

1 How does the brain represent the 2 semantic content of an image?

3 Huawei Xu^{1,2}, Ming Liu^{1,2,3*}, Delong Zhang^{1,2,3*}

*For correspondence:

lium@scnu.edu.cn (ML);
delong.zhang@m.scnu.edu.cn (DZ)

4 ¹Key Laboratory of Brain, Cognition and Education Sciences, Ministry of Education,
5 China; ²School of Psychology, Center for Studies of Psychological Application, and
6 Guangdong Key Laboratory of Mental Health and Cognitive Science, South China
7 Normal University, China; ³Plateau Brain Science Research Center, South China Normal
8 University/Tibet University, Guangzhou/Lhasa, China

10 **Abstract** Using deep neural networks (DNNs) as models to explore the biological brain is
11 controversial, which is mainly due to the impenetrability of DNNs. Inspired by neural style
12 transfer, we circumvented this problem by using deep features that were given a clear
13 meaning—the representation of the semantic content of an image. Using encoding models and
14 the representational similarity analysis, we quantitatively showed that the deep features which
15 represented the semantic content of an image mainly modulated the activity of voxels in the
16 early visual areas (V1, V2, and V3) and these features were essentially depictive but also
17 propositional. This result is in line with the core viewpoint of the grounded cognition to some
18 extent, which suggested that the representation of information in our brain is essentially
19 depictive and can implement symbolic functions naturally.

21 Introduction

22 Deep neural networks (DNNs) for image recognition provided an important tool for understand-
23 ing the nature of visual object recognition (*Cichy and Kaiser, 2019; Glaser et al., 2019; Kriegeskorte, 2015; Lindsay, 2020; Richards et al., 2019*). This is not only because DNNs are currently the
24 only models known to achieve near human-level performance in object recognition, but also be-
25 cause they have the properties such as the hierarchical organization and the parallel distributed
26 processing which are similar to the visual ventral stream—key circuits that underlie visual object
27 recognition (*DiCarlo and Cox, 2007; DiCarlo et al., 2012*). Using DNNs as computational models,
28 researchers found that DNNs could predict brain activity of visual processing across multiple hier-
29 archical levels at unprecedented accuracy for both macaque (*Cadena et al., 2019; Khaligh-Razavi
30 and Kriegeskorte, 2014; Yamins et al., 2014; Yamins and DiCarlo, 2016*) and human (*Eickenberg
31 et al., 2017; Guclu and van Gerven, 2015; Naselaris et al., 2015; Seeliger et al., 2018*) wherein later
32 layers in DNNs better predict higher areas of the visual ventral stream. The predictive power of
33 DNNs made “mind-reading” possible (*Horikawa and Kamitani, 2017; Shen et al., 2019; Wen et al.,
34 2017*) and promoted the integration of neuroscience and artificial intelligence (*Barrett et al., 2018;
35 Hassabis et al., 2017*).

36 Besides the predictive power, an ideal model should also possess the explanatory power, which
37 means that we should know how the model works (*Kay, 2018*). This is not the case of DNNs. DNNs
38 are essentially black boxes and we can not understand how the input data were transformed into
39 model output (*Rudin, 2019*). This is mainly due to the end-to-end learning and the huge number
40 of parameters in DNNs (the complex architectures of DNNs). For example, AlexNet has about 60
41 million self-learned parameters (*Krizhevsky et al., 2017*) and VGG16 has 138 million self-learned pa-
42

43 rameters (*Simonyan and Zisserman, 2014*). Even though we know the exact value of all parameters
44 for each input, we still can not understand what do these parameters really mean. So using DNNs
45 as models to explore the biological brain is something like replacing a black box with another, the
46 lack of explanatory made it controversial (*Cichy and Kaiser, 2019*). To open the black box and look
47 inside, researchers developed methods such as network dissection (*Zhou et al., 2019*) and visual-
48 ization (*Mahendran and Vedaldi, 2015; Nguyen et al., 2019; Olah et al., 2017; Yosinski et al., 2015;*
49 *Zeiler and Fergus, 2014*), and experimented with network architecture (*Kar et al., 2019*), learning
50 algorithm (*Han et al., 2019; Zhuang et al., 2021*), and input statistics (*Geirhos et al., 2018*). But none
51 of them can directly explain the meaning of the parameters (deep features) learned by DNNs.

52 However, an interesting application of DNNs may give us a hint about the meanings of some
53 deep features. Neural style transfer (NST) is a computer vision technique that allows us to render
54 the semantic content of an image in the style of another (*Jing et al., 2020; Gatys et al., 2016, 2017*).
55 Using NST, for example, we can blend a photo with van Gogh's "Sunflowers" to get a new image
56 which preserve the content of the photo but looks like if it was painted by van Gogh. According
57 to the seminal work of *Gatys et al. (2016)*, the implementation of the original NST algorithm was
58 based on a DNN optimized for object recognition—VGG19. This process took two images, a content
59 image and a style image. First, two images were fed into the pre-trained VGG19 model to extract
60 feature maps, respectively. Second, the feature maps of the conv4_2 layer of the content image
61 were selected as the semantic content representation. Third, the feature maps of the conv1_1 layer,
62 conv2_1 layer, conv3_1 layer, conv4_1 layer, and conv5_1 layer of the style image were selected to
63 compute the Gram matrix as the style representation. Last, through jointly minimizing the distance
64 of the feature representations of a white noise image from the content representation and the
65 style representation (feature inversion using the same VGG19 model), a new image was generated
66 which simultaneously match the content of the content image and the style of the style image.
67 The key to NST lies in the ability to extract representation from an image which explicitly separate
68 image content from style (*Gatys et al., 2016*).

69 In this study, we focused on the feature maps of the conv4_2 layer of the VGG19, which were
70 selected as the representation of the semantic content of an image in the original NST algorithm.
71 Although there was no clear explanation about why choose the layer conv4_2 as the semantic con-
72 tent representation of an image, NST was indeed effective and led to many successful applications
73 (e.g., Prisma). So it gave us an opportunity to explore the question of how does the brain repre-
74 sent the semantic content of an image. We used voxel-wise encoding models (*Kriegeskorte and*
75 *Douglas, 2019; Naselaris et al., 2011; van Gerven, 2017*) to answer this question, which could test
76 hypotheses about how information is represented in our brain. The results showed that, the deep
77 features, which represented the semantic content of an image, mainly modulated the activity of
78 voxels in the early visual areas (V1, V2, and V3). These semantics-related features mainly modu-
79 lated the voxels in the early visual areas rather than those in the higher visual areas naturally led
80 us to another question—what these features really are. For this question, we constructed encod-
81 ing models based on Gabor features which also modulated the activity of voxels in the early visual
82 areas (*Kay et al., 2008*) to compare it with encoding models using the deep features and used
83 representational similarity analysis (RSA, *Kriegeskorte, 2008; Kriegeskorte and Kievit, 2013; Nili*
84 *et al., 2014*) to explore the representational similarity between the representation of the semantic
85 content of an image and other representations such as the representation of semantics and the
86 representation of Gabor features. We found that these features were essentially depictive but also
87 propositional. It is in line with the core viewpoint of the grounded cognition (*Barsalou, 2008, 2010,*
88 *2020*) to some extent, which suggested that the representation of information in our brain was
89 essentially depictive and could implement symbolic functions naturally.

90 Results

91 The fMRI data we used was from *Horikawa and Kamitani (2017)*, which contained a training set
92 (subject viewed 1200 natural images), a testing set (subject viewed 50 natural images), and an im-

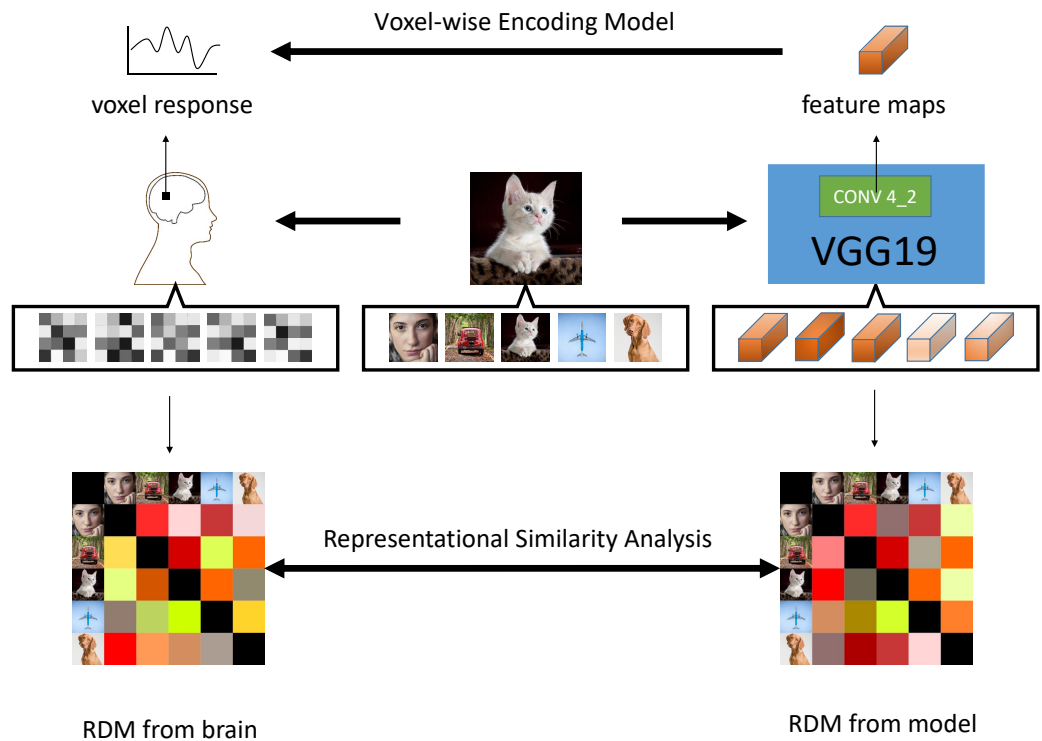


Figure 1. The Experimental Workflow: encoding models and representational similarity analysis.

93 agery set (subject imagined objects according to 50 nouns) for each of 5 subjects. In addition, 7
94 ROIs (V1, V2, V3, V4, LOC, FFA, and PPA) were identified for each subject (The fMRI data only con-
95 tained these ROIs). To explore how does the brain represent the semantic content of an image, we
96 extracted deep features which represented the semantic content of images from the conv4_2 layer
97 of a pre-trained VGG19 and constructed lasso-regularized linear models to predict voxel responses
98 from these features using the training sets for each voxel in each subject. Once models were fit-
99 ted, we used the testing sets and the imagery sets to evaluate models with Pearson's correlation
100 coefficient (r) for each voxel and decoding performance (identifying stimuli from measured brain
101 activity) for all survived voxels. After that we analysed survived models to find the relationship
102 between features and voxels. Because the models showed that these deep features mainly mod-
103 ulated the activity of voxels in the early visual areas, we then compared it with encoding models
104 using Gabor features—the low-level visual features which have been proved to tune simple cells
105 in primary visual cortex (*Hubel and Wiesel, 1962*) and modulate the activity of voxels in the early
106 visual areas (*Kay et al., 2008*). To further explore what the representation of the semantic content
107 of an image really is, we constructed three types of representational distance matrixs (RDMs)—the
108 RDMs from VGG19 (the conv5_4 layer and the fc2 layer), the RDMs from brain activity (7 ROIs for
109 each subject), and the RDMs directly from stimuli (Gabor features, silhouette, and semantics)—to
110 compare them with the RDM of the conv4_2 layer using the testing sets and the imagery sets (see
111 *Figure 1* and Methods and Materials).

112 **The deep features which represented the semantic content of an image mainly**
113 **modulated the activity of voxels in the early visual areas.**

114 Because only 4 models survived in the imagery sets (2 for S2, 1 for S3, and 1 for S5), the following
115 analysis mainly focused on the testing sets. The number of survived models (voxels) in the testing
116 sets was 201 of 4466 for S1, 356 of 4404 for S2, 789 of 4643 for S3, 701 of 4133 for S4, and 369
117 of 4370 for S5. The distribution of survived models in ROIs was different between 5 Subjects (*Fig-*

