

1 **Neural representations of social valence bias economic interpersonal choices**

2

3 **Authors:** Paloma Díaz-Gutiérrez¹, Juan E. Arco¹, Sonia Alguacil², Carlos González-
4 García³ & María Ruz¹

5 ¹Mind, Brain and Behavior Research Center, University of Granada, Spain

6 ²University Isabel I, Spain

7 ³Ghent University, Belgium

8

9

10 **Word count:** 10789

11

12

13

14

15 **Contact information:** Department of Experimental Psychology, University of Granada,
16 18071, Granada, Spain.

17 Telephone: +34 958 24 06 60. Facsimile number: +34 958 24 62 39

18 E-mail address: mruz@ugr.es (M. Ruz).

19

20 **Abstract**

21 Prior personal information is highly relevant during social interactions. Such knowledge
22 aids in the prediction of others, and it affects choices even when it is unrelated to actual
23 behaviour. In this investigation, we aimed to study the neural representation of positive
24 and negative personal expectations, how these impact subsequent choices, and the effect
25 of mismatches between expectations and encountered behaviour. We employed
26 functional Magnetic Resonance Imaging in combination with a version of the
27 Ultimatum Game (UG) where participants were provided with information about their
28 partners' moral traits previous to their fair or unfair offers. Univariate and multivariate
29 analyses revealed the implication of the supplementary motor area (SMA) and inferior
30 frontal gyrus (IFG) in the representation of expectations about the partners in the game.
31 Further, these regions also represented the valence of expectations, together with the
32 ventromedial prefrontal cortex (vmPFC). Importantly, the performance of multivariate
33 classifiers in these clusters correlated with a behavioural choice bias to accept more
34 offers following positive descriptions, highlighting the impact of the valence on the
35 expectations on participants' economic decisions. Altogether, our results suggest that
36 expectations based on social information guide future interpersonal decisions and that
37 the neural representation of such expectations in the vmPFC is related to their influence
38 on behaviour.

39

40 **1. Introduction**

41 Decision-making is a crucial constituent of our daily life. To make choices that best fit
42 our goals, we must rapidly weight different sources of information in an efficient
43 manner. An elegant approach to understand how we perform such weighting comes
44 from the framework of predictive coding (Friston, 2005), where optimal decision-
45 making combines sensory input (*evidence*) with predictions (*priors*; Schwarz et al.,
46 2016; Summerfield and De Lange, 2014). The role of these predictions has been
47 thoroughly examined in non-social decisions, where several studies have shown pre-
48 activation of target-related brain areas during the expectation period, prior to target
49 onset (e.g., Esterman and Yantis, 2010; González-García et al., 2016; Puri et al., 2009).
50 However, a large part of decisions involve social contexts, where we constantly engage
51 in interactions with others. Still, the role of expectations in such scenarios remains
52 unclear.

53

54 When making decisions in complex scenarios, people tend to choose more often and
55 faster the options that match their personal preferences (with higher personal value)
56 even when the objective task value of the different alternatives is similar (Lopez-Persem
57 et al., 2016). This leads to suboptimal decisions that do not properly consider potential
58 future outcomes (Fleming, Thomas, & Dolan, 2010). This is also the case for
59 interpersonal decisions, which can be biased by several sources of information at
60 different stages of processing (Díaz-Gutiérrez, Alguacil, & Ruz, 2017). For instance, in
61 the Ultimatum Game (UG; Güth, Schmittberger, & Schwarze, 1982; Moser, Gaertig, &
62 Ruz, 2014), participants receive monetary offers from game partners and decide
63 whether to accept them or not. Acceptance leads to both parts earning their split;
64 whereas no gains are earned after a rejection. Here, “rational” decisions from an
65 economic point of view should be of acceptance, since you can only earn money.
66 However, choices are strongly influenced by the fairness of the offer (how balanced
67 both halves of the split are). People often show high rejection rates towards unfair offers
68 (Sanfey, Rilling, Aronson, Nystrom, & Cohen, 2003), which has been explained in
69 terms of inequity-aversion tendencies (Fehr & Camerer, 2007) and punishment (Brañas-
70 Garza, Espín, Exadaktylos, & Herrmann, 2014). Others have emphasized the
71 importance of social norms, and how these impact the perception of fairness (Chang &
72 Sanfey, 2013). In these scenarios, the mechanisms underlying the processing of offers
73 depending on their fairness and participants’ subsequent responses have been

74 extensively studied. Here the role of the anterior cingulate cortex (ACC) and
75 supplementary motor area (SMA) stands out, concerning both fairness and people's
76 decisions (for a meta-analysis, see Gabay et al., 2014). Authors such as Sanfey et al.,
77 (2003) have shown the involvement of the anterior insula (aI) in fairness processing.
78 Also, Corradi-Dell'Acqua, Civai, Rumiati, & Fink (2013) differentiated its role from the
79 one of the medial prefrontal cortex (mPFC), which appears to be linked to emotional
80 self-related responses during interpersonal bargaining situations.

81

82 Despite the extensive and diverse studies in interpersonal games, it is largely unknown
83 how the brain represents socially relevant priors in these scenarios. Recent proposals
84 have tried to link predictive coding and the representation of social traits in relation to
85 social expectations (e.g., Tamir and Thornton, 2018). Several studies have described a
86 set of regions underlying the representation of knowledge that guides social predictions
87 in a broad context (termed Social Cognition Network; Frith & Frith, 2008), including
88 personal traits, stereotyping, semantic knowledge about people or inferences about
89 others and their mental states (Tamir & Thornton, 2018; Tamir, Thornton, Contreras, &
90 Mitchell, 2016). This network includes the temporoparietal junction (TPJ), superior
91 temporal sulcus (STS), precuneus (PC), anterior temporal lobes (ATL), amygdala and
92 the mPFC (Contreras et al., 2013; Frith, 2007; Frith and Frith, 2001; Mitchell et al.,
93 2008). These regions underlie processes such as Theory of Mind (ToM; Saxe and
94 Kanwisher, 2003). Similarly, in decisions in social contexts, the mPFC has been related
95 to expectations about others' behaviour (Corradi-Dell'Acqua, Turri, Kaufmann,
96 Clément, & Schwartz, 2015). Importantly, prior expectations during social decisions
97 also influence behaviour when they are not followed by their usual consequences. In
98 this line, different studies (Fouragnan et al., 2013; Ruz and Tudela, 2011) have observed
99 increased activation in brain areas associated with cognitive control, such as the ACC
100 and the aI when expectations about partners do not match their subsequent behaviour.
101 Similarly, Chang and Sanfey (2013) found a relationship between the deviation of the
102 expectations and increased activation in the aI, ACC and SMA. Specifically, in the UG,
103 an increase of activation in the dorsolateral PFC (dlPFC) and aI has been related to
104 participants' reaction to unfair offers (Knoch, Pascual-Leone, Meyer, Treyer, & Fehr,
105 2006; Sanfey et al., 2003), which has also been interpreted as a violation of what we
106 expect from others.

107

108 In addition to this, social expectations can also be based on the personal traits of others,
109 which are an essential component of social representations (Tamir & Thornton, 2018).
110 The priors that they generate relate to stereotypes and interact with perceptual processes
111 (Stolier & Freeman, 2016, 2017). These personality traits can be decomposed in three
112 different dimensions: rationality, social impact and, crucially to our investigation,
113 valence (positive vs. negative; Tamir and Thornton, 2018; Thornton and Mitchell,
114 2017). The representation of the character of others in association with positive or
115 negative information is an important source of bias in interpersonal decisions (Díaz-
116 Gutiérrez et al., 2017). For instance, Delgado et al. (2005), found that participants
117 trusted partners associated with positive moral traits more than those having negative
118 ones. Furthermore, a variety of studies employing the UG paradigm have observed that
119 participants tend to accept more offers from partners associated with positive
120 descriptions, compared to negative ones (Gaertig, Moser, Alguacil, & Ruz, 2012). This
121 tendency is steeper when participants navigate uncertain scenarios (Ruz, Moser, &
122 Webster, 2011). Moreover, in this context, the use of high-density
123 electroencephalography (EEG) has shown that negative descriptions of partners lead to
124 a higher amplitude of the medial frontal negativity (MFN; associated with the
125 evaluation of outcomes, Hajcak et al., 2006; Yeung and Sanfey, 2004) when decisions
126 are made (Moser et al., 2014). These data indicate how, regardless of fairness, people
127 evaluate offers as more negative when they come from a disagreeable partner. Such
128 knowledge about personal traits has been suggested to be integrated by the mPFC (Van
129 Overwalle, 2009). For example, this area increases its coupling with other regions
130 responding to specific traits (Hassabis et al., 2014), and shows heightened activation
131 when a partner's behaviour violates previous trait implications (Ma et al., 2012).

132

133 Nonetheless, despite the key relevance of valence in psychological theories and its
134 marked impact on social decision-making, it is not well understood how valence is
135 represented at the neural level and its effect on subsequent choices (Barrett & Bliss-
136 Moreau, 2009). Results of a recent meta-analysis (Lindquist, Satpute, Wager, Weber, &
137 Barrett, 2015) provide evidence of a general recruitment of a set of regions for valenced
138 versus neutral information, including the bilateral aI, the ventral and dorsal portions of
139 the mPFC (vm/dmPFC), the dorsal ACC, SMA, and lateral PFC. Lindquist et al. (2015)
140 found that the vmPFC/ACC was more frequently activated in positive vs. negative than

141 in positive vs. neutral contrasts, which could indicate that these regions represent
142 valence information along a single bipolar dimension.

143

144 Taking all this into account, in the current functional Magnetic Resonance Imaging
145 (fMRI) study, we employed a modified version of the UG (Gaertig et al., 2012) to
146 investigate how socially relevant priors represented by the valence of personal
147 descriptions of partners bias interpersonal economic choices. First, we aimed to study
148 which neural regions code for the generation and maintenance of positive and negative
149 expectations about other people. In a second step, we assessed how these expectations
150 bias decisions. We expected to find specific neural representations underlying the
151 expectations about the partners, with different patterns depending on the valence of
152 these predictions (Lindquist et al., 2015). Specifically, we hypothesized that these
153 patterns would be represented in regions related to social cognition and priors in
154 decision-making (Contreras et al., 2012; González-García et al., 2016; Saxe &
155 Kanwisher, 2003). Last, we intended to ascertain which neural mechanisms were
156 engaged when there is a mismatch between personal expectations and the partners'
157 behaviour. We predicted that control-related areas would be engaged when the valenced
158 description was not congruent with the subsequent partner's behaviour.

159

160 **2. Methods**

161 **2.1. Participants**

162 Twenty-four volunteers were recruited from the University of Granada ($M = 21.08$, SD
163 $= 2.92$, 12 men), matching the sample size employed in Moser et al. (2014), who
164 implemented the same version of the task for electroencephalography (EEG). This
165 sample is similar to previous fMRI studies using the UG (Chang and Sanfey, 2013;
166 Grecucci, Giorgetta, Bonini & Sanfey, 2013). All participants were right-handed with
167 normal or corrected vision and received economic remuneration (20-25 Euros,
168 proportionally to their acceptance rates). Participants signed a consent form approved by
169 the Ethics Committee of the University of Granada.

