

*Research Articles: Behavioral/Cognitive*

## Representational organization of novel task sets during proactive encoding

<https://doi.org/10.1523/JNEUROSCI.0725-19.2019>

**Cite as:** J. Neurosci 2019; 10.1523/JNEUROSCI.0725-19.2019

Received: 1 April 2019

Revised: 19 July 2019

Accepted: 13 August 2019

---

*This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.*

**Alerts:** Sign up at [www.jneurosci.org/alerts](http://www.jneurosci.org/alerts) to receive customized email alerts when the fully formatted version of this article is published.

1 **Representational organization of novel task sets during proactive**  
2 **encoding.**

3 **Abbreviated title:** Representational structure for novel instructions.

4 Ana F. Palenciano<sup>1</sup>, Carlos González-García<sup>2</sup>, Juan E. Arco<sup>1</sup>, Luiz Pessoa<sup>3</sup> & María  
5 Ruz<sup>1</sup>

6 <sup>1</sup> *Mind, Brain, and Behavior Research Center (CIMCYC), University of Granada, 18011,*  
7 *Granada, Spain.*

8 <sup>2</sup> *Department of Experimental Psychology, Ghent University, 9000, Ghent, Belgium.*

9 <sup>3</sup> *Psychology Department, University of Maryland, 20742, Maryland, United States of America*

10

11 **Corresponding author:** María Ruz. E-mail address: [mruz@ugr.es](mailto:mruz@ugr.es) (M. Ruz).

12 **Number of pages:** 34

13 **Number of figures:** 6

14 **Number of tables:** 3

15 **Number of words.** Abstract: 241. Introduction: 622. Discussion: 1464.

16 **Financial interests or conflicts of interest:** none declared.

17 **Acknowledgments:** This work was supported by the Spanish Ministry of Science and  
18 Innovation (PSI2016-78236-P) and the Spanish Education, Culture and Sports Ministry  
19 (FPU2014/04271 and EST16/00772 to A.F.P.). This research is part of A.F.P.'s activities for the  
20 Psychology Graduate Program of the University of Granada. We are grateful to Srikanth  
21 Padmala for his valuable help during the planning and implementation of the different fMRI  
22 data analysis employed in the current experiment.

23 **Abstract**

24 Recent multivariate analyses of brain data have boosted our understanding of the organizational  
25 principles that shape neural coding. However, most of this progress has focused on perceptual  
26 visual regions (Connolly et al., 2012), whereas far less is known about the organization of more  
27 abstract, action-oriented representations. In this study, we focused on humans' remarkable  
28 ability to turn novel instructions into actions. While previous research shows that instruction  
29 encoding is tightly linked to proactive activations in fronto-parietal brain regions, little is known  
30 about the structure that orchestrates such anticipatory representation. We collected fMRI data  
31 while participants (both males and females) followed novel complex verbal rules that varied  
32 across control-related variables (integrating within/across stimuli dimensions, response  
33 complexity, target category) and reward expectations. Using Representational Similarity  
34 Analysis (Kriegeskorte et al., 2008) we explored where in the brain these variables explained  
35 the organization of novel task encoding, and whether motivation modulated these  
36 representational spaces. Instruction representations in the lateral prefrontal cortex were  
37 structured by the three control-related variables, while intraparietal sulcus encoded response  
38 complexity and the fusiform gyrus and precuneus organized its activity according to the relevant  
39 stimulus category. Reward exerted a general effect, increasing the representational similarity  
40 among different instructions, which was robustly correlated with behavioral improvements.  
41 Overall, our results highlight the flexibility of proactive task encoding, governed by distinct  
42 representational organizations in specific brain regions. They also stress the variability of  
43 motivation-control interactions, which appear to be highly dependent on task attributes such as  
44 complexity or novelty.

45 **Significance Statement**

46 In comparison with other primates, humans display a remarkable success in novel task contexts  
47 thanks to our ability to transform instructions into effective actions. This skill is associated with  
48 proactive task-set reconfigurations in fronto-parietal cortices. It remains yet unknown, however,  
49 *how* the brain encodes in anticipation the flexible, rich repertoire of novel tasks that we can

50 achieve. Here we explored cognitive control and motivation-related variables that might  
51 orchestrate the representational space for novel instructions. Our results showed that different  
52 dimensions become relevant for task prospective encoding depending on the brain region, and  
53 that the lateral prefrontal cortex simultaneously organized task representations following  
54 different control-related variables. Motivation exerted a general modulation upon this process,  
55 diminishing rather than increasing distances among instruction representations.

**56 Introduction**

57 Humans quickly learn from instructions which elements are relevant in a context and their  
58 respective appropriate actions. These parameters are encoded proactively in our brain in an  
59 action-based code (Brass, Liefoghe, Braem, & De Houwer, 2017; Cole, Braver, & Meiran,  
60 2017), preparing our perceptual and motor systems in advance (Cole, Laurent, & Stocco, 2013)  
61 and facilitating success in novel environments. Instructed behavior is thus critical to avoid less  
62 effective and slow trial-and-error learning, and also enables the social transmission of task  
63 procedures. There is scarce knowledge, however, about how the informational and motivational  
64 content of novel instructions organizes neural activity in a proactive manner.

65 Behavioral results support the role of proactive control (Braver, 2012) on instructed action (e.g.  
66 Liefoghe, Wenke, & De Houwer, 2012; see also Cole, Patrick, & Braver, 2018; Duncan et al.,  
67 2008; Luria, 1966). Recently, neuroimaging studies have revealed a link between novel  
68 instruction preparation and the fronto-parietal (FP) network (e.g. Cole, Bagic, Kass, &  
69 Schneider, 2010; Hartstra, Kühn, Verguts, & Brass, 2011; Palenciano, González-García, Arco,  
70 & Ruz, 2018). The middle (MFG) and inferior (IFG) frontal gyri, and the inferior frontal sulcus  
71 (IFS), together with the intraparietal sulcus (IPS), encode novel instruction content both in  
72 multivoxel activity patterns (Bourguignon, Braem, Hartstra, De Houwer, & Brass, 2018;  
73 González-García, Arco, Palenciano, Ramírez, & Ruz, 2017; Muhle-Karbe, Duncan, De Baene,  
74 Mitchell, & Brass, 2017) and distributed functional connectivity (Cole, Laurent, et al., 2013).  
75 Crucially, the fidelity of information encoding is linked to the intention to implement the  
76 instruction (versus mere memorization demands; Bourguignon et al., 2018; Muhle-Karbe et al.,  
77 2017) and it is also closely related to the efficiency of behavior (Cole, Ito, & Braver, 2016;  
78 González-García et al., 2017). Nonetheless, while current studies have mainly focused on  
79 decoding the upcoming target category (González-García et al., 2017; Muhle-Karbe et al.,  
80 2017), the wider organizational structure that shapes anticipatory task representation remains  
81 unknown. To study the relevant dimensions organizing novel instruction encoding, we selected  
82 three variables known to be relevant for proactive control.

83 Task preparation consists of a two-step process (Rubinstein et al., 2001), composed first by an  
84 abstract goal reconfiguration and second by the activation of specific stimulus-response  
85 contingencies (De Baene & Brass, 2014; Muhle-Karbe, Andres, & Brass, 2014). Our study  
86 exploited these two phases. First, in relation to the high-level task goal setting, we manipulated  
87 the integration of information within or across feature dimensions of stimuli (Rigotti et al.,  
88 2013), a variable traditionally linked to task complexity and top-down attention (e.g. Treisman  
89 & Gelade, 1980). Second, the stimulus-response reconfiguration process was manipulated by  
90 the response set complexity, requiring single or sequential motor responses. Moreover, to  
91 explore stimuli-specific preparatory mechanisms previously documented (e.g. González-García,  
92 Mas-Herrero, de Diego-Balaguer, & Ruz, 2016; Sakai & Passingham, 2003, 2006), we also  
93 manipulated the relevant target category.

94 Finally, cognitive control and motivation maintain an intricate relationship during task  
95 preparation (Pessoa, 2009, 2017). Reward expectation boosts cue-locked activity across the FP  
96 network (Parro, Dixon, & Christoff, 2017), and it has been recently linked to stronger  
97 anticipatory rule encoding (Etzel, Cole, Zacks, Kay, & Braver, 2016). Nonetheless,  
98 contradictory findings have also been found (Wisniewski, Forstmann, & Brass, 2018), and a  
99 comprehensive characterization of this interaction in complex, novel scenarios is still pending.  
100 Consequently, we included economic incentives in our paradigm and assessed the nature of  
101 their effect on instruction preparation. By varying these four variables (dimension integration,  
102 response-set complexity, target category, and reward), we built a set of novel, verbal  
103 instructions that were followed by healthy participants while functional magnetic imaging  
104 (fMRI) data were collected. Using Representation Similarity Analysis (RSA; Kriegeskorte,  
105 Mur, & Bandettini, 2008), we assessed the extent to which each of our control-related variables  
106 organized instruction encoding, as well as the effect of motivation upon this organization.

107 **Materials and methods**

108 *Participants*

109 Thirty-six students from the University of Granada completed the experimental paradigm inside  
110 an MRI scanner (16 women, mean age = 22.97 years, SD = 3.32 years). All of them were right-  
111 handed, with normal or corrected-to-normal vision, and native Spanish speakers. In exchange  
112 for their participation, they received between 20 and 40€, depending on their performance on  
113 the rewarded trials (see below). They all signed a consent form approved by the Ethics  
114 Committee of the University of Granada. Four participants were later excluded due to excess of  
115 head movement (> 3mm) or poor performance (<70% of correct responses).

116 *Apparatus, stimuli, and procedure*

117 For the experiment, we built a set of 192 different novel verbal instructions. Each instruction  
118 referred to two independent conditions about faces or food items that could be met or not by the  
119 upcoming grids, and their associated responses (e.g.: “*If there are two women and an additional*  
120 *sad person, press A; if not, press L*”). The conditions in the instructions referred to several  
121 dimensions of the stimuli: gender (*woman, man*), race (*black, white*), emotion (*happy, sad*) and  
122 size (*big, small*) of faces, or kind (*fruit, vegetable*), color (*green, yellow*), form (*round,*  
123 *elongated*) and size (*big, small*) of food items.

124 Instructions were created by manipulating in an orthogonal manner (1) the ***Integration of***  
125 ***stimuli dimensions*** (within vs. across dimensions), (2) the ***Response set*** required (single vs.  
126 sequential) and (3) the ***Category*** of the relevant stimuli that they referred to (faces vs. food). For  
127 example, the instruction “*If there is a woman and there is a man, press A; if not, press L*”  
128 involves within-dimension integration (i.e., gender), requires a single response (a left –“A”– or  
129 a right –“L”– index button press) and is face-related. On the other hand, “*If there is a fruit and a*  
130 *small food item, press AL; if not, press LA*” requires across-dimension integration (the type of  
131 food and its size), demands a sequence of two button presses to respond and is food-related.

132 Instructions referred to either 2, 3 or 4 stimuli of the target grid. Equivalent trials were created  
133 for the different levels of these three variables.

134 In addition, we included *Motivation* as another variable: half of the instructions were associated  
135 with the possibility of receiving an economic reward if responses were fast and accurate while  
136 the other half were non-rewarded. To do so, we split our 192 instructions into two equivalent  
137 sets in terms of the manipulations of the other independent variables, and also regarding the  
138 specific attributes specified (e.g., the same number of instructions referring to happy faces in  
139 both groups). We counterbalanced across participants the assignment of these two halves to the  
140 rewarded and non-rewarded conditions. The reward status of each trial was indicated by a cue  
141 consisting on either a plus (+) or a cross (x) sign, in either silhouette or filled in black. We  
142 counterbalanced across participants whether they should attend to the shape (plus vs. cross) or  
143 the appearance (contour vs. filled sign) to obtain the reward information. This way, each  
144 participant had two different cues indicating each motivation condition, preventing a one-to-one  
145 mapping between reward expectation and visual cue identity, which otherwise could generate  
146 spurious confounds in further analysis.

147 For each instruction, we created two grids of stimuli, one that fulfilled the conditions instructed,  
148 and another one that did not. We counterbalanced them so that individual participants saw only  
149 one of the two instruction-grid pairings. All grids were unique combinations of images of 4  
150 faces and 4 food items, which were pseudo-randomly selected from a pool of 32 pictures,  
151 composed by 16 faces pictures (8 different identities, half of them women and half men, half  
152 with happy expression and half with sad ones, half white and half black, appearing each of them  
153 in large and small sizes), extracted from the NimStim database (Tottenham et al., 2009), and 16  
154 food pictures (8 different items, half of them vegetables and half fruits, half in green color and  
155 half in yellow, half with a round shape and half elongated, appearing each of them in large and  
156 small sizes) obtained from available sources on the internet (all of them with Creative  
157 Commons license). Upon target presentation, the responses required were always one or two  
158 sequential button presses, performed with the left (“A”) and/or right (“L”) index. The sequence



159 of trial events is depicted in Figure 1. Each trial started with a jittered fixation point ( $0.5^\circ$ ), with  
160 a duration that ranged from 4500 to 7500ms, in steps of 500ms (mean = 5750ms). Then, a  
161 reward cue was presented ( $1.5^\circ$ ; 2000ms), followed by the instruction ( $25.75^\circ$ ; 2500ms). Next a  
162 second jittered fixation appeared (with the same characteristics as the previous one), and the  
163 target grid ( $21^\circ$ ) was presented for 2500ms, where participants were required to respond.  
164 Afterward, a feedback symbol was presented ( $1.65^\circ$ ; 500ms), indicating whether the participant  
165 had earned money in that trial (with a Euro symbol), whether the response was correct but no  
166 money was achieved (tick symbol) or whether the response was incorrect (cross symbol).