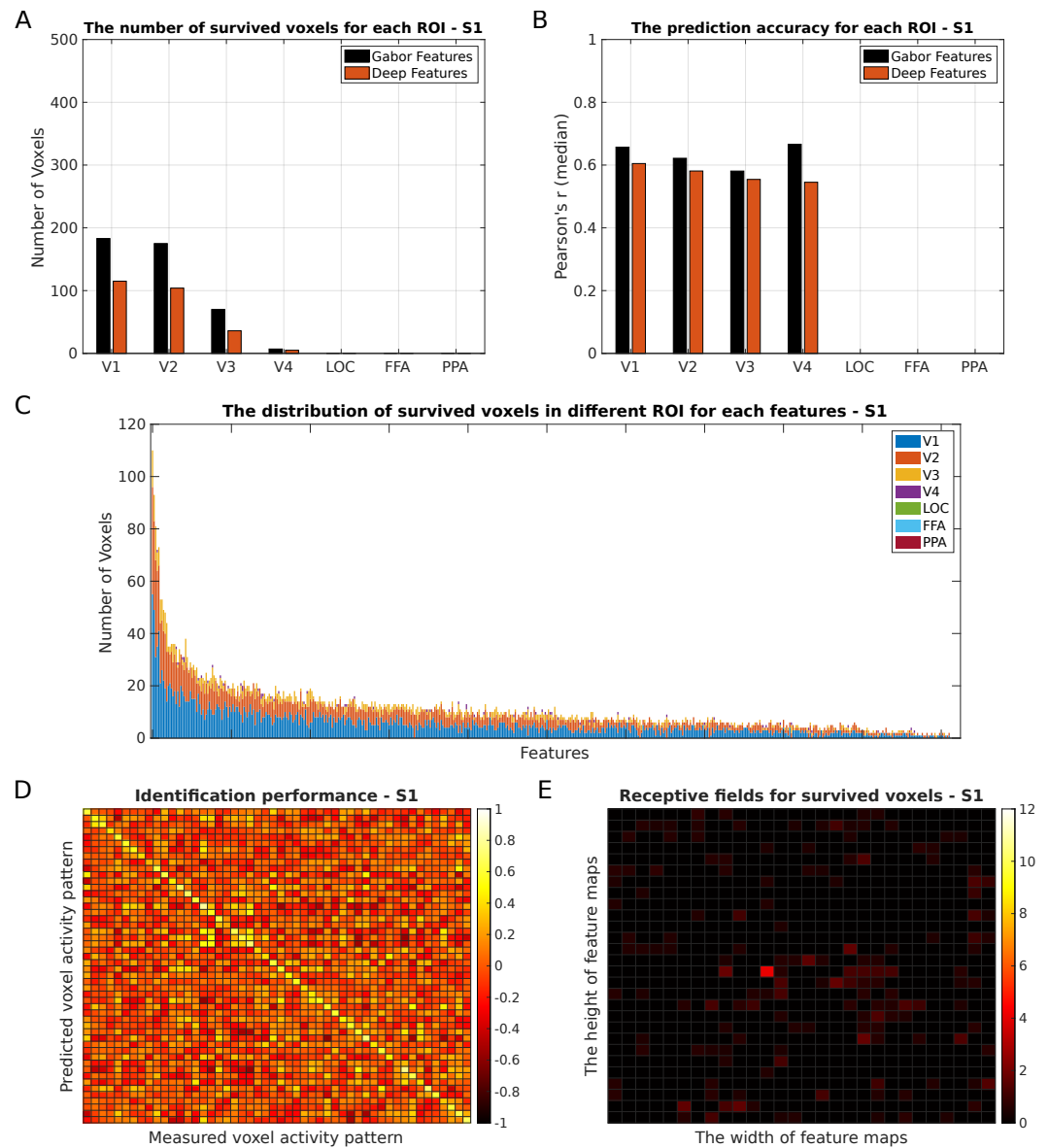


Figure 2. The result of encoding models based on the deep features for S1. (A) The number of survived voxels for each ROI. (B) The prediction accuracy for each ROI (the median of Pearson's correlation coefficients of all survived voxels in each ROI). (C) The distribution of survived voxels in different ROI for each features. The features were ranked according to the number of voxels its related. (D) The decoding performance for S1 (identifying stimuli from measured brain activity using the testing set). (E) The distribution of receptive fields for survived voxels on the feature map. The value of each location equaled the sum of prediction accuracy (r) for all survived voxels located in that location.

Figure 2-Figure supplement 1. The result of encoding models based on the deep features for S2.

Figure 2-Figure supplement 2. The result of encoding models based on the deep features for S3.

Figure 2-Figure supplement 3. The result of encoding models based on the deep features for S4.

Figure 2-Figure supplement 4. The result of encoding models based on the deep features for S5.

Table 1. Top-5 features (Index)

Feature Order	Subject1	Subject2	Subject3	Subject4	Subject5
1	383	250	250	383	383
2	60	60	383	60	250
3	289	383	60	250	60
4	250	289	355	448	482
5	355	198	289	355	289

118 *ure 2.A* for S1, *Figure 2-Figure* supplement for other subjects). There were only 4 ROIs left in S1: V1
119 (115), V2 (104), V3 (36), and V4 (5). The number of ROIs left in S5 was 6: V1 (223), V2 (185), V3 (53),
120 V4 (14), LOC (1), and PPA (7). All the 7 ROIs left in S2, S3, and S4. For S2, V1 remained 210 voxels;
121 V2 remained 164 voxels; V3 remained 54 voxels; V4 remained 4 voxels; LOC remained 11 voxels;
122 FFA remained 4 voxels; PPA remained 1 voxels. For S3, V1 remained 334 voxels; V2 remained 356
123 voxels; V3 remained 180 voxels; V4 remained 58 voxels; LOC remained 48 voxels; FFA remained
124 43 voxels; PPA remained 20 voxels. For S4, V1 remained 328 voxels; V2 remained 308 voxels; V3
125 remained 177 voxels; V4 remained 28 voxels; LOC remained 20 voxels; FFA remained 37 voxels;
126 PPA remained 17 voxels (Because some voxels simultaneously belonged to two different ROIs, the
127 sum of the number of all the voxels in different ROIs may be larger than the total number of the
128 survived voxels for each subject).

129 Another measurement for prediction accuracy in different ROIs is the median of Pearson's cor-
130 relation coefficients of all survived voxels in each ROI (*Figure 2.B* for S1, *Figure 2-Figure* supplement
131 for other subjects). The prediction accuracies of the V1, V2, V3, and V4 were 0.60, 0.58, 0.55, and
132 0.55 for S1; The prediction accuracies of the V1, V2, V3, V4, LOC, FFA, and PPA were 0.64, 0.58, 0.57,
133 0.52, 0.54, 0.52, and 0.50 for S2; The prediction accuracies of the V1, V2, V3, V4, LOC, FFA, and
134 PPA were 0.68, 0.65, 0.63, 0.58, 0.57, 0.58, and 0.57 for S3; The prediction accuracies of the V1, V2,
135 V3, V4, LOC, FFA, and PPA were 0.69, 0.65, 0.61, 0.57, 0.55, 0.56, and 0.54 for S4; The prediction
136 accuracies of the V1, V2, V3, V4, LOC, and PPA were 0.63, 0.60, 0.58, 0.54, 0.50, and 0.56 for S5.

137 Because of individual differences in brain structure and function, the pattern of prediction accu-
138 racy across ROIs was different among subjects. But we still observed some clear common trends:
139 the features of the conv4_2 layer of the VGG19, which were selected as the semantic content rep-
140 resentation of an image in the NST algorithm, mainly modulated the activity of voxels in the early
141 visual areas (V1, V2, and V3). First, most of the survived voxels located in the early visual areas for
142 each subject, and the number of survived voxels in other ROIs (V4, LOC, FFA, and PPA) are rela-
143 tively few or just zero; Second, the prediction accuracy for early visual areas were slightly higher
144 than other ROIs.

145 The survived models could be used to decode stimuli from the measured brain activity—image
146 identification using the testing sets. The identification accuracies of 5 subjects (*Figure 2.D*) were
147 92% (46/50), 90% (45/50), 100%, 100%, and 92% (46/50). After checked all the identification errors,
148 we founded that there were some common mistakes among different subjects. All the 4 images
149 (No.17, No.19, No.41, and No.44) that were incorrectly identified by the encoding models of S5
150 were also incorrectly identified in S2, and three of them (No.19, No.41 and No.44) were incorrectly
151 identified in S1 too. The encoding models of S1 made the same mistake as the models of S2, which
152 identified the No.41 image as the No.42 image. And the encoding models of S2 made the same
153 mistake as the models of S5, which identified the No.44 image as the No.26 image and the No.17
154 image as the No.22 image (For copyright reasons, we can not show the actual images).

155 Because Lasso regression enables feature selection, the survived models also described the
156 relationship between features (X) and voxel responses (y) through regression coefficients. From
157 the perspective of voxels, we calculated the number of features each ROI related (median) and
158 found no common trend among subjects. From the perspective of features, we calculated the

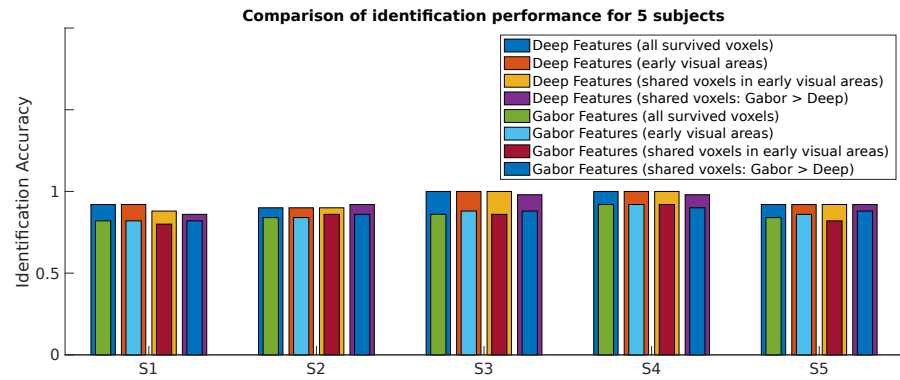


Figure 3. The comparison of identification performance for 5 subjects.

159 number of voxels each feature related and analyzed the location distribution of these voxels in
160 different ROIs. After ranked features according to the number of voxels its related, we found that
161 the deep features were mainly related to the voxels in the early visual areas (**Figure 2.C** for S1, **Fig-**
162 **ure 2-Figure** supplement for other subjects). And we calculated Pearson's correlation coefficients
163 of ranked feature index between each subject pair to examine if there were similar patterns among
164 subjects. The result showed that there was no significant correlation. But if only considered the
165 top-5 features, we found that most of the features were same among subjects (**Table 1**).

166 The features (X) of each survived model corresponded to a spatial location on the feature map
167 (we only used the features from one spatial location of the feature map as X , see **Methods** and
168 **Materials**), which could be seen as the center of population receptive field of the related voxel
169 (y). So we could visualize the distribution of receptive fields of survived voxels for each subject
170 (**Figure 2.E** for S1, **Figure 2-Figure** supplement for other subjects). The result showed that survived
171 voxels distributed widely on the feature map and there was a slight trend that some of voxels
172 clustered near the center of the feature map.

173 **The deep features may contained more information about stimuli than Gabor fea-** 174 **tures.**

175 We also constructed encoding models based on Gabor features and compared it with encoding
176 models using the deep features. Like the encoding models using the deep features, few voxels sur-
177 vived in the imagery sets when encoding models used gabor features (1 for S2 and 1 for S4). So we
178 only compared two different types of encoding models using the test sets. From the perspective
179 of individual voxel, there were more voxels survived in the early visual areas with encoding models
180 based on Gabor features for all subjects(**Figure 2.A** for S1, **Figure 2-Figure** supplement for other
181 subjects). And the prediction accuracy of the early visual areas (the median of Pearson's correla-
182 tion coefficients of all survived voxels in each ROI) was also higher for all subjects when encoding
183 models used Gabor features (**Figure 2.B** for S1, **Figure 2-Figure** supplement for other subjects). It
184 implied that, relative to the deep features, Gabor features were preferentially represented by the
185 early visual areas. From the perspective of activity pattern of voxels (decoding performance), how-
186 ever, the identification performances were better for all subjects when encoding models used the
187 deep features (**Figure 3**).

188 The better identification performance of encoding models using the deep features could be
189 due to the survived voxels in the higher visual areas, so we excluded survived voxels not in the
190 early visual areas for both models and compared identification performance again. The results
191 showed that, for all subjects, there were more survived voxels when encoding models used Gabor
192 features (**Figure 4.B** and **Figure 4.C** for S1, **Figure 4-Figure** supplement for other subjects) but the
193 identification performances were still better when encoding models used the deep features (**Fig-**

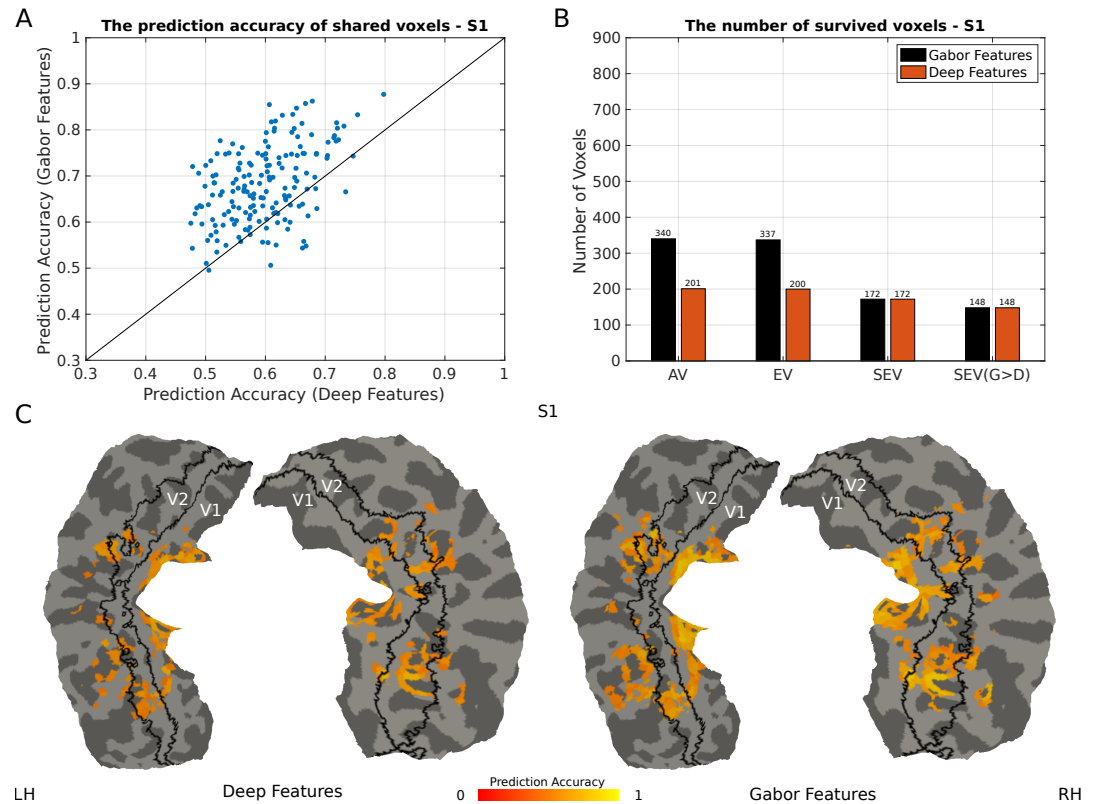


Figure 4. The comparison of two types of encoding models for S1. (A) The prediction accuracy of shared voxels (voxels simultaneously survived in the early visual areas for both models) for S1. (B) The number of survived voxels for S1. AV = all survived voxels, EV = all survived voxels in the early visual areas, SEV = voxels simultaneously survived in the early visual areas for both models, SEV(G>D) = voxels which were better predicted by encoding models using Gabor features in SEV. (C) Prediction accuracy of survived voxels in the early visual areas mapped on the occipital surface. Surface reconstruction and flattening were performed using FreeSurfer (<https://surfer.nmr.mgh.harvard.edu/>).