170

171 **2.2. Apparatus and stimuli**

172 We employed 16 adjectives used in previous studies (Gaertig et al., 2012; Moser et al.,
173 2014; Ruz et al., 2011; see Table 1) as trait-valenced descriptions of the game

174 proposers, extracted from the Spanish translation of the Affective Norms for English
175 Words database (ANEW; Redondo et al., 2007). Half of the adjectives were positive (M
176 = 7.65 valence, SD = 0.43), and the other half were negative (M = 2.3 valence, SD =
177 0.67). All words were matched in arousal (M = 5.69, SD = 0.76), number of letters (M =
178 6.19, SD = 1.42) and frequency of use (M = 20.19, SD = 18.47). In addition, we
179 employed numbers from 1 to 9 (two in each trial) in black colour to represent different
180 monetary offers. Stimuli were controlled and presented by E-Prime software (Schneider,
181 Eschman, & Zuccolotto, 2002). Inside the scanner, the task was projected on a screen
182 visible to participants through a set of mirrors placed on the radiofrequency coil.

183

184 **2.3. Task and procedure**

185 To add credibility to the interpersonal game setting, participants were told that they
186 were about to receive offers made by real participants in a study of a previous
187 collaboration with a foreign university. Furthermore, to engage participants in the game
188 as a real social scenario, prior to the scanner they performed two tasks in which they
189 had to make economic offers that would be used for other participants in future studies.
190 In one of the tasks, participants acted as proposers, filling a questionnaire where they
191 had to make offers for 16 different unknown partners, who would be involved in future
192 experimental games. Here, they had to split 10 Euros into two parts, one for themselves
193 and the other for their partners. Additionally, in a second task, they played a short
194 version of the Dictator Game (Kahneman, Knetsch, & Thaler, 1986), where they
195 decided how to divide another 10 Euros between themselves and an anonymous partner,
196 who would have a merely passive role concerning the output of the offer. Moreover,
197 participants were told that the offers that they were about to see in the scanner were
198 each provided by a different partner who previously performed the same tasks as they
199 did before the scanner, and therefore, the offers were real examples of other
200 participants' responses when acting as proposers. Participants were informed that each
201 offer would be preceded by a word that had been obtained as an output from a series of
202 personality and social questionnaires filled by their partners and, therefore, that these
203 adjectives described them in some way (see Table 1). Choices made by participants had
204 an influence in their final payment, as it actually varied (20-25 Euros) according to their
205 choices during the game in the scanner. In a post-scanning informal debriefing session,
206 none of the participants reported suspicions regarding the background story of this

207 procedure, which has also been used successfully in other settings (e.g. Correa,
208 Alguacil, Ciria, Jiménez, & Ruz, 2020; Correa et al., 2017).

209

210 In the scanner, participants played the role of the responder in a modified UG (e.g.,
211 Gaertig et al., 2012), deciding whether to accept or reject monetary offers made by
212 different partners (proposers). If they accepted the offer, both parts earned their
213 respective splits, whereas if they rejected it, neither of them earned money from that
214 exchange. Offers consisted of splits of 10 Euros, which could be fair (5/5, 4/6) or unfair
215 (3/7, 2/8, 1/9). The number presented at the left on the screen was always the amount of
216 money given to the participant, and the one on the right side was the one proposed by
217 the partners for themselves.

218

219

Table 1. List of adjectives employed in the task (Gaertig et al., 2012).

Positive words	Negative words
Friend	Criminal
Generous	Cruel
Honest	Disloyal
Honourable	False
Humble	Guilty
Kind	Hostile
Loyal	Selfish
Warm	Traitor

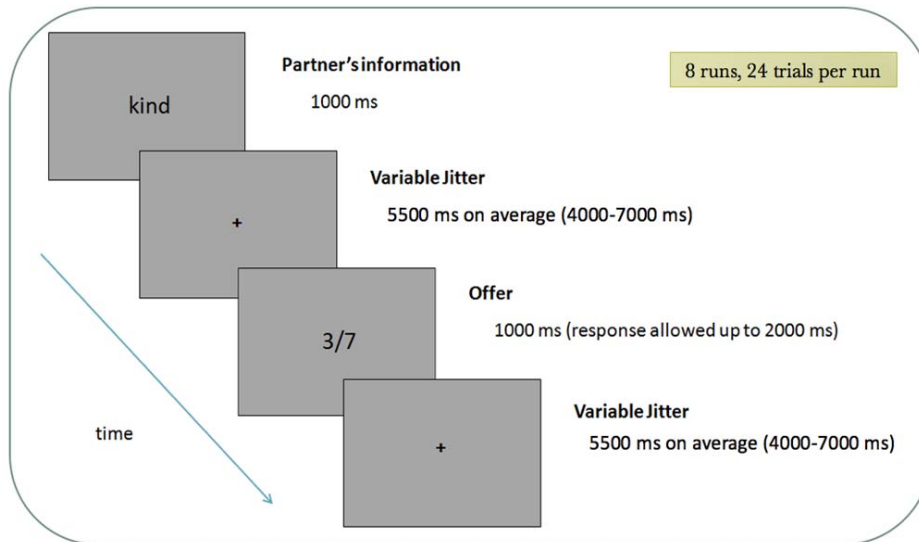
220

221

222 Personal information about the partners was included as adjectives with different
223 valence. A third of these descriptions was positive, another third negative, and the last
224 third was neutral, represented by text indicating the absence of information about that
225 partner ("no test"). The valence of the adjectives was orthogonal to the fair-unfair nature
226 of the offer. The order of the offers and adjectives was randomized, and each type of
227 personal information (positive, negative, no information) preceded each offer equally
228 within and across runs. Decision-response associations were counterbalanced between
229 participants.

230

231 Participants performed a total of 192 trials, arranged in 8 runs (24 trials per run). In each
232 run, a start cue of 6 s was followed by 24 trials. Each trial (see Figure 1) started with an
233 adjective for 1 s (mean = 2.98°), preceding a jittered interval lasting 5.5 s on average (4-
234 7 s, +/0.76°). Then, the offer appeared for 0.5 s (1.87°), followed by a second jittered
235 interval (mean = 5.5 s; 4-7 s, +/0.76°). Overall, each run lasted 5.1 minutes and the
236 whole task 41 minutes approximately.
237



238
239 **Figure 1.** Sequence of events in a trial. The task varied the Valence of the partner's information (Positive,
240 Negative, No information) and the Fairness of the offer (Fair/Unfair), which were manipulated
241 orthogonally in the design.

242

243 **2.4. Image acquisition and preprocessing**

244 MRI images were acquired using a Siemens Magnetom TrioTim 3T scanner, located at
245 the Mind, Brain and Behavior Research Center in Granada. Functional images were
246 obtained with a T2*-weighted echo-planar imaging (EPI) sequence, with a TR of 2000
247 ms. Thirty-two descendent slices with a thickness of 3.5 mm (20% gap) were extracted
248 (TE = 30 ms, flip angle = 80°, voxel size of 3.5 mm³). The sequence was divided into 8
249 runs, consisting of 166 volumes each. After the functional sessions, a structural image
250 of each participant with a high-resolution T1-weighted sequence (TR = 1900 ms; TE =
251 2.52 ms; flip angle = 9°, voxel size of 1 mm³) was acquired.

252

253 Data were preprocessed with SPM12 software
254 (<http://www.fil.ion.ucl.ac.uk/spm/software/spm12/>). The first three volumes of each run

255 were discarded to allow the signal to stabilize. Images were realigned and unwarped to
256 correct for head motion, followed by slice-timing correction. Afterwards, T1 images
257 were coregistered with the realigned functional images. Then, functional images were
258 spatially normalized according to the standard Montreal Neurological Institute (MNI)
259 template and smoothed employing an 8 mm Gaussian kernel. Low-frequency artefacts
260 were removed using a 128 high-pass filter. Data for multivariate analyses was only
261 head-motion and slice-time corrected and coregistered.

262

263 **2.5. Univariate analyses**

264 First-level analyses were conducted for each participant, following a General Linear
265 Model in SPM12. We employed an event-related design, where activity was modelled
266 using regressors for each valence type of adjective and for the offers. The estimated
267 model included three regressors for the Words (positive, negative, no information) and
268 six for the Offers (Fair offers_Positive, Fair offers_Negative, Fair offers_Neutral,
269 Unfair offers_Positive, Unfair offers_Negative, Unfair offers_Neutral). Note that since
270 decisions were made when the offers appeared, and that responses (choices) showed a
271 strong dependency on offer fairness, offer fairness and decisions cannot be modelled
272 separately. Given our research questions, we modelled the offer events considering their
273 fairness regardless of participants' choices. Regressors were convolved with a standard
274 hemodynamic response, with adjectives modelled with their duration (1 s + jitter), and
275 offers modelled as events with zero duration. This temporal difference is accounted by
276 the fact that the words describing the partners trigger preparatory processes, which
277 extend in time (e.g. Bode & Haynes, 2009; Di Russo et al., 2017; González-García,
278 Arco, Palenciano, Ramírez, & Ruz, 2017; González-García et al., 2016; Sakai, 2008),
279 whereas the processing of the offers ends shortly after with the response of each trial
280 (see Moser et al., 2014). In addition, the orthogonal manipulation of these variables in
281 the design avoided covariance confounds between word cues and target offers.

282

283 At the second level of analysis, *t*-tests were conducted for comparisons related to the
284 presence of expectations (information about the partner > no information), the valence
285 of the information (positive > negative, negative > positive) and the fairness of the offer
286 (fair > unfair, unfair > fair). We also carried out contrasts for congruence effects
287 between the events, where we had congruent (positive descriptions followed by fair
288 offers, negative descriptions followed by unfair offers) and incongruent trials (positive

289 descriptions followed by unfair offers, negative descriptions followed by fair offers). To
290 control for false positives at the group level, we employed permutations tests with
291 statistical non-parametric mapping (SnPM13, <http://warwick.ac.uk/snpm>) and 5000
292 permutations. We performed cluster-wise inference on the resulting voxels with a
293 cluster-forming threshold of 0.001, which was later used to obtain significant clusters
294 (FWE corrected at $p < 0.05$).

