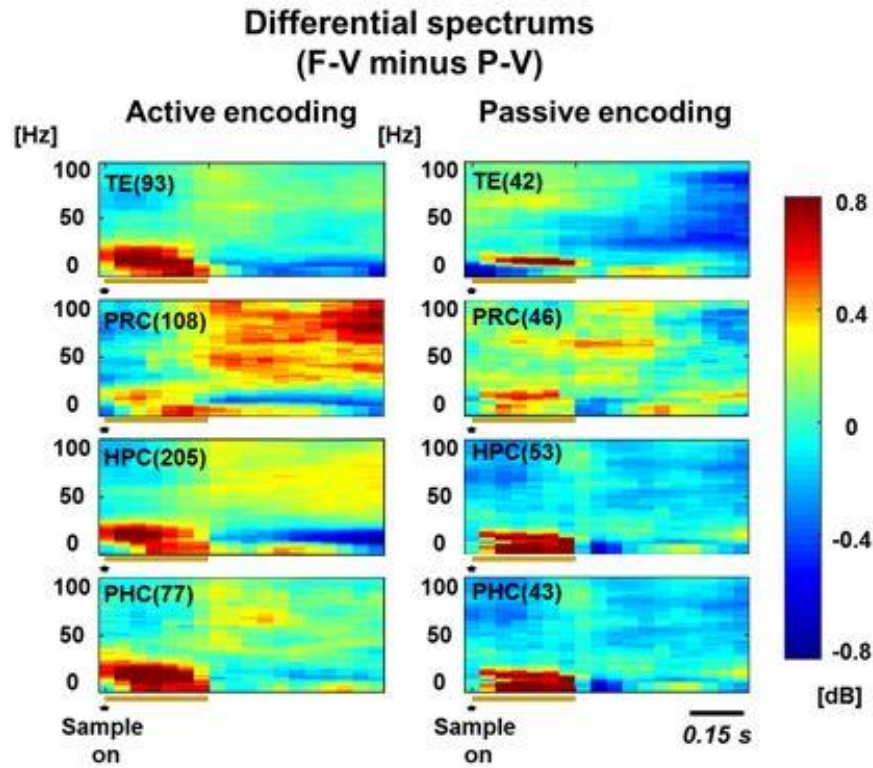
657 the same and different location pairs in the F-V and P-V conditions of the (C) active-encoding

658 and (D) passive-encoding tasks in each brain region. The correlation coefficients between the

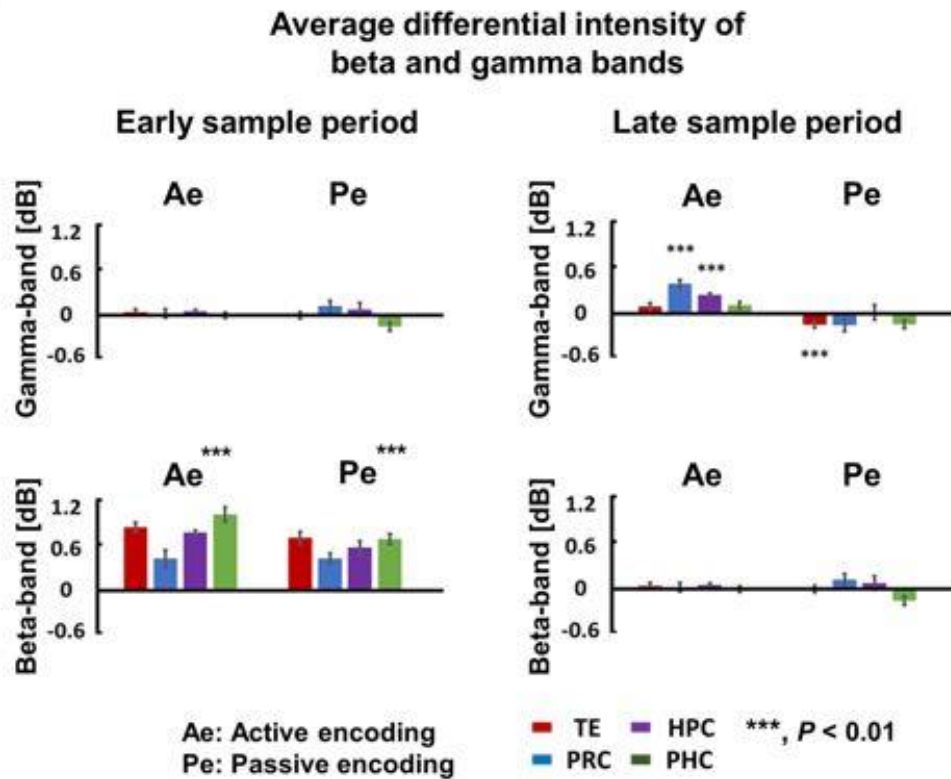
659 population vectors for the trial-types with the same items (i.e., a diagonal line of each small  
660 matrix sorted by the locations, blue pixels in Figs. 5A&B) were excluded from this analysis.