Figure 4-Figure supplement 1. The comparison of two types of encoding models for S2.

Figure 4-Figure supplement 2. The comparison of two types of encoding models for S3.

Figure 4-Figure supplement 3. The comparison of two types of encoding models for S4.

Figure 4-Figure supplement 4. The comparison of two types of encoding models for S5.

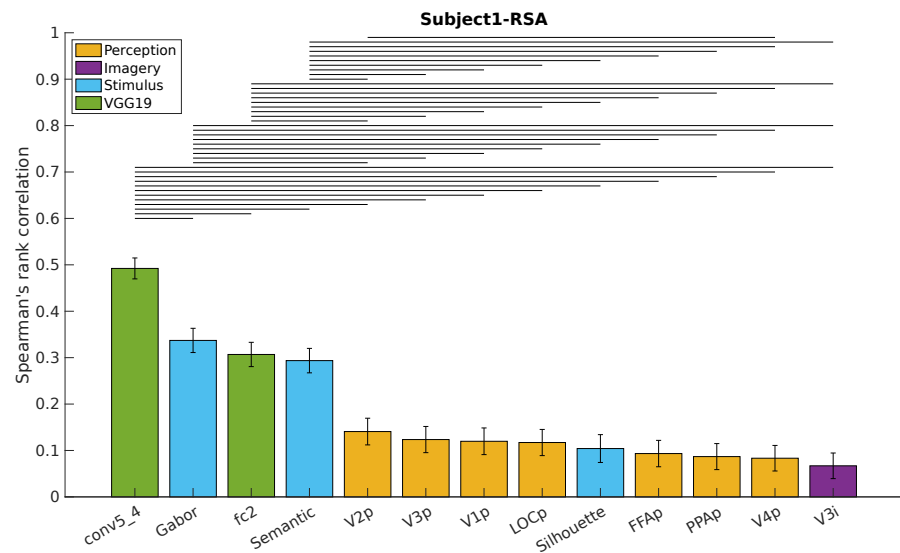


Figure 5. The result of RSA for S1.

Figure 5-Figure supplement 1. The result of RSA for S2.

Figure 5-Figure supplement 2. The result of RSA for S3.

Figure 5-Figure supplement 3. The result of RSA for S4.

Figure 5-Figure supplement 4. The result of RSA for S5.

194 **ure 3).** Then we chosed voxels simultaneously survived in the early visual areas for both models
195 and found that the identification performances were better when encoding models used the deep
196 features for all subjects (**Figure 3**). Further analyses on these voxels, we observed a common trend
197 among all subjects—the prediction accuracies of shared voxels showed a positive correlation be-
198 tween two types of models and most of voxels were better predicted when encoding models used
199 Gabor features (**Figure 4.A** for S1, **Figure 4-Figure supplement** for other subjects). Finally, we only
200 chosed voxels which were better predicted by encoding models using Gabor features from shared
201 voxels to compare the identification performance between two types of models. The result did
202 not change—the identification performances were better when encoding models used the deep
203 features for all subjects (**Figure 3**), which may suggest that there was additional information in the
204 deep features.

205 **The representation of the semantic content of an image did related to the seman-**
206 **tics of the image and also preserved visual details of the image to some extent.**

207 In accordance with the analysis using encoding models, the RSA was also based on individuals
208 (**Figure 5** for S1, **Figure 5-Figure supplement** for other subjects). The RDMs from the pre-trained
209 VGG19 were significantly correlated with the RDM of the layer conv4_2. For all subjects, the RDM
210 of the layer conv5_4 (the last convolutional layer) showed the strongest correlation ($r_s = 0.49$) and
211 the RDM of the layer fc2 (the last fully connected layer before SoftMax layer) was in the second
212 echelon among all candidate RDMs ($r_s = 0.31$). This was a reasonable result given that all three
213 RDMs derived from the same model (VGG19) and the layer conv5_4 was more similar to the layer
214 conv4_2 than the layer fc2.

215 The RDMs from the stimuli were significantly correlated with the RDM of the layer conv4_2, too.
216 For all subjects, the RDM of Gabor features ($r_s = 0.34$) and the RDM of semantics ($r_s = 0.29$) were
217 both in the position of the second echelon. It implied that the representation of the layer conv4_2
218 did relate to the semantics of stimuli and also preserved visual details of stimuli to some extent.
219 The RDM of silhouette was in the position of the lowest echelon for all subjects ($r_s = 0.10$). Because

220 the silhouette of an object provided a limited description of the specific shape of the object, the
221 representation of the layer conv4_2 also related to the specific shape of stimuli.

222 The situation of the RDMs from brain activity was more complex. For each subject, all 7 RDMs
223 from the testing sets were significantly correlated with the RDM of the layer conv4_2. Although
224 there were some individual differences about the relative position of these RDMs, the RDMs from
225 the early visual areas which were in the second or third echelon roughly showed stronger correla-
226 tion than those from the higher visual areas. In contrast, there were few RDMs from the imagery
227 sets significantly related to the RDM of the layer conv4_2 (V3 for S1, PPA for S2, PPA for S3, V3, V4,
228 LOC, and FFA for S5). And all of these RDMs were in the position of the lowest echelon. This result
229 was in line with the result of encoding models to some extent.

230 Discussion

231 The impenetrability of DNNs reduced the explanatory power of studies which used DNNs as com-
232 putational models to explore the biological brain. Inspired by NST, we circumvented this problem
233 by using deep features that were given clear meaning—the representation of the semantic con-
234 tent of an image. Using encoding models, we quantitatively showed that the deep features which
235 represented the semantic content of an image mainly modulated the activity of voxels in the early
236 visual areas. It was a surprise that the semantics-related features mainly modulated the voxels
237 in the early visual areas rather than those in the higher visual areas. Then we compared encod-
238 ing models using the deep features with encoding models using Gabor features which have been
239 proved to modulate the activity of voxels in the early visual areas (*Kay et al., 2008*) and found that
240 the deep features may contained more information about stimuli than Gabor features. These re-
241 sults naturally led us to another question: what the representation of the semantic content of an
242 image really is? The result of RSA showed that, the representation of the semantic content of an
243 image did related to the semantics of the image and also preserved visual details of the image to
244 some extent. It implied that the representation of the semantic content of an image might be a
245 hybrid form—both in propositional format and depictive format.

246 How could the format of a representation be both propositional and depictive? The question
247 of how the information is represented in our brain had been discussed for many years, which
248 was known as the imagery debate (*Pearson and Kosslyn, 2015*). At the heart of the debate was
249 whether all information is represented in a symbolic, propositional format. Convergent evidence
250 from empirical studies of mental imagery suggested that information can be represented in a
251 pictorial, depictive format (*Pearson, 2019*). The existence of the depictive format of information
252 ended the imagery debate but also raised new questions: how many formats can the brain use
253 and what is the relationship between these formats and the propositional format? With the devel-
254 opment of theories of grounded cognition, the dominant position of the propositional format in
255 cognition is being challenged. From the perspective of grounded cognition, there were no amodal
256 symbols in our brain that were independent of the modal representation and all cognitive phe-
257 nomena were ultimately grounded in modal simulations, bodily states, and situated action, which
258 was supported by many researches on perception, memory, language, thought, social cognition,
259 and development (*Barsalou, 2008, 2010, 2020*). This view emphasized the key role of the depic-
260 tive, modality-specific representation in cognition and denied the independent existence of the
261 symbolic, propositional representation, which was clearly articulated by Comenius from several
262 hundred years ago—"things are essential, words only accidental; things are the body, words but
263 the garment; things are the kernel, words the shell and husk. Both should be presented to the intel-
264 lect at the same time, but particularly the things, since they are as much objects of understanding
265 as is language" (*Paivio, 2007*).

266 From this view, the representation of information in our brain is essentially depictive and can
267 implement symbolic functions naturally. This is in line with our result to some extent. On the
268 one hand, the representation of the semantic content of an image (the feature maps of the layer

conv4_2) was essentially depictive. This was because the feature maps extracted from the convolutional layer of the VGG19 naturally preserved the topology of the original image. Besides, The result of RSA also showed that it preserved visual details of the image. On the other hand, this representation did reflect the semantics of the image in some degree. To the best of our knowledge, it is the first time that the existence of such representation in our brain is quantitatively proved, at least in the early visual areas.

Unlike the previous study (*Naselaris et al., 2015*), we did not observe that encoding models which were trained using the perceptual data could successfully predict voxel responses from the imagery data. This could be due to different experimental tasks. In the study of *Naselaris et al. (2015)*, subjects were asked to imagine particular artworks, such as “Betty” by Gerhard Richter and “Horse Bath” by Odd Nerredum. In the study of *Horikawa and Kamitani (2017)*, which provided data for this paper, subjects were asked to imagine as many object images as possible from concrete nouns, such as leopard and swan. The difference between two tasks was whether the imagery had a particular content. For example, when you were asked to draw your cat or dog, what you drew must be a particular cat or dog; but when you were asked to draw a cat or dog, you could draw any cat or dog, even Hello Kitty or Snoopy. Because of the individualization and arbitrariness of the imagery in our study, it seems reasonable that our result was not consistent with the previous study and could not address the issue of the relationship between perception and imagery.

In addition to the obvious individual divergences in encoding mechanisms, our result showed that there were some common mechanisms among subjects (common mistakes in image identification and similarity of the top 5 features). In contrast to the symbolic, propositional representations, the depictive, modality-specific representations of information were grounded in the modalities, the body, and the environment. So they were highly personal and changed from time to time. This was a key difference between the grounded cognition and traditional cognitive theories and could be used to explain individual divergences in cognition. Meanwhile, we did share a common physical and environmental basis, which was also reflected in cognitive process and made communication possible. This may explain the existence of the common mechanisms.

In our study, we quantitatively showed that the deep features which represented the semantic content of an image mainly modulated the activity of voxels in the early visual areas and these features were essentially depictive but also propositional. This result implied that some depictive representation of an object in our brain can naturally reflect semantics of the object to some extent and this phenomena can be found in the early visual areas, which provided empirical evidence to the core viewpoint of the grounded cognition. In fact, there was another theory also addressing the relationship between the propositional representation and the depictive representation of information in our brain—the dual coding theory (*Clark and Paivio, 1991*), which emphasized the beneficial effects of the depictive representation of information on cognition (concreteness) and suggested that the two types of representations are independent from each other in our brain (Paivio believed that there were two distinct subsystems in our brain specialized for dealing with different types of representations). Both theories admitted the association between the two types of representations but disagreed with each other about the relationship between the two types of representations. From our result, we tend to support the monistic view of the two types of representations. But we also noticed that our result only involved the early visual areas which were in the early stages of the visual ventral stream. So How does the depictive representation of an object change along the visual ventral stream to make object recognition possible—whether the independent propositional representation will eventually appear—is not clear. This question needs further studies.

315 **Methods and Materials**

316 **Data**

317 We used the fmri data that was originally published in *Horikawa and Kamitani (2017)*, which can be
318 downloaded from <https://github.com/KamitaniLab/GenericObjectDecoding>. The data was obtained
319 from two fMRI experiments for each of 5 subjects: an image presentation experiment and an im-
320 agery experiment. There were two sessions in the image presentation experiment—the training
321 session and the testing session. In the training session, subjects viewed 1200 images from 150
322 categories (8 images from each category) as each image presented once. In the testing session,
323 subjects viewed 50 images from 50 categories (one image from each category) as each image pre-
324 sented 35 times. In the imagery experiment, subjects were asked to imagine about 50 nouns from
325 50 categories (one noun from each category) as each noun presented 10 times. The categories
326 used in the imagery experiment were the same as those in the testing session, which were not
327 used in the training session. All the images (1250 natural images with the resolution of 512×512×3)
328 and the corresponding categories (200) were collected from ImageNet (*Deng et al., 2009*). In addi-
329 tion, there were standard retinotopic mapping experiment and functional localizer experiment to
330 identify lower visual areas (V1, V2, V3, and V4) and higher visual areas (LOC, FFA and PPA) for each
331 participant. The details of the experimental design, MRI acquisition protocol, and preprocessing of
332 the fMRI data could be found in *Horikawa and Kamitani (2017)*.

333 Before further analysis, we averaged the repeated trials in the testing session and the imagery
334 experiment. First, we standardized the fMRI data from the training session. The mean and stan-
335 dard deviation of the training set were then used to standardize the testing sets (from the test-
336 ing session) and the imagery sets (from the imagery experiment). After that, we performed trial-
337 averaging for the testing sets and the imagery sets to improve the signal-noise ratio. Because of
338 trial-averaging, there were statistical difference between the training set and the other two. So we
339 rescaled the averaged data by a factor of \sqrt{n} where n is the number of trials averaged (*Shen et al.,*
340 *2019*).

341 **Encoding Models Based On Deep Features**

342 To probe how does the brain represent the semantic content of an image, we used voxel-wise
343 encoding models. Such models were constructed separately for each voxel in each individual, so
344 our analysis was also individual-based. There were two steps to construct voxel-wise encoding
345 models (*Naselaris et al., 2011*): the first step was a nonlinear transformation from a stimulus space
346 to a feature space; the second step was a linear transformation from the feature space to a voxel
347 space.

348 As the first step, we got deep features which represented the semantic content of an image. We
349 employed the pre-trained VGG19 model based on the open source machine learning framework of
350 PyTorch (*Paszke et al., 2019*) to extract the feature maps of the conv4_2 layer as the image content
351 representation. After image preprocessing (such as image scaling and cropping, more details can
352 be found from https://pytorch.org/hub/pytorch_vision_vgg/), images from the training data session
353 and the testing data session were fed into the VGG19 model and the feature maps of the conv4_2
354 layer were extracted, respectively. As a result, the size of the training feature maps was [1200, 512,
355 28, 28] ([the number of images, the number of kernels, the height of feature map, the width of
356 feature map]), and the size of the testing feature maps was [50, 512, 28, 28].