295

296 **2.6 Multivariate analyses**

297 We performed MVPA to examine the brain areas representing the valence of the
298 expectations, that is, the regions containing information about whether the partners were
299 described with positive vs. negative adjectives. To this end, we performed a whole-brain
300 searchlight (Kriegeskorte et al., 2006) on the realigned images (prior to normalization).
301 We employed The Decoding Toolbox (TDT; Hebart et al., 2015), to create 12-mm
302 radius spheres, where linear support vector machine classifiers ($C=1$; Pereira et al.,
303 2009) were trained and tested using a leave-one-out cross-validation scheme, employing
304 the data from the 8 scanning runs (training was performed with data from 7 runs and
305 tested in the remaining run, in an iterative fashion). We used a Least-Squares Separate
306 model (LSS; Turner, 2010) to reduce collinearity between regressors (Abdulrahman &
307 Henson, 2016; Arco et al., 2018). This approach fits the standard hemodynamic
308 response to two regressors: one for the current event of a trial (positive/negative
309 adjective) and a second one for all the remaining events and trials. As in the previous
310 analyses, adjective regressors were modelled with their duration (1 s + jitter) and offers
311 with zero duration. Consequently, the output of this model was one beta image per
312 event (total = 128 images, 64 for each type of adjective, 112 for training and 16 for
313 testing in each iteration). Afterwards, at the group level, non-parametrical statistical
314 analyses were performed on the resulting accuracy maps following the method proposed
315 by Stelzer et al. (2013) for MVPA data. We permuted the labels and trained the
316 classifier 100 times for each participant. The resulting maps were then normalized to an
317 MNI space. Afterwards, we randomly picked one of these maps per each participant and
318 averaged them, obtaining a map of group accuracies. This procedure was repeated
319 50000 times, building an empirical chance distribution for each voxel position and
320 selecting the 50th greatest value, which corresponds to the threshold that marks the
321 statistical significance. Only the voxels that surpassed this were considered significant.

322 The resulting map was FWE corrected at 0.05, computing previously the cluster size
323 that matched this value from the clusters obtained in the empirical distribution.

324

325 Importantly, the valence of the description influenced acceptance rates, which could
326 generate potential confounds in the previous decoding. The association between hand
327 and decision (left/right, acceptance/rejection) was fully counterbalanced across
328 participants, but remained constant for each of them. Therefore, the classifier could use
329 response information (accept vs. reject) when decoding valence. To clarify this issue,
330 we performed a response classification at the offer period (following the same
331 procedure as for the valence decoding). Then, we ran a conjunction analysis, computing
332 the intersection between valence and response group maps to examine whether the
333 regions containing relevant information about the valence were the same as those
334 representing the participants' decisions (accept vs. reject). Moreover, to test additionally
335 the potential overlap between the neural representations of participants' decisions and
336 the valence of the expectations about the partners, we performed a cross-classification
337 analysis (Kaplan, Man, & Greening, 2015) between these two domains. Following again
338 the same classification procedure described above in this section, we trained the
339 classifier with the participants' responses to the offers (accept vs. reject) and tested it on
340 the valence of the partner's descriptions (positive vs. negative).

341

342 **2.7. Relationship between decoding accuracy and choices**

343 To examine the extent to which the fidelity of representation of (positive vs. negative)
344 personal priors relates to the decisions made by participants, we performed a correlation
345 analysis between an individual bias index and mean decoding accuracy values from
346 each significant cluster in the MVPA described above. To obtain this behavioural index,
347 for each participant we subtracted the average acceptance rate following negative
348 descriptions from the average acceptance rate after positive descriptions (regardless of
349 the nature of the offer). For each subject, we performed a one-tailed (right) Spearman's
350 correlation between the behavioural index and the decoding accuracy from each
351 significant cluster (Bonferroni-corrected for multiple comparisons). To further ascertain
352 that participants' motor responses were not contaminating this link between valence
353 representation and interpersonal choices, we ran an additional correlation analysis
354 following the same approach, this time to examine the link between valence' decoding
355 results and the response made by participants (acceptance or rejection of the offer).

356 Therefore, for each participant, we calculated their average acceptance rate in general,
357 regardless of the valence of the expectation and the fairness of the offers.

358

359 **3. Results**

360 **3.1. Behavioural data**

361 Acceptance rates (AR) and reaction times (RTs) were analysed in a Repeated Measures
362 ANOVA, with Offers (fair/unfair) and Valence of the descriptions (positive, negative,
363 neutral) as factors. The Greenhouse-Geisser correction was applied whenever the
364 sphericity assumption was violated.

365 **3.1.1. Acceptance rates**

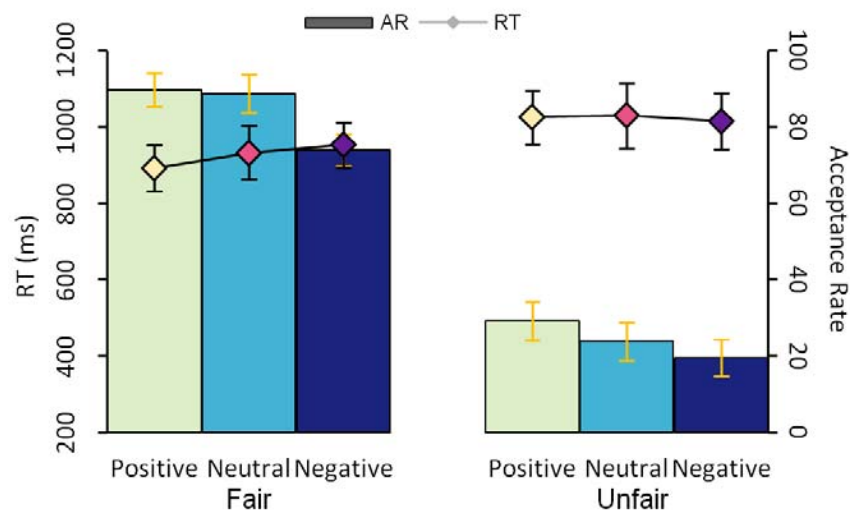
366 Participants responded on 100% of the trials. Data showed (see Figure 2) a main effect
367 of Offer $F_{1,23} = 74.50, p < .001, \eta_p^2 = .764$, where fair offers were accepted more often
368 ($M = 84.09\%$; $SD = 22.10$) than unfair ones ($M = 24.18\%$; $SD = 24.10$). Valence was
369 also significant, $F_{2,22} = 13.735, p = .001, \eta_p^2 = .374$. Participants accepted more offers
370 when they were preceded by a positive description of the partner ($M = 59.39\%$; $SD =$
371 23.09), than when there was no information ($M = 56.31\%$; $SD = 21.89$) or when this
372 was negative ($M = 46.70\%$; $SD = 24.33$). Planned comparisons revealed that these
373 differences were significant between all pairs (all $p_s < .05$). Finally, the Offer X Valence
374 interaction was also significant, $F_{2,22} = 4.262, p = .033, \eta_p^2 = .156$. Planned
375 comparisons showed that for fair offers, there were differences between all comparisons
376 ($p_s = .002$) except between positive and neutral information ($p = .399$), whereas for
377 unfair offers, there was no difference in acceptance rates between negative and neutral
378 information ($p = .074$) but there was for the rest of the pairwise comparisons: $p_s < .01$

379

380 **3.1.2. Reaction times**

381 Results showed (see Figure 2) a main effect of Offer $F_{1,23} = 22.489, p < .001, \eta_p^2 = .494$,
382 where participants took longer to respond to unfair ($M = 1023.53$ ms; $SD = 373.10$ ms)
383 than to fair offers ($M = 925.62$ ms; $SD = 309.57$ ms). Neither Valence, $F_{2,22} = 1.05, p =$
384 $.341$, or its interaction with Fairness, $F_{2,22} = 1.956, p = .168$ were significant. In
385 addition, to measure the influence of expectations on participant's responses (see Ruz et
386 al., 2011), we ran an ANOVA where we included the valence of the descriptions and the
387 decision (accept, reject) made to the offers. Here, we did not find any effect of Valence,
388 $F < 1$, but we found significant effects of Decision, $F_{1,23} = 5.519, p = .028, \eta_p^2 = .194$,

389 since participants were faster to accept ($M = 951.37$ ms; $SD = 356.01$ ms) than to reject
390 the offers ($M = 988.97$ ms; $SD = 316.91$ ms). Interestingly, data showed an interaction
391 Valence X Decision, $F_{2,22} = 4.23$, $p = .025$, $\eta_p^2 = .155$, replicating previous findings
392 (Gaertig et al., 2012; Ruz et al., 2011). Planned comparisons indicated that these
393 differences in RT for responses took place only after positive, $F_{1,23} = 13.997$, $p = .001$,
394 $\eta_p^2 = .378$ (Accept: $M = 927.60$ ms, $SD = 297.37$ ms; Reject: $M = 993.91$ ms, $SD =$
395 335.52 ms), and neutral descriptions, $F_{1,23} = 4.504$, $p = .045$, $\eta_p^2 = .165$ (Accept: $M =$
396 955.8 ms, $SD = 304.96$ ms; Reject: $M = 987.80$ ms, $SD = 328.48$ ms), but not for
397 negative descriptions, $F < 1$.



398
399 **Figure 2.** Acceptance Rates (AR, bars) and reaction times (RT, lines) to fair and unfair offers preceded
400 by positive, negative and neutral descriptions of the partner (error bars represent S.E.M).

401

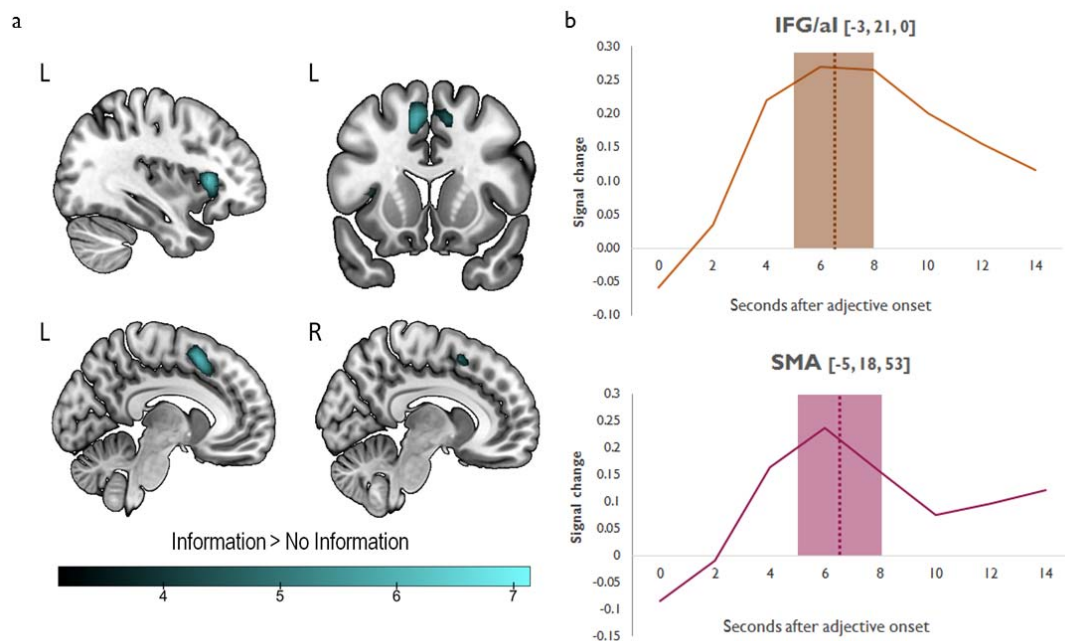
402 3.2. Neuroimaging data

403 3.2.1. Univariate results

404 *Expectations*

405 During the presentation of the description and the time interval that followed, that is,
406 when participants had personal information to **generate expectations** [(Positive
407 adjective & Negative adjective) > No Information], we observed a cluster of activity
408 (see Figure 3a) in the left dorsal aI ($k = 109$; $-33, 21, 4$) and bilateral Supplementary
409 Motor Cortex (SMA; $k = 138$; $-8, 11, 53$; see Fig. 3). Additionally, the right inferior
410 parietal lobe (right IPL) showed higher activity ($k = 264$; $55, -35, 53$) for **positive**

411 **descriptions** compared to negative ones. No cluster surpassed the statistical threshold
412 ($p > 0.05$) for the opposite contrast.