167 Before being scanned, participants completed a behavioral practice session. They received  
168 indications about how to perform the task, as well as details on how rewards would be  
169 administered, emphasizing that both accurate and fast responses were needed to accumulate  
170 money for a maximum of 40€. Specifically, they were informed that they would receive 20€ for  
171 their time and that the rest of the compensation would depend on their performance on rewarded  
172 trials: the initial extra increases would be easier to earn while approaching the upper limit of the  
173 payment would require a higher accuracy rate. Then, they performed a simple discrimination  
174 task with the different reward cues, and after that, they practiced the instruction-following task,  
175 completing one block of 32 trials. Practice instructions were drawn from a separate set (which  
176 was equivalent in all the parameters specified above) and were not employed in the MRI  
177 experiment, to maintain trial novelty. Participants repeated the practice block as many times as  
178 needed to obtain an accuracy rate above 75% (on average, participants performed the practice  
179 block 1.75 times). Once this phase was completed, the experimental paradigm was performed  
180 inside the scanner. This was composed by the full 192 instructions set, presented in six different  
181 runs (32 trials each). All runs included an equal number of face and food-related, single and  
182 sequential responses, within and across-dimension integration and rewarded and non-rewarded  
183 instructions. Overall, participants spent 90 minutes approximately inside the MRI scanner.

184 *Experimental Design and behavioral statistical analysis*

185 Our task was built following a 4-way factorial design, in which the following within-subjects  
186 independent variables were orthogonally manipulated: (1) Dimension integration; (2) Response  
187 set complexity; (3) Target category and (4) Reward.

188 We conducted an a priori power analysis to compute sample size. Using the PANGAEA software  
189 (<https://jakewestfall.shinyapps.io/pangea/>), we calculated the minimum number of participants  
190 to detect a behavioral two-way interaction term (i.e., between reward and any other proactive  
191 control-related variable), assuming a medium effect size (Cohen's  $d = .3$ ).

192 We used IBM SPSS Statistics v20 software to analyze accuracy and reaction time data. We  
193 conducted two repeated-measures ANOVAs, specifying four factors corresponding to our  
194 independent variables. To explore significant interaction terms, we carried out further post hoc  
195 tests, using a Bonferroni correction for multiple comparisons.

#### 196 *fMRI preprocessing*

197 MRI data were acquired using a 3-Tesla Siemens Trio scanner located at the Mind, Brain, and  
198 Behavior Research Center (CIMCYC, University of Granada, Spain). Functional images were  
199 collected employing a T2\* Echo Planar Imaging (EPI) sequence (TR = 2210ms, TE = 23ms, flip  
200 angle = 70°). Each volume consisted of 40 slices, obtained in descending order, with 2.3mm of  
201 thickness (gap = 20%, voxel size = 3mm<sup>3</sup>). A total of 1716 volumes were obtained, in 6 runs of  
202 286 volumes each. We also acquired a high-resolution anatomical T1-weighted image (192  
203 slices of 1mm, TR = 2500ms, TE = 3.69ms, flip angle = 7°, voxel size = 1mm<sup>3</sup>).

204 The functional images were preprocessed and analyzed with SPM12  
205 (<http://www.fil.ion.ucl.ac.uk/spm/software/spm12/>), with the exception of single-trial parameter  
206 estimation (see *RSA section*), which was conducted on AFNI. After discarding the first four  
207 volumes of each run to allow for stabilization of the signal, the images were spatially realigned  
208 and slice-time corrected. Then, the participants' structural T1 image, which had been  
209 coregistered with the EPI volumes, was segmented to obtain the transformation matrices needed  
210 to normalize the functional images to the MNI space. Finally, they were smoothed with an 8mm

211 FWHM Gaussian kernel. The full preprocessing pipeline was completed before conducting the  
212 univariate analysis, while only realigned and slice-timing corrected images were employed for  
213 the multivariate tests (see next section). In the latter, normalization and smoothing were  
214 performed after the individual-level analysis, following the same strategy as above.

### 215 *fMRI statistical analysis*

#### 216 *Control univariate analysis*

217 We first conducted a univariate standard GLM, modelling each of the sixteen combinations of  
218 our variables (for example: within-dimension integration/simple response required/faces-  
219 related/ rewarded) and specifying two regressors per trial: one for the encoding phase (from the  
220 reward cue until the end of the instruction), and another for the implementation stage  
221 (encompassing the target grid presentation and until the end of the feedback cue). All regressors  
222 were convolved with the canonical hemodynamic response function. We also added error trials  
223 and six motion parameters as nuisance regressors, and a high-pass filter of 128s to avoid low-  
224 frequency noise.

225 The rationale of this analysis was to check the effect of motivation during the encoding of novel  
226 instructions with the aim of ensuring that our manipulation successfully generated typical  
227 reward-related patterns of activation (Parro et al., 2017). This was done by performing *t*-tests at  
228 the individual (first) level, contrasting rewarded versus non-rewarded encoding regressors, and  
229 carrying these statistical maps to a group one-sample *t*-test. The result was cluster-wise FWE-  
230 corrected for multiple comparison at  $P < .05$  (from an initial threshold of  $P < .001$  and  $k = 10$ ).  
231 With this approach, we obtained one large cluster that extended across multiple brain regions.  
232 To obtain smaller, anatomically coherent clusters, we employed a stricter threshold (uncorrected  
233 cluster-forming threshold of  $P < .0001$ , with the corresponding FWE correction at  $P < .05$ ), as  
234 done previously (e.g. Dumontheil et al., 2011; Palenciano et al., 2018).

#### 235 *Representational Similarity Analyses*

236 We conducted a series of multivariate RSAs, following a two-step approach. First, we analyzed

237 whole-brain data, using a searchlight approach, to find regions encoding novel instructions  
238 according to each of our three control-related variables. Second, we used the significant areas as  
239 Regions Of Interest (ROIs) and focused on them to explore the effect of reward on their  
240 representational geometry.

241 *Whole-brain model-based RSA.* We first studied whether the representational structure of novel  
242 instructions was explained by three variables related to cognitive control preparation: dimension  
243 integration, response set complexity and target category. Importantly, we specifically wanted to  
244 explore this during the initial encoding stage, where proactive task-set reconfiguration takes  
245 place. To do so, we first obtained trial-by-trial estimations of our signal, following a Least-  
246 Square-Sum approach (LSS; Turner, 2010) to ensure the smallest possible collinearity among  
247 regressors (Arco, González-García, Díaz-Gutiérrez, Ramírez, & Ruz, 2018). We generated and  
248 estimated one separate model per trial, in which we defined: (1) a regressor isolating the  
249 encoding phase of the individual trial of interest; (2) a second regressor containing the rest of  
250 trials (encoding phase) of the same condition; (3) thirty-one additional regressors encompassing  
251 the rest of conditions at the encoding and implementation phases (as in the GLM specified  
252 above), and (4) nuisance regressors (movement, errors). To do so, we employed AFNI's function  
253 3dLSS ([https://afni.nimh.nih.gov/pub/dist/doc/program\\_help/3dLSS.html](https://afni.nimh.nih.gov/pub/dist/doc/program_help/3dLSS.html)). Once the trial-wise  
254 parameter images were obtained, the rest of the RSA was performed with The Decoding  
255 Toolbox (Hebart, Görden, & Haynes, 2014).

256 In our analysis, we compared three theoretical models of representational organization (one per  
257 preparation-related independent variable) with the empirical one, built from spatially distributed  
258 activity patterns. To do so, we employed a spherical searchlight (radius: 4 voxels) and applied it  
259 to the whole brain (Kriegeskorte, Goebel, & Bandettini, 2006). First, we built three theoretical  
260 representational dissimilarity matrices (RDM, Fig. 2a), which captured the expected  
261 dissimilarity (represented with 0s and 1s) between pairs of trials, according to the corresponding  
262 variables of interest. For example, in the Category RDM, dissimilarity is expected to be minimal  
263 within pairs of trials that refer either to faces or to food, while maximal between pairs of trials

264 referring to different target categories. Then, in each iteration of the searchlight, we generated a  
265 neural RDM, using a measure of distance based on Pearson correlation. Specifically, we  
266 extracted the corresponding single-trial beta values of the voxels involved, correlated each pair  
267 of the trials' activity patterns, and subtracted that value from 1. Afterwards, this neural RDM  
268 was Spearman-correlated with the theoretical ones (Fig. 2c), and the coefficients were  
269 normalized with Fisher's  $z$  transformation and assigned to the central voxel of the searchlight  
270 sphere. Importantly, both theoretical and neural matrices were built trial-wise (i.e., not  
271 averaging within conditions), and thus, were fully symmetrical with a diagonal of 0s.  
272 Consequently, only the lower triangle of the matrices, excluding the diagonal, was included in  
273 the correlation to avoid inflated positive results (Ritchie, Bracci, & Op de Beeck, 2017). After  
274 iterating the searchlight across the whole brain, we obtained three maps per participant  
275 representing how well the representational geometry in different regions matched the one  
276 expected by each of our three theoretical models.

277 Statistical significance was assessed non-parametrically via permutation testing, as proposed by  
278 Stelzer, Chen, & Turner (Stelzer, Chen, & Turner, 2013). We first performed 100 permutations  
279 at the individual level, where trial labels were randomly shifted and the whole analysis was  
280 repeated. Then, at the group level, we resampled 50,000 times one of the permuted maps of each  
281 subject and averaged them. The resulting bootstrapped group maps were used to build a voxel-  
282 wise null distribution of correlation values, which was used to extract the correlation coefficient  
283 coinciding with a probability of 0.001 of the right-tailed area of the distribution (i.e., linked to a  
284  $p \leq .001$ ) of each individual voxel. The group map of the results was then thresholded using  
285 these values. From the bootstrapped maps we also built a null distribution of cluster sizes  
286 (Stelzer, Chen, & Turner, 2013), which determined the probability of each cluster extent under  
287 the null distribution. We used this to assign the corresponding  $P$  value to the surviving clusters  
288 of the group results map, and FWE-corrected ( $P < .05$ ) them to control for multiple  
289 comparisons.

290 We performed a further conjunction test to find areas sharing the three representational

291 organizational schemes. To do so, we thresholded ( $P < .05$ , FWE corrected) and binarized the  
292 three maps from the previous step, and obtained the overlapping voxels (Nichols, Brett,  
293 Andersson, Wager, & Poline, 2005).

294 Importantly, the RSA results could be influenced by other variables statistically related to our  
295 manipulations (Popov, Ostarek, & Tenison, 2018), such as instructions' length and speed of  
296 responses, which differed slightly between conditions. To examine their influence on the results,  
297 we performed an additional multiple regression analysis taking both variables into account. We  
298 built two different RDMs (see Fig. 2.b) in which each cell contained the absolute difference in  
299 the number of letters (instruction's length RDM) or reaction time (response speed RDM),  
300 respectively, between specific pairs of instructions. We then used them as regressors together  
301 with the three proactive control-related RDMs, predicting the neural pattern of dissimilarities in  
302 each iteration of a searchlight. The regressors were built vectorizing the lower triangle of the  
303 RDM, excluding the diagonal values. It is important to note that there were small but still  
304 significant correlations among some of the regressors included in the analysis. Specifically,  
305 dimension integration correlated with instruction length and RT, and target category did so with  
306 instruction length. To assess the impact of these correlations on the regression estimation, we  
307 computed Variance Inflation Factors (Mumford, Poline, & Poldrack, 2015), an index of the  
308 regressors' collinearity. For our five models, and in all the participants, VIF were always below  
309 1.1 (being 5 a typical cutoff above which the estimation would be compromised; Mumford et  
310 al., 2015). Thus, even despite the relationship among variables, the results of our main analyses  
311 are still meaningful. The corresponding beta weight maps obtained showed the regions where  
312 the effect of our variables of interest remained significant even when instruction's length and  
313 response speed were included.

314 Finally, even when the distance measure employed to build the neural RDMs (i.e., Pearson  
315 correlation) is insensitive to differences in mean signal intensity between conditions, differences  
316 in signal variance could be affecting it (Walther et al., 2016). For that reason, these analyses as  
317 well as the reward-related tests (see below), were repeated after a z-normalization of the

318 multivoxel activity patterns, ensuring equal mean (0) and standard deviation (1) across all pairs  
319 of trials. The results thus obtained did not differ from the initial non-normalized ones, so we do  
320 not report them here.

321 *ROI-based RSA.* The previous analysis identified brain areas encoding instructions according  
322 to each one of three proactive control variables, separately. We next ran ROI analyses to further  
323 explore the role of the three variables for task coding in these regions. Specifically, we  
324 estimated the extent to which each of the manipulated control variables explained the neural  
325 organization in the ROIs identified in the previous analysis. We followed a Leave-One-Subject-  
326 Out (LOSO) cross-validation procedure (Esterman, Tamber-Rosenau, Chiu, & Yantis, 2010),  
327 using the searchlight maps obtained before. First, we identified regions sensitive to each of the  
328 three models for each participant, running a group level *t*-test with the corresponding maps from  
329 the rest of the sample, i.e., excluding their own data. Significant clusters showing consistency  
330 across all LOSO iterations were selected as ROIs, and inverse normalized to the participants'  
331 native space. In a second step, we estimated the ROIs RDMs and correlated them with the three  
332 models RDMs. Importantly, thanks to the LOSO procedure we avoided circularity in the  
333 analysis, as independent data was employed to select the ROIs and to compute de correlations  
334 with the models. The correlation coefficients (for each participant, one per ROI and model)  
335 were then introduced in a repeated measures ANOVA, with ROI and Model as factors, and the  
336 interaction term was examined to detect heterogeneity in task encoding organization across  
337 regions (Reverberi, Gorgen, & Haynes, 2012). Interactions were further characterized by one  
338 sample *t*-tests, in order to determine which models had an effect on the different regions studied.  
339 Whenever the normality assumption was not met (assessed with the Saphiro-Wilk test), we  
340 employed Wilcoxon signed-rank tests instead. All *P* values were Bonferroni-corrected for  
341 multiple comparisons, adjusting them to the number of ROIs explored.