357 As the second step, we constructed linear regression models to predict the brain activity evoked
358 by an image from the features which represented the semantic content of the same image. For
359 each voxel, the model can be expressed by

$$y = X\beta + \epsilon \quad (1)$$

360 Where y is a measured voxel response and X is features which elicited the response of the voxel.

361 Just as each neuron has its own receptive field, each voxel has its own population receptive field
362 (*Dumoulin and Wandell, 2008*). It means that a voxel only responds to the features in its population

363 receptive field. So there was no need to put all the features into the model (The number of features
364 for each image is 401408, It will pose a problem known as the curse of dimensionality if we put all
365 the features into the model). And the features we extracted from the VGG19 model were naturally
366 organized into feature maps that preserved the topology of stimuli. For example, the size of the
367 feature maps for each image was $512 \times 28 \times 28$. It could be seen as 784 (28×28) spatial locations with
368 512 features at each spatial location. The spatial arrangement of 784 locations (28×28) preserved
369 the topology of the original image (512×512). According to the model of the population receptive
370 field, receptive fields are center-surround organized and features at the center of a receptive field
371 make the largest contribution to the activity measured in the voxel (*St-Yves and Naselaris, 2018*).
372 So we only used 512 features at one of the 784 locations as X for each voxel (To simplify the model
373 and reduce computation time, we ignored the surround of receptive fields and assumed that each
374 location is a candidate for the center of a receptive field). To find the best center for each voxel,
375 we constructed separate linear regression models for each location/voxel combination.

376 We used the training data to fit models. For each model, X was the features from one of the
377 784 locations represented by a 1200×512 matrix (1200 images), y was a measured voxel response
378 represented by a 1200×1 matrix. It was reasonable to assume that a voxel only responded to a
379 fraction of the X . So we estimated regression coefficients of each model using lasso regression:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad (2)$$

380 where lambda is a complexity parameter that controls the amount of regularization. To accelerate
381 model fitting, we used the function “lasso_gpu” from a MATLAB package developed by *Mohr and*
382 *Ruge (2020)*, which can be efficiently implemented in parallel on a GPU. The optimal value of lambda
383 was selected from a lambda sequence with 10 lambdas: $2^{-2}, 2^{-3}, 2^{-4}, \dots, 2^{-11}$ (the first lambda was
384 set using the function “calculate_lambda_start” from the same MATLAB package and the lambda
385 range was set to guarantee that the last lambda is not chosen as the best one). We chose the best
386 lambda and the best location for each voxel using five-fold cross-validation with the coefficient of
387 determination (R^2).

388 **Model Evaluation**

389 Once fitted, encoding models were evaluated using the testing sets and the imagery sets, respec-
390 tively. For each voxel, we defined the model’s prediction accuracy as the Pearson’s correlation
391 coefficient (r) between the measured voxel response and the response predicted by the model.
392 The significance of the correlation was assessed by a permutation test with 10000 permutations
393 (Bonferroni correction for the number of voxels, and $p < 0.05$). For each ROI, we used the number
394 of survived voxels within the ROI and the median r as measurements for prediction accuracy.

395 Another measurement for prediction accuracy was decoding performance—identifying stimuli
396 from measured brain activity (*Kay et al., 2008*). First, we used survived models to predict the voxel
397 activity pattern from the test feature maps for each of the 50 stimuli. Second, we calculated the
398 Pearson’s correlation coefficient between the predicted voxel activity pattern and the measured
399 voxel activity pattern for each stimulus/stimulus combination. The stimulus whose predicted voxel
400 activity pattern was most correlated with the measured voxel activity pattern of itself was regarded
401 as correct decoding. We defined the identification accuracy as the percentage of stimuli that are
402 correctly identified from the testing data or the imagery data.

403 **Model Analysis**

404 After model evaluation, we explored the relationship between features and voxels through regres-
405 sion coefficients of survived models. Because of the lasso regression, some of regression coeffi-
406 cients in each model were set to zero automatically. The nonzero coefficients indicated that which
407 features were related to the activity of a voxel and how are they related. For each subject, we de-
408 scribed the relationship between features and voxels from two complementary perspective—the

409 perspective of features and the perspective of voxels. From the perspective of features, we used
410 the number of voxels each feature related and the location distribution of these voxels in differ-
411 ent ROIs to analyse the property of different features. From the perspective of voxels, we used
412 the number of features each voxel related and the location of the voxel to analyse the property of
413 different ROIs.

414 Furthermore, the survived models also provided information about the population receptive
415 fields of voxels. For each model, the X was selected from one of the 784 locations (28×28). The
416 location of the X could be seen as the center of the population receptive field of the correspond-
417 ing voxel on the feature maps. We used heatmap to visualize and explore the distribution of the
418 receptive fields of survived voxels for each subject.

419 **Encoding Models Based On Gabor Features**

420 To compare with encoding models using the deep features, we also constructed encoding models
421 based on Gabor features. We got gabor features of stimuli according to the method of *Li et al.*
422 (2018): Firstly, a Gabor Wavelet Pyramid (GWP, *Gaziv, 2021*) model was used to get original Gabor
423 features from stimuli (six spatial frequencies: 1, 2, 4, 8, 16, and 32 cycles/FOV; eight orientations: 0° ,
424 22.5° , 45° , ..., and 157.5° ; and two phases: 0° and 90° ; The FOV covered full of a image and all images
425 were downsampled to 128×128 pixels); Secondly, the absolute values of the projections of each
426 quadrature wavelet pair were averaged to get the contrast energies; Thirdly, the contrast energies
427 were normalized to linearize the relationship between contrast energies and voxel responses (each
428 contrast energy was divided by the sum of the contrast energy and the median of all contrast
429 energies in the training set which were at the same position and the same orientation); Fourthly,
430 the normalized contrast energies at the same position (eight orientations) were averaged to reduce
431 the dimension of features; Lastly, the average luminance of stimuli were also added into Gabor
432 features. As a result, each stimulus had 1366 features. After that we used the same method as
433 described above to construct linear regression models from Gabor features to voxel responses.

434 **Representational Similarity Analysis**

435 To probe the representation property of the layer conv4_2, we used RSA, which characterized the
436 representational geometry of a set of stimuli in a brain region or computational model by repre-
437 sentational distance matrix (RDM) and compared RDMs to explore the representational similarity
438 between different brain regions or brain regions and computation models. Three types of RDMs
439 (the candidate RDMs) were constructed to compare with the RDM of the layer conv4_2 (the refer-
440 ence RDM).

441 The first type were RDMs derived from the pre-trained VGG19. There were two RDMs—the
442 RDM of the layer conv5_4 and the RDM of the layer fc2. As the RDM of the layer conv4_2, the
443 two candidate RDMs were constructed using the corresponding feature maps extracted from the
444 VGG19 by the 50 images from the testing set. We selected the correlation distance (1 minus the
445 linear correlation between each pair of feature maps) as the measurement of representational
446 dissimilarity to construct each RDM.

447 The second type were RDMs derived from brain activity. We used the measured brain activity
448 from the testing sets and the imagery sets to constructed RDMs, respectively. Because the cate-
449 gories of the stimuli (nouns) in the imagery data were the same as the categories of the stimuli
450 (images) in the test data, the RDMs from the imagery data could be treated as RDMs using the
451 same set of stimuli—the 50 images from the test data. There were 14 RDMs for each subject, 7
452 RDMs from the test data for each ROI and 7 RDMs from the imagery data for each ROI (V1, V2, V3,
453 V4, LOC, FFA, and PPA).

454 The third type were RDMs derived from stimuli (50 images from the testing set) directly. There
455 were three RDMs—the RDM of Gabor features, the RDM of silhouette, and the RDM of semantics.
456 The RDM of Gabor features was constructed using the Gabor features of images extracted from
457 the GWP model. To construct the RDM of silhouette, we converted images to silhouettes (binary

458 images in which each figure pixel is 0 and each background pixel is 1) and calculated the correlation
459 distance between each pair of silhouettes. To construct the RDM of semantics, we calculated the
460 semantic distance between each pair of images. We used the function “path_similarity” from the
461 Natural Language Toolkit library (*Bird et al., 2009*) to calculate how similar two categories of images
462 are (semantic similarity), which based on the WordNet (*Miller, 1995*). The score returned from the
463 function was in range 0 to 1, so we converted the semantic similarity to the semantic distance by
464 subtracting the score from 1.

465 In total 19 candidate RDMs were constructed to compare with the reference RDM for each
466 subject and the Spearman’s rank correlation coefficient was selected to measure the similarity be-
467 tween each candidate RDM and the reference RDM. After that, we performed statistical inference
468 to answer two questions—whether a candidate RDM and the reference RDM were significantly
469 correlated (by permutation test with 10000 permutations) and whether the correlation between
470 a candidate RDM and the reference RDM was significantly different from the correlation between
471 another candidate RDM and the reference RDM (by bootstrap test with 1000 replications). For each
472 test, FDR was applied for multiple comparison correction (*Benjamini and Hochberg, 1995*). All cal-
473 culations were done using MATLAB 2020a and Python 3.7 on a Linux (Ubuntu 18.04 LTS) desktop
474 with a Geforce GTX 1660 Ti graphics card (6 Gb of VRAM).

475 Acknowledgments

476 This work was supported by the National Natural Science Foundation of China (No. 31600907).

477 References

- 478 **Barrett DGT**, Morcos AS, Macke JH. Analyzing biological and artificial neural networks: challenges with oppor-
479 tunities for synergy? . 2018 Oct; .
- 480 **Barsalou LW**. Grounded Cognition. *Annual Review of Psychology*. 2008 jan; 59(1):617–645. doi: [10.1146/an-](https://doi.org/10.1146/annurev.psych.59.103006.093639)
481 [nurev.psych.59.103006.093639](https://doi.org/10.1146/annurev.psych.59.103006.093639).
- 482 **Barsalou LW**. Grounded Cognition: Past, Present, and Future. *Topics in Cognitive Science*. 2010 sep; 2(4):716–
483 724. doi: [10.1111/j.1756-8765.2010.01115.x](https://doi.org/10.1111/j.1756-8765.2010.01115.x).
- 484 **Barsalou LW**. Challenges and Opportunities for Grounding Cognition. *Journal of Cognition*. 2020; 3(1). doi:
485 [10.5334/joc.116](https://doi.org/10.5334/joc.116).
- 486 **Benjamini Y**, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple
487 Testing. *Journal of the Royal Statistical Society: Series B (Methodological)*. 1995 jan; 57(1):289–300. doi:
488 [10.1111/j.2517-6161.1995.tb02031.x](https://doi.org/10.1111/j.2517-6161.1995.tb02031.x).
- 489 **Bird S**, Klein E, Loper E. *Natural Language Processing with Python*. O’Reilly UK Ltd.; 2009. [https:](https://www.ebook.de/de/product/8149545/steven_bird_ewan_klein_edward_loper_natural_language_processing_with_python.html)
490 [//www.ebook.de/de/product/8149545/steven_bird_ewan_klein_edward_loper_natural_language_](https://www.ebook.de/de/product/8149545/steven_bird_ewan_klein_edward_loper_natural_language_processing_with_python.html)
491 [processing_with_python.html](https://www.ebook.de/de/product/8149545/steven_bird_ewan_klein_edward_loper_natural_language_processing_with_python.html).
- 492 **Cadena SA**, Denfield GH, Walker EY, Gatys LA, Tolia AS, Bethge M, Ecker AS. Deep convolutional models
493 improve predictions of macaque V1 responses to natural images. *PLOS Computational Biology*. 2019 apr;
494 15(4):e1006897. doi: [10.1371/journal.pcbi.1006897](https://doi.org/10.1371/journal.pcbi.1006897).
- 495 **Cichy RM**, Kaiser D. Deep Neural Networks as Scientific Models. *Trends in Cognitive Sciences*. 2019 apr;
496 23(4):305–317. doi: [10.1016/j.tics.2019.01.009](https://doi.org/10.1016/j.tics.2019.01.009).
- 497 **Clark JM**, Paivio A. Dual coding theory and education. *Educational Psychology Review*. 1991 sep; 3(3):149–210.
498 doi: [10.1007/bf01320076](https://doi.org/10.1007/bf01320076).
- 499 **Deng J**, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. ImageNet: A large-scale hierarchical image database. In: *2009*
500 *IEEE Conference on Computer Vision and Pattern Recognition* IEEE; 2009. doi: [10.1109/cvpr.2009.5206848](https://doi.org/10.1109/cvpr.2009.5206848).
- 501 **DiCarlo JJ**, Cox DD. Untangling invariant object recognition. *Trends in Cognitive Sciences*. 2007 aug; 11(8):333–
502 341. doi: [10.1016/j.tics.2007.06.010](https://doi.org/10.1016/j.tics.2007.06.010).
- 503 **DiCarlo JJ**, Zoccolan D, Rust NC. How Does the Brain Solve Visual Object Recognition? *Neuron*. 2012 feb;
504 73(3):415–434. doi: [10.1016/j.neuron.2012.01.010](https://doi.org/10.1016/j.neuron.2012.01.010).