413 **Figure 3. a)** Univariate results during the expectation period. Scales reflect peaks of significant t-values
414 ($p < .05$, FWE-corrected for multiple comparisons). **b)** Time course of activation in the IFG/aI (-3, 21, 0;
415 *top*) and SMA (-5, 18, 53; *bottom*) clusters obtained from the conjunction analysis. From these regions,
416 we extracted the signal change values related to the processing of personal information minus the average
417 during the neutral condition, time-locked to the adjective onset. The shaded areas show the variable time
418 window during which the offer could appear (5-8 s after the adjective onset) whereas the dotted lines
419 show its average (6.5 after the adjective onset)
420

421 During *offer processing*, the previous presentation of **personal information** about the
422 partner [(Offer_Pos & Offer_Neg > Offer_Neu)] yielded again significant activity
423 involving the bilateral dorsal aI and right SMA ($k = 23349$; -33, 21, 4).

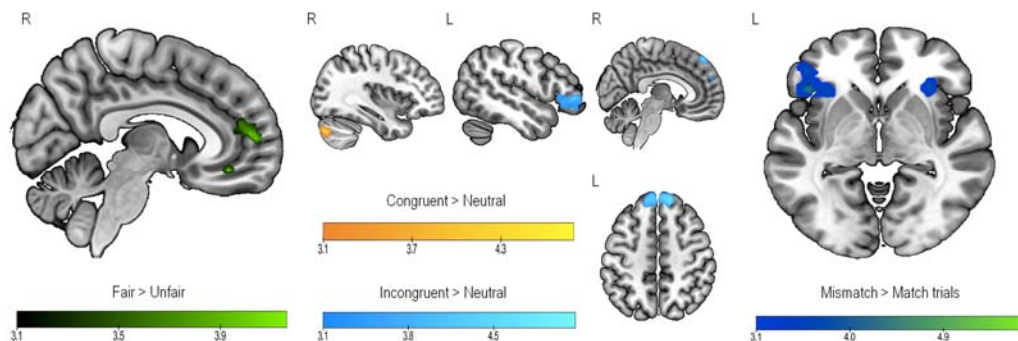
424 To check whether the regions related to personal information were the same during the
425 presentation of the valenced adjectives and during the presentation of the offer (positive
426 and negative > neutral in both cases), we ran a conjunction analysis with the regions
427 significant in both contrasts (Nichols, Brett, Andersson, Wager, & Poline, 2005).
428 Similar to each contrast individually, we observed two clusters: one in the left IFG/aI (k
429 = 93; -3, 21, 0) and one involving bilateral SMA ($k = 126$; -5, 18, 53), suggesting that
430 both areas increased their activation during the expectation and offer stages (see Figure
431 3b).

432

433 *Offer fairness*

434 **Fair offers** (Fair > Unfair) generated activity (see Figure 4) in the right medial frontal
435 gyrus (mFG) and ACC ($k = 171$; 6, 39, -14), while the opposite contrast (unfair > fair)
436 did not yield any significant clusters ($p > 0.05$). Furthermore, we examined neural
437 responses depending on whether previous expectations were matched or not by the
438 nature (fair vs. unfair) of the offer. Here, **congruence** (see Figure 4) between
439 expectations and offer (Congruent > Neutral) showed a cluster of activity in right
440 cerebellum (right Crus; $k = 153$; 17, -88, -32). Conversely, **incongruence** (see Figure 4)
441 between expectations and offer (Incongruent > Neutral) yielded activations in the right
442 medial Superior Frontal Gyrus (mSFG) and its lateral portion bilaterally ($k = 401$; 13,
443 39, 56), as well as in left IFG ($k = 177$; -54, 39, 0). Lastly, regarding general conflict
444 effects, a comparison between **mismatch** (incongruent) vs. **match** (congruent) trials
445 showed clusters of bilateral activity in the IFG/aI ($k = 232$; -43, 25, -11/ $k = 140$; 34,
446 35, 4; see Figure 4).

447



448

449 **Figure 4.** Univariate results for the offer. Scales reflect peaks of significant t-values ($p < .05$, FWE-
450 corrected for multiple comparisons).

451

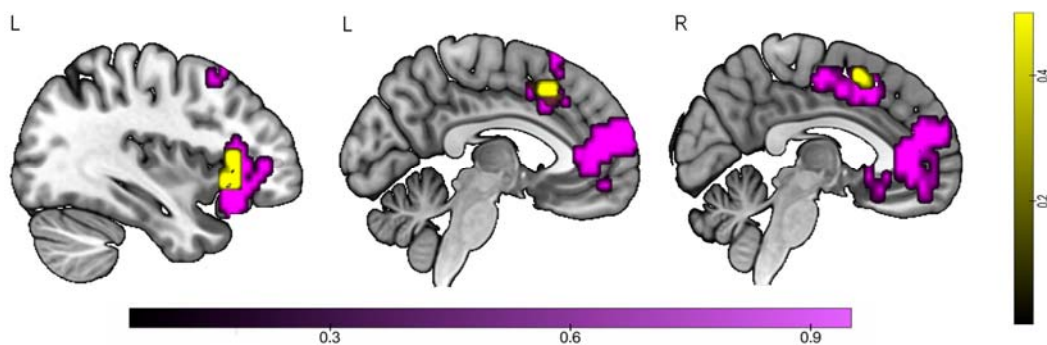
452 **3.2.2. Multivariate results**

453 *Valence of expectations' classification*

454 Expectations about the partners (positive vs. negative information) showed distinct
455 patterns of neural activity in a cluster including the left inferior and middle frontal gyrus
456 (IFG/MFG) and aI ($k = 319$; -46.5, 28, -32.2), the bilateral ventromedial prefrontal
457 cortex (vmPFC) and ACC ($k = 483$; 6, 21, -19.6), and the bilateral middle cingulate
458 cortex (MCC) and SMA ($k = 339$; -4.5, 14, 35; see Figure 5).

459

460 Although the same comparisons (positive vs. negative) in univariate GLM only yielded
461 a significant cluster activation in the IPL for positive > negative expectations, we ran a
462 conjunction analysis (Nichols, Brett, Andersson, Wager & Poline, 2005) to test whether
463 the regions that increased their activation during the presentation of the adjectives
464 (positive & negative > neutral) were similar to those that contained relevant information
465 about the valence (as reflected by multivariate results). For this, we computed the
466 intersection between the group maps from both contrasts. Results showed two clusters
467 (see Figure 5): one in the left IFG/aI ($k = 56$; $-36, 25, 0$) and one involving the bilateral
468 SMA ($k = 69$; $-8, 18, 46$).
469



470
471 **Figure 5.** Multivariate results (violet). Different neural patterns for the valence of the adjective (positive
472 vs. negative) during the expectation stage. Scales reflect corrected p-values (<.05). Significant regions in
473 both univariate and multivariate analyses are highlighted in yellow.
474

475 Moreover, the valence of the partners' descriptions influenced participants' choices,
476 where they accepted more offers after positive than negative descriptions. As explained
477 in the methods section (2.6 Multivariate analyses), information about participants'
478 responses might be employed to decode the valence of partners' descriptions. To
479 examine whether the regions containing relevant information about the valence were the
480 same as those representing the participants' decisions (accept vs. reject), we performed
481 a response classification at the offer period and ran a conjunction analysis. Here, we
482 observed that only a cluster in the bilateral SMA ($k = 95$; $-1, 7, 48$) resulted significant
483 for both classification analyses. Additionally, we carried out a cross-classification
484 analysis (Kaplan et al., 2015) to examine the overlap between the neural representations
485 of participants' choices and the valence of partners' descriptions. In this case, that a
486 classifier trained with response data is not able to decode valence category accurately
487 would suggest that the neural codes underlying valence and response classifications are

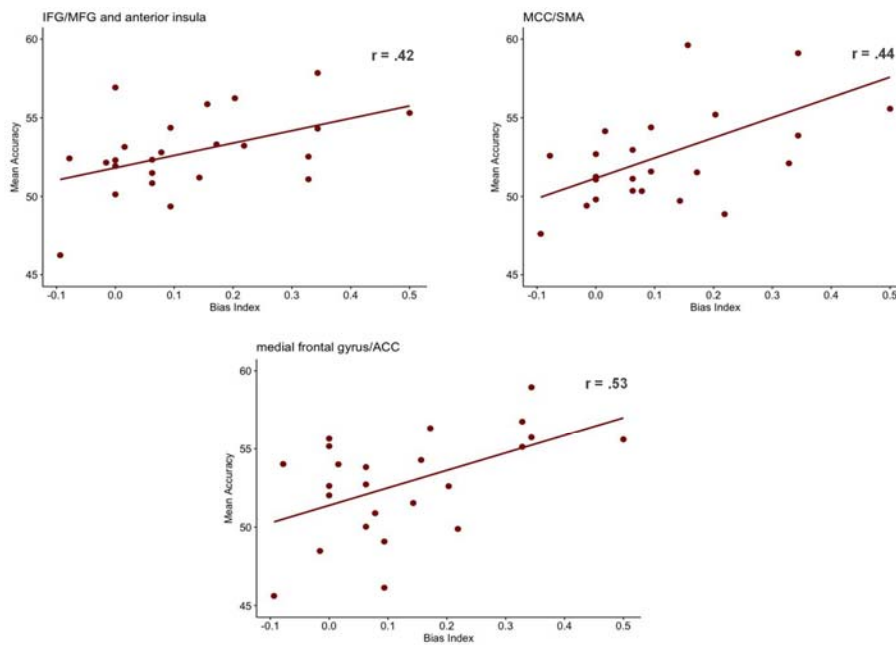
488 different and, therefore, that the valence decoding results are not explained by
489 participants' responses. Results from this analysis showed that cross-decoding was only
490 possible from bilateral SMA extending to left parietal lobe ($k = 671$; -1, -11, 45), as well
491 as from a cluster in left cerebellum extending to lingual and fusiform gyri ($k = 381$; -18,
492 -60, -15). This indicates that classification of valence in IFG/aI and vmPFC/ACC
493 cannot be explained by the patterns related to participants' responses.