342 Additionally, we aimed to extrapolate our findings to regions consistently found in the literature  
343 during both practiced (e.g. Woolgar, Hampshire, Thompson, & Duncan, 2011) and novel (e.g.  
344 González-García et al., 2017) task preparation, and in general, when demanding cognitive

345 processing is deployed (Duncan, 2010). This set of brain areas belong to the Multiple Demand  
346 Network (MDN; Duncan, 2010), which includes the bilateral RLPFC, MFG, IFS, anterior  
347 insula/frontal operculum (aIfO) area, IPS, anterior cingulate cortex (ACC) and pre-  
348 supplementary area (preSMA). To assess the organization of novel task encoding across this  
349 MDN, we employed functionally derived masks of its nodes (from Fedorenko, Duncan, &  
350 Kanwisher, 2013; template available at <http://imaging.mrc-cbu.cam.ac.uk/imaging/MDsystem>),  
351 inverse normalized them to the participants' native space, and followed the same ROI-approach  
352 as above, extracting each ROI RDM and correlating it with the models' matrices. Again,  
353 correlation coefficients were entered into a repeated measures ANOVA with ROI and Model as  
354 factors, interactions were examined, and finally, a series of one-sample *t*-tests (or Wilcoxon  
355 signed-rank test when normality was violated) were conducted.

356 *Analysis of reward-related effects on RSA results.* A final goal of our study was to assess  
357 whether the representational space of novel instructions was affected by motivation. Our initial  
358 hypothesis was that reward would polarize the representational geometry, enhancing the effect  
359 of our control-related variables at structuring rule encoding. In other words, and taking as an  
360 example the target category variable, we assessed whether reward expectations would increase  
361 the distance between representations of instructions referring to different stimulus categories (in  
362 extension to the other variables, indicated as *different-condition dissimilarity*), while decreasing  
363 the distance among those referring to same target category (*same-condition dissimilarity*). Our  
364 second, alternative hypothesis was that reward would exert a general effect, globally increasing  
365 the distances among instruction representations, independently of the other variables  
366 manipulated. In this sense, we expected that both *different* and *same-condition dissimilarity*  
367 would be increased in rewarded trials, in comparison with non-rewarded ones. The two  
368 possibilities would be compatible with previous findings showing that reward expectancy  
369 enhances rule decodability (Etzel et al., 2016).

370 To test these two hypotheses, we run ROI analyses (Fig. 2d) for each of our control-related  
371 variables, focusing on the regions that resulted statistically significant in the main RSA. To do



372 so, at the individual level and for each variable, we first ran a searchlight and generated four  
373 whole-brain maps containing dissimilarity values among: (1) same-condition rewarded trials;  
374 (2) different-conditions rewarded trials; (3) same-condition non-rewarded trials; and (4)  
375 different-conditions non-rewarded trials. These values were the result of averaging and  
376 normalizing (with the Fisher transformation) the pertinent cells of the neural RDM (see Fig. 2c  
377 for an example) in each searchlight iteration. The maps thus obtained were normalized to the  
378 MNI space, so we could extract participants' mean dissimilarities for each of our ROIs using  
379 MarsBar (Brett, Anton, Valabregue, & Poline, 2002). After that, and for each ROI and variable,  
380 we conducted two Wilcoxon signed-rank tests (Nili et al., 2014). First, to assess our main  
381 hypothesis, we tested whether  $(\text{DifferentCond.}_{\text{Rewarded}} - \text{SameCond.}_{\text{Rewarded}}) >$   
382  $(\text{DifferentCond.}_{\text{NonRewarded}} - \text{SameCond.}_{\text{NonRewarded}})$ . To explore the second possible hypothesis, we  
383 collapsed across same and different conditions, and tested if  $(\text{DifferentCond.}_{\text{Rewarded}} +$   
384  $\text{SameCond.}_{\text{Rewarded}})/2 - (\text{DifferentCond.}_{\text{NonRewarded}} + \text{SameCond.}_{\text{NonRewarded}})/2$  was greater than 0  
385 (Fig 2c). In both analyses, we corrected for multiple comparisons (number of ROIs being tested)  
386 with an FWE threshold of  $P < .05$ .

387 Last, to investigate the relevance for behavior of the effect of motivation on representational  
388 structure, we correlated this effect with behavioral data. Specifically, for each participant, we  
389 computed the average decrease in dissimilarity and in the inverse efficiency scores (IES;  
390 Townsend & Ashby, 1978) linked to rewarded trials (in comparison with non-rewarded ones).  
391 The IES was employed in this analysis to take into account, simultaneously, improvements in  
392 accuracy and response speed. As we performed as many correlations as ROIs assessed in this  
393 analysis, we again controlled for multiple comparisons with an FWE threshold of  $P < .05$ .

394 Additionally, to explore the possibility of motivation exerting an effect during the subsequent  
395 implementation of instructions, we also ran the analyses detailed above with beta images  
396 obtained from this stage.

397 *MVPA-based assessment of reward effects.*

398 Finally, to further connect our results with previous findings, we performed multivoxel pattern

399 analysis (MVPA) to explore the effect of reward on decoding precisions (Etzel et al., 2016). We  
400 decoded the two conditions of each of our three control-related variables, training three binary  
401 classifiers: one for distinguishing between within versus across-dimension integration  
402 instructions, other for single versus sequential response requirements, and the last one for faces  
403 and food-related trials. This was done separately for rewarded and non-rewarded trials. Again,  
404 we used non-normalized and unsmoothed trial-wise beta images from the encoding stage. As we  
405 aimed to detect any region with reward-related increases in task decodability, we performed the  
406 MVPA in a whole brain fashion, using searchlight (instead of biasing the results using ROIs  
407 resulting from the RSA). In each searchlight iteration, we followed a leave one-run-out cross-  
408 validation approach, training a linear support-vector machine classifier (C=1; Pereira, Mitchell,  
409 & Botvinick, 2009) with five of our six runs, and testing it with the remaining one, in an  
410 iterative fashion. Then, for each of our variables, we subtracted the accuracy map obtained from  
411 non-rewarded trials to the map from rewarded ones, and then normalized and smoothed these  
412 images, to conduct an above zero one-sample *t*-test at the group level. This way, we assessed the  
413 benefits in classification precision associated with reward.

## 414 **Results**

### 415 *Behavioral results*

416 We analyzed RT and accuracy data separately, conducting two repeated measures ANOVA with  
417 four factors, corresponding to the four variables manipulated: dimension integration (within vs.  
418 across), response set complexity (single vs. sequential), category (faces vs. food items) and  
419 motivation (rewarded vs. non-rewarded). Importantly, the main effect of motivation was  
420 statistically significant on both accuracy ( $F_{1,31} = 4.97, P < .05, \eta_p^2 = .14$ ) and RT ( $F_{1,31} = 6.52,$   
421  $P < .05, \eta_p^2 = .17$ ) data, with more accurate (rewarded:  $M = 0.85, SD = 0.11$ ; non-rewarded:  $M$   
422  $= 0.83, SD = 0.12$ ) and faster (rewarded:  $M = 1.16, SD = 0.21$ ; non-rewarded:  $M = 1.20, SD =$   
423  $0.20$ ) responses on the rewarded condition (see Fig. 3). This indicates that participants made use

424 of reward cues and the economic incentives had the expected effect on behavior, improving its  
425 efficiency

426 In addition, accuracy data showed a main effect of dimension integration ( $F_{1,31} = 9.24, P < .05,$   
427  $\eta_p^2 = .23$ ), with better performance when within-dimension integration was required (within  
428 dimension:  $M = .86, SD = 0.13$ ; across dimensions:  $M = .83, SD = 0.12$ ), and a significant  
429 three-way interaction of category, response set complexity and dimension integration ( $F_{1,31} =$   
430  $4.46, P = .043, \eta_p^2 = .13$ ). Even despite the lack of hypothesis regarding an interaction at this  
431 level, we performed post hoc pair-wise comparisons, which revealed that the interaction was  
432 driven by less robust ( $P > .05$ ) differences among within and across-dimensions trials that  
433 required a single response and was food-related (while, in the rest of combinations of  
434 independent variables, this difference was significant).

435 On the other hand, RT results also showed a main effect of dimension integration ( $F_{1,31} = 61.81,$   
436  $P < .001, \eta_p^2 = .67$ ) in the same direction as above (within-dimension:  $M = 1.12, SD = 0.17$ ;  
437 across-dimensions:  $M = 1.24, SD = 0.2$ ), and a main effect of category ( $F_{1,31} = 74.89, P < .001,$   
438  $\eta_p^2 = .71$ ), with faster responses to food-related instructions (faces:  $M = 1.23, SD = 0.21$ ; food  
439 items:  $M = 1.14, SD = 0.19$ ). Neither the effect of response set complexity (accuracy:  $F_{1,31} =$   
440  $0.31, P = .579, \eta_p^2 = .01$ ; reaction time:  $F_{1,31} = 0.21, P = .653, \eta_p^2 = .01$ ) nor any other ANOVA  
441 term resulted significant in the behavioral measures (main effect of Category on accuracy:  $F_{1,31}$   
442  $= 3.23, P = .082, \eta_p^2 = .094$ ; all interactions terms, except the ones stated above,  $P > .100$ ).

443 *Univariate results: reward-related activations during instruction encoding.*

444 We first assessed mean activity during novel instruction encoding, comparing rewarded against  
445 non-rewarded trials. To do so, we performed a univariate GLM, defining regressors for each  
446 combination of variables (e.g.: within-dimension integration, single response, face-related  
447 rewarded trials), separately for the encoding and the implementation stages. A group level *t*-test  
448 showed that, in accordance with our expectations and previous literature (Parro et al., 2017), the  
449 basal ganglia and fronto-parietal cortices were more active for rewarded than non-rewarded

450 instruction encoding. We observed peaks of activation (see Fig. 4) in the bilateral inferior  
451 frontal junction (IFJ), premotor and supplementary motor areas (left: [-33, 5, 26],  $z = 5.07$ ,  $k =$   
452 489; right: [33, 2, 59],  $z = 4.79$ ,  $k = 572$ ), cingulate cortex ([-9, 5, 32],  $z = 5.48$ ,  $k = 20$ ),  
453 bilateral IPS extending into the precuneus (left: [-18, -64, 35],  $z = 4.77$ ,  $k = 357$ ; right: [33, -52,  
454 53],  $z = 4.36$ ,  $k = 324$ ), the accumbens, ventral portion of the caudate and thalamus ([12, -22,  
455 20],  $z = 5.13$ ,  $k = 1176$ ), inferior temporal gyrus ([48, -58, -13],  $z = 4.48$ ,  $k = 52$ ), occipital  
456 cortex ([30, -61, -25],  $z = 5$ ,  $k = 1364$ ) and midbrain ([0, -31, -4],  $z = 5.19$ ,  $k = 255$ ). Thus,  
457 regions involved in reward processing (Haber & Knutson, 2009), as well as in cognitive control  
458 paradigms with monetary incentive manipulations (e.g. Engelmann, 2009), were engaged by our  
459 task, indicating the success of the reward manipulation.

460 *Model-based RSA results: instruction encoding structured by proactive-control variables.*

461 We aimed to identify regions whose organization during task encoding was explained by  
462 dimension integration, response set complexity and target category. With that purpose, we  
463 employed an RSA (Kriegeskorte, Mur, & Bandettini, 2008) to compare the representational  
464 dissimilarity matrices (RDMs) found in neural data during the encoding stage with theoretical  
465 RDMs corresponding to the three proactive control-related variables (see Fig. 2). In neural  
466 RDMs, each cell contained the dissimilarity ( $1 - \text{Pearson correlation}$ ) between the multivariate  
467 patterns of activation of two trials. In the theoretical RDMs, cells contained dissimilarities (1:  
468 maximal, 0: minimal) that we would expect if a certain variable organized encoding (i.e.: for  
469 target category, all faces-related trials would be minimally dissimilar, while face and food-  
470 related trials would be maximally dissimilar). Using searchlight (Kriegeskorte et al., 2006), we  
471 Spearman-correlated neural and theoretical RDMs across the brain and obtained maps showing  
472 how well these three variables captured the representational space of different areas. The  
473 modality of **dimension integration** (Fig. 5a) only had a significant effect on rule encoding at  
474 the left MFG and IFG, incurring into the IFS ([-51, 20, 26],  $k = 642$ ). **Response set complexity**  
475 (Fig. 5b), on the other hand, organized task representations on a wide cluster including the  
476 bilateral IFG, premotor, supplementary and primary motor cortices, somatosensory area, middle

477 temporal gyrus and superior and inferior parietal lobe extending along the IPS ([-42, -31, 44], k  
478 = 8583) and in the left parahippocampal cortex ([-18, -40, -1], k = 301). Finally, in the case of  
479 the **target category** RSA (Fig. 5c), significant correlations were found in an extensive cluster  
480 on the left hemisphere covering the IFG incurring into the IFJ, the fusiform gyrus, the temporo-  
481 parietal junction (TPJ), the inferior and middle temporal gyrus and the precuneus ([-39, -67, 17],  
482 k = 5581). On the right hemisphere, the analysis was also significant on the right middle  
483 temporal gyrus and TPJ ([39, -58, 23], k = 442) and the IFG ([42, 26, 14, k = 295]. Finally, the  
484 medial superior frontal gyrus ([-9, 53, 26], k = 377) was also involved.

