The Journal of Neuroscience

Vicarious reinforcement learning signals when instructing others

Matthew Apps, University of Oxford
Elise Lesage, NIDA
Narender Ramnani, Royal Holloway University of London

Commercial Interest: No

1

2

### Vicarious reinforcement learning signals when instructing others

Apps, M.A.J[1,2,4]., Lesage, E[3,4]., & Ramnani, N[4].

5

[1] Nuffield Department of Clinical Neuroscience, University of Oxford, Oxford, UK

[2] Department of Experimental Psychology, University of Oxford, Oxford, UK

[3] Currently at: NIDA IRP, National Institutes of Health, Baltimore, US

[4] Department of Psychology, Royal Holloway, University of London, UK.

10

11 Word Count: Intro (500), Methods and Results (5745) and Discussion (1495)

12 Corresponding author email: matthew.apps@ndcn.ox.ac.uk

13

14

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32    Abstract

33    *Reinforcement learning (RL) theory posits that learning is driven by discrepancies between the*

34    *predicted and actual outcomes of actions (prediction errors, PEs). In social environments, learning is*

35    *often guided by similar RL mechanisms. For example, teachers monitor the actions of students and*

36    *provide feedback to them. This feedback evokes PEs in students that guide their learning. We report*

37    *the first study that investigates the neural mechanisms that underpin these processes.  Neurons in*

38    *the Anterior Cingulate Cortex (ACC) signal PEs when learning from the outcomes of one's own*

39    *actions, but also signal information when outcomes are received by others. Does a teacher's ACC*

40    *signal PEs when monitoring a student's learning? Using fMRI, we studied brain activity in human*

41    *subjects (teachers) as they taught a confederate (student) action-outcome associations by providing*

42    *positive or negative feedback. We examined activity time-locked to the students' responses, when*

43    *teachers infer student predictions and know actual outcomes. We fitted a RL-based computational*

44    *model to the behaviour of the student to characterise their learning, and examined whether a*

45    *teacher's ACC signals when the student's predictions were wrong. In line with our hypothesis, activity*

46    *in the teacher's ACC covaried with the PE values in the model. Additionally, activity in the teacher's*

47    *insula and ventromedial prefrontal cortex covaried with the predicted value according to the student.*

48    *Our findings highlight that the ACC signals prediction errors vicariously for others' erroneous*

49    *predictions, when monitoring and instructing their learning. These results suggest that RL*

50    *mechanisms, processed vicariously, may underpin and facilitate teaching behaviours.*

51

52

53

54

55

56

57

58

59

60     Introduction

61     In reinforcement learning (RL) theory, learning is driven by prediction errors (PEs) (Sutton and Barto,

62     1998), which occur when the outcome of an action is discrepant from that which is predicted. A

63     wealth of research has found neurons that signal PEs when the outcomes of one's own actions are

64     unexpected (Rushworth et al., 2009). However, learning rarely occurs in a social vacuum. Often the

65     learning of 'students'  is guided by feedback provided by a 'teacher'. Such instructed learning is

66     thought to be fundamental for the transmission of abstract, complex information between humans

67     (Hoppitt et al., 2008). However, to date, there is no understanding of the neural or computational

68     mechanisms that underpin teaching behaviours (Stanley and Adolphs, 2013; Gariépy et al., 2014;

69     Ruff and Fehr, 2014). Does the brain of a teacher process the learning of a student under the

70     computational principles of RL theory?

71

72     The Anterior Cingulate Cortex (ACC) is well known for its role in social behaviour (Singer et al., 2004;

73     Ruff and Fehr, 2014). Lesions to the ACC disrupt the processing of social stimuli (Hadland et al., 2003;

74     Rudebeck et al., 2006), neurons in the ACC are sensitive to rewarding stimuli that others will receive

75     (Chang et al., 2013) and neuroimaging studies have shown that the ACC processes predictions about

76     the value of others' actions (Behrens et al., 2008; Jones et al., 2011; Zhu et al., 2012; Apps et al.,

77     2013b; Boorman et al., 2013; Apps and Ramnani, 2014). In contrast, theories of ACC function suggest

78     that it processes PEs relating to the outcomes of one's own decisions, in a manner that conforms to

79     RL principles (Silvetti et al. in press; Amiez et al., 2005; Alexander and Brown, 2011; Hayden et al.,

80     2011; Kennerley et al., 2011).

81     How can these viewpoints be reconciled? It has been claimed that the ACC gyrus (ACCg) processes

82     social information, but the computational principles that it instantiates parallel those of the adjacent

83     ACC regions (Apps et al., 2013a). That is, the ACCg processes PEs about others' actions. However, no

84     previous study has examined whether PEs are processed in the ACCg when monitoring,

85     understanding and guiding the learning of others.