494

495 *Correlation between decoding accuracy and the bias index*

496 To explore how much influence the valence of the adjectives had on choices, we
497 correlated the mean decoding accuracies (positive vs. negative) for each significant
498 cluster in the MVPA with the behavioural bias index for each participant. This analysis
499 yielded significant positive correlations between the decoding accuracy for the
500 descriptions' valence and the behavioural bias in all 3 significant clusters (see Figure 6):
501 the left IFG/MFG and aI ($r = .42$; $p = .02$), bilateral vmPFC/ACC ($r = .44$; $p = .015$),
502 and the left MCC/SMA ($r = .53$; $p = .0038$). Hence, the better the activation patterns in
503 these regions discriminated between the valence of the partners' information, the larger
504 the effect of valenced information on subsequent behavioural choices. A second
505 correlation control analysis showed that this link was not contaminated by participants'
506 motor responses, since there was no correlation between any of the ROIs mean
507 accuracies and general acceptance rate per participant (all $ps > .39$), which supports the
508 specificity of the link between valenced expectations and choices.

509



510

511 **Figure 6.** Scatter plots showing significant correlations between mean decoding accuracies in each cluster
512 and the behavioural index. IFG: Inferior frontal gyrus. MFG: Middle frontal gyrus. ACC: Anterior
513 Cingulate Cortex. MCC: Middle Cingulate Cortex. SMA: Supplementary Motor Area.

514

515 **4. Discussion**

516 Our study investigated the neural basis of social valenced expectations during an
517 interpersonal UG. Results revealed that social information about other people bias
518 subsequent economic choices, as well as it increases activity in the anterior insula and
519 SMA. Furthermore, decoding analysis allowed to observe that these areas, together with
520 the vmPFC, represent the content of such expectations. Notably, the better this
521 information is represented in these regions, the more biased are participants to employ
522 such knowledge when making their economic decisions.

523

524 The UG employed showed a clear behavioural effect of interpersonal expectations,
525 where positive descriptions of others led to higher acceptance rates compared to
526 negative ones. Additionally, the impact of the expectations was reflected on the speed of
527 choices, where people needed more time to reject offers after positive (or neutral)
528 expectations. This pattern indicates that participants integrate social information in their
529 decision-making process, showing a tendency to process offers as fairer when the
530 partner is described positively. Further, this data replicates previous results (Gaertig et
531 al., 2012; Moser et al., 2014; Ruz et al., 2011), emphasizing the role of expectations

532 (Sanfey, 2009) and valenced morality in decision-making (Barrett & Bliss-Moreau,
533 2009). Overall, the behavioural pattern of choices observed supports the utility of the
534 experimental paradigm to induce interpersonal valenced expectations about others that
535 bias subsequent choices made to the same set of objective behaviour (offers made by
536 partners).

537

538 Several regions increased their activation when participants held in mind social
539 expectations about game partners. This information engaged the SMA and the dorsal aI,
540 which were also active at the offer stage. These are regions have been previously related
541 to preparation processes (Brass & von Cramon, 2004), as well as sustained (Dosenbach,
542 Fair, Cohen, Schlaggar, & Petersen, 2008; Palenciano, González-García, Arco, & Ruz,
543 2019) and transient (Menon & Uddin, 2010; Sridharan, Levitin, & Menon, 2008) top-
544 down control, in paradigms where participants use cue-related information to perform
545 tasks of different nature on subsequent targets. In previous studies using the UG, these
546 regions have been linked to response to unfairness (Gabay et al., 2014). In addition,
547 previous work has related aI activation with the rejection of unfair offers (Sanfey et al.,
548 2003). In the current context, these areas may be involved in using the interpersonal
549 information contained in the cue to guide or bias the action towards a certain choice,
550 according to the valence of the expectation. However, univariate contrasts between the
551 words containing positive vs. negative information, in stark contrast with behavioural
552 outcomes, showed effects restricted on a cluster in the IPL. This region has been related
553 to the simulation of others' action in shared representations (Van Overwalle, 2009), and
554 a part of our cluster it is included in the TPJ (e.g., Scholz et al., 2009), which plays a
555 main role in ToM (Saxe & Kanwisher, 2003). The increase of activation in this region
556 for positive expectations could indicate a higher reliance on positive descriptions by the
557 ToM processes involved in our task. This fits with the pattern found in RTs where only
558 positive expectations speeded acceptance choices, whereas negative descriptions did not
559 speed rejections. Further research will be needed to replicate this imbalance of
560 information and to better understand the nature of the underlying brain processes.

561

562 Importantly, the use of a multivariate classification analysis (MVPA) unveiled the brain
563 regions that contain differential patterns for positive vs. negative expectations about
564 partners. This is especially relevant since previous work has indicated how valence
565 differences at a neural level are particularly hard to observe (Lindquist et al., 2015).

566 These areas included the SMA/MCC, IFG/MPFC and vmPFC/ACC. There was no
567 difference in RT between positive and negative conditions (see Behavioural data,
568 section 3.1.), which rules out the possibility that the classifier was mistakenly
569 discriminating faster vs. slower conditions.

570

571 The relevance of the SMA in social scenarios has been reported previously (Chang &
572 Sanfey, 2013). These authors observed a relationship between the activity in this area
573 and the deviation of previous expectations. Moreover, Lindquist et al. (2015) linked this
574 region to the unspecific representation of valence. Our conjunction analysis shows that
575 part of the SMA increases its activity during the expectation period and also shows
576 different patterns depending on the valence of the expectation. This data suggests that
577 the SMA has a role in general preparation but it also contains specific fine information
578 relevant to the task. In addition, we observe partial overlapping activation with the
579 response classification, which suggests that this region also contains some information
580 about participants' responses. The MCC, on the other hand, has been associated with an
581 increase of the efficiency in decision-making, being involved in the anticipation and
582 consequent expectations of outcomes in a variety of non-social tasks (Vogt, 2016).
583 Further, it has also been related to the prediction and monitoring of outcomes in social
584 decisions (Apps, Lockwood, & Balsters, 2013), and it may play a similar role in our
585 study.

586

587 On the other hand, the patterns of activity in a lateral prefrontal cortex cluster (LPFC),
588 including the IFG and MPFC, also discriminated the valence of the expectations.
589 Interestingly, these areas were part of a large cluster that also increased their activation
590 during the maintenance of social information, as revealed by univariate results. In non-
591 social paradigms, the LPFC has been related to working memory maintenance (Morgan,
592 Jackson, Van Koningsbruggen, Shapiro, & Linden, 2013; Sala, Rämä, & Courtney,
593 2003) and other forms of cognitive control (e.g., Reverberi et al., 2012). The IFG
594 specifically has also been associated with the selection of semantic information
595 (Jefferies, 2013; Wagner, Paré-Blagoev, Clark, & Poldrack, 2001), and it is also
596 involved in the expectation to perform different non-social tasks employing verbal
597 material (e.g., González-García et al., 2017; Sakai and Passingham, 2006). Notably, our
598 results extend this role to a social context (see also Filkowski et al., 2016; Thye et al.,
599 2018; Van Overwalle, 2009), where verbal information is used to generate positive or

600 negative expectations about game partners, by showing that the pattern of activity in this
601 frontal region differs depending on the nature of the information used to predict the
602 proximal behaviour of others.

603

604 On the other hand, the vmPFC/ACC did not increase its overall activation during the
605 expectation period but contained patterns related to the valence of the predictions.
606 Crucially, this area overlaps with the region isolated in the meta-analysis by Lindquist et
607 al., (2015), where they linked its activity with a bipolar representation of valence. On a
608 broader context, this region is part of the social cognition network, associated with
609 mentalizing processes (Koster-Hale & Saxe, 2013; Tamir et al., 2016), and behaviour
610 guided by social cues, along with the ACC. Previous studies relate the mPFC with
611 predictions about others' desires (Corradi-Dell'Acqua et al., 2015), and priors during
612 valued decisions (Lopez-Persem et al., 2016). Additionally, Van Overwalle (2009)
613 linked this region to the integration of personal traits, and it has been extensively
614 associated with the representation of intentions as well (Haynes et al., 2007).

615

616 The association between a brain region and a given behaviour is strengthened when a
617 link can be observed between the fidelity of a pattern of activity and the behavioural
618 outcome studied (Naselaris, Kay, Nishimoto, & Gallant, 2011; Tong & Pratte, 2012).
619 To find this evidence we obtained, for each participant, a bias index representing how
620 much the valence of the personal information influenced their choices and correlated
621 this index with the accuracy of the classifier in disentangling the patterns generated by
622 positive vs. negative words. We observed a positive correlation between these two
623 factors in the three clusters sensitive to the valence of expectations. Thus, the better the
624 classifier distinguished between descriptions of different valence, the more people
625 tended to accept offers preceded by positive compared to negative descriptions. These
626 results strongly suggest that these valenced representations were used to weight the
627 posterior acceptance or rejection decisions to the same set of objective offers, biasing
628 behaviour. Importantly, additional control correlation analysis evidenced that this
629 finding was not contaminated by participants' responses.

630

631 We could also observe the effect of expectations by studying the brain activity
632 generated by offers that matched or mismatched them, that is, fair and unfair offers
633 preceded by descriptions of the same or opposing valence. Here we found cerebellum

634 activity when fair offers were preceded by positive descriptions and unfair ones
635 followed negative adjectives. This region is associated with prediction in a variety of
636 contexts, such as language (Lesage, Hansen, & Miall, 2017; Pleger & Timmann, 2018)
637 and also social cognition (Van Overwalle, Baetens, Mariën, & Vandekerckhove, 2014),
638 among others. In social scenarios, where people frequently anticipate others' needs or
639 actions, the understanding of the role of the cerebellum in predictions is particularly
640 relevant (Sokolov, Miall, & Ivry, 2017). Although previous studies (Berthoz, 2002)
641 found increased activity in the cerebellum when predictions (social norms) were
642 violated, we observed the opposite. Hence, our data suggest that in the current context
643 the cerebellum may signal when predictions are matched by social observations.
644 Conversely, when predictions are not met, we observed activation in the IPFC,
645 specifically the IFG and aI. In this contexts, the IFG has been associated with semantic
646 cueing (González-García et al., 2016), semantic control (Jefferies, 2013) and emotional
647 regulation during social decisions (Grecucci, Giorgetta, Bonini, & Sanfey, 2013).
648 Conversely, the aI has been linked to responses to unfair offers, which represent a
649 violation of social norms (Corradi-Dell'Acqua et al., 2013). This agrees with the
650 incompatibility we observe here between previous expectations and actual events.
651 Altogether, this data also supports the relevance of expectations when participants face
652 the outcome of an interaction. At this point, they may need to suppress the previous
653 information to act in accordance with the offer.

654

655 Although it was not the main goal of this work, we also examined brain responses to the
656 fairness of the offer. While previous work has shown activation in areas such as aI,
657 cingulate cortex and mPFC in reaction to unfair offers (Corradi-Dell'Acqua et al., 2013;
658 Gabay et al., 2014), we observed higher activation in ACC/mPFC when participants
659 faced fair (vs. unfair) offers. In this line, the mPFC has been linked to the monitoring of
660 emotional reactions in bargaining scenarios (Corradi-Dell'Acqua et al., 2013), and its
661 involvement could represent the positive outcome related to fair offers, in line with
662 previous work associating the mPFC with value assessment of outcomes (Amodio &
663 Frith, 2006). The ACC, on the other hand, has been related to the proposal of fair offers
664 due to strategic motives (Chen, Chen, Kuo, Kan, & Yang, 2017), suggesting a role of
665 this area in computing reward. This, in turn, would be in line with our results of the
666 fairness of the offer, where the ACC could be relevant to signal their rewarding
667 outcomes.