485 As instructions' length and speed of responses varied among some of our variables, we  
486 performed an additional multiple regression analysis, in which we included our three theoretical  
487 models, an RDM based on dissimilarities in length, and another one based on RT as regressors.  
488 Importantly, the multiple regression statistical model was examined to detect an excess of  
489 collinearity which could have impaired the interpretability of these results. We computed the  
490 VIF for all the regressors and across our whole sample of participants, and all of were under 1.1,  
491 an index of good estimability of regression weights. The beta maps (one per model) obtained  
492 after iterating the analysis in a searchlight procedure ensured that the variance linked to our  
493 RSA models was not misattributed due to differences in instruction length or speed of  
494 responses. Importantly, the results obtained this way were very similar to the ones extracted  
495 with the standard approach, identifying the same clusters than before.

496 We also conducted a **conjunction analysis** to assess the overlap among regions common to the  
497 three organizational schemes. Only the left IFG and IFJ resulted significant in this test (Fig. 6).

498 *LOSO-based ROI analysis: assessing confluence of models within regions.*

499 The previous analyses left unexplained the extent to which each of the brain areas isolated by  
500 RDM analyses reflected in their organization the three manipulated variables. Furthermore, the  
501 conservative correction for multiple comparisons used in the searchlight could overshadow this  
502 effect elsewhere in the brain. To shed some light upon this issue, we employed a more sensitive  
503 ROI analysis, together with a LOSO approach to avoid double dipping when selecting regions.

504 All the clusters identified in the main group results (Fig. 5) were consistently found across all  
505 participants with the LOSO approach, with the exception of the medial superior frontal gyrus  
506 under the category model, which was absent in four subjects and thus not included in the  
507 analysis. The correlations of the ROIs' RDMs and the three models' matrices were analyzed  
508 with a repeated measures ANOVA, in which we found a significant interaction of ROI and  
509 Model ( $F_{12, 348} = 6.050$ ,  $P < .001$ ,  $\eta_p^2 = .173$ ), evidencing variability in instruction coding  
510 structure across regions. We then ran one sample  $t$ -tests or Wilcoxon signed-rank tests  
511 (depending on data distribution) to assess model performance in each ROI (see Table 1). The  
512 general pattern obtained replicated the searchlight results: the model which originally identified  
513 each specific ROI in the searchlight was the one explaining most robustly its encoding activity.  
514 Further, in almost all the regions, we did not find enough evidence supporting the effect of the  
515 remaining variables. Converging with the previous analyses, the left IFG identified with the  
516 dimension integration model was also significantly correlated with response set complexity and  
517 category. Similarly, the left IFG cluster found in the category RSA was correlated with the  
518 dimension integration model too. In addition, this confluence of models analysis revealed that  
519 the response set model was also significant in the category-related cluster involving the left  
520 fusiform and precuneus (see Table 1).

521 *ROI analysis spanning Multiple Demand Network regions.*

522 Following a similar strategy as above, we also examined task encoding organization across the  
523 regions comprising the MD network. We extracted each MD region's RDM and correlated it  
524 with our three models' RDM, and then entered the correlation coefficients into a repeated  
525 measures ANOVA. Again, a significant ROI\*Model interaction was found ( $F_{20, 620} = 2.168$ ,  $P$   
526  $= .002$ ,  $\eta_p^2 = .065$ ). To assess which models significantly structured activations across MD  
527 ROIs, we conducted one-sample  $t$ -tests or Wilcoxon signed-rank tests when data were not  
528 normally distributed (see Table 2).

529 Only a subset of MD network regions encoded instructions consistently according to any of the  
530 proactive control variables, and all of them were located on the left hemisphere and in the LPFC

531 and parietal cortex. The findings were, however, consistent with the searchlight and ROI-related  
532 results presented so far. The three variables exerted an effect on different left lateral prefrontal  
533 sections: dimension integration and response complexity on the IFG; dimension integration and  
534 target category on the more dorsal MFG; and finally, category on the RLPFC. Response  
535 complexity was the attribute which most robustly captured representational organization in the  
536 IPS.

537 *Effects of reward on representational geometry.*

538 We then explored the effects of motivation in each of the ROIs encoding different attributes of  
539 the instructions (Fig. 5), assessing two possible mechanisms that could underlie the behavioral  
540 improvements linked to reward (Fig. 2). On the one hand, we tested whether reward made our  
541 variables more efficient in sharpening the representational space (Fig. 2d, Hypothesis 1), In  
542 other words, and taking as an example the target category variable, we assessed whether reward  
543 expectations would increase the distance between representations of instructions referring to  
544 different stimulus categories (in extension to the other variables, indicated as *different-condition*  
545 *dissimilarity*), while decreasing the distance among those referring to same target category  
546 (*same-condition dissimilarity*). On the other, we tested the alternative possibility that  
547 dissimilarities would be, in general, greater in the rewarded trials (Fig 2d, Hypothesis 2),  
548 regardless of the variables manipulated (i.e., regardless of the pair of instructions being same or  
549 different-condition). This could reflect a mechanism for making rule representations more  
550 distinguishable among each other, and also, it would be compatible with the increase in rule  
551 decoding accuracy that has been linked to motivation in previous reports (Etzel et al., 2016). With  
552 that purpose, we extracted, for each region, the average dissimilarity among pairs of instructions  
553 pertaining to the same and different conditions, separately for rewarded and non-rewarded trials.  
554 We then used Wilcoxon signed-rank tests (Nili et al., 2014) to check whether the difference  
555 between different-condition and same-condition trials was larger in the rewarded than in the  
556 non-rewarded condition, and also, whether the mean dissimilarity (collapsing across same and  
557 different-condition) was increased by motivation.

558 In the first case, no reward-related differences were observed for any of the instruction-related  
559 variables (all  $P$ s  $>.1$ ). It is important to note, however, that these results (as most of the findings  
560 presented in this study) are anchored to the instruction's encoding stage, in which proactive  
561 control configuration takes place. To explore the possibility that the hypothesized interaction  
562 shaped neural activations during the later implementation phase (more related to reactive  
563 control; Braver, 2012; Palenciano, González-García, Arco, & Ruz, 2018), we conducted a  
564 further test employing beta images from this epoch. However, and again, the expected effect  
565 was not significant for any of the ROIs examined (all  $P$ s  $>.1$ ).

566 When addressing the second hypothesis, surprisingly, we found the opposite pattern: reward  
567 systematically decreased the dissimilarity values in all the ROIs evaluated (all  $P$ s  $< .05$ , see  
568 Table 2). To test the behavioral relevance of this finding we correlated, across our participants,  
569 the average decrease in dissimilarities associated with reward, with the benefit of motivation on  
570 performance (IES; Townsend & Ashby, 1978). We found that in fact, the decrease in  
571 representational distances due to reward was significantly correlated with the motivation-related  
572 improvements in behavioral performance. Furthermore, this seemed to be a quite robust effect,  
573 being present in all of the ROIs included in the analysis (see Table 3 for further details).

#### 574 *MVPA results*

575 We finally aimed to explore the effect of reward directly on decoding accuracies, employing  
576 MVPA (Haxby, Connolly, & Guntupalli, 2014), as it has been previously reported during rule  
577 encoding in a classic, repetitive task-switching setting (Etzel et al., 2016). We discriminated  
578 between the two conditions of each instruction-related variable (i.e., one among faces and food-  
579 related trials, other for single versus sequential response requirements, and a last one for within  
580 versus across-dimension integration instructions) separately for rewarded and non-rewarded  
581 trials. We trained and tested our classifiers across the whole brain using searchlight and  
582 obtained, as a result, an accuracy map for each motivation condition and variable. Nonetheless,  
583 while classification was above chance in different brain regions for the three variables, we did  
584 not detect any differences in accuracies between rewarded and non-rewarded trials, as no cluster



585 survived at the group-level the *t*-test assessing above zero differences between the two  
586 motivation conditions.

### 587 **Discussion**

588 In the present study, we aimed to characterize the representational space for novel instructions  
589 during their proactive preparation. We assessed whether variables linked to proactive control  
590 organized encoding activity patterns and whether this structure was affected by reward  
591 expectations. Our results portrayed a complex landscape, where different organizational  
592 principles governed instruction encoding in FP cortices and lower-level perceptual and motor  
593 areas.

594 The left IFG/IFJ reflected the most complex and overarching representational structure, with  
595 activity patterns structured by dimension integration, response complexity and target category.  
596 Robust evidence supports the role of the IFJ in task-set reconfiguration (Brass, Derrfuss,  
597 Forstmann, & Cramon, 2005) in practiced (e.g. Woolgar, Hampshire, Thompson, & Duncan,  
598 2011) and novel contexts (e.g. González-García et al., 2016; Muhle-Karbe et al., 2017),  
599 orchestrating neural dynamics during attentional selection (e.g. Baldauf & Desimone, 2014).  
600 This region seems to be involved in task-set maintenance (Sakai, 2008), selecting task-relevant  
601 information represented in perceptual regions (Cole, Reynolds, et al., 2013; Miller & Cohen,  
602 2001). The current study advances our knowledge about the structure underlying *how*  
603 information is coded during novel instruction encoding, and stresses the diversity of task  
604 parameters that orchestrate task encoding in the IFG/IFJ. Such a complex, multidimensional  
605 representational space (Rigotti et al., 2013) could be key to support the richness and flexibility  
606 of human behavior in novel environments. This perspective qualifies recent research, based on  
607 MVPA, that highlights the compositionality characterizing representations held in the IFG  
608 (Cole, Laurent, et al., 2013; Deraeve, Vassena, & Alexander, 2019; Reverberi, Görden, &  
609 Haynes, 2012), by which complex tasks are coded by combining their simpler constituent  
610 elements.

611 The IPS also encoded novel rules proactively, but now according to response complexity. While

612 this is quite consistent with previous studies linking the parietal cortex to action preparation, it  
613 is worth noticing the distinction found in our data between parietal and prefrontal regions, a  
614 finding further confirmed with a more sensitive ROI analysis. Dimension integration, the  
615 variable manipulated to appeal to a higher-level task goal representation, had an effect only on  
616 LPFC, while the IPS was linked to the more specific response-set complexity (De Baene &  
617 Brass, 2014; Rubinstein et al., 2001). The frequent coactivation of IFG/IFJ and IPS in demanding  
618 paradigms (Duncan, 2010) had complicated the identification of their separate contributions.  
619 The differential pattern we observed is highly relevant to disentangle their proactive role.  
620 Interestingly, the emerging picture portrays the IFG/IFJ and the IPS collaborating during novel  
621 task representation, with the former maintaining overarching representations of all relevant  
622 variables, and the latter activating the relevant stimulus-response contingencies (see also Muhle-  
623 Karbe et al., 2014). The use of RSA in our paradigm provides a deeper understanding of this  
624 process, emphasizing that the proposed two-stage preparatory mechanism also guides task-set  
625 encoding in FP cortices. In this sense, variables key for abstract goal or specific S-R settings  
626 become relevant differentially depending on the region.

627 Additional medial and lateral frontal cortices also participate in the FP network and are  
628 frequently recruited during task preparation (Duncan, 2010). Consequently, we also examined  
629 instruction coding in these MD regions. Our findings highlighted other LPFC areas reflecting  
630 target category (both the RLPFC and MFG) and dimension integration (MFG). The overall  
631 pattern of results obtained both with whole-brain and with ROI approaches reflects high  
632 heterogeneity within the FP network in general, and in the LPFC in particular, in terms of the  
633 attributes structuring task-set representation. In contrast, we did not obtain evidence supporting  
634 proactive task-set encoding in the ACC/preSMA and the alFO regions. This finding fits with the  
635 subdivision of the FP network into two differentiated components: one anchored in the LPFC  
636 and IPS, and a second one composed by the ACC and the alFO (Dosenbach et al., 2007;  
637 Palenciano et al., 2018). In line with our results, anticipatory task coding has been  
638 predominantly found in regions from the former rather than in the latter (Crittenden, Mitchell, &

639 Duncan, 2016). Ultimately, the variability found within the FP control network during proactive  
640 novel task setting (Palenciano et al., 2018), with different processes and representational  
641 formats being combined, could be key to maximize flexibility.

642 Fronto-parietal cortices were not the sole brain regions encoding novel instruction parameters.  
643 Activity in fusiform gyri was organized according to target category, whereas patterns in  
644 somatomotor cortices reflected response complexity. While these regions are not associated *per*  
645 *se* with proactive control, their involvement reflects that their representational geometry is tuned  
646 in an anticipatory fashion by relevant task parameters conveyed by instructions. It is important  
647 to stress that all the results discussed were locked to instruction encoding, where no target  
648 stimuli had been presented, neither any specific motor response could have been prepared.  
649 These findings suggest that FP areas exert a bias in posterior cortices, according to the content  
650 of instructions. Supporting this, increments of mean activity (Esterman & Yantis, 2010) and  
651 target-specific information encoding (e.g. Stokes, Thompson, Nobre, & Duncan, 2009) have  
652 been reported in perceptual and motor regions during preparation. Importantly, these changes  
653 have been linked to boosts in functional connectivity between the FP and posterior cortices  
654 (González-García et al., 2016; Sakai & Passingham, 2006). In direct relation to our findings, a  
655 recent study showed that the representational organization in regions along the visual pathway is  
656 dynamically adapted to task demands (Nastase et al., 2017). Our current results add to these  
657 findings by showing that representational space tuning could be a mechanism of preparatory  
658 bias, which could reflect predictive coding principles where iterative loops of feedback and  
659 feedforward communication shape cognition (Friston, 2005).