- 505 **Dumoulin SO**, Wandell BA. Population receptive field estimates in human visual cortex. *NeuroImage*. 2008
506 jan; 39(2):647–660. doi: [10.1016/j.neuroimage.2007.09.034](https://doi.org/10.1016/j.neuroimage.2007.09.034).
- 507 **Eickenberg M**, Gramfort A, Varoquaux G, Thirion B. Seeing it all: Convolutional network layers map the function
508 of the human visual system. *NeuroImage*. 2017 may; 152:184–194. doi: [10.1016/j.neuroimage.2016.10.001](https://doi.org/10.1016/j.neuroimage.2016.10.001).
- 509 **Gatys LA**, Ecker AS, Bethge M. Image Style Transfer Using Convolutional Neural Networks. In: *2016 IEEE Con-*
510 *ference on Computer Vision and Pattern Recognition (CVPR)* IEEE; 2016. doi: [10.1109/cvpr.2016.265](https://doi.org/10.1109/cvpr.2016.265).
- 511 **Gatys LA**, Ecker AS, Bethge M. Texture and art with deep neural networks. *Current Opinion in Neurobiology*.
512 2017 oct; 46:178–186. doi: [10.1016/j.conb.2017.08.019](https://doi.org/10.1016/j.conb.2017.08.019).
- 513 **Gaziv G**, Gabor Wavelet Pyramid ([https://www.mathworks.com/matlabcentral/fileexchange/60088-gabor-](https://www.mathworks.com/matlabcentral/fileexchange/60088-gabor-wavelet-pyramid)
514 [wavelet-pyramid](https://www.mathworks.com/matlabcentral/fileexchange/60088-gabor-wavelet-pyramid)); 2021. MATLAB Central File Exchange.
- 515 **Geirhos R**, Rubisch P, Michaelis C, Bethge M, Wichmann FA, Brendel W. ImageNet-trained CNNs are biased
516 towards texture; increasing shape bias improves accuracy and robustness. . 2018 Nov; .
- 517 **van Gerven MAJ**. A primer on encoding models in sensory neuroscience. *Journal of Mathematical Psychology*.
518 2017 feb; 76:172–183. doi: [10.1016/j.jmp.2016.06.009](https://doi.org/10.1016/j.jmp.2016.06.009).
- 519 **Glaser JJ**, Benjamin AS, Farhoodi R, Kording KP. The roles of supervised machine learning in systems neuro-
520 science. *Progress in Neurobiology*. 2019 apr; 175:126–137. doi: [10.1016/j.pneurobio.2019.01.008](https://doi.org/10.1016/j.pneurobio.2019.01.008).
- 521 **Guclu U**, van Gerven MAJ. Deep Neural Networks Reveal a Gradient in the Complexity of Neural Rep-
522 resentations across the Ventral Stream. *Journal of Neuroscience*. 2015 jul; 35(27):10005–10014. doi:
523 [10.1523/jneurosci.5023-14.2015](https://doi.org/10.1523/jneurosci.5023-14.2015).
- 524 **Han K**, Wen H, Shi J, Lu KH, Zhang Y, Fu D, Liu Z. Variational autoencoder: An unsupervised model
525 for encoding and decoding fMRI activity in visual cortex. *NeuroImage*. 2019 sep; 198:125–136. doi:
526 [10.1016/j.neuroimage.2019.05.039](https://doi.org/10.1016/j.neuroimage.2019.05.039).
- 527 **Hassabis D**, Kumaran D, Summerfield C, Botvinick M. Neuroscience-Inspired Artificial Intelligence. *Neuron*.
528 2017 jul; 95(2):245–258. doi: [10.1016/j.neuron.2017.06.011](https://doi.org/10.1016/j.neuron.2017.06.011).
- 529 **Horikawa T**, Kamitani Y. Generic decoding of seen and imagined objects using hierarchical visual features.
530 *Nature Communications*. 2017 may; 8(1). doi: [10.1038/ncomms15037](https://doi.org/10.1038/ncomms15037).
- 531 **Hubel DH**, Wiesel TN. Receptive fields, binocular interaction and functional architecture in the cat's visual
532 cortex. *The Journal of Physiology*. 1962 jan; 160(1):106–154. doi: [10.1113/jphysiol.1962.sp006837](https://doi.org/10.1113/jphysiol.1962.sp006837).
- 533 **Jing Y**, Yang Y, Feng Z, Ye J, Yu Y, Song M. Neural Style Transfer: A Review. *IEEE Transactions on Visualization*
534 *and Computer Graphics*. 2020 nov; 26(11):3365–3385. doi: [10.1109/tvcg.2019.2921336](https://doi.org/10.1109/tvcg.2019.2921336).
- 535 **Kar K**, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ. Evidence that recurrent circuits are critical to the ventral
536 stream's execution of core object recognition behavior. *Nature Neuroscience*. 2019 apr; 22(6):974–983. doi:
537 [10.1038/s41593-019-0392-5](https://doi.org/10.1038/s41593-019-0392-5).
- 538 **Kay KN**. Principles for models of neural information processing. *NeuroImage*. 2018 oct; 180:101–109. doi:
539 [10.1016/j.neuroimage.2017.08.016](https://doi.org/10.1016/j.neuroimage.2017.08.016).
- 540 **Kay KN**, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008
541 mar; 452(7185):352–355. doi: [10.1038/nature06713](https://doi.org/10.1038/nature06713).
- 542 **Khaligh-Razavi SM**, Kriegeskorte N. Deep Supervised, but Not Unsupervised, Models May Explain IT
543 Cortical Representation. *PLoS Computational Biology*. 2014 nov; 10(11):e1003915. doi: [10.1371/jour-](https://doi.org/10.1371/journal.pcbi.1003915)
544 [nal.pcbi.1003915](https://doi.org/10.1371/journal.pcbi.1003915).
- 545 **Kriegeskorte N**. Representational similarity analysis – connecting the branches of systems neuroscience. *Fron-*
546 *tiers in Systems Neuroscience*. 2008; doi: [10.3389/neuro.06.004.2008](https://doi.org/10.3389/neuro.06.004.2008).
- 547 **Kriegeskorte N**. Deep Neural Networks: A New Framework for Modeling Biological Vision and Brain Infor-
548 mation Processing. *Annual Review of Vision Science*. 2015 nov; 1(1):417–446. doi: [10.1146/annurev-vision-](https://doi.org/10.1146/annurev-vision-082114-035447)
549 [082114-035447](https://doi.org/10.1146/annurev-vision-082114-035447).
- 550 **Kriegeskorte N**, Douglas PK. Interpreting encoding and decoding models. *Current Opinion in Neurobiology*.
551 2019 apr; 55:167–179. doi: [10.1016/j.conb.2019.04.002](https://doi.org/10.1016/j.conb.2019.04.002).

- 552 **Kriegeskorte N**, Kievit RA. Representational geometry: integrating cognition, computation, and the brain.
553 Trends in Cognitive Sciences. 2013 aug; 17(8):401–412. doi: [10.1016/j.tics.2013.06.007](https://doi.org/10.1016/j.tics.2013.06.007).
- 554 **Krizhevsky A**, Sutskever I, Hinton GE. ImageNet classification with deep convolutional neural networks. Com-
555 munications of the ACM. 2017 may; 60(6):84–90. doi: [10.1145/3065386](https://doi.org/10.1145/3065386).
- 556 **Li C**, Xu J, Liu B. Decoding natural images from evoked brain activities using encoding models with invertible
557 mapping. Neural Networks. 2018 sep; 105:227–235. doi: [10.1016/j.neunet.2018.05.010](https://doi.org/10.1016/j.neunet.2018.05.010).
- 558 **Lindsay GW**. Convolutional Neural Networks as a Model of the Visual System: Past, Present, and Future. Journal
559 of Cognitive Neuroscience. 2020 feb; p. 1–15. doi: [10.1162/jocn_a_01544](https://doi.org/10.1162/jocn_a_01544).
- 560 **Mahendran A**, Vedaldi A. Understanding deep image representations by inverting them. In: *2015 IEEE Confer-*
561 *ence on Computer Vision and Pattern Recognition (CVPR)* IEEE; 2015. doi: [10.1109/cvpr.2015.7299155](https://doi.org/10.1109/cvpr.2015.7299155).
- 562 **Miller GA**. WordNet: A Lexical Database for English. Communications of the ACM. 1995 nov; 38(11):39–41. doi:
563 [10.1145/219717.219748](https://doi.org/10.1145/219717.219748).
- 564 **Mohr H**, Ruge H. Fast Estimation of L1-Regularized Linear Models in the Mass-Univariate Setting. Neuroinfor-
565 matics. 2020 sep; doi: [10.1007/s12021-020-09489-1](https://doi.org/10.1007/s12021-020-09489-1).
- 566 **Naselaris T**, Kay KN, Nishimoto S, Gallant JL. Encoding and decoding in fMRI. NeuroImage. 2011 may; 56(2):400–
567 410. doi: [10.1016/j.neuroimage.2010.07.073](https://doi.org/10.1016/j.neuroimage.2010.07.073).
- 568 **Naselaris T**, Olman CA, Stansbury DE, Ugurbil K, Gallant JL. A voxel-wise encoding model for early vi-
569 sual areas decodes mental images of remembered scenes. NeuroImage. 2015 jan; 105:215–228. doi:
570 [10.1016/j.neuroimage.2014.10.018](https://doi.org/10.1016/j.neuroimage.2014.10.018).
- 571 **Nguyen A**, Yosinski J, Clune J. Understanding Neural Networks via Feature Visualization: A Survey. In: *Explain-*
572 *able AI: Interpreting, Explaining and Visualizing Deep Learning* Springer International Publishing; 2019.p. 55–76.
573 doi: [10.1007/978-3-030-28954-6_4](https://doi.org/10.1007/978-3-030-28954-6_4).
- 574 **Nili H**, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. A Toolbox for Representational Similarity
575 Analysis. PLoS Computational Biology. 2014 apr; 10(4):e1003553. doi: [10.1371/journal.pcbi.1003553](https://doi.org/10.1371/journal.pcbi.1003553).
- 576 **Olah C**, Mordvintsev A, Schubert L. Feature Visualization. Distill. 2017 nov; 2(11). doi: [10.23915/distill.00007](https://doi.org/10.23915/distill.00007).
- 577 **Paivio A**. Mind and its evolution : a dual coding theoretical approach. Mahwah, N.J: L. Erlbaum Associates;
578 2007.
- 579 **Paszke A**, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison
580 A, Köpf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, et al. PyTorch: An
581 Imperative Style, High-Performance Deep Learning Library. . 2019 Dec; .
- 582 **Pearson J**. The human imagination: the cognitive neuroscience of visual mental imagery. Nature Reviews
583 Neuroscience. 2019 aug; 20(10):624–634. doi: [10.1038/s41583-019-0202-9](https://doi.org/10.1038/s41583-019-0202-9).
- 584 **Pearson J**, Kosslyn SM. The heterogeneity of mental representation: Ending the imagery debate. Proceedings
585 of the National Academy of Sciences. 2015 jul; 112(33):10089–10092. doi: [10.1073/pnas.1504933112](https://doi.org/10.1073/pnas.1504933112).
- 586 **Richards BA**, Lillicrap TP, Beaudoin P, Bengio Y, Bogacz R, Christensen A, Clopath C, Costa RP, de Berker A,
587 Ganguli S, Gillon CJ, Hafner D, Kepecs A, Kriegeskorte N, Latham P, Lindsay GW, Miller KD, Naud R, Pack CC,
588 Poirazi P, et al. A deep learning framework for neuroscience. Nature Neuroscience. 2019 oct; 22(11):1761–
589 1770. doi: [10.1038/s41593-019-0520-2](https://doi.org/10.1038/s41593-019-0520-2).
- 590 **Rudin C**. Stop explaining black box machine learning models for high stakes decisions and use interpretable
591 models instead. Nature Machine Intelligence. 2019 may; 1(5):206–215. doi: [10.1038/s42256-019-0048-x](https://doi.org/10.1038/s42256-019-0048-x).
- 592 **Seeliger K**, Fritsche M, Güçlü U, Schoenmakers S, Schöffelen JM, Bosch SE, van Gerven MAJ. Convolutional
593 neural network-based encoding and decoding of visual object recognition in space and time. NeuroImage.
594 2018 oct; 180:253–266. doi: [10.1016/j.neuroimage.2017.07.018](https://doi.org/10.1016/j.neuroimage.2017.07.018).
- 595 **Shen G**, Dwivedi K, Majima K, Horikawa T, Kamitani Y. End-to-End Deep Image Reconstruction From Human
596 Brain Activity. Frontiers in Computational Neuroscience. 2019 apr; 13. doi: [10.3389/fncom.2019.00021](https://doi.org/10.3389/fncom.2019.00021).
- 597 **Simonyan K**, Zisserman A. Very Deep Convolutional Networks for Large-Scale Image Recognition. . 2014 Sep; .

- 598 **St-Yves G**, Naselaris T. The feature-weighted receptive field: an interpretable encoding model for complex
599 feature spaces. *NeuroImage*. 2018 oct; 180:188–202. doi: [10.1016/j.neuroimage.2017.06.035](https://doi.org/10.1016/j.neuroimage.2017.06.035).
- 600 **Wen H**, Shi J, Zhang Y, Lu KH, Cao J, Liu Z. Neural Encoding and Decoding with Deep Learning for Dynamic
601 Natural Vision. *Cerebral Cortex*. 2017 oct; 28(12):4136–4160. doi: [10.1093/cercor/bhx268](https://doi.org/10.1093/cercor/bhx268).
- 602 **Yamins DLK**, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-optimized hierarchical models
603 predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*. 2014 may;
604 111(23):8619–8624. doi: [10.1073/pnas.1403112111](https://doi.org/10.1073/pnas.1403112111).
- 605 **Yamins DLK**, DiCarlo JJ. Using goal-driven deep learning models to understand sensory cortex. *Nature Neuro-*
606 *science*. 2016 feb; 19(3):356–365. doi: [10.1038/nn.4244](https://doi.org/10.1038/nn.4244).
- 607 **Yosinski J**, Clune J, Nguyen A, Fuchs T, Lipson H. Understanding Neural Networks Through Deep Visualization.
608 . 2015 Jun; .
- 609 **Zeiler MD**, Fergus R. Visualizing and Understanding Convolutional Networks. In: *Computer Vision – ECCV 2014*
610 Springer International Publishing; 2014.p. 818–833. doi: [10.1007/978-3-319-10590-1_53](https://doi.org/10.1007/978-3-319-10590-1_53).
- 611 **Zhou B**, Bau D, Oliva A, Torralba A. Interpreting Deep Visual Representations via Network Dissec-
612 tion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2019 sep; 41(9):2131–2145. doi:
613 [10.1109/tpami.2018.2858759](https://doi.org/10.1109/tpami.2018.2858759).
- 614 **Zhuang C**, Yan S, Nayebi A, Schrimpf M, Frank MC, DiCarlo JJ, Yamins DLK. Unsupervised neural network models
615 of the ventral visual stream. *Proceedings of the National Academy of Sciences*. 2021 jan; 118(3):e2014196118.
616 doi: [10.1073/pnas.2014196118](https://doi.org/10.1073/pnas.2014196118).