668

669 Our study has certain limitations, which should be addressed in future investigations.
670 First, the optimal procedure to perform multivariate analyses and avoid response-related
671 confounds is to counterbalance response options for each participant (Todd, Nystrom, &
672 Cohen, 2013). In the current experiment, however, the association between hand and
673 response was counterbalanced at the group but not the individual level. Thus, our
674 valence-related classifications could have been affected by the response patterns linked
675 to acceptance and rejection choices. To rule this out, we performed an additional
676 conjunction analysis, which showed that only a small portion of the SMA cluster was
677 common to both contrasts. Also, we observed that patterns in part of this region
678 overlapped between participants' decisions and the valence of their expectations. These
679 results suggest that the SMA represents both events with similar codes, although it
680 could also be the case that findings in this region are due to confounds from
681 participants' responses. In further support of the relevance of the representation of the
682 valence in the bias observed in decisions, an additional control analysis showed that the
683 performance of the classifier for the valence decoding was only related to a specific
684 behavioural bias resulting from the valence of the expectation, but not with the response
685 itself. Therefore, our data highlight that the fidelity of the valence representation in
686 IFG/aI and vmPFC is associated with the extent to which the partners' descriptions
687 modulate participants' decisions.

688

689 Further, it may be argued that the influence of partners' moral information could be due
690 to alterations in participants' mood after reading these descriptions, rather than because
691 the generation of expectations about their likely behaviour. Although we cannot deny
692 completely this possibility, our findings show a specific link between participants'
693 behavioural bias and the neural representation of partners' social information, which
694 would not be in line with an explanation related to general mood fluctuations.
695 Alternatively, following previous work on affective priming and conflict (Dignath,
696 Eder, Steinhauser, & Kiesel, 2020; Fritz & Dreisbach, 2013), adjectives could act as
697 affective primes (Bush et al., 2018). Although we cannot completely rule out this
698 possibility, previous results suggest otherwise. Gaertig et al. (2012) carried out an
699 experiment without the social cover story to test this alternative explanation. Here, the
700 same words failed to trigger valence bias in choices. This indicates that, rather than an
701 automatic priming effect triggered by the adjectives, it is the association between these

702 and the character of the partners which impacted participants' decisions. An additional
703 concern relates to the ecological validity of our study, which is limited by the context of
704 fMRI scanning in a single location. However, we increased the credibility of the social
705 scenario by means of instructions and a cover story, where we recreated an actual
706 delayed interaction between participants of different studies, and where actual earnings
707 were contingent on the choices made during the game. In fact, none of the participants
708 showed signs of susceptibility about the underlying nature of the study when informally
709 debriefed at the end of the session. Nonetheless, participants could have approached the
710 task in various ways, engaging in the social context differently. Thus, we believe that
711 including a more detailed and structured debriefing where this and other points are
712 addressed should be included in future studies. Moreover, another step forward would
713 be to assess participants' personality and prosocial tendencies, since individual
714 predispositions can also influence these dynamics (Díaz-Gutiérrez et al., 2017). Futures
715 studies could use some form of virtual reality during scanning (Mueller et al., 2012)
716 together with more complex verbal descriptions of others to examine whether similar
717 brain regions represent this content and the way this is structured, perhaps employing
718 neuroimaging methods with higher ecological validity (e.g. Pinti et al., 2018).
719 Additionally, another interesting research question would be to find if there is a sort of
720 "common valence space" for the two stages of the paradigm. That is, to find out if there
721 is shared information underlying the valence of the adjective (positive/negative) but also
722 the "pleasantness" of the offer (fair-positive, unfair-negative). A future study designed
723 to employ cross-classification decoding approaches (Kaplan et al., 2015) between the
724 expectation and the evidence game periods with temporally precise methods such as
725 electroencephalography could offer valuable information on this respect.

726

727 **Acknowledgments**

728 This work was supported through grants by the Spanish Ministry of Science and
729 Innovation (PSI2013-45567-P and PSI2016-78236-P to M.R.), the Spanish Ministry of
730 Education, Culture and Sports (FPU2014/04272 to P.D.G.) and the University of
731 Granada, through a "Contratos puente" scholarship to P.D.G.

732

733 **References**

734 Abdulrahman, H., & Henson, R. N. (2016). Effect of trial-to-trial variability on optimal

- 735 event-related fMRI design: Implications for Beta-series correlation and multi-voxel
736 pattern analysis. *NeuroImage*, *125*, 756–766.
737 <https://doi.org/10.1016/j.neuroimage.2015.11.009>
- 738 Amodio, D. M., & Frith, C. D. (2006). Meeting of minds: The medial frontal cortex and
739 social cognition. *Nature Reviews Neuroscience*, *7*(4), 268–277.
740 <https://doi.org/10.1038/nrn1884>
- 741 Apps, M. A. J., Lockwood, P. L., & Balsters, J. H. (2013). The role of the midcingulate
742 cortex in monitoring others' decisions. *Frontiers in Neuroscience*, *7*, 251.
743 <https://doi.org/10.3389/fnins.2013.00251>
- 744 Arco, J. E., González-García, C., Díaz-Gutiérrez, P., Ramírez, J., & Ruz, M. (2018).
745 Influence of activation pattern estimates and statistical significance tests in fMRI
746 decoding analysis.
- 747 Barrett, L. F., & Bliss-Moreau, E. (2009). Affect as a psychological primitive.
748 *Advances in Experimental Social Psychology*, *41*(08), 167–218.
749 [https://doi.org/10.1016/S0065-2601\(08\)00404-8](https://doi.org/10.1016/S0065-2601(08)00404-8).Affect
- 750 Berthoz, S. (2002). An fMRI study of intentional and unintentional (embarrassing)
751 violations of social norms. *Brain*, *125*(8), 1696–1708.
752 <https://doi.org/10.1093/brain/awf190>
- 753 Bode, S., & Haynes, J. D. (2009). Decoding sequential stages of task preparation in the
754 human brain. *NeuroImage*, *45*(2), 606–613.
755 <https://doi.org/10.1016/j.neuroimage.2008.11.031>
- 756 Brañas-Garza, P., Espín, A. M., Exadaktylos, F., & Herrmann, B. (2014). Fair and
757 unfair punishers coexist in the Ultimatum Game. *Scientific Reports*, *4*, 6025.
758 <https://doi.org/10.1038/srep06025>
- 759 Brass, M., & von Cramon, D. Y. (2004). Decomposing Components of Task
760 Preparation with Functional Magnetic Resonance Imaging. *Journal of Cognitive*
761 *Neuroscience*, *16*(4), 609–620. <https://doi.org/10.1162/089892904323057335>
- 762 Bush, K. A., Gardner, J., Privratsky, A., Chung, M. H., James, G. A., & Kilts, C. D.
763 (2018). Brain states that encode perceived emotion are reproducible but their
764 classification accuracy is stimulus-dependent. *Frontiers in Human Neuroscience*,
765 *12*(July), 1–15. <https://doi.org/10.3389/fnhum.2018.00262>
- 766 Chang, L. J., & Sanfey, A. G. (2013). Great expectations: Neural computations

- 767 underlying the use of social norms in decision-making. *Social Cognitive and*
768 *Affective Neuroscience*, 8(3), 277–284. <https://doi.org/10.1093/scan/nsr094>
- 769 Chen, Y., Chen, Y., Kuo, W., Kan, K., & Yang, C. C. (2017). Strategic Motives Drive
770 Proposers to Offer Fairly in Ultimatum Games: An fMRI Study. *Scientific Reports*,
771 7(527), 1–11. <https://doi.org/10.1038/s41598-017-00608-8>
- 772 Contreras, J. M., Banaji, M. R., & Mitchell, J. P. (2012). Dissociable neural correlates
773 of stereotypes and other forms of semantic knowledge. *Social Cognitive and*
774 *Affective Neuroscience*, 7(7), 764–770. <https://doi.org/10.1093/scan/nsr053>
- 775 Contreras, J. M., Banaji, M. R., & Mitchell, J. P. (2013). Multivoxel Patterns in
776 Fusiform Face Area Differentiate Faces by Sex and Race. *PLoS ONE*, 8(7),
777 e69684. <https://doi.org/10.1371/journal.pone.0069684>
- 778 Corradi-Dell’Acqua, C., Civai, C., Rumiati, R. I., & Fink, G. R. (2013). Disentangling
779 self- and fairness-related neural mechanisms involved in the ultimatum game: An
780 fMRI study. *Social Cognitive and Affective Neuroscience*, 8(4), 424–431.
781 <https://doi.org/10.1093/scan/nss014>
- 782 Corradi-Dell’Acqua, C., Turri, F., Kaufmann, L., Clément, F., & Schwartz, S. (2015).
783 How the brain predicts people’s behavior in relation to rules and desires. Evidence
784 of a medio-prefrontal dissociation. *Cortex*, 70, 21–34.
785 <https://doi.org/10.1016/j.cortex.2015.02.011>
- 786 Correa, A., Alguacil, S., Ciria, L. F., Jiménez, A., & Ruz, M. (2020). Circadian rhythms
787 and decision-making: a review and new evidence from electroencephalography.
788 *Chronobiology International*. <https://doi.org/10.1080/07420528.2020.1715421>
- 789 Correa, A., Ruiz-Herrera, N., Ruz, M., Tonetti, L., Martoni, M., Fabbri, M., & Natale,
790 V. (2017). Economic decision-making in morning/evening-type people as a
791 function of time of day. *Chronobiology International*, 34(2), 139–147.
792 <https://doi.org/10.1080/07420528.2016.1246455>
- 793 Delgado, M. R., Frank, R. H., & Phelps, E. A. (2005). Perceptions of moral character
794 modulate the neural systems of reward during the trust game. *Nature*
795 *Neuroscience*, 8(11), 1611–1618. <https://doi.org/10.1038/nn1575>
- 796 Di Russo, F., Berchicci, M., Bozzacchi, C., Perri, R. L., Pitzalis, S., & Spinelli, D.
797 (2017). Beyond the “Bereitschaftspotential”: Action preparation behind cognitive
798 functions. *Neuroscience and Biobehavioral Reviews*, 78(April), 57–81.