86     Using fMRI, for the first time, we examine whether activity in the brain of a teacher can be

87     characterised by the computational principles of RL theory when monitoring and guiding the trial

88     and error learning of a student. We examined activity in subjects ('teacher') whose role was to teach

89     action-outcome contingencies to a confederate ('student') by monitoring their responses and

90     providing positive and negative feedback. Teachers had pre-learnt the correct associations and

91     therefore knew the actual value of each action. In addition, they could also model and simulate the

92     students' prediction of each outcome. Thus, the teachers could process a PE at the time of students'

93     actions. We fitted a RL-based computational model to student's behaviour, and tested the

94     hypothesis that activity in the ACCg of teachers would covary with PEs from the model at the time of

95     students' actions.

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

116 **Methods**

117 *Subjects*

118 Sixteen healthy right-handed participants were screened for neurological, psychiatric and

119 psychological disorders (ages 18-32; 10 female). One subject failed to complete the whole scanning

120 session and was excluded from the analyses. Each subject was paired with one of three confederate

121 participants, who they believed were a naïve participant. All participants gave written informed

122 consent. The studies were approved by the Royal Holloway, University of London Psychology

123 Department Ethics Committee and conformed to the regulations set out in the CUBIC MRI Rules of

124 Operation. The subjects were not paid for their participation but were given the incentive of

125 receiving a picture of their brain for taking part. The subjects were informed that the other

126 participant performing the task with them (the confederate) was being paid £5 for their participation

127 since they were not being scanned, but that this payment was unrelated to task performance.

128

129 *Task design*

130 Subjects performed a task in which they acted as a 'teacher' providing a 'student' (confederate) with

131 positive or negative feedback. The student learned the associations between a set of 10 arbitrary

132 instruction cues and one of four responses on a keypad.  The teacher had pre-learnt the same

133 associations one day prior to scanning, and was therefore able to determine whether an action

134 chosen for a particular visual cue was correct or incorrect. The teacher's task was to determine

135 whether the student's actions were correct or incorrect and then use a keypad of their own to

136 deliver this feedback to the student.

137

138 During the training the teacher was required to learn the arbitrary stimulus-response associations

139 between ten instruction cues (coloured shapes that gave no indication of which response was

140 correct) and one of four motor responses by trial and error (fig.1). That is, there was only one correct

141 response for each instruction cue ensuring that learning the correct association for one instruction

142 cue was not informative as to the correct associations for any other instruction cue. There were 100

143 trials in total, with ten presentations of each instruction cue. The instruction cues were presented in

144 two blocks, five instruction cues in the first 50 trials and five in the last 50 trials. The cues were

145 pseudorandomly presented, in a predefined sequence (see fig.1). A correct response was indicated

146 by the presence of a picture of a one pound coin at time of the feedback screen and an incorrect

147 response by a crossed out one pound coin. If the subjects did not respond within 750ms of the

148 trigger cue, feedback was displayed as "missed".

149

150 During the scanning session the teacher, monitored the student's responses and provided them with

151 feedback. The student learnt exactly the same associations as the teacher had learnt during the

152 training session, with trials presented in the same order. The teachers were also informed of the

153 identical nature of the trial structure. To maintain experimental control, we deceived teachers as to

154 the nature of the student. Whilst the teachers believed they were performing the task with another

155 genuine participant, the responses they saw were computer-generated and modelled on the

156 behaviour of a participant in the pilot training session. The students were drawn from one of three

157 confederates. This approach was necessary in order to maintain control over the performance of the

158 third-person, such that the behaviour of the other person was consistent across participants.

159

160 During the teaching task the teachers saw two sets of information that were not presented to the

161 student. Firstly, on one screen, the teachers were reminded of the correct association on each trial,

162 before the student made a response (fig.1). This eliminated the possibility that trials would be lost,

163 or that the student's learning would be compromised by poor performance of the teacher, as a

164 result of the teacher's failing to recall the correct association for each stimulus that they had learned

165 in the previous session. It also ensured that participants were able to register the discrepancy

166 between the student's prediction and the actual value of their action, a key component of our

167 hypotheses.

168

169 *Procedure*

170 *Training session*

171 Teachers were trained in two phases one day prior to scanning. In the first phase, the teacher was

172 seated in front of a monitor, with a response keypad. This first phase of the training was designed to

173 ensure that all teachers had learnt all the stimulus-response associations through trial-and-error. All

174 teachers made at least two consecutive correct responses for the last two presentations of each

175 instruction cue. All teachers had therefore learnt the correct associations for each stimulus. This

176 enabled them to act effectively as a teacher during the scanning session.