660 Crucially, the structure of information encoded by all these regions was sensitive to trial-wise  
661 motivational states. Surprisingly, reward expectation diminished the dissimilarities between the  
662 representations of the instructions although preserving the organizational scheme found in each  
663 area. Based on recent findings of increased task decodability (Etzel et al., 2016), we had  
664 hypothesized that reward would either polarize the representational structure or increase the  
665 representational distances overall. Results were, however, in the opposite direction, even when

666 our reward manipulation was successful at boosting performance and also increased activity in  
667 control and reward-related regions (Parro et al., 2017). Most importantly, decreases in  
668 dissimilarities were also robustly correlated with behavioral improvements. Taking into account  
669 that additional analysis employing MVPA and using data from the implementation stage  
670 corroborated these results, their implication must be thoughtfully considered. One possibility is  
671 that the decrease in dissimilarities is generated by a general boost of reward in signal-to-noise  
672 ratio. Although our results persisted after normalizing data across trials, a reward-related  
673 reduction of multivariate noise pattern could still be possible, and it could benefit task coding in  
674 the absence of the hypothesized RSA results. However, the MVPA did not reveal improved task  
675 classification accuracy in the rewarded condition, and thus this interpretation remains uncertain.  
676 Alternatively, motivation could have influenced task coding in ways that our searchlight  
677 procedure was not sensitive to. That would be the case if reward affected the spatial distribution  
678 of information: as ROIs were defined by size-fixed searchlight spheres, and were equal in  
679 rewarded and non-rewarded conditions, an effect like that would remain shadowed. Finally, the  
680 task complexity could also be key. In less demanding situations such as repetitive task switching  
681 (Etzel et al., 2016), reward could directly sharpen task encoding representations. In novel  
682 environments, however, motivation could exert a more general effect at the process level -  
683 instead of at the representational one. It could increase the efficiency of task reconfiguration  
684 (Braem & Egner, 2018), as indexed by the improvements in behavior, while the specific rule  
685 representations would remain equally structured. Nonetheless, more research is needed to  
686 properly characterize the intricate interactions among proactive control and motivation (Pessoa,  
687 2017) in rich task environments, more akin to daily life situations.

688 The current study entails some limitations that constrain the scope of our findings and call for  
689 further research. On the one hand, the nature of our paradigm demanded the selection of a few  
690 instruction-organizing variables. Some other dimensions, critical for anticipatory encoding, may  
691 have been left unaddressed. Furthermore, non-linear combinations of variables could add to the  
692 organization principles governing control regions (Rigotti et al., 2013). Considering an

693 increasing number of plausible models in more complex and/or naturalistic scenarios, together  
694 with data-driven methods such as multidimensional scaling or component analysis, will  
695 complement our results. On the other hand, our main dependent variable (fMRI hemodynamic  
696 signal) provided spatially precise, but temporal impoverished data. Temporally resolved  
697 techniques, such as electroencephalography or magnetoencephalography, could be key to unveil  
698 the temporal dynamics of the representational patterns.

699 Overall, our findings provide novel insights on how verbal complex novel instructions organize  
700 proactive brain activations. The emerging picture departs from pure localizationist approaches  
701 where brain regions carry fixed information about concrete cognitive processes. Rather, the  
702 different dimensions relevant for efficient instructed action shape brain activity across an  
703 extended set of areas, flexibly structuring encoding activity according to the relevant task  
704 parameters.

## 705 **References**

706 Arco, J. E., González-García, C., Díaz-Gutiérrez, P., Ramírez, J., & Ruz, M. (2018). Influence  
707 of activation pattern estimates and statistical significance tests in fMRI decoding analysis.  
708 <https://doi.org/10.1101/344549>

709 Baldauf, D., & Desimone, R. (2014). Neural Mechanisms of Object-Based Attention. *Science*,  
710 *344*(6182), 424–427. <https://doi.org/10.1126/science.1247003>

711 Bourguignon, N. J., Braem, S., Hartstra, E., De Houwer, J., & Brass, M. (2018). Encoding of  
712 Novel Verbal Instructions for Prospective Action in the Lateral Prefrontal Cortex:  
713 Evidence from Univariate and Multivariate Functional Magnetic Resonance Imaging  
714 Analysis. *Journal of Cognitive Neuroscience*, 1–15. <https://doi.org/10.1162/jocn>

715 Braem, S., & Egner, T. (2018). Getting a Grip on Cognitive Flexibility. *Current Directions in*  
716 *Psychological Science*, *27*(6), 470–476. <https://doi.org/10.1177/0963721418787475>

717 Brass, M., Derrfuss, J., Forstmann, B., & Cramon, D. Y. von. (2005). The role of the inferior  
718 frontal junction area in cognitive control. *Trends in Cognitive Sciences*, *9*(7), 314–316.

- 719 <https://doi.org/10.1016/J.TICS.2005.05.001>
- 720 Brass, M., Liefoghe, B., Braem, S., & De Houwer, J. (2017). Following new task instructions:  
721 Evidence for a dissociation between knowing and doing. *Neuroscience & Biobehavioral*  
722 *Reviews*, *81*, 16–28. <https://doi.org/10.1016/J.NEUBIOREV.2017.02.012>
- 723 Braver, T. S. (2012). The variable nature of cognitive control: A dual mechanisms framework.  
724 *Trends in Cognitive Sciences*, *16*(2), 106–113. <https://doi.org/10.1016/j.tics.2011.12.010>
- 725 Brett, M., Anton, J., Valabregue, R., & Poline, J. (2002). Region of interest analysis using the  
726 MarsBar toolbox for SPM 99. *Neuroimage*, *16*, 99. Retrieved from [http://www.mrc-](http://www.mrc-cbu.cam.ac.uk/Imaging/marsbar.html)  
727 [cbu.cam.ac.uk/Imaging/marsbar.html](http://www.mrc-cbu.cam.ac.uk/Imaging/marsbar.html)
- 728 Cole, M. W., Bagic, A., Kass, R., & Schneider, W. (2010). Prefrontal Dynamics Underlying  
729 Rapid Instructed Task Learning Reverse with Practice. *Journal of Neuroscience*, *30*(42),  
730 14245–14254. <https://doi.org/10.1523/JNEUROSCI.1662-10.2010>
- 731 Cole, M. W., Braver, T. S., & Meiran, N. (2017). The task novelty paradox: Flexible control of  
732 inflexible neural pathways during rapid instructed task learning. *Neuroscience and*  
733 *Biobehavioral Reviews*. <https://doi.org/10.1016/j.neubiorev.2017.02.009>
- 734 Cole, M. W., Ito, T., & Braver, T. S. (2016). The Behavioral Relevance of Task Information in  
735 Human Prefrontal Cortex. *Cerebral Cortex*, *26*(6), 2497–2505.  
736 <https://doi.org/10.1093/cercor/bhv072>
- 737 Cole, M. W., Laurent, P., & Stocco, A. (2013). Rapid instructed task learning: A new window  
738 into the human brain's unique capacity for flexible cognitive control. *Cognitive, Affective*  
739 *and Behavioral Neuroscience*, *13*(1), 1–22. <https://doi.org/10.3758/s13415-012-0125-7>
- 740 Cole, M. W., Patrick, L. M., & Braver, T. S. (2018). A role for proactive control in rapid  
741 instructed task learning. *Acta Psychologica*, *184*, 20–30.  
742 <https://doi.org/10.1016/J.ACTPSY.2017.06.004>
- 743 Cole, M. W., Reynolds, J. R., Power, J. D., Repovs, G., Anticevic, A., & Braver, T. S. (2013).  
744 Multi-task connectivity reveals flexible hubs for adaptive task control. *Nature*

- 745        *Neuroscience*, 16(9), 1348–1355. <https://doi.org/10.1038/nn.3470>
- 746    Crittenden, B. M., Mitchell, D. J., & Duncan, J. (2016). Task Encoding across the Multiple  
747        Demand Cortex Is Consistent with a Frontoparietal and Cingulo-Opercular Dual Networks  
748        Distinction. *The Journal of Neuroscience : The Official Journal of the Society for*  
749        *Neuroscience*, 36(23), 6147–6155. <https://doi.org/10.1523/JNEUROSCI.4590-15.2016>
- 750    De Baene, W., & Brass, M. (2014). Dissociating strategy-dependent and independent  
751        components in task preparation. *Neuropsychologia*, 62, 331–340.  
752        <https://doi.org/10.1016/j.neuropsychologia.2014.04.015>
- 753    Deraeve, J., Vassena, E., & Alexander, W. (2019). Conjunction or co-activation? A multi-level  
754        MVPA approach to task set representations. *BioRxiv*, 521385.  
755        <https://doi.org/10.1101/521385>
- 756    Dosenbach, N. U. F., Fair, D. A., Miezin, F. M., Cohen, A. L., Wenger, K. K., Dosenbach, R.  
757        A. T., ... Petersen, S. E. (2007). Distinct brain networks for adaptive and stable task  
758        control in humans. *Proceedings of the National Academy of Sciences*, 104(26), 11073–  
759        11078. <https://doi.org/10.1073/pnas.0704320104>
- 760    Dumontheil, I., Thompson, R., & Duncan, J. (2011). Assembly and Use of New Task Rules in  
761        Fronto-parietal Cortex. *Journal of Cognitive Neuroscience*, 23(1), 168–182.  
762        <https://doi.org/10.1162/jocn.2010.21439>
- 763    Duncan, J. (2010). The multiple-demand (MD) system of the primate brain: mental programs  
764        for intelligent behaviour. *Trends in Cognitive Sciences*, 14(4), 172–179.  
765        <https://doi.org/10.1016/j.tics.2010.01.004>
- 766    Duncan, J., Parr, A., Woolgar, A., Thompson, R., Bright, P., Cox, S., ... Nimmo-Smith, I.  
767        (2008). Goal Neglect and Spearman's g: Competing Parts of a Complex Task. *Journal of*  
768        *Experimental Psychology: General*, 137(1), 131–148. [https://doi.org/10.1037/0096-](https://doi.org/10.1037/0096-3445.137.1.131)  
769        [3445.137.1.131](https://doi.org/10.1037/0096-3445.137.1.131)
- 770    Engelmann, J. B., Damaraju, E., Padmala, S., & Pessoa, L. (2009). Combined effects of

- 771 attention and motivation on visual task performance: Transient and sustained motivational  
772 effects. *Frontiers in Human Neuroscience*, 3. <https://doi.org/10.3389/neuro.09.004.2009>
- 773 Esterman, M., Tamber-Rosenau, B. J., Chiu, Y.-C., & Yantis, S. (2010). Avoiding non-  
774 independence in fMRI data analysis: Leave one subject out. *NeuroImage*, 50(2), 572–576.  
775 <https://doi.org/10.1016/J.NEUROIMAGE.2009.10.092>
- 776 Esterman, M., & Yantis, S. (2010). Perceptual Expectation Evokes Category-Selective Cortical  
777 Activity. *Cerebral Cortex*, 20(5), 1245–1253. <https://doi.org/10.1093/cercor/bhp188>
- 778 Etzel, J. A., Cole, M. W., Zacks, J. M., Kay, K. N., & Braver, T. S. (2016). Reward Motivation  
779 Enhances Task Coding in Frontoparietal Cortex. *Cerebral Cortex*, 26(4), 1647–1659.  
780 <https://doi.org/10.1093/cercor/bhu327>
- 781 Fedorenko, E., Duncan, J., & Kanwisher, N. (2013). Broad domain generality in focal regions of  
782 frontal and parietal cortex. *Proceedings of the National Academy of Sciences of the United*  
783 *States of America*, 110(41), 16616–16621. <https://doi.org/10.1073/pnas.1315235110>
- 784 Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the Royal*  
785 *Society B: Biological Sciences*, 360(1456), 815–836.  
786 <https://doi.org/10.1098/rstb.2005.1622>
- 787 González-García, C., Arco, J. E., Palenciano, A. F., Ramírez, J., & Ruz, M. (2017). Encoding,  
788 preparation and implementation of novel complex verbal instructions. *NeuroImage*, 148,  
789 264–273. <https://doi.org/10.1016/J.NEUROIMAGE.2017.01.037>
- 790 González-García, C., Mas-Herrero, E., de Diego-Balaguer, R., & Ruz, M. (2016). Task-specific  
791 preparatory neural activations in low-interference contexts. *Brain Structure and Function*,  
792 221(8), 3997–4006. <https://doi.org/10.1007/s00429-015-1141-5>
- 793 Haber, S. N., & Knutson, B. (2009). The Reward Circuit: Linking Primate Anatomy and Human  
794 Imaging. *Neuropsychopharmacology*, 35(10), 4–26. <https://doi.org/10.1038/npp.2009.129>
- 795 Hartstra, E., Kühn, S., Verguts, T., & Brass, M. (2011). The implementation of verbal  
796 instructions: An fMRI study. *Human Brain Mapping*, 32(11), 1811–1824.