- 799 <https://doi.org/10.1016/j.neubiorev.2017.04.019>
- 800 Díaz-Gutiérrez, P., Alguacil, S., & Ruz, M. (2017). Bias and control in social decision-
801 making. In A. Ibáñez, L. Sedeño, & A. Gacriá (Eds.), *Neuroscience and Social*
802 *Science: The Missing Link* (pp. 47–68). Cham, Switzerland: Springer.
803 <https://doi.org/10.1007/978-3-319-68421-5>
- 804 Dignath, D., Eder, A. B., Steinhauser, M., & Kiesel, A. (2020). Conflict monitoring and
805 the affective-signaling hypothesis—An integrative review. *Psychonomic Bulletin*
806 *and Review*, 27(2), 193–216. <https://doi.org/10.3758/s13423-019-01668-9>
- 807 Dosenbach, N. U. F., Fair, D. A., Cohen, A. L., Schlaggar, B. L., & Petersen, S. E.
808 (2008). A dual-networks architecture of top-down control. *Trends in Cognitive*
809 *Sciences*, 12(3), 99–105. <https://doi.org/10.1016/j.tics.2008.01.001>
- 810 Esterman, M., & Yantis, S. (2010). Perceptual expectation evokes category-selective
811 cortical activity. *Cerebral Cortex*, 20(5), 1245–1253.
812 <https://doi.org/10.1093/cercor/bhp188>
- 813 Fehr, E., & Camerer, C. F. (2007). Social neuroeconomics: the neural circuitry of social
814 preferences. *Trends in Cognitive Sciences*, 11(10), 419–427.
815 <https://doi.org/10.1016/j.tics.2007.09.002>
- 816 Filkowski, M. M., Anderson, I. W., & Haas, B. W. (2016). Trying to trust: Brain
817 activity during interpersonal social attitude change. *Cognitive, Affective and*
818 *Behavioral Neuroscience*, 16(2), 325–338. [https://doi.org/10.3758/s13415-015-](https://doi.org/10.3758/s13415-015-0393-0)
819 [0393-0](https://doi.org/10.3758/s13415-015-0393-0)
- 820 Fleming, S. M., Thomas, C. L., & Dolan, R. J. (2010). Overcoming status quo bias in
821 the human brain. *Proceedings of the National Academy of Sciences*, 107(13),
822 6005–6009. <https://doi.org/10.1073/pnas.0910380107>
- 823 Fouragnan, E., Chierchia, G., Greiner, S., Neveu, R., Avesani, P., & Coricelli, G.
824 (2013). Reputational Priors Magnify Striatal Responses to Violations of Trust. *The*
825 *Journal of Neuroscience*, 33(8), 3602–3611.
826 <https://doi.org/10.1523/JNEUROSCI.3086-12.2013>
- 827 Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the*
828 *Royal Society B: Biological Sciences*, 360(1456), 815–836.
829 <https://doi.org/10.1098/rstb.2005.1622>
- 830 Frith, C. D. (2007). The social brain? *Phil. Trans. R. Soc. B*, 362(10), 671–678.

- 831 <https://doi.org/10.1098/rstb.2006.2003>
- 832 Frith, C. D., & Frith, U. (2008). Implicit and Explicit Processes in Social Cognition.
833 *Neuron*, 60(3), 503–510. <https://doi.org/10.1016/j.neuron.2008.10.032>
- 834 Frith, U., & Frith, C. (2001). The biological basis of social interaction. *American*
835 *Psychological Society*, 10(5), 151–155.
836 <https://doi.org/https://doi.org/10.1111/1467-8721.00137>
- 837 Fritz, J., & Dreisbach, G. (2013). Conflicts as aversive signals: Conflict priming
838 increases negative judgments for neutral stimuli. *Cognitive, Affective and*
839 *Behavioral Neuroscience*, 13(2), 311–317. [https://doi.org/10.3758/s13415-012-](https://doi.org/10.3758/s13415-012-0147-1)
840 [0147-1](https://doi.org/10.3758/s13415-012-0147-1)
- 841 Gabay, A. S., Radua, J., Kempton, M. J., & Mehta, M. A. (2014). The Ultimatum Game
842 and the brain: A meta-analysis of neuroimaging studies. *Neuroscience and*
843 *Biobehavioral Reviews*, 47, 549–558.
844 <https://doi.org/10.1016/j.neubiorev.2014.10.014>
- 845 Gaertig, C., Moser, A., Alguacil, S., & Ruz, M. (2012). Social information and
846 economic decision-making in the ultimatum game. *Frontiers in Neuroscience*,
847 6:103. <https://doi.org/10.3389/fnins.2012.00103>
- 848 González-García, C., Arco, J. E., Palenciano, A. F., Ramírez, J., & Ruz, M. (2017).
849 Encoding, preparation and implementation of novel complex verbal instructions.
850 *NeuroImage*, 148(January), 264–273.
851 <https://doi.org/10.1016/j.neuroimage.2017.01.037>
- 852 González-García, C., Mas-Herrero, E., de Diego-Balaguer, R., & Ruz, M. (2016). Task-
853 specific preparatory neural activations in low-interference contexts. *Brain*
854 *Structure and Function*, 221(8), 3997–4006. [https://doi.org/10.1007/s00429-015-](https://doi.org/10.1007/s00429-015-1141-5)
855 [1141-5](https://doi.org/10.1007/s00429-015-1141-5)
- 856 Grecucci, A., Giorgetta, C., Bonini, N., & Sanfey, A. G. (2013). Reappraising social
857 emotions: the role of inferior frontal gyrus, temporo-parietal junction and insula in
858 interpersonal emotion regulation. *Frontiers in Human Neuroscience*, 7:523.
859 <https://doi.org/10.3389/fnhum.2013.00523>
- 860 Güth, W., Schmittberger, R., & Schwarze, B. (1982). An experimental analysis of
861 ultimatum bargaining. *Journal of Economic Behavior and Organization*, 3(4),
862 367–388. [https://doi.org/10.1016/0167-2681\(82\)90011-7](https://doi.org/10.1016/0167-2681(82)90011-7)

- 863 Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2006). The feedback-related
864 negativity reflects the binary evaluation of good versus bad outcomes. *Biological*
865 *Psychology*, *71*(2), 148–154. <https://doi.org/10.1016/j.biopsycho.2005.04.001>
- 866 Hassabis, D., Spreng, R. N., Rusu, A. A., Robbins, C. A., Mar, R. A., & Schacter, D. L.
867 (2014). Imagine all the people: How the brain creates and uses personality models
868 to predict behavior. *Cerebral Cortex*, *24*(8), 1979–1987.
869 <https://doi.org/10.1093/cercor/bht042>
- 870 Haynes, J., Sakai, K., Rees, G., Gilbert, S., Frith, C., & Passingham, R. E. (2007).
871 Reading Hidden Intentions in the Human Brain. *Current Biology*, *17*(4), 323–328.
872 <https://doi.org/10.1016/j.cub.2006.11.072>
- 873 Hebart, M. N., Gorgen, K., & Haynes, J.-D. (2015). The Decoding Toolbox (TDT): a
874 versatile software package for multivariate analyses of functional imaging data.
875 *Frontiers in Neuroinformatics*, *8*, 88. <https://doi.org/10.3389/fninf.2014.00088>
- 876 Jefferies, E. (2013). The neural basis of semantic cognition: Converging evidence from
877 neuropsychology, neuroimaging and TMS. *Cortex*, *49*(3), 611–625.
878 <https://doi.org/10.1016/j.cortex.2012.10.008>
- 879 Kahneman, D., Knetsch, J. L., & Thaler, R. H. (1986). Fairness and Assumptions of
880 Economics. *Journal of Business*, *59*(4), S285-300.
- 881 Kaplan, J. T., Man, K., & Greening, S. G. (2015). Multivariate cross-classification:
882 applying machine learning techniques to characterize abstraction in neural
883 representations. *Frontiers in Human Neuroscience*, *9*, 151.
884 <https://doi.org/10.3389/fnhum.2015.00151>
- 885 Knoch, D., Pascual-Leone, A., Meyer, K., Treyer, V., & Fehr, E. (2006). Diminishing
886 reciprocal fairness by disrupting the right prefrontal cortex. *Science*, *314*(5800),
887 829–832. <https://doi.org/10.1126/science.1129156>
- 888 Koster-Hale, J., & Saxe, R. (2013). Theory of Mind: A Neural Prediction Problem.
889 *Neuron*, *79*(5), 836–848. <https://doi.org/10.1016/j.neuron.2013.08.020>
- 890 Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional
891 brain mapping. *Proceedings of the National Academy of Sciences of the United*
892 *States of America*, *103*, 3863–3868. <https://doi.org/10.1073/pnas.0600244103>
- 893 Lesage, E., Hansen, P. C., & Miall, R. C. (2017). Right Lateral Cerebellum Represents
894 Linguistic Predictability. *The Journal of Neuroscience*, *37*(26), 6231–6241.

- 895 <https://doi.org/10.1523/JNEUROSCI.3203-16.2017>
- 896 Lindquist, K. A., Satpute, A. B., Wager, T. D., Weber, J., & Barrett, L. F. (2015). The
897 Brain Basis of Positive and Negative Affect: Evidence from a Meta-Analysis of the
898 Human Neuroimaging Literature. *Cerebral Cortex*, *26*(5), 1910–1922.
899 <https://doi.org/10.1093/cercor/bhv001>
- 900 Lopez-Persem, A., Domenech, P., & Pessiglione, M. (2016). How prior preferences
901 determine decision-making frames and biases in the human brain. *eLife*, *5*:
902 *e20317*. <https://doi.org/10.7554/eLife.20317>
- 903 Ma, N., Vandekerckhove, M., Baetens, K., Overwalle, F. Van, Seurinck, R., & Fias, W.
904 (2012). Inconsistencies in spontaneous and intentional trait inferences. *Social*
905 *Cognitive and Affective Neuroscience*, *7*(8), 937–950.
906 <https://doi.org/10.1093/scan/nsr064>
- 907 Menon, V., & Uddin, L. Q. (2010). Saliency, switching, attention and control: a
908 network model of insula function. *Brain Structure and Function*, *214*(5–6), 655–
909 667. <https://doi.org/10.1007/s00429-010-0262-0>
- 910 Mitchell, T. M., Shinkareva, S. V., Carlson, A., Chang, K.-M., Malave, V. L., Mason,
911 R. A., & Just, M. A. (2008). Predicting Human Brain Activity Associated with the
912 Meanings of Nouns. *Science*, *320*, 1191–1195.
913 <https://doi.org/10.1126/science.1152876>
- 914 Morgan, H. M., Jackson, M. C., Van Koningsbruggen, M. G., Shapiro, K. L., & Linden,
915 D. E. J. (2013). Frontal and parietal theta burst TMS impairs working memory for
916 visual-spatial conjunctions. *Brain Stimulation*, *6*(2), 122–129.
917 <https://doi.org/10.1016/j.brs.2012.03.001>
- 918 Moser, A., Gaertig, C., & Ruz, M. (2014). Social information and personal interests
919 modulate neural activity during economic decision-making. *Frontiers in Human*
920 *Neuroscience*, *8*: 31. <https://doi.org/10.3389/fnhum.2014.00031>
- 921 Mueller, C., Luehrs, M., Baecke, S., Adolf, D., Luetzkendorf, R., Luchtman, M., &
922 Bernarding, J. (2012). Building virtual reality fMRI paradigms: A framework for
923 presenting immersive virtual environments. *Journal of Neuroscience Methods*,
924 *209*(2), 290–298. <https://doi.org/10.1016/j.jneumeth.2012.06.025>
- 925 Naselaris, T., Kay, K. N., Nishimoto, S., & Gallant, J. L. (2011). Encoding and
926 decoding in fMRI. *NeuroImage*, *56*(2), 400–410.