177

178 In the second part of the training session, the teacher was required to become familiar with their

179 role as a teacher, and therefore the task that they would perform in the scanner. During this session

180 the participant lay supine within a mock MRI scanner and provided positive and negative feedback

181 to the experimenter outside the mock scanner. They practised this role with the experimenter (see

182 scanning session below) such that they became familiar with the task they would perform during the

183 scanning session but were not teaching the student any information about associations that the

184 student would need to learn in the scanning session – i.e. they learnt how to teach a student,

185 without teaching a student stimulus-response associations that would be later used during scanning.

186 In this part of the training, exactly the same setup was used as during scanning, but with the

187 experimenter taking the place of the student and only a reduced number of trials (20) were used. It

188 is important to note that given the requirement to maintain control of responses of the

189 experimenter across subjects, the actions of the experimenter, as with the actual student, were

190 actually a set of pre-programmed computer-controlled responses.

191

192 *Scanning session*

193 Before the teacher entered the scanner they were shown the student sitting in the MRI control

194 room, in front of the monitor with a response keypad. The corner of the student's screen was

195 covered, allowing information to be presented to the teacher inside the scanner that the student

196 was not presented with (see trial structure below for more details). Crucially the teacher was made

197 aware that they would have access to information in the corner of the screen that was not able to be

198 seen by the student.

199

200 By obscuring that corner of only the student's screen (and not the teacher's screen) it was also

201 possible to present the teacher's trigger cue and response to them without the student being able to

202 observe this information. Hence, the teacher was also aware that the only feedback displayed to the

203 student was that of a pound coin or a pound coin with a cross through it at the time of the final

204 feedback.  If the teacher failed to accurately indicate whether the response of the student was

205 correct or incorrect, then the words "no feedback" were presented on the screen to the teacher and

206 the student. This strategy ensured that teachers believed that the student was learning from the

207 feedback that they were providing and ensured that they performed the task accurately. The

208 teacher believed that the student was responding to the trials in real-time, but in fact the trials were

209 computer-controlled, and the profile of responses were based on those of a participant during a

210 previous pilot experiment. This participant was chosen due to a fast learning rate (see behavioural

211 modelling below) and also as they missed only three trials. These trials were also shown to the

212 teacher, thus ensuring that the pre-programmed behaviour of the student seemed genuine to the

213 teacher. At the end of the scanning session the participants were asked standard debriefing

214 questions, as used in previous studies (Apps et al., 2012; Apps et al., 2013b; Apps and Ramnani,

215   2014), to ensure that they had maintained a full belief in the deception throughout the experiment.

216   Specifically, we asked four yes/no questions. (1) Are you surprised to read that you were deceived

217   on the task (yes/no) ? (2) Did you believe that the responses that you were observing were those of

218   the other person (yes/no)? (3)Did you believe the other person was learning the correct responses

219   from your feedback (yes/no)? (4) Did you believe that the other person was learning the correct

220   responses for the different shapes for the first time? (yes/no). A 'no' response on question one or a

221   'yes' response on questions two to four would have led to exclusion from the experiment.

222   *Trial structure (see fig.1).*

223   The teachers' trials consisted of an instruction cue (one of the ten that they had learnt associations

224   for during training), immediately followed by the cue indicating the correct button (which reminded

225   the teacher only – and not the student - of the correct association for that instruction cue), a student

226   trigger cue and response (indicating to the teacher which response the student had made), a teacher

227   trigger cue (to which the teacher pressed one button on a keypad for a correct student response and

228   another for an incorrect student response – cued by the presence of a pound or coin or a crossed

229   out pound coin switching pseudorandomly from left to right across trials) and then the feedback

230   (indicating to the student whether the response was correct or incorrect).

231

232   Computational Modelling

233   *Behavioural Modelling*

234   The behaviour of the student was modelled using a simple Rescorla-Wagner (R-W) based

235   reinforcement learning algorithm (Rescorla and Wagner, 1972) which has been extensively used to

236   examine the behavioural and neural basis of arbitrary visuomotor associations (Dayan and Balleine,

237   2002; Schultz, 2006; Brovelli et al., 2008; Dayan and Daw, 2008). This model also bears considerable

238   similarity to recent, influential models of ACC function (Silvetti et al., in press; Alexander and Brown,

239   2011). As the aim of this study was to examine brain activity in teachers, we maintained

240   experimental control by ensuring that all subjects observed the same learning behaviour exhibited

241   by the student. This requirement did not allow us to make comparisons between different

242   computational models of behaviour, as model comparison cannot be meaningfully applied to a

243   single subject's data. However, given the extensive use of the R-W model for associative learning

244   tasks similar to that used here (Dayan and Daw, 2008), and the fact that most recent computational

245   models of ACC function that we know of are underpinned by the same principles as a R-W model

246   (Silvetti et al. in press), this approach was more than sufficient for meeting the aims of this study.