- 797 <https://doi.org/10.1002/hbm.21152>
- 798 Haxby, J. V., Connolly, A. C., & Guntupalli, J. S. (2014). Decoding Neural Representational  
799 Spaces Using Multivariate Pattern Analysis. *Annual Review of Neuroscience*, *37*(1), 435–  
800 456. <https://doi.org/10.1146/annurev-neuro-062012-170325>
- 801 Hebart, M. N., Gorgen, K., & Haynes, J.-D. (2014). The Decoding Toolbox (TDT): a versatile  
802 software package for multivariate analyses of functional imaging data. *Frontiers in*  
803 *Neuroinformatics*, *8*, 88. <https://doi.org/10.3389/fninf.2014.00088>
- 804 Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional brain  
805 mapping. *Proceedings of the National Academy of Sciences*, *103*(10), 3863–3868.  
806 <https://doi.org/10.1073/pnas.0600244103>
- 807 Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity analysis –  
808 connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, *2*, 4.  
809 <https://doi.org/10.3389/neuro.06.004.2008>
- 810 Liefoghe, B., Wenke, D., & De Houwer, J. (2012). Instruction-based task-rule congruency  
811 effects. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *38*(5),  
812 1325–1335. <https://doi.org/10.1037/a0028148>
- 813 Luria, A. R. (1966). *Higher Cortical Functions in Man*. Boston, MA: Springer US.  
814 <https://doi.org/10.1007/978-1-4684-7741-2>
- 815 Miller, E. K., & Cohen, J. D. (2001). An Integrative Theory of Prefrontal Cortex Function.  
816 *Annual Review of Neuroscience*, *24*(1), 167–202.  
817 <https://doi.org/10.1146/annurev.neuro.24.1.167>
- 818 Muhle-Karbe, P. S., Andres, M., & Brass, M. (2014). Transcranial magnetic stimulation  
819 dissociates prefrontal and parietal contributions to task preparation. *The Journal of*  
820 *Neuroscience : The Official Journal of the Society for Neuroscience*, *34*(37), 12481–  
821 12489. <https://doi.org/10.1523/JNEUROSCI.4931-13.2014>
- 822 Muhle-Karbe, P. S., Duncan, J., De Baene, W., Mitchell, D. J., & Brass, M. (2017). Neural

- 823 Coding for Instruction-Based Task Sets in Human Frontoparietal and Visual Cortex.  
824 *Cerebral Cortex*, 27(3), 1891–1905. <https://doi.org/10.1093/cercor/bhw032>
- 825 Mumford, J. A., Poline, J.-B., & Poldrack, R. A. (2015). Orthogonalization of Regressors in  
826 fMRI Models. *PLOS ONE*, 10(4), e0126255. <https://doi.org/10.1371/journal.pone.0126255>
- 827 Nastase, S. A., Connolly, A. C., Oosterhof, N. N., Halchenko, Y. O., Guntupalli, J. S., Visconti  
828 di Oleggio Castello, M., ... Haxby, J. V. (2017). Attention Selectively Reshapes the  
829 Geometry of Distributed Semantic Representation. *Cerebral Cortex*, 27(8), 4277–4291.  
830 <https://doi.org/10.1093/cercor/bhx138>
- 831 Nichols, T., Brett, M., Andersson, J., Wager, T., & Poline, J. B. (2005). Valid conjunction  
832 inference with the minimum statistic. *NeuroImage*, 25(3), 653–660.  
833 <https://doi.org/10.1016/j.neuroimage.2004.12.005>
- 834 Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., & Kriegeskorte, N. (2014). A  
835 Toolbox for Representational Similarity Analysis. *PLoS Computational Biology*, 10(4),  
836 e1003553. <https://doi.org/10.1371/journal.pcbi.1003553>
- 837 Palenciano, A. F., González-García, C., Arco, J. E., & Ruz, M. (2018). Transient and Sustained  
838 Control Mechanisms Supporting Novel Instructed Behavior. *Cerebral Cortex*.  
839 <https://doi.org/10.1093/cercor/bhy273>
- 840 Parro, C., Dixon, M. L., & Christoff, K. (2017). The Neural Basis of Motivational Influences on  
841 Cognitive Control. <https://doi.org/10.1101/113126>
- 842 Pereira, F., Mitchell, T., & Botvinick, M. (2009). Machine learning classifiers and fMRI: a  
843 tutorial overview. *NeuroImage*, 45(1 Suppl), S199-209.  
844 <https://doi.org/10.1016/j.neuroimage.2008.11.007>
- 845 Pessoa, L. (2009). How do emotion and motivation direct executive control? *Trends in*  
846 *Cognitive Sciences*, 13(4), 160–166. <https://doi.org/10.1016/j.tics.2009.01.006>
- 847 Pessoa, L. (2017). Cognitive-motivational interactions: Beyond boxes-and-arrows models of the  
848 mind-brain. *Motivation Science*, 3(3), 287–303. <https://doi.org/10.1037/mot0000074>

- 849 Popov, V., Ostarek, M., & Tenison, C. (2018). Practices and pitfalls in inferring neural  
850 representations. *NeuroImage*, *174*, 340–351.  
851 <https://doi.org/10.1016/J.NEUROIMAGE.2018.03.041>
- 852 Reverberi, C., Gorgen, K., & Haynes, J.-D. (2012). Distributed Representations of Rule Identity  
853 and Rule Order in Human Frontal Cortex and Striatum. *Journal of Neuroscience*, *32*(48),  
854 17420–17430. <https://doi.org/10.1523/JNEUROSCI.2344-12.2012>
- 855 Reverberi, C., Gorgen, K., & Haynes, J.-D. (2012). Compositionality of Rule Representations in  
856 Human Prefrontal Cortex. *Cerebral Cortex*, *22*(6), 1237–1246.  
857 <https://doi.org/10.1093/cercor/bhr200>
- 858 Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller, E. K., & Fusi, S.  
859 (2013). The importance of mixed selectivity in complex cognitive tasks. *Nature*,  
860 *497*(7451), 585–590. <https://doi.org/10.1038/nature12160>
- 861 Ritchie, J. B., Bracci, S., & Op de Beeck, H. (2017). Avoiding illusory effects in  
862 representational similarity analysis: What (not) to do with the diagonal. *NeuroImage*, *148*,  
863 197–200. <https://doi.org/10.1016/j.neuroimage.2016.12.079>
- 864 Rubinstein, J. S., Meyer, D. E., Evans, J. E., Allport, A., Carr, T., Kieras, D., ... Stemberg, S.  
865 (2001). Executive Control of Cognitive Processes in Task Switching Federal Aviation  
866 Administration. *Journal of Experimental Psychology: Human Perception and*  
867 *Performance*, *27*(4), 763–797. <https://doi.org/10.1037/0096-1523.27.4.763>
- 868 Sakai, K. (2008). Task Set and Prefrontal Cortex. *Annual Review of Neuroscience*, *31*(1), 219–  
869 245. <https://doi.org/10.1146/annurev.neuro.31.060407.125642>
- 870 Sakai, K., & Passingham, R. E. (2003). Prefrontal interactions reflect future task operations.  
871 *Nature Neuroscience*, *6*(1), 75–81. <https://doi.org/10.1038/nn987>
- 872 Sakai, K., & Passingham, R. E. (2006). Prefrontal set activity predicts rule-specific neural  
873 processing during subsequent cognitive performance. *The Journal of Neuroscience : The*  
874 *Official Journal of the Society for Neuroscience*, *26*(4), 1211–1218.

- 875 <https://doi.org/10.1523/JNEUROSCI.3887-05.2006>
- 876 Stelzer, J., Chen, Y., & Turner, R. (2013a). Statistical inference and multiple testing correction  
877 in classification-based multi-voxel pattern analysis (MVPA): Random permutations and  
878 cluster size control. *NeuroImage*, *65*, 69–82.  
879 <https://doi.org/10.1016/J.NEUROIMAGE.2012.09.063>
- 880 Stelzer, J., Chen, Y., & Turner, R. (2013b). Statistical inference and multiple testing correction  
881 in classification-based multi-voxel pattern analysis (MVPA): Random permutations and  
882 cluster size control. *NeuroImage*, *65*, 69–82.  
883 <https://doi.org/10.1016/j.neuroimage.2012.09.063>
- 884 Stokes, M., Thompson, R., Nobre, A. C., & Duncan, J. (2009). Shape-specific preparatory  
885 activity mediates attention to targets in human visual cortex. *Proceedings of the National*  
886 *Academy of Sciences*, *106*(46), 19569–19574. <https://doi.org/10.1073/pnas.0905306106>
- 887 Tottenham, N., Tanaka, J. W., Leon, A. C., McCarry, T., Nurse, M., Hare, T. A., ... Nelson, C.  
888 (2009). The NimStim set of facial expressions: Judgments from untrained research  
889 participants. *Psychiatry Research*, *168*(3), 242–249.  
890 <https://doi.org/10.1016/j.psychres.2008.05.006>
- 891 Townsend, J., & Ashby, G. (1978). Methods of modeling capacity in simple processing  
892 systems. In J. Castellan & F. Restle (Eds.), *Cognitive theory* (Vol. 3, pp. 200–239).  
893 Hillsdale, N.J: Erlbaum. Retrieved from  
894 <https://labs.psych.ucsb.edu/ashby/gregory/publications/281>
- 895 Treisman, A. M., & Gelade, G. (1980). A feature-integration theory of attention. *Cognitive*  
896 *Psychology*, *12*(1), 97–136. [https://doi.org/10.1016/0010-0285\(80\)90005-5](https://doi.org/10.1016/0010-0285(80)90005-5)
- 897 Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., & Diedrichsen, J. (2016). Reliability  
898 of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*, *137*, 188–200.  
899 <https://doi.org/10.1016/j.neuroimage.2015.12.012>
- 900 Wisniewski, D., Forstmann, B., & Brass, M. (2018). How exerting control over outcomes

901 affects the neural coding of tasks and outcomes. *BioRxiv*. <https://doi.org/10.1101/375642>

902 Woolgar, A., Hampshire, A., Thompson, R., & Duncan, J. (2011). Adaptive Coding of Task-  
903 Relevant Information in Human Frontoparietal Cortex. *Journal of Neuroscience*, *31*(41),  
904 14592–14599. <https://doi.org/10.1523/JNEUROSCI.2616-11.2011>

## 905 **Figure legend**

906 **Fig. 1:** Sequence of events in a single trial.

907 **Fig. 2:** Main analysis procedure. (a) Theoretical Representational Dissimilarity Matrices  
908 (RDMs) employed in the Representational Similarity Analysis (RSA). Within/Across-D. stands  
909 for within-dimension and across-dimension integration, while Single/Sequential R. stands for  
910 single response and sequential response. (b) RDMs capturing differences in instruction length  
911 (number of letters) and reaction time, included in a multiple regression analysis together with  
912 matrices shown in (a) to control for the effect of these two variables. (c) Following a searchlight  
913 approach, we extracted the neural RDM at each brain location and compared it – via Spearman  
914 correlation – with our three theoretical RDMs. As a result, we obtained three whole-brain  
915 correlation maps, one per model. (d) To assess the effect of motivation, for each region  
916 significant in (c) we extracted the neural RDMs from rewarded (R+) and non-rewarded (NR)  
917 trials. To study potential interactions of reward expectation and the corresponding model  
918 variable (Hypothesis 1), we averaged the dissimilarity values among same-condition and  
919 different-condition trials and tested if the subtraction among these two values was higher in the  
920 rewarded condition (using Wilcoxon signed-rank test). We also checked for a general increase in  
921 dissimilarities associated to reward (Hypothesis 2). *Note:* All matrices in the figure were  
922 simplified for visualization purposes by averaging cells within conditions. The matrices shown  
923 in (b) were further averaged across the sample. In (d), matrices display only one task variable  
924 (collapsing between the remaining two) to highlight the analysis logic. In all the analyses,  
925 however, trial-wise and single subject matrices were employed.

926 **Fig. 3:** Behavioral data. Violin plots showing correct responses (a) and Reaction Time (b) data  
927 for each condition, in rewarded and non-rewarded trials.

928 **Fig. 4:** Regions showing greater activity during the encoding of rewarded compared to non-  
929 rewarded instructions. Abbreviations stand for Nucleus Accumbens (N. Acc), inferior frontal  
930 junction (IFJ), premotor cortex (PMC), supplementary motor cortex (SMA), pre-supplementary  
931 motor cortex (preSMA) and intraparietal sulcus (IPS).

932 **Fig. 5:** Model-based RSA searchlight results for the three models (a-c) and render image  
933 showing the overlap among them (d). *Note:* Identical sections were employed to display the  
934 results across models.

935 **Fig. 6:** Conjunction analysis results.