- 927 <https://doi.org/10.1016/j.neuroimage.2010.07.073>
- 928 Nichols, T., Brett, M., Andersson, J., Wager, T., & Poline, J. B. (2005). Valid
929 conjunction inference with the minimum statistic. *NeuroImage*, 25(3), 653–660.
930 <https://doi.org/10.1016/j.neuroimage.2004.12.005>
- 931 Palenciano, A. F., González-García, C., Arco, J. E., & Ruz, M. (2019). Transient and
932 Sustained Control Mechanisms Supporting Novel Instructed Behavior. *Cerebral*
933 *Cortex*, 29(9), 3948–3960. <https://doi.org/10.1093/cercor/bhy273>
- 934 Pereira, F., Mitchell, T., & Botvinick, M. (2009). Machine learning classifiers and fMRI:
935 A tutorial overview. *NeuroImage*, 45(1), S199–S209.
936 <https://doi.org/10.1016/j.neuroimage.2008.11.007>.Machine
- 937 Pinti, P., Tachtsidis, I., Hamilton, A., Hirsch, J., Aichelburg, C., Gilbert, S., & Burgess,
938 P. W. (2018). The present and future use of functional near-infrared spectroscopy
939 (fNIRS) for cognitive neuroscience. *Annals of the New York Academy of Sciences*,
940 1–25. <https://doi.org/10.1111/nyas.13948>
- 941 Pleger, B., & Timmann, D. (2018). The role of the human cerebellum in linguistic
942 prediction, word generation and verbal working memory: evidence from brain
943 imaging, non-invasive cerebellar stimulation and lesion studies. *Neuropsychologia*.
944 <https://doi.org/10.1016/j.neuropsychologia.2018.03.012>
- 945 Puri, A. M., Wojciulik, E., & Ranganath, C. (2009). Category expectation modulates
946 baseline and stimulus-evoked activity in human inferotemporal cortex. *Brain*
947 *Research*, 1301, 89–99. <https://doi.org/10.1016/j.brainres.2009.08.085>
- 948 Redondo, J., Fraga, I., Padrón, I., & Comesaña, M. (2007). The Spanish adaptation of
949 anew (Affective Norms for English Words). *Behavior Research Methods*, 39(3),
950 600–605. <https://doi.org/10.3758/BF03193031>
- 951 Reverberi, C., Görgen, K., & Haynes, J.-D. (2012). Compositionality of Rule
952 Representations in Human Prefrontal Cortex. *Cerebral Cortex*, 22(6), 1237–1246.
953 <https://doi.org/10.1093/cercor/bhr200>
- 954 Ruz, M., Moser, A., & Webster, K. (2011). Social expectations bias decision-making in
955 uncertain inter-personal situations. *PLoS ONE*, 6(2): e157.
956 <https://doi.org/10.1371/journal.pone.0015762>
- 957 Ruz, M., & Tudela, P. (2011). Emotional conflict in interpersonal interactions.
958 *NeuroImage*, 54(2), 1685–1691. <https://doi.org/10.1016/j.neuroimage.2010.08.039>

- 959 Sakai, K., & Passingham, R. E. (2006). Prefrontal Set Activity Predicts Rule-Specific
960 Neural Processing during Subsequent Cognitive Performance. *Journal of*
961 *Neuroscience*, 26(4), 1211–1218. [https://doi.org/10.1523/JNEUROSCI.3887-](https://doi.org/10.1523/JNEUROSCI.3887-05.2006)
962 05.2006
- 963 Sakai, Katsuyuki. (2008). Task Set and Prefrontal Cortex. *Annual Review of*
964 *Neuroscience*, 31(1), 219–245.
965 <https://doi.org/10.1146/annurev.neuro.31.060407.125642>
- 966 Sala, J. B., Rämä, P., & Courtney, S. M. (2003). Functional topography of a distributed
967 neural system for spatial and nonspatial information maintenance in working
968 memory. *Neuropsychologia*, 41(3), 341–356. [https://doi.org/10.1016/S0028-](https://doi.org/10.1016/S0028-3932(02)00166-5)
969 3932(02)00166-5
- 970 Sanfey, A. G. (2009). Expectations and social decision-making: Biasing effects of prior
971 knowledge on Ultimatum responses. *Mind and Society*, 8(1), 93–107.
972 <https://doi.org/10.1007/s11299-009-0053-6>
- 973 Sanfey, A. G., Rilling, J. K., Aronson, J. A., Nystrom, L. E., & Cohen, J. D. (2003). The
974 neural basis of economic decision-making in the Ultimatum Game. *Science*, 300,
975 1755–1758. <https://doi.org/10.1126/science.1082976>
- 976 Saxe, R., & Kanwisher, N. (2003). People thinking about thinking peopleThe role of the
977 temporo-parietal junction in “theory of mind.” *NeuroImage*, 19(4), 1835–1842.
978 [https://doi.org/10.1016/S1053-8119\(03\)00230-1](https://doi.org/10.1016/S1053-8119(03)00230-1)
- 979 Schneider, W., Eschman, A., & Zuccolotto, A. (2002). E-Prime user’s guide. Pittsburgh:
980 Psychology Software Tools Inc.
- 981 Scholz, J., Triantafyllou, C., Whitfield-Gabrieli, S., Brown, E. N., & Saxe, R. (2009).
982 Distinct regions of right temporo-parietal junction are selective for theory of mind
983 and exogenous attention. *PLoS ONE*, 4(3).
984 <https://doi.org/10.1371/journal.pone.0004869>
- 985 Schwarz, K. A., Pfister, R., & Büchel, C. (2016). Rethinking Explicit Expectations:
986 Connecting Placebos, Social Cognition, and Contextual Perception. *Trends in*
987 *Cognitive Sciences*, 20(6), 469–480. <https://doi.org/10.1016/j.tics.2016.04.001>
- 988 Sokolov, A. A., Miall, R. C., & Ivry, R. B. (2017). The Cerebellum: Adaptive
989 Prediction for Movement and Cognition. *Trends in Cognitive Sciences*, 21(5), 313–
990 332. <https://doi.org/10.1016/j.tics.2017.02.005>

- 991 Sridharan, D., Levitin, D. J., & Menon, V. (2008). A critical role for the right fronto-
992 insular cortex in switching between central-executive and default-mode networks.
993 *Proceedings of the National Academy of Sciences*, *105*(34), 12569–12574.
994 <https://doi.org/10.1073/pnas.0800005105>
- 995 Stelzer, J., Chen, Y., & Turner, R. (2013). Statistical inference and multiple testing
996 correction in classification-based multi-voxel pattern analysis (MVPA): Random
997 permutations and cluster size control. *NeuroImage*, *65*, 69–82.
998 <https://doi.org/10.1016/j.neuroimage.2012.09.063>
- 999 Stolier, R. M., & Freeman, J. B. (2016). Neural pattern similarity reveals the inherent
1000 intersection of social categories. *Nature Neuroscience*, *19*(6), 795–797.
1001 <https://doi.org/10.1038/nn.4296>
- 1002 Stolier, R. M., & Freeman, J. B. (2017). A Neural Mechanism of Social Categorization.
1003 *The Journal of Neuroscience*, *37*(23), 5711–5721.
1004 <https://doi.org/10.1523/JNEUROSCI.3334-16.2017>
- 1005 Summerfield, C., & De Lange, F. P. (2014). Expectation in perceptual decision making:
1006 Neural and computational mechanisms. *Nature Reviews Neuroscience*, *15*(11),
1007 745–756. <https://doi.org/10.1038/nrn3838>
- 1008 Tamir, D. I., & Thornton, M. A. (2018). Modeling the Predictive Social Mind. *Trends in*
1009 *Cognitive Sciences*, *22*(3), 201–212. <https://doi.org/10.1016/j.tics.2017.12.005>
- 1010 Tamir, D. I., Thornton, M. A., Contreras, J. M., & Mitchell, J. P. (2016). Neural
1011 evidence that three dimensions organize mental state representation: Rationality,
1012 social impact, and valence. *Proceedings of the National Academy of Sciences of*
1013 *the United States of America*, *113*(1), 194–199.
1014 <https://doi.org/10.1073/pnas.1511905112>
- 1015 Thornton, M. A., & Mitchell, J. P. (2017). Theories of Person Perception Predict
1016 Patterns of Neural Activity During Mentalizing. *Cerebral Cortex*, 1–16.
1017 <https://doi.org/10.1093/cercor/bhx216>
- 1018 Thyne, M. D., Murdaugh, D. L., & Kana, R. K. (2018). Brain Mechanisms Underlying
1019 Reading the Mind from Eyes, Voice, and Actions. *Neuroscience*, *374*, 172–186.
1020 <https://doi.org/10.1016/j.neuroscience.2018.01.045>
- 1021 Todd, M. T., Nystrom, L. E., & Cohen, J. D. (2013). Confounds in multivariate pattern
1022 analysis: Theory and rule representation case study. *NeuroImage*, *77*, 157–165.

- 1023 <https://doi.org/10.1016/j.neuroimage.2013.03.039>
- 1024 Tong, F., & Pratte, M. S. (2012). Decoding Patterns of Human Brain Activity. *Annual*
1025 *Review of Psychology*, 63(1), 483–509. [https://doi.org/10.1146/annurev-psych-](https://doi.org/10.1146/annurev-psych-120710-100412)
1026 [120710-100412](https://doi.org/10.1146/annurev-psych-120710-100412)
- 1027 Turner, B. (2010). Comparison of methods for the use of pattern classification on rapid
1028 event-related fMRI data. Poster session presented at the Annual Meeting of the
1029 Society for Neuroscience, San Diego, CA.
- 1030 Van Overwalle, F. (2009). Social cognition and the brain: A meta-analysis. *Human*
1031 *Brain Mapping*, 30(3), 829–858. <https://doi.org/10.1002/hbm.20547>
- 1032 Van Overwalle, F., Baetens, K., Mariën, P., & Vandekerckhove, M. (2014). Social
1033 cognition and the cerebellum: A meta-analysis of over 350 fMRI studies.
1034 *NeuroImage*, 86, 554–572. <https://doi.org/10.1016/j.neuroimage.2013.09.033>
- 1035 Vogt, B. A. (2016). Midcingulate cortex: Structure, connections, homologies, functions
1036 and diseases. *Journal of Chemical Neuroanatomy*, 74, 28–46.
1037 <https://doi.org/10.1016/j.jchemneu.2016.01.010>
- 1038 Wagner, A. D., Paré-Blagoev, E. J., Clark, J., & Poldrack, R. A. (2001). Recovering
1039 meaning: Left prefrontal cortex guides controlled semantic retrieval. *Neuron*,
1040 31(2), 329–338. [https://doi.org/10.1016/S0896-6273\(01\)00359-2](https://doi.org/10.1016/S0896-6273(01)00359-2)
- 1041 Yeung, N., & Sanfey, A. G. (2004). Independent Coding of Reward Magnitude and
1042 Valence in the Human Brain. *Journal of Neuroscience*, 24(28), 6258–6264.
1043 <https://doi.org/10.1523/JNEUROSCI.4537-03.2004>
- 1044