247

248    The R-W model assumes that the associative value of an action (or stimulus) changes once new

249    information reveals that the actual outcome of a decision is different from the predicted outcome

250    (Rescorla and Wagner, 1972). Thus, on each trial, an action has a predicted associative value, that is

251    updated by a prediction error signal when the outcome reveals that this prediction is erroneous. The

252    evolution of the associative values for each action are given by:

253

254    *(1)*

$$V_{a(n+1)} = V_{a(n)} + \eta \, x \, \delta$$

255    Where:

256    *(2)* $$\delta = \lambda_a - V_{a(n)}$$

257

258    In both (1) and (2), *n* is the trial number, *a* = 1 ….*k* with *k* representing the available actions and η is

259    the learning rate. The asymptotic value (λ) of a correct action is greater than 0, but is a free

260    parameter that is estimated, and is 0 for an incorrect response. A prediction error is therefore the

261    student's prediction of its associative value ($V_{a(n)}$) subtracted from the actual value of the action ($\lambda$)

262    known by the teacher. We instructed the students (and teachers on the first day) that 1 of the four

263    finger movements could be correct for each instruction cue stimulus. Importantly, this also ensured

264    that learning the correct association for one instruction cue was not informative as to the correct

265    associations for any other instruction cue. Thus the associative values of actions for one instruction

266    cue were not informative as to the value of an action for another instruction cue. The initial

267    associative strength of each action for each stimulus was set to λ/4, given the equiprobability of

268    each of the four actions being correct.

269

270

271    *Model estimation*

272    To model the action selection process of the student we transformed the associative values into

273    probabilities using the softmax equation. This method is a standard approach used in reinforcement

274    learning theory (Sutton and Barto, 1981). The probability of the action chosen by a subject is given

275    by:

276

277    *(3)*

$$P_a(n) = \frac{exp(\beta \, V_a(n))}{\sum_{a} exp(\beta \, V_a(n))}$$

9

278

279   This equation converts the associative values of the action chosen by a subject to a probability

280   $(P_a(n))$. The coefficient β represents the stochasticity (or temperature) of the student's behaviour

281   (i.e. the sensitivity to the value of each option). A high β (greater than 1) causes all actions to be

282   nearly equiprobable, with a low β amplifying the differences in associative values. These two

283   algorithms were used to model action selection by the student over time. The associative value the

284   student placed on the chosen action ($V_a(n)$) was then updated in the R-W model, based on the

285   feedback.

286

287   Crucially, in this study, the feedback was provided by a teacher (the subject being scanned). As the

288   teacher had expert knowledge of all the associations –and was informed of the correct action on

289   each trial- they knew the asymptotic value (λ) of each action chosen by the student. In this

290   experiment, an aim was to examine whether the teacher modelled the learning of the student. It

291   was therefore assumed that to instruct the student, the teacher would have to calculate the

292   discrepancy between the student's prediction of the outcome ($Va_{(n)}$) and the asymptotic value (λ)

293   of the action chosen by the student. This asymptotic value would be known only by the teacher

294   whilst the student would still be learning. Only when the student has learnt the correct stimulus-

295   response associations for each cue would there be no discrepancy between the asymptotic value

296   known by the teacher and the prediction made by the student. The aim of the teacher was therefore

297   to provide the student with appropriate feedback to minimise the discrepancy between their own

298   expert knowledge and predictions made by the student.

299

300   Within the R-W model and the softmax algorithm there are free parameters which need to be

301   estimated. To identify the optimal set of free parameters for the student's behaviour (given the

302   teacher's feedback), the learning rate, the stochasticity parameter β and the asymptotic value λ

303   were varied. The output of the softmax algorithm is a series of probabilities, based on the values of

304   each of these parameters and the actions chosen by the student. By varying the parameters, the

305   probabilities output by the softmax algorithm differ. To select the parameters that best fitted the

306   student's behavioural data (given the teacher's feedback) a maximum likelihood approach was used.