661

A



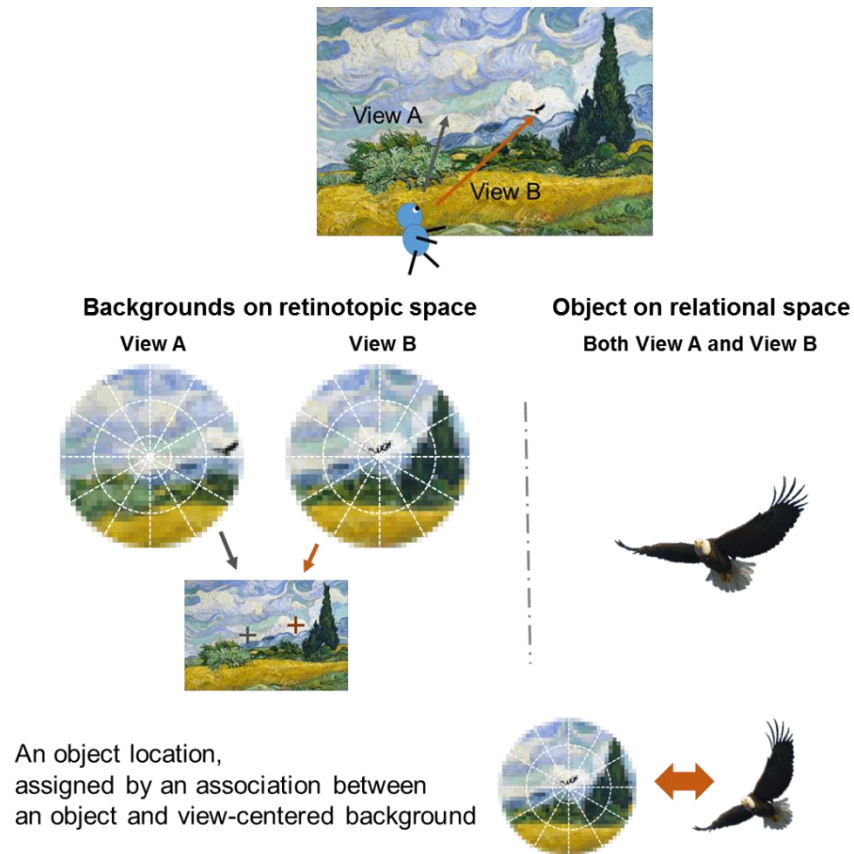
B



663 **Fig. 6. The difference spectrums of the local field potential activities between the F-V and**  
664 **P-V conditions in the active-encoding and passive-encoding tasks**

665 (A) The average difference local field potential (LFP) spectrums between the F-V and P-V  
666 conditions of active-encoding and passive-encoding tasks during the sample period ( 0:1000  
667 msec after sample on). Raw spectrums from different recording sites (indicated in parentheses)  
668 were used for this analysis. Average intensity of each frequency during the baseline period  
669 (600:0 msec before sample on) was subtracted at the corresponding frequency. (B) The average  
670 difference value of beta-band (1-25 Hz) and gamma-band (30-80 Hz) intensity during early  
671 sample (0:300 msec after sample on) and late sample (350:800 msec after sample on) periods  
672 (see Fig. 6A) in each task and recording region.

673



674

675 **Fig. 7. Parallel scene processing on the retinotopic and relational spaces.**

676 *Top*, Assume that a subject was in wheat field and viewing the valley. An eagle was in  
677 parafoveal vision of the subject in view A, while it was in the subject's foveal vision in View B.

678 *Middle*, When the subject attended the eagle either voluntarily or involuntarily, the eagle would  
679 be selected as an object from the retinotopic image and processed on the relational space (*right*)

680 regardless of its original retinotopic. Conversely, background images would be automatically  
681 captured and processed on the retinotopic space, which specify a location of the view point in the

682 scene accordingly (*left*). *Bottom*, The location of the object in the scene would be assigned by an  
683 associated information between the view-centered background and the object. This model