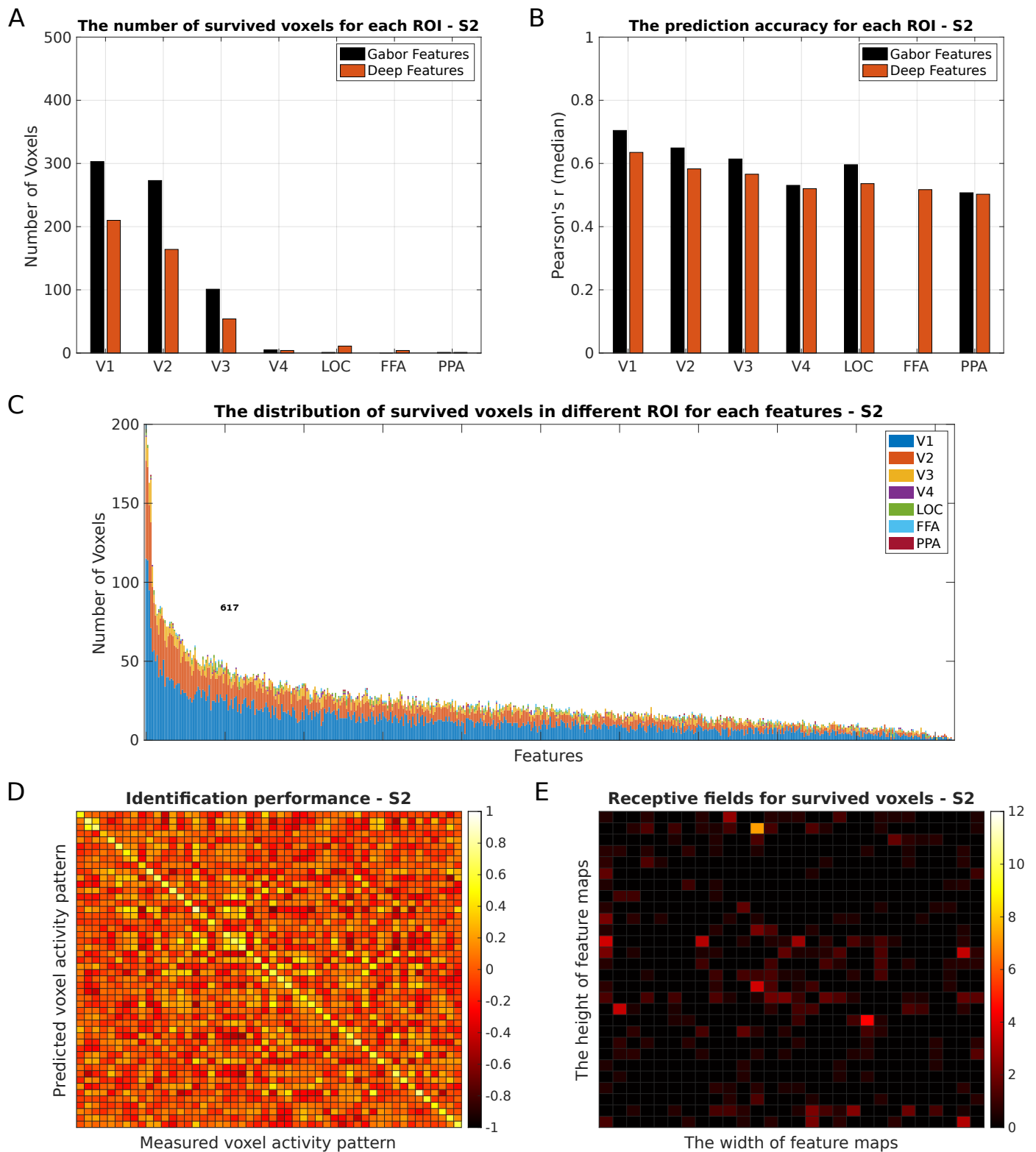


Figure 2-Figure supplement 1. The result of encoding models based on the deep features for S2.

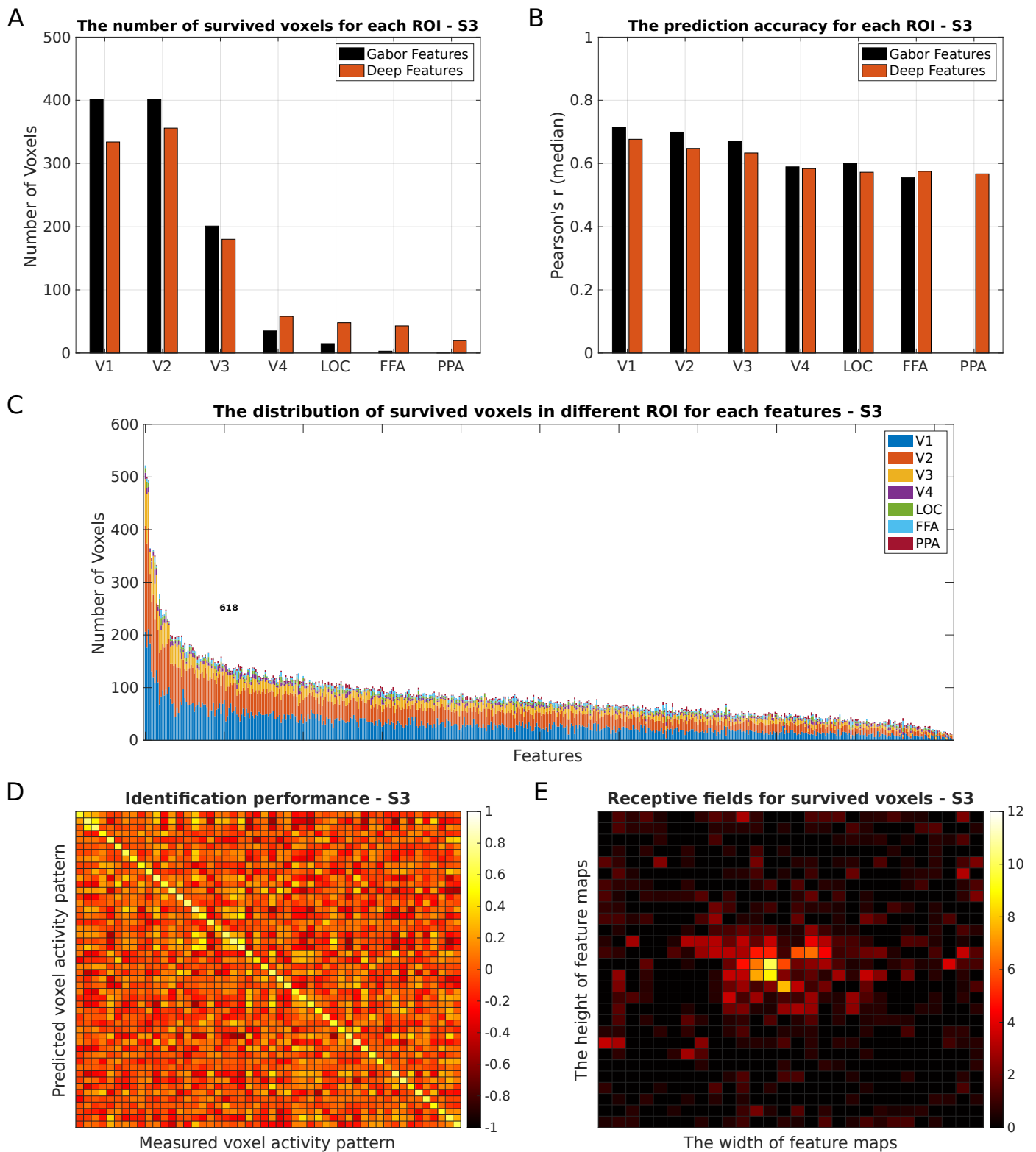


Figure 2-Figure supplement 2. The result of encoding models based on the deep features for S3.

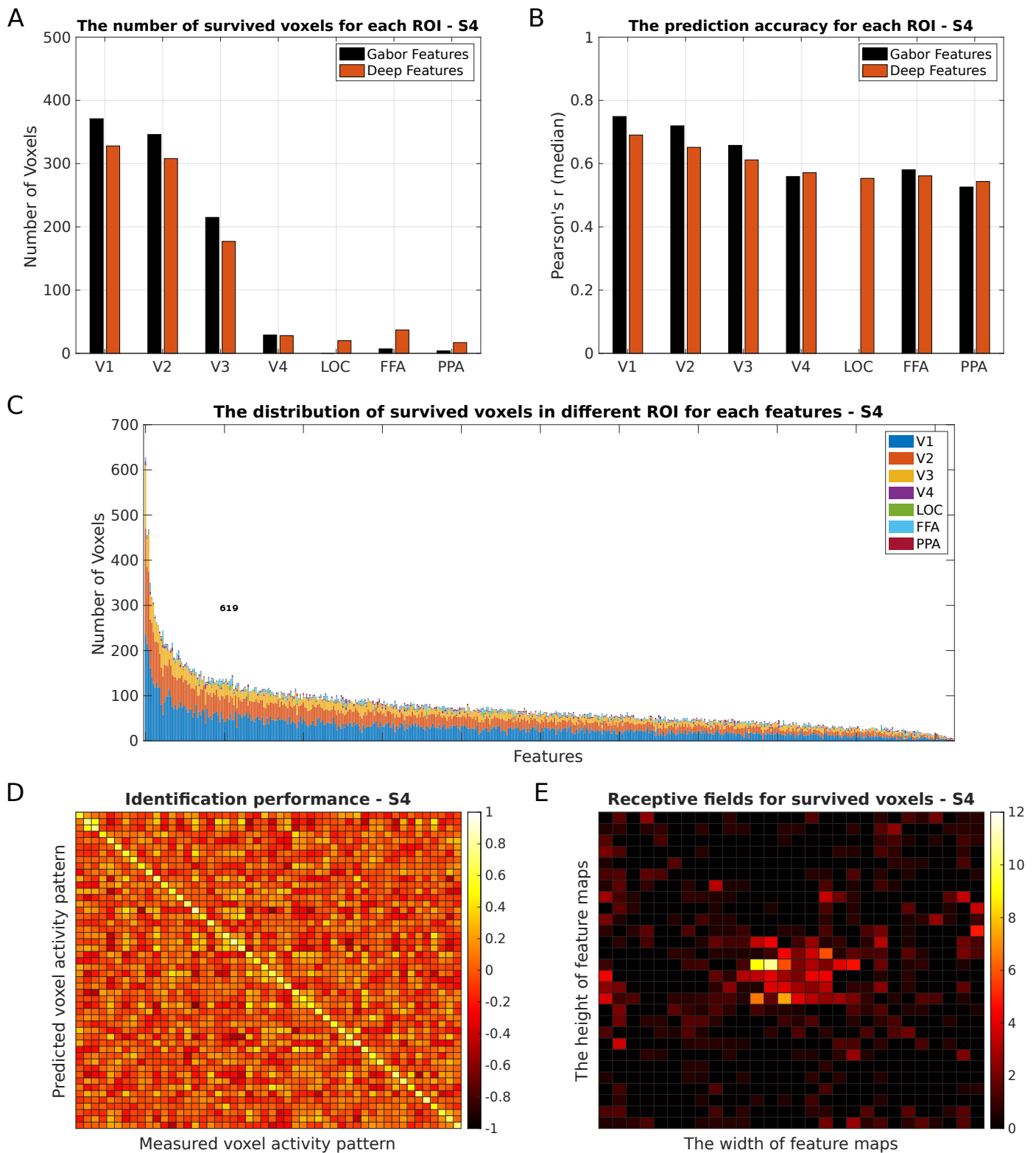


Figure 2-Figure supplement 3. The result of encoding models based on the deep features for S4.

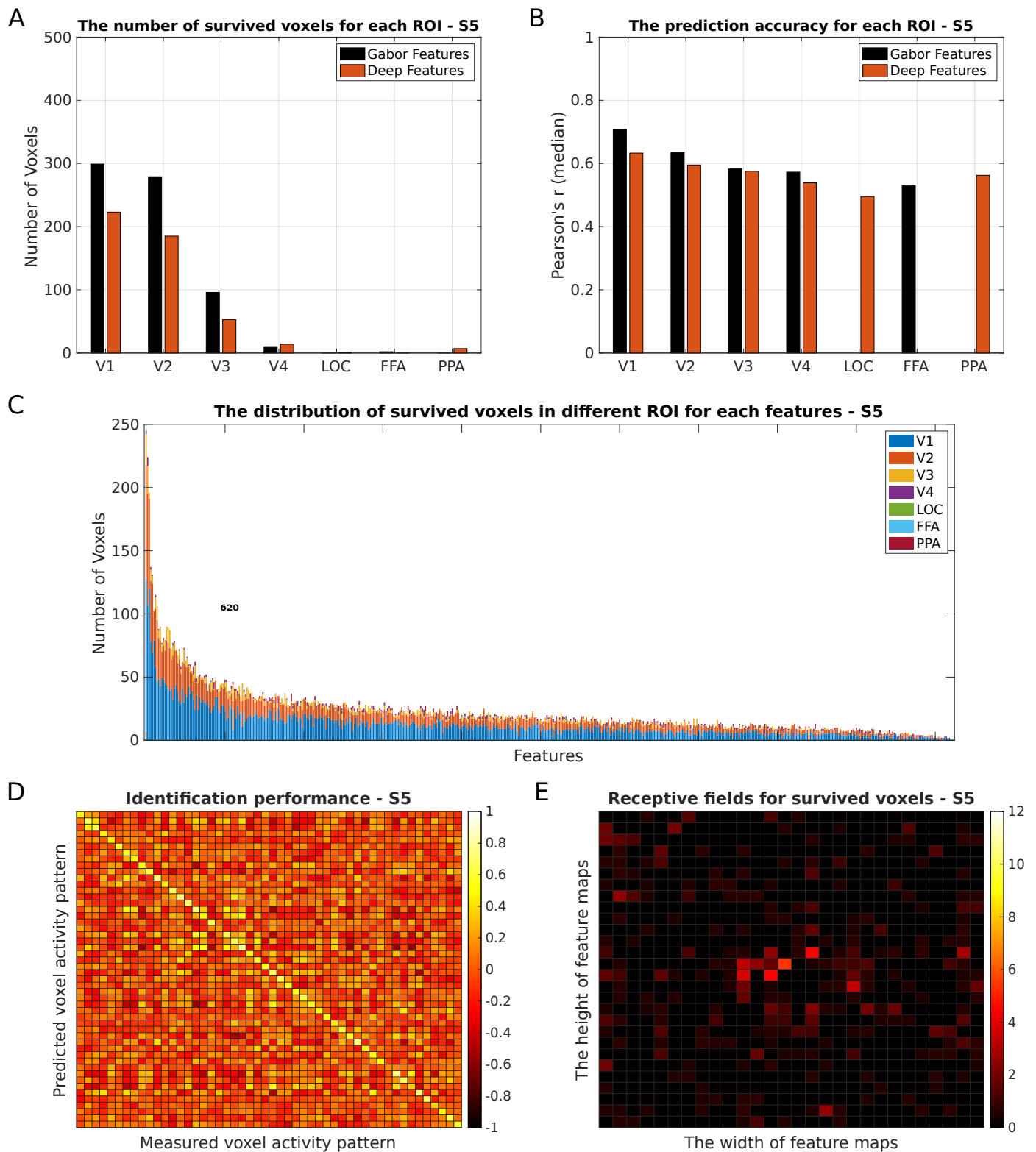


Figure 2-Figure supplement 4. The result of encoding models based on the deep features for S5.