307   By using a maximum likelihood algorithm it was possible to maximise the probabilities of the actions

308   chosen by the student and identify the values of each of the parameters that produced them. The

309   learning rate η was varied between 0 and 1in steps of 0.05, β between 0 and 5 in steps of 0.1 and λ

310   between 0 and 5 in steps of 0.1.  The likelihood of the chosen actions were found using:

311

312 *(3)* $$L = \sum_n \ln P_a(n)$$

313 where the likelihood of each set of parameters (L) is determined by the log of probability of the

314 performed action ($P_a(n)$) of the student at trial *n,* according to the model. If the model perfectly

315 predicts the actions, the probability of every chosen action would = 1 and L would be 0. As the

316 probabilities become less than 1 the log-likelihood L assumes negative values. The best fitting

317 parameters were then selected using:

318

319 *(4)* $$\theta' = \arg\max \theta \ (L)$$

320

321 This identified the set of parameters for which L was closest to 0 i.e. the best fitting parameter set.

322 Where $\theta$ is the parameter set and L is the log-likelihood. Importantly, in this study, the student's

323 data was computer controlled and thus every teacher observed the same responses of the student.

324 Variations in these parameters could therefore only be explained by changes in the feedback, i.e. if

325 the teacher failed to give the student feedback on a particular trial. If this happened, then those

326 trials were removed from the modelling and likewise, data at the time of the student's response on

327 those trials was removed from the fMRI analysis. The maximum likelihood approach revealed that

328 for the behaviour of the student, the best fitting parameters were a λ of 1, a learning rate η of 0.95

329 and a β values ranging from 2.3 to 2.7- reflecting the apparent differences in stochasticity of the

330 behaviour given the teacher's feedback (see fig.1). Importantly, we used the behaviour of a

331 participant from a pilot experiment as the 'student' behaviour. This student had a high learning rate

332 (0.95) and thus, this ensured that any effects we observed in the ACCg could not be accounted for by

333 teachers learning the learning rate of the student, as in Behrens et al. (2008).

334

335 *Apparatus*

336 Subjects lay supine in an MRI scanner (3T Siemens Trio, CUBIC, Royal Holloway, University of

337 London) with the fingers of the right hand positioned on an MRI-compatible response box. Stimuli

338 were projected onto a screen behind the subject and viewed in a mirror positioned above the

339 subjects face. Presentation software (Neurobehavioral Systems, Inc., USA) was used for

340 experimental control (stimulus presentation and response collection). A custom-built parallel port

341 interface connected to the Presentation PC received transistor-transistor logic (TTL) pulse inputs

342 from the response keypad. It also received TTL pulses from the MRI scanner at the onset of each

343 volume acquisition, allowing events in the experiment to become precisely synchronized with the

344 onset of each scan. The timings of all events in the experiment were sampled accurately,

345 continuously and simultaneously (independently of Presentation) at a frequency of 1 kHz using an

346 A/D 1401 unit (Cambridge Electronic Design, UK). Spike2 software was used to create a temporal

347 record of these events. Reaction times were calculated off-line, and event timings were prepared for

348 subsequent general linear model (GLM) analysis of fMRI data (see event definition and modelling

349 below).

350 *Functional Imaging and analysis*

351 *Data Acquisition*

352

353 Scans were acquired on a Siemens Trio 3T scanner. T1-weighted structural images were acquired at

354 a resolution of 1×1×1 mm using an MPRAGE sequence. 1016 EPI scans were acquired from each

355 participant. 38 slices were acquired in an ascending manner, at an oblique angle (≈30˚) to the AC-PC

356 line to decrease the impact of susceptibility artefact in subgenual cortex (Deichmann et al., 2003). A

357 voxel size of 3×3×3 mm (20% slice gap, 0.6 mm) was used; TR=3s, TE=32, flip angle=85°. The

358 functional sequence lasted 51 minutes. Immediately following the functional sequence, phase and

359 magnitude maps were collected using a GRE field map sequence ($TE_1$ = 5.19ms, $TE_2$ = 7.65ms).

360

361 *Image Preprocessing*

362 Scans were pre-processed using SPM8 ([www.fil.ion.ucl.ac.uk/spm](www.fil.ion.ucl.ac.uk/spm)). The EPI images from each

363 subject were corrected for distortions caused by susceptibility-induced field inhomogeneities using

364 the FieldMap toolbox (Andersson et al., 2001). This approach corrects for both static distortions and

365 changes in these distortions attributable to head motion (Hutton et al., 2002). The static distortions

366 were calculated using the phase and magnitude field maps acquired after the EPI sequence. The EPI

367 images were then realigned, and coregistered to the subject's own anatomical image. The structural

368 image was processed using a unified segmentation procedure combining segmentation, bias

369 correction, and spatial normalization to the MNI template (Ashburner and Friston, 2005); the same

370 normalization parameters were then used to normalize the EPI images. Lastly, a Gaussian kernel of 8

371 mm FWHM was applied to spatially smooth the images in order to conform to the assumptions of

372 the GLM implemented in SPM8.