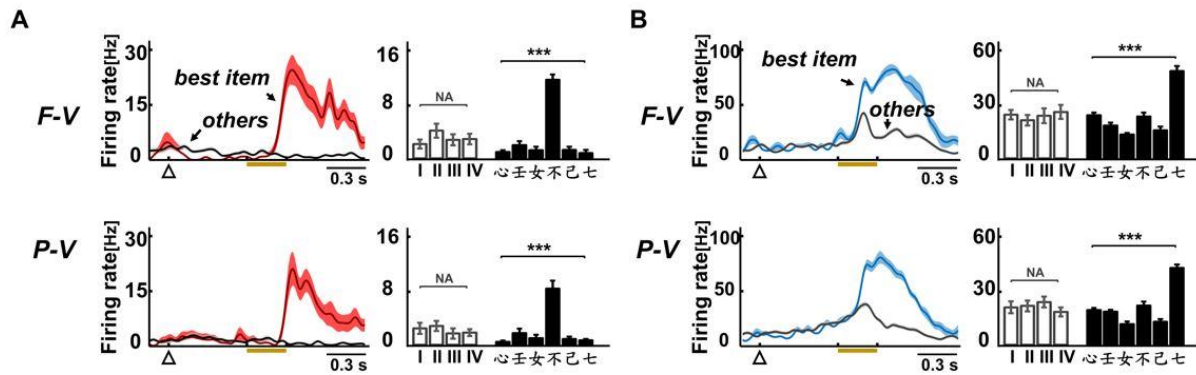
684 hypothesizes that first person's perspective of a scene containing objects depends on the parallel

685 visual processing on the retinotopic and relational spaces, and their association. The original  
686 painting is titled Wheat Field with Cypresses by Vincent Willem van Gogh.

687

688 **Supplementary Figures**

689



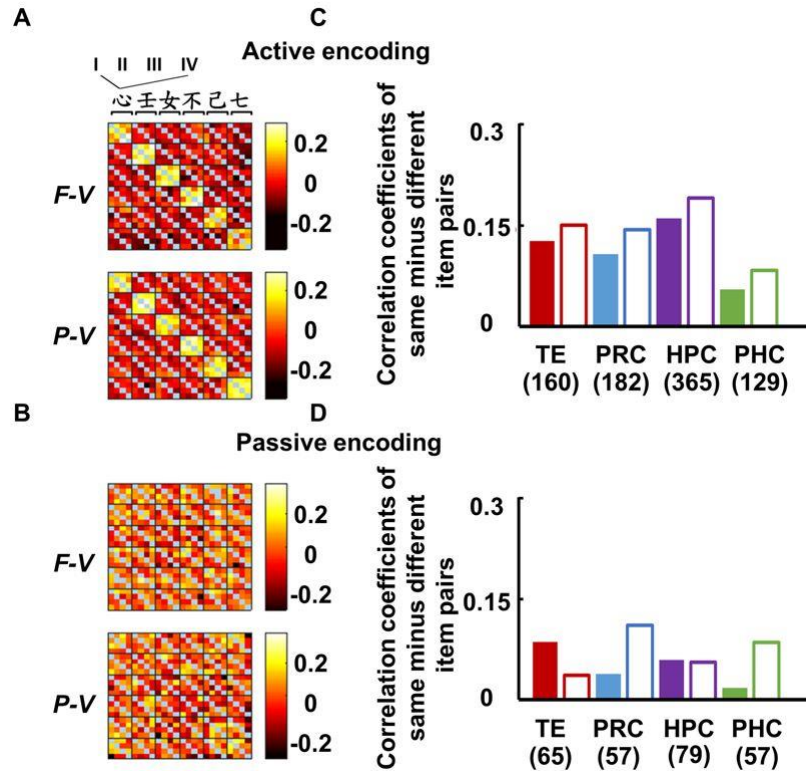
690

691 **Fig. S1. Responses of the item-selective cells in the active-encoding task**

692 (A) Example of the item-selective cells from TE in the F-V and P-V conditions of the active-  
693 encoding task. (Left) Spike-density functions (SDFs) (sigma = 20 ms) indicating the firing rates  
694 under two conditions (best item and the average of other five items). (Right) Bar graph indicating  
695 the mean firing rate during sample period (80-1000 ms after sample on) under each location and  
696 each item. (B) Example of the item-selective cells from PRC in the F-V and P-V conditions of  
697 the active-encoding task.

698





699

700 **Fig. S2. Item effects at population level**

701 (A-B) Correlation coefficients of each pair out of the full 24 (six items  $\times$  four locations) \* 24 (six  
702 items  $\times$  four locations) population vectors in the HPC under the F-V and P-V conditions of the  
703 (A) active-encoding and (B) passive-encoding tasks. Correlation coefficients of dummy data sets  
704 with item labels randomly shuffled (n=1000) were subtracted from the raw correlation  
705 coefficients. All recorded neurons from HPC were used in this analysis. Pearson's linear  
706 correlation coefficient. (C-D) Difference value between the mean correlation coefficient under  
707 the same and different item pairs in the F-V and P-V conditions of the (C) active-encoding and  
708 (D) passive-encoding tasks in each brain region. The correlation coefficients between the  
709 population vectors for the trial-types with the same items (i.e., a diagonal line of each small  
710 matrix sorted by the items, blue pixels in Figs. S2A&B) were excluded from this analysis.