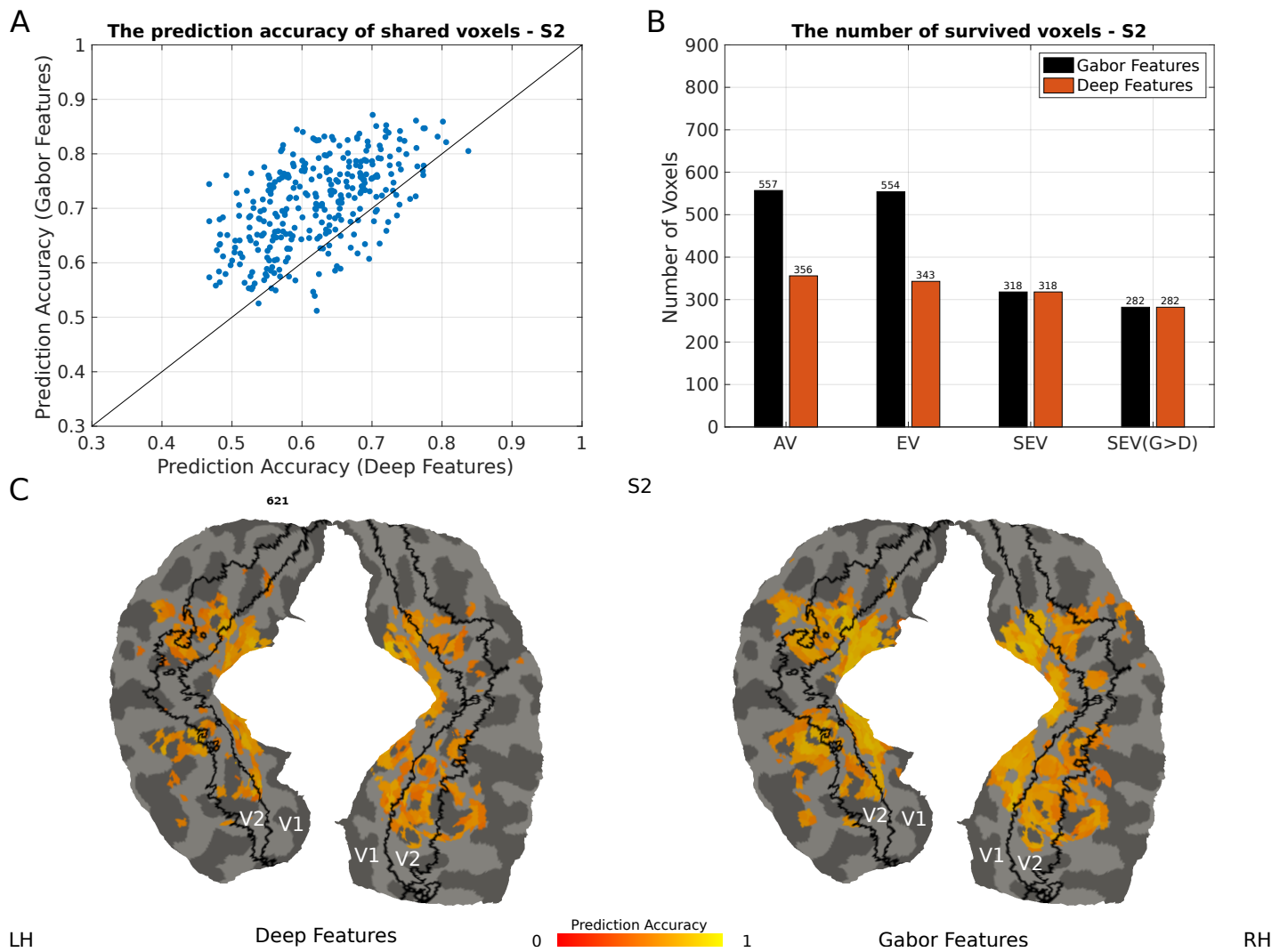


Figure 4-Figure supplement 1. The comparison of two types of encoding models for S2.

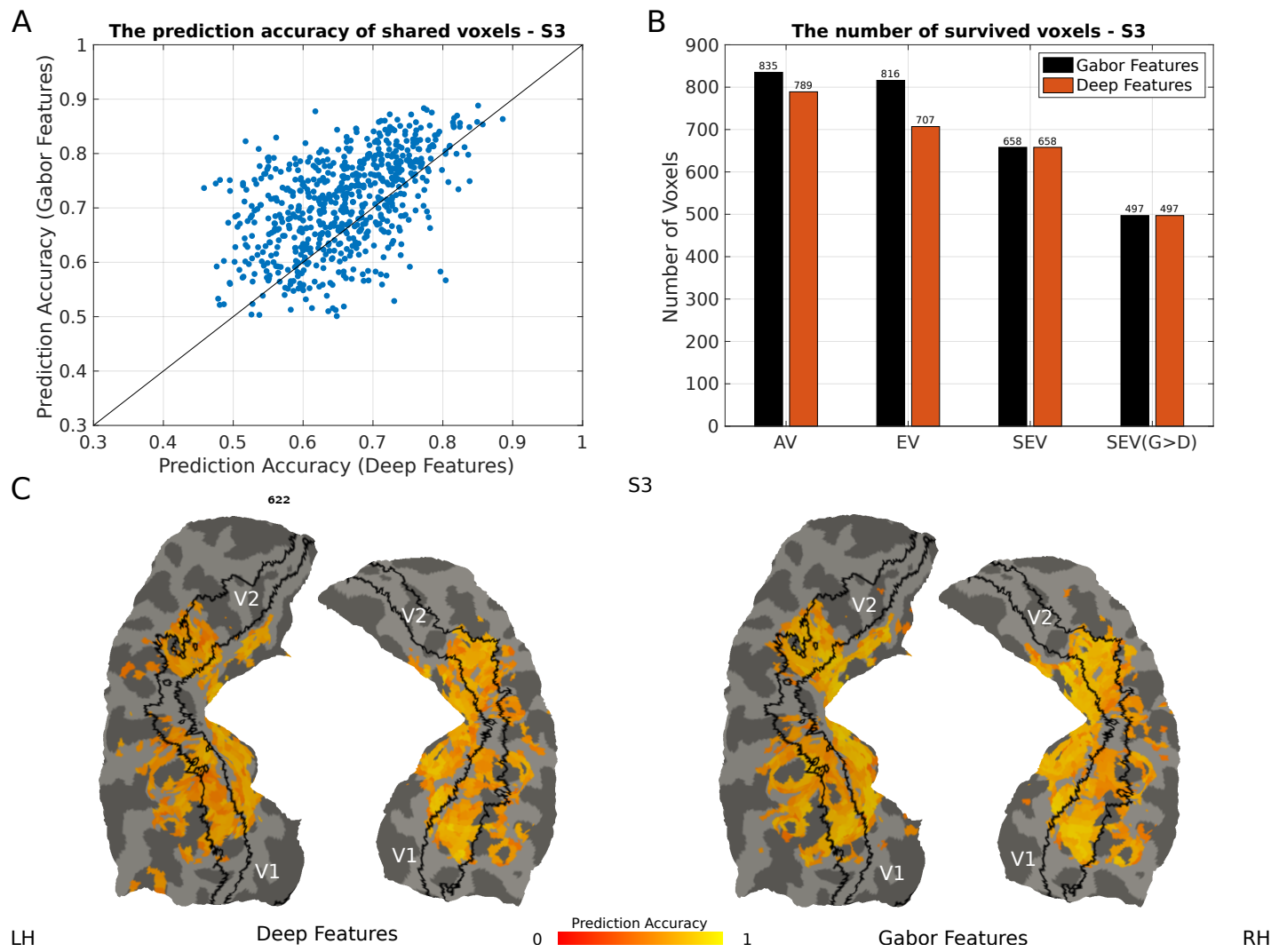


Figure 4-Figure supplement 2. The comparison of two types of encoding models for S3.

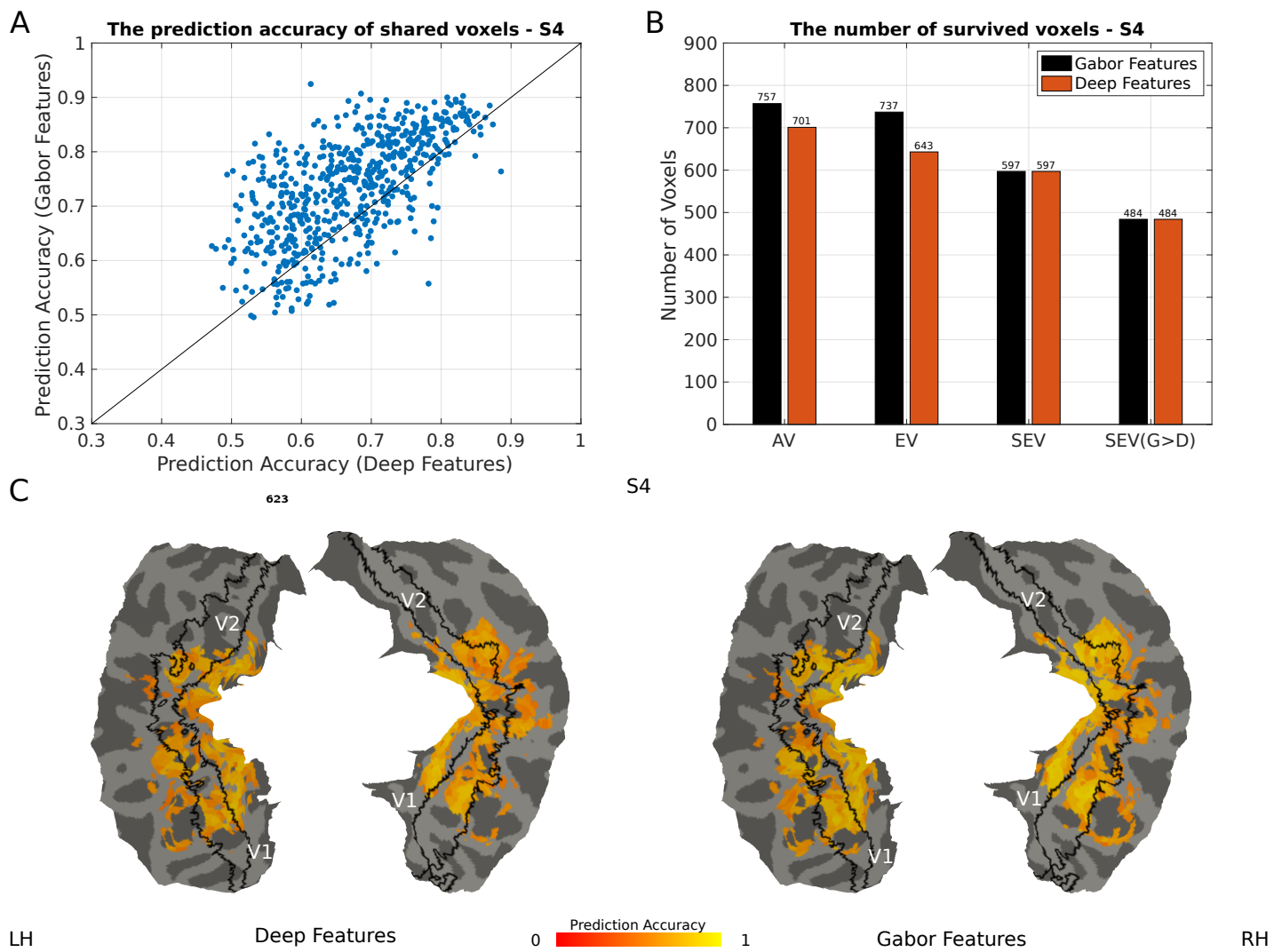


Figure 4-Figure supplement 3. The comparison of two types of encoding models for S4.