373

374 *Event definition and modelling (Student response)*

375 Multiple GLMs analyses were performed to investigate activity time-locked to the teacher's

376 observation of the student's response. These were performed to ensure that activations identified

377 could only be accounted for by the uniquely explained variance of a parameter in the R-W model.

378 Although each of the GLMs differed from the others, they shared several common properties. Each

379    GLM contained regressors modelling the instruction cue, the student response cue, the teacher

380    trigger cue and the feedback cue. Regressors were constructed for each of these events by

381    convolving the event timings with the canonical Heamodynamic Response Function (HRF). The

382    effects of head motion were modelled in the analysis by including the six parameters of head motion

383    acquired during preprocessing as covariates of no interest. In addition to these regressors defined

384    for the event types, each GLM also contained regressors which were first order parametric

385    modulations of the student response cue event. These modulators scaled the amplitude of the HRF

386    in line with either the $\lambda_a$, $V_a$ or $\delta$ parameters from the Rescorla-Wagner algorithm. The values of

387    these parameters corresponded to the teacher's valuation ($\lambda_a$, the actual value of the action); the

388    student's prediction ($V_a$, the student's prediction of the value) and the prediction error ($\delta$, the

389    discrepancy between the student's prediction and the actual value) respectively. The prediction

390    error could of course only be coded by the teacher at the time of the student's action, as the student

391    would not have known the actual value of the action when they are learning. When a trial was

392    missed by the student or when teachers delivered erroneous feedback or failed to respond, these

393    parameters were all assigned a value of zero. Two sets of analyses were conducted in this study to

394    examine responses at the time of the student's response:

395

396    (1) Nine separate GLMs were created in which the values of one of $\lambda$, $V_a$, and $\delta$ were used as first-

397    order parametric modulators of the student response cues. These models enabled areas of the brain

398    in which the BOLD response varied in the manner predicted by one of the parameters to be

399    identified (see paragraph below). However, due to correlations between the values of these

400    parameters in the R-W model and correlations due to these parameters being time-locked to the

401    same event on each trial, additional analyses were required.

402    To examine activity that covaried with the prediction error parameter, we created three GLMs. The

403    first contained only the values of the $\delta$ parameter as a parametric modulation of the student

404    response cues. The second contained $\lambda$ as a parametric modulator, with the values of the $\delta$

405    parametric modulator orthogonalised with respect to the values $\lambda$. The third contained $V_a$ as a

406    parametric modulator, with the values of the $\delta$ parametric modulator orthogonalized with respect to

407    the values of the $V_a$ parameter. Voxels were only considered if they were significant in an F-contrast

408    in all three of these GLMs. This approach was then repeated for the $\lambda$ and $V_a$ parameters. Thus, nine

409    GLMs were constructed to examine activity which varied with the values from the parameters of the

410    R-W model. It is important to note that typically one would orthogonalise the parameter of interest

411    with respect to both of the other parameters, in one GLM. However, this was not possible in the

412    present study, because the prediction error parameter is a product of the other two parameters in

413   the R-W model. Thus, orthogonalizing the prediction error (δ) parameter with respect to both of the

414   other parameters in this model would have removed most of the variance that could be explained.

415   The approach we have used provides a statistically conservative way to ensure that any variance

416   that could be explained by the PE parameter is not due to its correlations with the student's

417   prediction parameter or the actual value (the teacher's valuation).

418   (2) To control for other possible responses in the ACC at the time of the student's response, we

419   created a GLM that contained alternative control parameters that varied with other plausible

420   responses which were not components of the R-W model.

421

422   The hypothesis of this study was that the ACC would signal a PE at the time of another's action. In

423   the R-W model these PEs are 'signed', such that during learning a negative outcome results in a

424   negative PE signal and a positive outcome results in a positive PE. However, it is notable that there is

425   empirical data that suggests that neurons in the ACC, and models of ACC function, have found both

426   signed and unsigned PEs in the ACC (Alexander and Brown, 2011; Kennerley et al., 2011; Matsumoto

427   et al., 2007). It was therefore crucial that we test the possibility that PEs in the ACC reflect not

428   classical PE signals, as found in dopamine neurons in the midbrain, but may reflect 'unsigned' PEs

429   that simply code for the magnitude of a PE and not whether it is positive or negative. We therefore

430   created an unsigned PE parameter, that covaried with the magnitude of δ but was always positive.