936 **Tables**937 **Table 1.** Effect of the three models on the LOSO-estimated ROIs.

| Original model          | ROI                               | Model tested | T value | Z value | P value |
|-------------------------|-----------------------------------|--------------|---------|---------|---------|
| Dimension integration   | Left IFG                          | Dim.         | 3.354   |         | .008    |
|                         |                                   | Resp.        | 3.292   |         | .009    |
|                         |                                   | Cat.         | 3.635   |         | .004    |
| Response set complexity | Left IPS                          | Dim.         | 0.614   |         | 1       |
|                         |                                   | Resp.        | 5.351   |         | < .001  |
|                         |                                   | Cat.         |         | 1.975   | .163    |
|                         | Motor cortices, left LPFC         | Dim.         | 2.478   |         | .067    |
|                         |                                   | Resp.        | 3.647   |         | .004    |
|                         |                                   | Cat.         | 1.166   |         | .886    |
| Target category         | Left fusiform gyrus and precuneus | Dim.         | 0.476   |         | 1       |
|                         |                                   | Resp.        | 3.463   |         | .006    |
|                         |                                   | Cat.         | 5.466   |         | < .001  |
|                         | Left IFG                          | Dim.         | 2.832   |         | .029    |
|                         |                                   | Resp.        |         | 0.699   | .242    |
|                         |                                   | Cat.         | 4.930   |         | < .001  |
|                         | Right MTG                         | Dim.         |         | -0.144  | .557    |
|                         |                                   | Resp.        |         | -1.008  | .843    |
|                         |                                   | Cat.         |         | 2.859   | .002    |
| Right IFG               | Dim.                              |              | 1.275   | .101    |         |
|                         | Resp.                             |              | -0.206  | .582    |         |
|                         | Cat.                              |              | 3.085   | .001    |         |

938 *Note:* P values displayed are Bonferroni-corrected for multiple comparisons. Abbreviations  
939 stand for inferior frontal gyrus (IFG), intraparietal sulcus (IPS), and middle temporal gyrus  
940 (MTG), Dimension integration model (Dim.), Response complexity model (Resp.) and Target  
941 Category (Cat.).

942

943 **Table 2.** Effect of the three models on the MD network ROIs.

| ROI         | Model | T val  | Z val  | P value |
|-------------|-------|--------|--------|---------|
| ACC/preSMA  | Dim.  |        | 0.645  | 1       |
|             | Resp. |        | 1.673  | .115    |
|             | Cat.  | -0.026 |        | 1       |
| Left RLPFC  | Dim.  |        | 1.019  | .571    |
|             | Resp. |        | 0.346  | .365    |
|             | Cat.  |        | 2.665  | .023    |
| Left IFS    | Dim.  | 3.644  |        | .005    |
|             | Resp. | 4.423  |        | < .001  |
|             | Cat.  |        | 2.328  | .058    |
| Left MFG    | Dim.  |        | 2.739  | .014    |
|             | Resp. |        | 0.870  | .754    |
|             | Cat.  | 4.298  |        | .002    |
| Left alfo   | Dim.  | 0.667  |        | 1       |
|             | Resp. |        | 1.206  | .228    |
|             | Cat.  |        | 2.197  | .060    |
| Left IPS    | Dim.  | 1.617  |        | .638    |
|             | Resp. |        | 2.814  | .025    |
|             | Cat.  | 2.639  |        | .071    |
| Right RLPFC | Dim.  |        | 0.365  | 1       |
|             | Resp. | 1.460  |        | .849    |
|             | Cat.  | 0.861  |        | 1       |
| Right IFS   | Dim.  | 2.220  |        | .186    |
|             | Resp. |        | 1.599  | .211    |
|             | Cat.  |        | -0.626 | 1       |
| Right MFG   | Dim.  | 2.311  |        | .152    |
|             | Resp. | 1.294  |        | 1       |
|             | Cat.  | 2.042  |        | .273    |
| Right alfo  | Dim.  | 0.023  |        | 1       |
|             | Resp. |        | 1.299  | .280    |
|             | Cat.  | 1.352  |        | 1       |
| Right IPS   | Dim.  |        | 1.262  | .548    |
|             | Resp. |        | 1.842  | .330    |
|             | Cat.  |        | -0.701 | 1       |

944 *Note:* P values displayed are Bonferroni-corrected for multiple comparisons. Abbreviations  
 945 stand for anterior cingulate cortex (ACC), presupplementary motor area (preSMA), rostralateral  
 946 prefrontal cortex (RLPFC), inferior frontal sulcus (IFS), middle frontal gyrus (MFG), anterior  
 947 insula/frontal operculum area (alfo), intraparietal sulcus (IPS), Dimension integration model  
 948 (Dim.), Response complexity model (Resp.) and Target Category (Cat.).

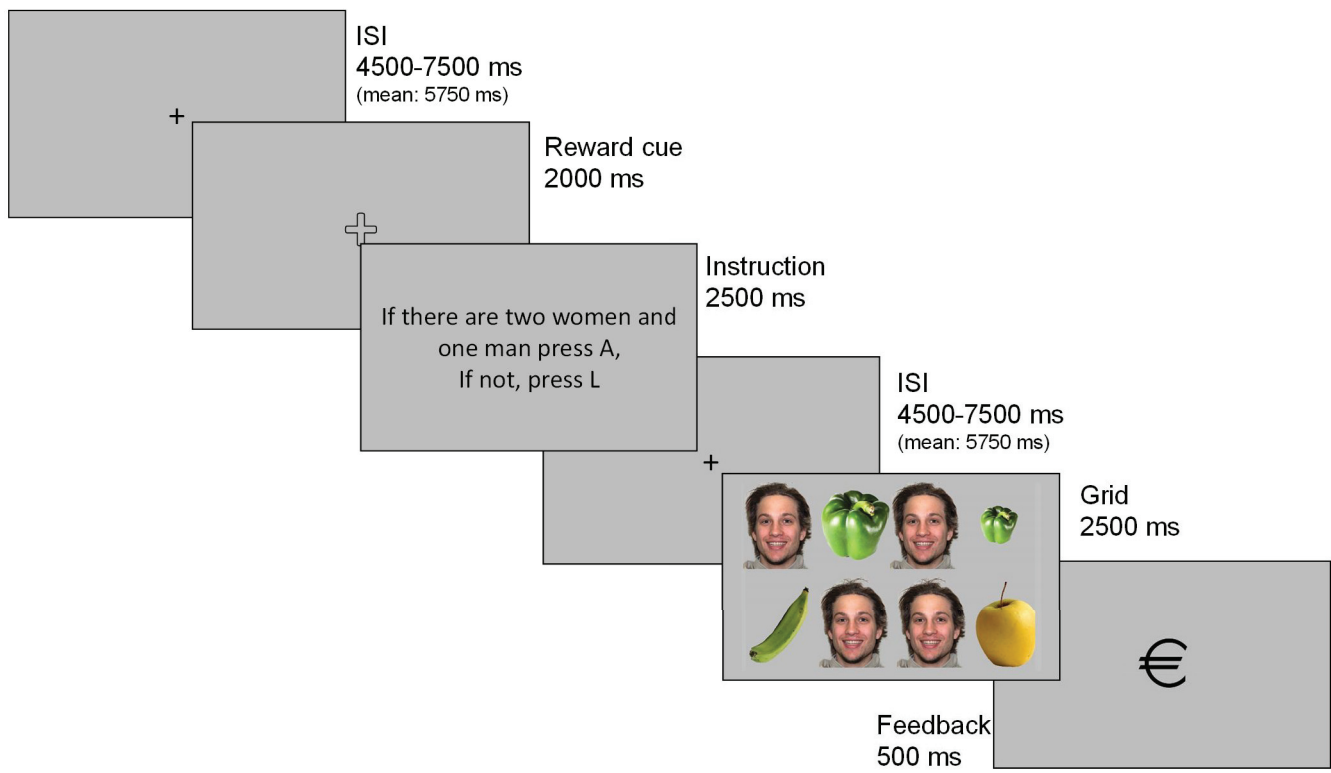
949



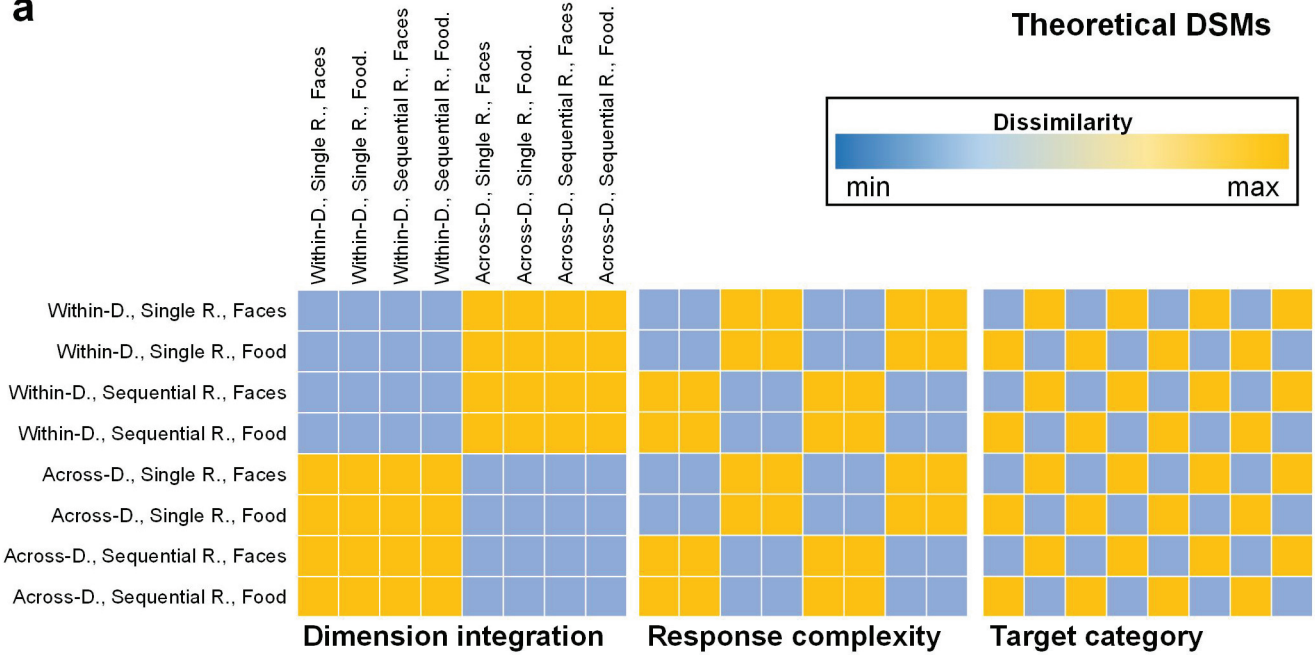
950 **Table 3.** Effect of reward on dissimilarity values and correlation with behavioral improvement.

| 951 | <b>ROI</b>                                     | <b>Effect of reward on<br/>dissimilarity values</b> | <b>Correlation<br/>RSA - behavior</b> |
|-----|--|---|---------------------------------------|
| 952 | <i>Task set complexity</i>                     |   |                                       |
|     | Left IFG/IFJ                                   | Z = -3.005*   | r = 0.515*                            |
| 953 | <i>Response set complexity</i>                 |   |                                       |
|     | M1 / PM / SMA /<br>IPS                         | Z = -3.712*   | r = 0.565*                            |
| 954 | Left PHC                                       | Z = -3.712*   | r = 0.558*                            |
| 955 | <i>Target category</i>                         |   |                                       |
|     | Left fusiform<br>gyrus/ precuneus /<br>IFG/IFJ | Z = -3.712*   | r = 0.543*                            |
| 956 | Right MTG/TPJ                                  | Z = -4.419*   | r = 0.495*                            |
| 957 | Right IFG                                      | Z = -3.712*   | r = 0.533*                            |
| 958 | Medial SFG                                     | Z = -2.652*   | r = 0.482*                            |

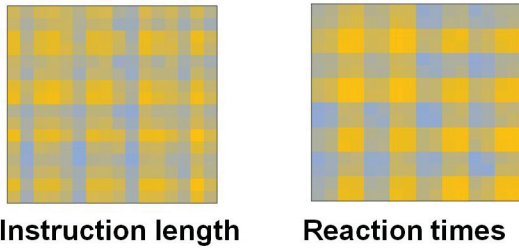
959 *Note:* The asterisks indicate significance at  $P < .05$  on the Wilcoxon paired-sample signed-rank  
 960 test (middle column) or Pearson correlation coefficient (left column). In the last case, multiple  
 961 comparisons were controlled with an FWE criterion. Abbreviations stand for inferior frontal  
 962 gyrus (IFG), inferior frontal junction (IFJ), primary motor cortex (M1), premotor cortex (PM)  
 963 supplementary motor area (SMA), parahippocampal cortex (PHC), middle temporal gyrus  
 964 (MTG), temporoparietal junction (TPJ) and superior frontal gyrus (SFG).