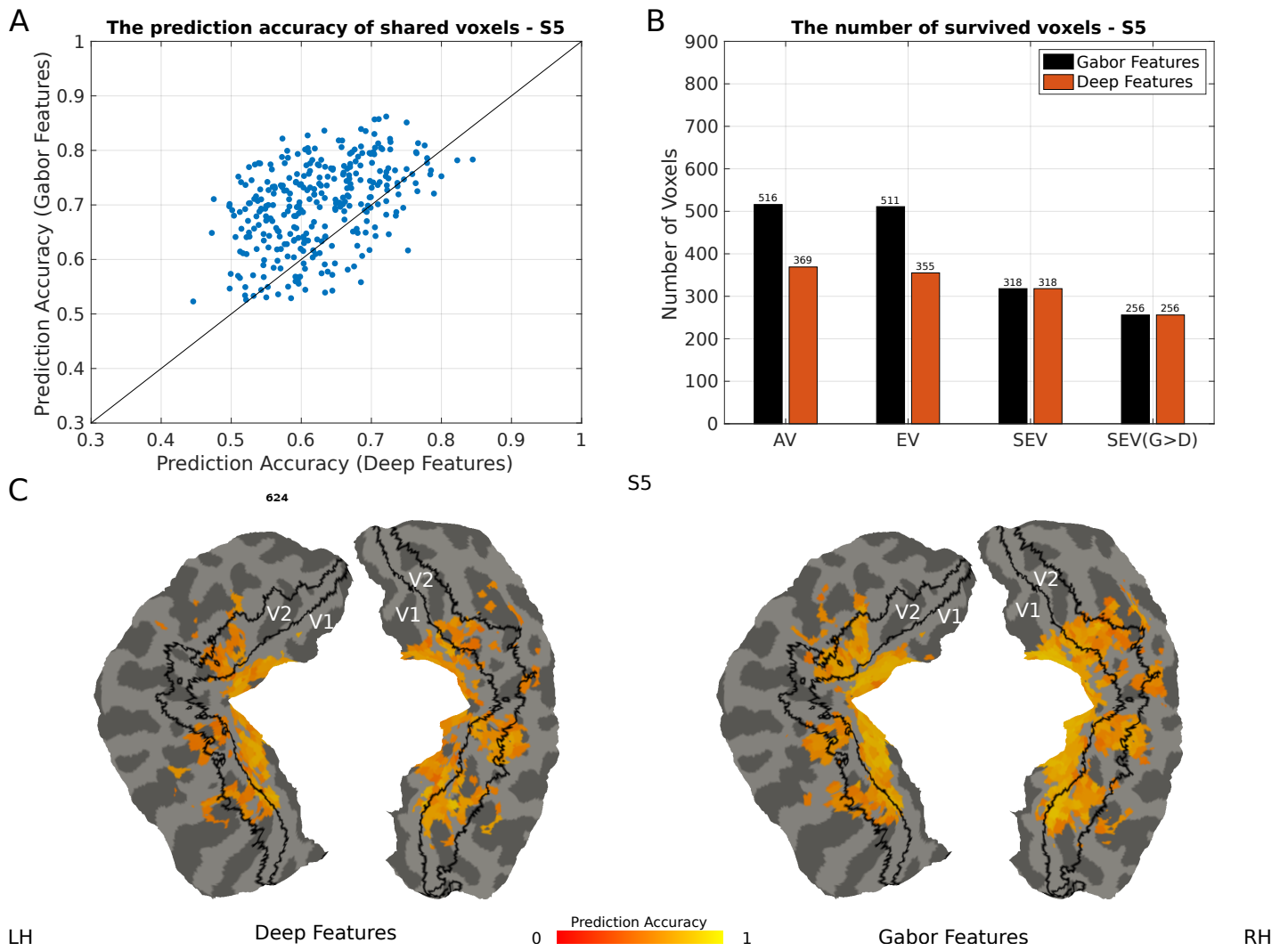


Figure 4-Figure supplement 4. The comparison of two types of encoding models for S5.

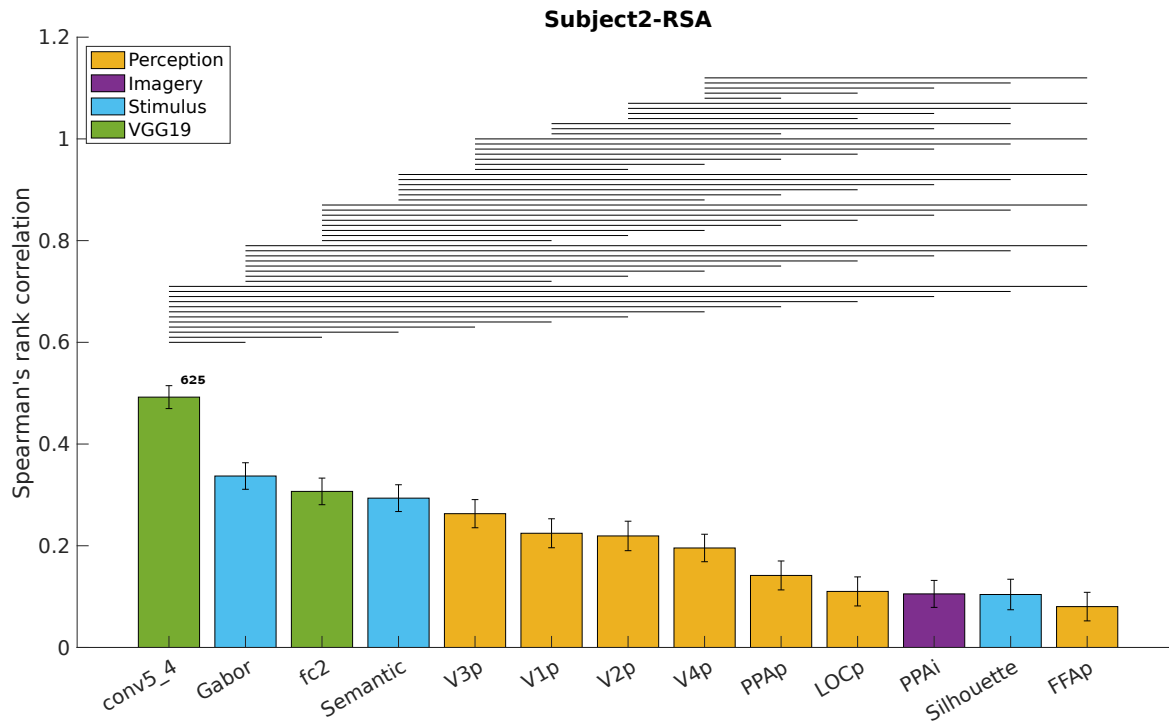


Figure 5-Figure supplement 1. The result of RSA for S2.

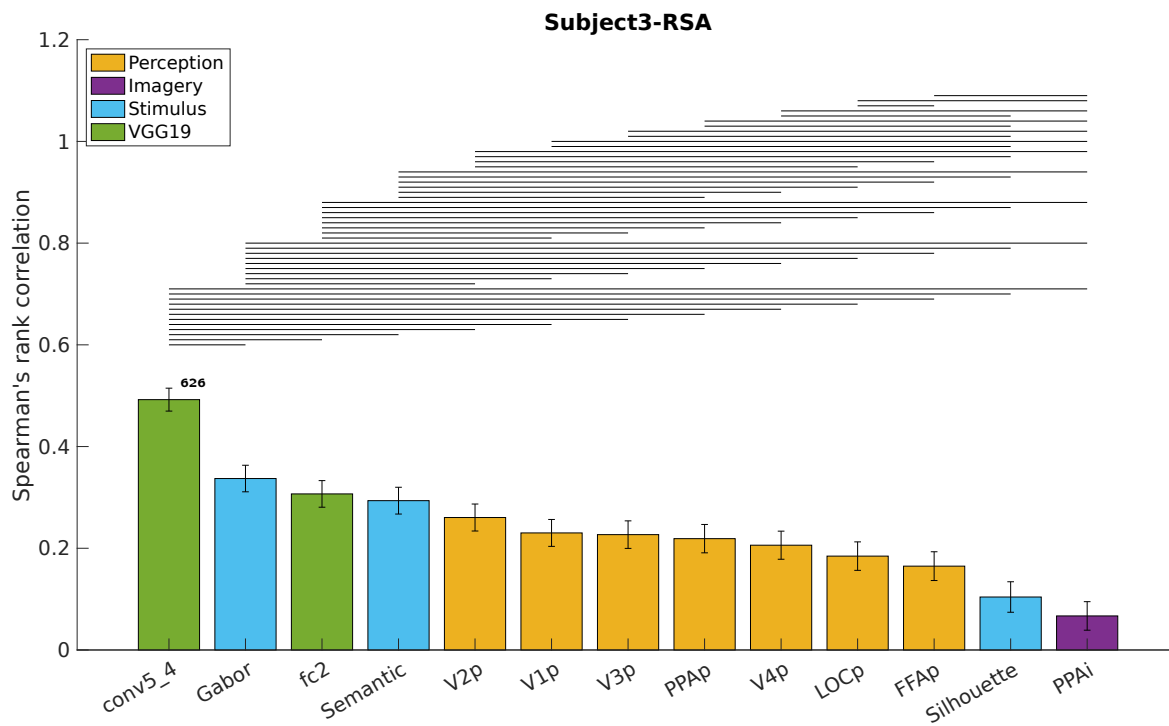


Figure 5-Figure supplement 2. The result of RSA for S3.

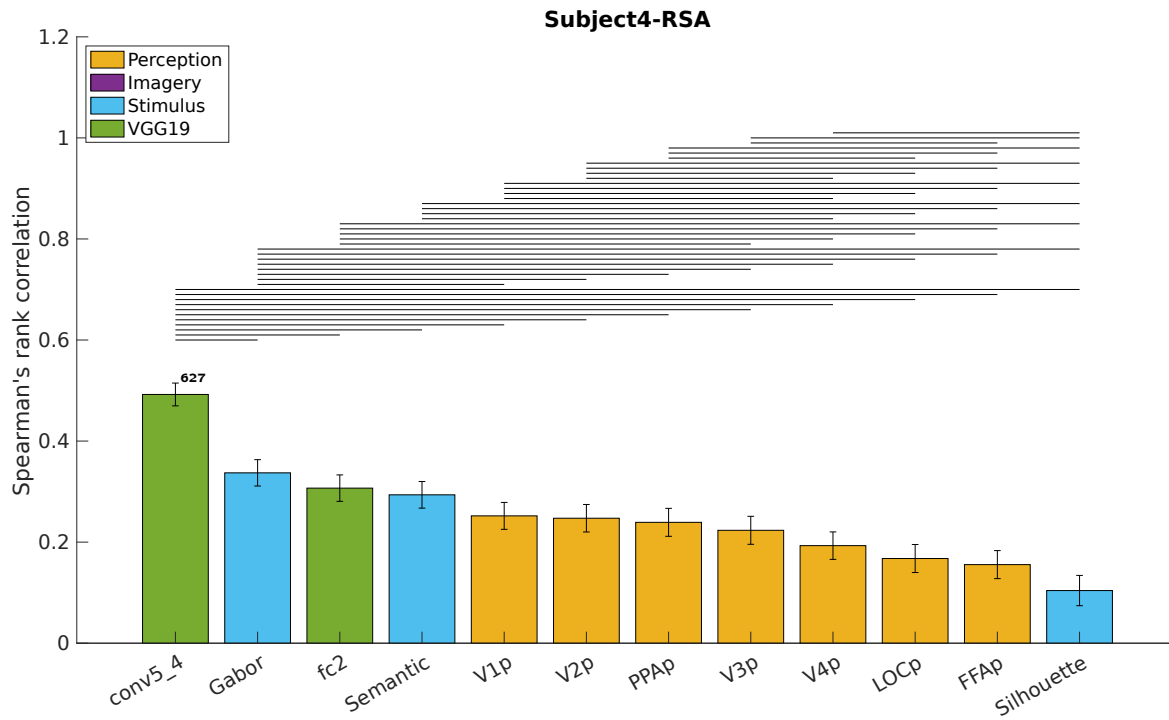


Figure 5-Figure supplement 3. The result of RSA for S4.

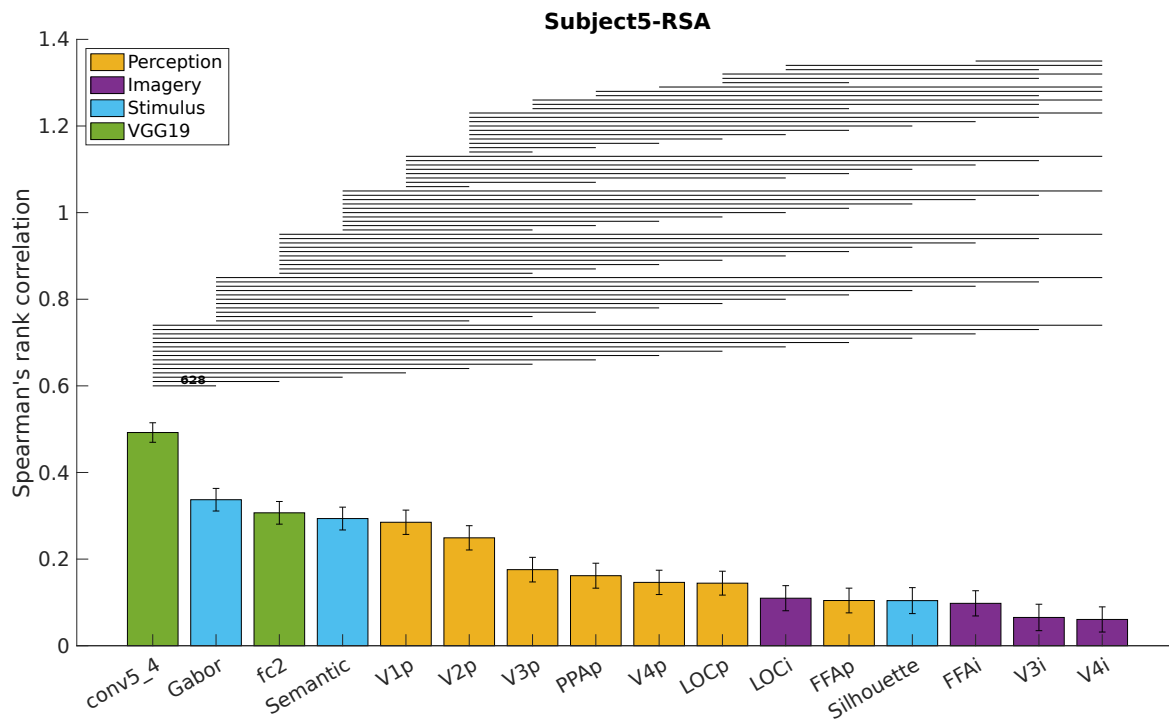


Figure 5-Figure supplement 4. The result of RSA for S5.