431

432   Classical error detection accounts of the ACC suggest that the region has a generalised role in

433   processing errors in information processing (Carter et al., 1998; Bush et al., 2000; Holroyd et al.,

434   2004; Yeung and Nieuwenhuis, 2009), including the processing of errors which are elicited by the

435   actions of others (Somerville et al., 2006; Shane et al., 2008; Yoshida et al., 2012). It is therefore

436   possible that the ACC might have exhibited an unsigned and uniform magnitude signal whenever the

437   student performed an incorrect action. To test this possibility we created a parameter that took on a

438   value of 1 whenever the student performed an incorrect action and 0 when there was no error.

439

440   The error detection and unsigned prediction error parameters were fitted to the responses of the

441   student and included in a GLM. In this GLM the parameters were not orthogonalized with respect to

442   each other, allowing them to compete to explain variance. This allowed us to determine which

443   parameter best explained activity in the ACCg at the time of the student's response. T-tests were

444   then conducted between them to test which parameter best explained activity in a given voxel.

445

446

447

448

449

450    *Outcome event*

451    In addition to the main analysis, we examined activity at the time of the outcome event. We used

452    the same strategy as that employed to examine activity at the time of the student's response,

453    namely to fit the parameters from the model to the time of the outcome events.

454

455    *Examining activity at the time of the teacher's response*

456    Whilst our design enabled us to examine activity at the time of the teacher's response, it was

457    suboptimal for asking questions about differences in how one's own compared to others actions are

458    processed in the brain. Thus, we did not compare activity between the student and teacher motor

459    events nor examine covariations with the BOLD response with parameter from the RW model at the

460    time of the teacher's response. However, other studies have used tasks specifically designed to

461    tackle such issues, which have nicely characterised responses in the brain comparing performing or

462    observing actions (Burke et al., 2010; Ramnani and Miall, 2004).

463

464    *Second-Level analysis*

465    Random effects analyses (Full-Factorial ANOVA) were applied to determine voxels significantly

466    different at the group level. SPM{t} images from all subjects at the first-level were entered into

467    second-level full factorial design matrices. T-contrasts and F-contrasts were conducted in each of the

468    GLMs. These contrasts identified voxels in which activity varied parametrically in the manner

469    predicted by the parameters in the R-W model. Separate corrections for multiple comparison were

470    used for the ACCg and the whole brain. To examine activity across the whole brain, FDR correction

471    was applied. In contrast, activity in the ACCg was corrected for by using an 80% probability mask of

472    the ACCg (see 'Anatomical Localization' below).

473

474    For the second set of analyses examining alternative models of ACC activity, the T-contrasts between

475    the prediction error parameter and the control parameters were examined at a lower threshold. This

476    was necessary due to the high covariance between each of these parameters. For these contrasts a

477    threshold of $P<0.01$, uncorrected for multiple comparisons, was employed.

478

479    It was possible that there may be individual differences in activity at the time of the student's

480    response, based on teacher's own learning history. To test this we input the learning rates from the

481    R-W model, which were estimated on the choices of the teacher in the initial training session, as

482    covariates of interest at the time of student's response.

483    *Anatomical Localization*

484    To test our hypothesis, we used an 80% probability anatomical masks of the ACCg. To create each

485    mask, subject-specific masks of the ACCg were constructed in FSL (http://www.fmrib.ox.ac.uk/fsl/).

486    Although the cytoarchitectonic boundaries of the ACC have no corresponding gross anatomical

487    landmarks, we defined the anatomical boundaries based on the location of these boundaries in

488    previous literature investigating cingulate cytoarchitecture (Vogt et al., 1995). To define the

489    posterior border of the midcingulate cortex, we used a boundary defined by a plane perpendicular

490    to the AC-PC line that lay 22 mm posterior to the anterior commissure (Vogt et al., 1995). We

491    included all voxels that lay within the ACCg extending anterior to this border, including subgenual

492    cingulate cortex. The final ACCg mask included only voxels which were within the ACCg in 80% of our

493    subjects. Importantly, this mask was of the ACCg only and did not extend into the adjacent sulcus.

494

495

496

**Results**

*Behavioural Results*

The teacher's task was to monitor the student's responses, determine whether the response was correct or incorrect, and deliver this as feedback to the student. The student's responses, unbeknown to the teachers, were computer-controlled replays of a real subject's responses during a pilot experiment, and included trials in which the student missed three trials (included such that the student's responses seemed realistic) and thus, teachers were required to respond on 97 trials. Teachers correctly gave feedback to the student on 95.2% (SD ± 2.9; range: 91-99%) of trials, indicating that all teachers understood the correct association for each stimulus and also understood whether the student's responses were correct or incorrect. In addition, responses to a standardised set of questions, revealed that none of the participants were aware of the nature of the deception. Thus, participants believed they were instructing another participants, and they were highly accurate at doing so.