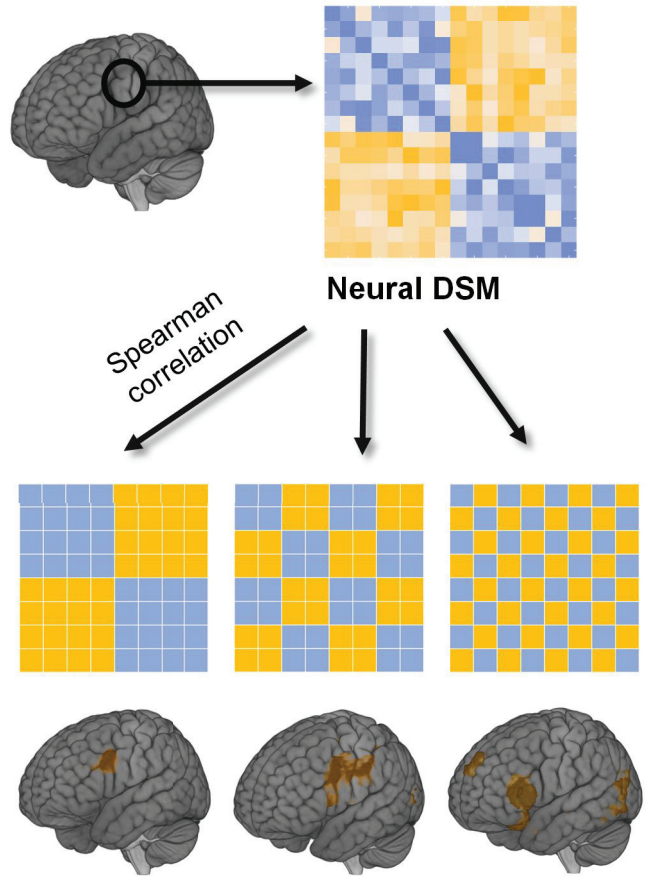
**a**



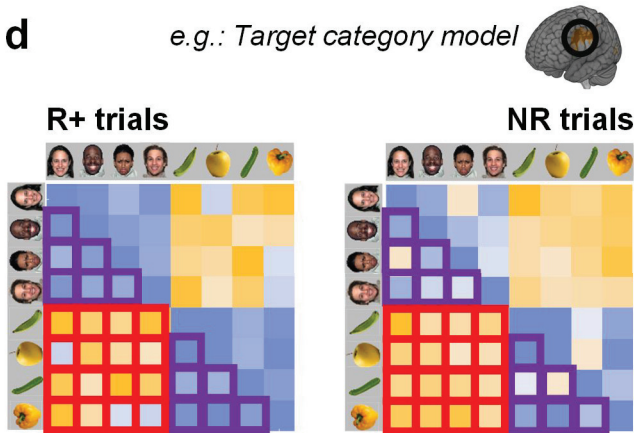
**b**



**c**

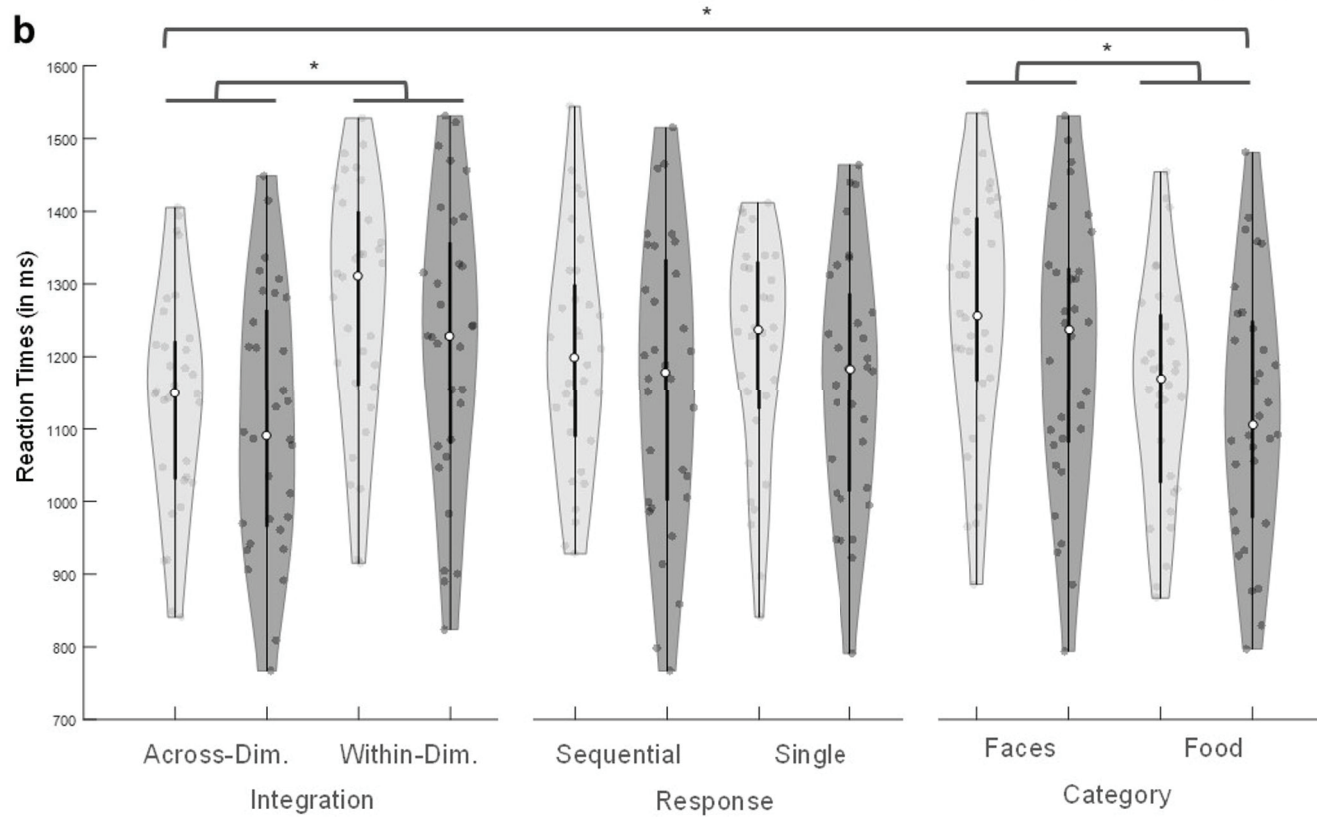
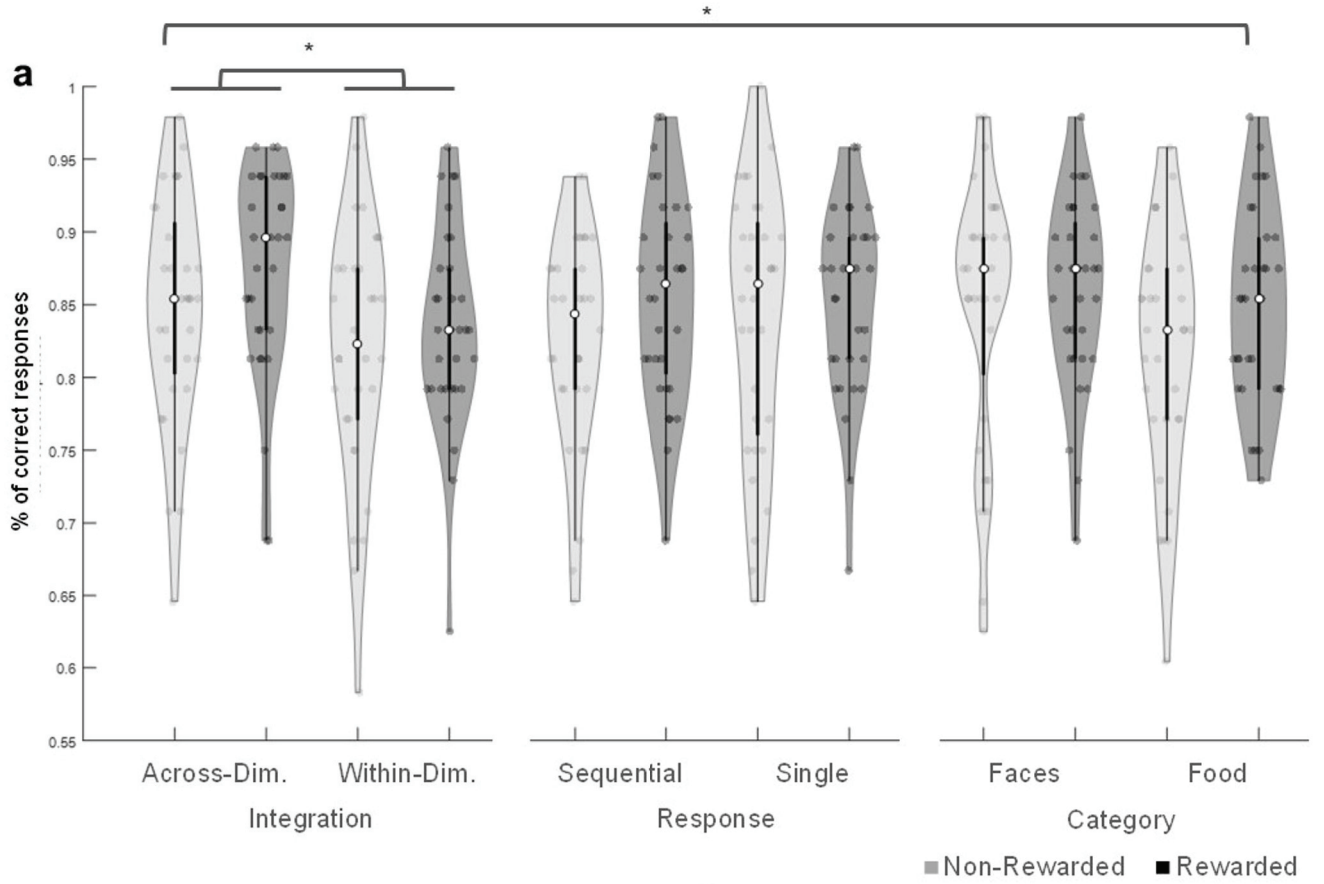


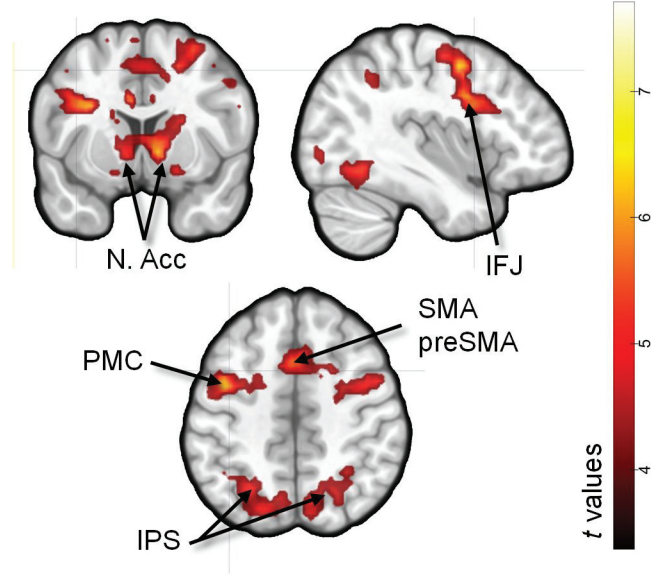
**d**

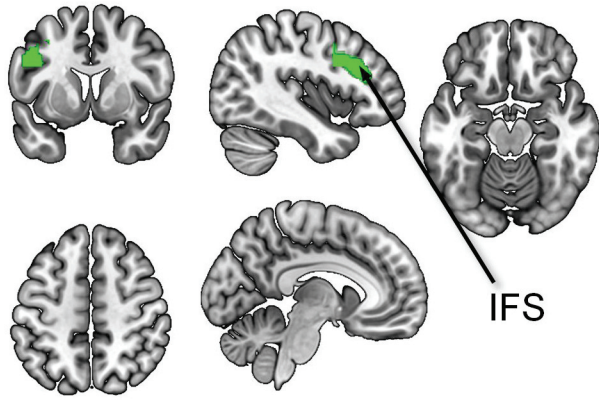


Hypothesis 1:  
 $(\text{Same Cond}_{R+} - \text{Different Cond}_{R+}) > (\text{Same Cond}_{NR} - \text{Different Cond}_{NR})$

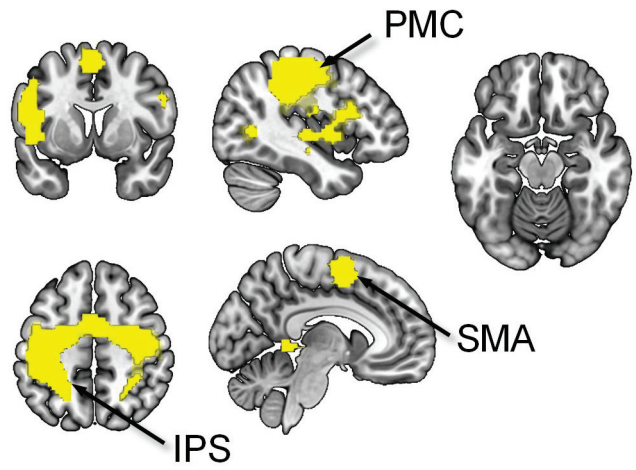
Hypothesis 2:  
 $(\text{Same Cond}_{R+} + \text{Different Cond}_{R+}) > (\text{Same Cond}_{NR} + \text{Different Cond}_{NR})$



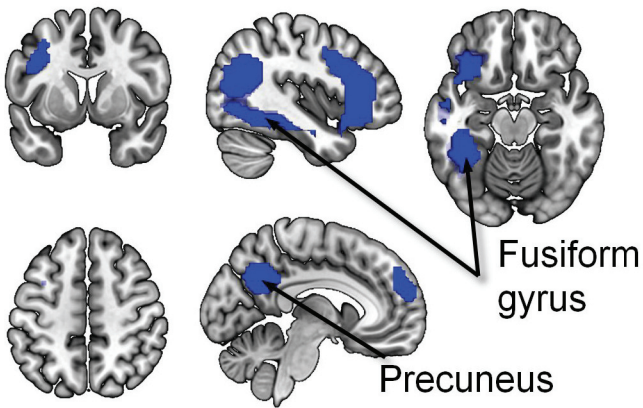




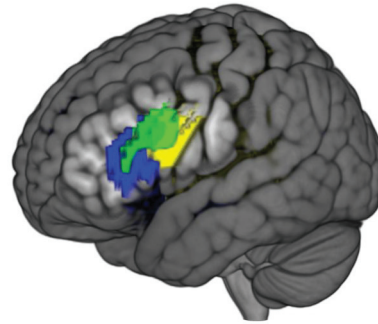
(a) Task complexity model



(b) Response complexity model



(c) Category model



(d) Model overlap