*Imaging results*

*Student's response*

The main aim of this experiment was to examine activity in the brain of a teacher when they monitor the responses of a student. We tested the hypothesis that the ACCg would signal the discrepancy between a student's prediction and the actual outcome known by a teacher – a student prediction error (PE). In line with the hypothesis, activity was found in the ACCg (fig.2), putatively in midcingulate area 24a'/24b', which varied significantly with the PE ($\delta$) parameter of the R-W model (MNI coordinates (x,y,z) 2, 30, 12; Z = 3.17; p < 0.005 svc). Activity in this area was also better explained by the signed R-W PE parameter than by an unsigned PE parameter, or by a parameter in which simple response errors (see methods) were modelled (p > 0.01 uncorrected). No other region in the ACC, even at a reduced threshold, showed a significant covariation with the PE parameter (p > 0.01 uncorrected). No portion of the ACC showed a significant effect of either the unsigned parameter or the parameter which modelled every erroneous response of the student, even at a reduced threshold (p > 0.01). No region of the ACC showed a significant effect of the student prediction parameter, or the actual value known by the teacher (p > 0.01). No other brain area significantly varied with the prediction error parameter when correcting for multiple comparisons (p < 0.05 FDR). At a reduced threshold, activity in an area consistent with the location of the Ventral Tegmental Area (VTA) and the head of the caudate nucleus covaried with the PE parameter from the R-W model (P<0.005 uncorrected).

531

532 *Simulating the student prediction*

533 At the time of the student's response, the predicted value according to the student could be

534 modelled by the teacher. We examined whether activity in the brain of the teacher time-locked to

535 the student's action covaried with the student's prediction parameter ($V_{a(n)}$ ). Activity which varied

536 significantly with this parameter was found in a portion of the Ventromedial Prefrontal cortex

537 (VmPFC; -14, 32, -10, Z = 5.06, p < 0.05 FDR, putatively BA 32) and in the right short insular gyrus (48,

538 -4, -2, Z = 4.08 FDR, putatively area Idg; fig.3). These were the only regions in which the unique

539 variance could be accounted for significantly by the predicted value according to the student.

540

541 *The Teacher's valuation*

542 At the time of the student's action, the teacher knew the actual value of the student's choice. We

543 examined activity time-locked to the student's choice that covaried with the actual value of the

544 chosen action. Activity which varied statistically with this parameter was found in the Superior

545 Frontal Sulcus (SFS) bordering BAs 8,9 and 9/46 (-20, 32, 46; Z = 5.06, p < 0.05 FDR) and Posterior

546 Cingulate Cortex (PCC; -14, -52, 32; Z = 5.57, p < 0.05 FDR) putatively in BA. These were the only

547 regions in which the variance could be uniquely and significantly accounted for by the actual value of

548 the action known by the teacher.

549 Individual differences in the brains of teachers

550 To test whether activity at the time of the student's response varied depending on the teacher's

551 own learning history, we examined whether activity covaried with the learning rates of the teachers

552 in the initial training session. No areas of the brain covaried significantly when correcting for multiple

553 comparisons. However, at a reduced threshold (p < 0.001 uncorrected) we found activity in the three

554 regions, including regions that also responded to the teacher's valuation in bilateral SFS (MNI 26, 0,

555 42; Z = 4.4; -34, -2, 40;  Z = 3.87), and in the PCC (MNI -14, -22, 34; Z = 3.59), as well as in the intra-

556 parietal sulcus (MNI -44, -38, 50; Z = 4.05). However, these results should be interpreted with

557 caution, given the low sample size for exploring individual differences and that the results are

558 reported at an uncorrected threshold.

559

560

561

562

563    Outcome events

564    In addition to the main analysis, we also examined activity time-locked to the outcome event.

565    Activity was not found to covary with any of the parameters from the model at the time of the

566    outcome when correcting for multiple comparisons. However, activity was found to covary with PE

567    parameter from the model in several areas, Cerebellar Lobule VI (MNI -20, -38, 34, Z = 4.05), VmPFC

568    (MNI 10, 54, 12, Z = 3.92), Hippocampus (MNI 36, -12, -20), and the left temporal pole (MNI -56, -10,

569    -24; Z = 3.58), but only at a reduced threshold (p < 0.001 uncorrected).

570

571

572

573

574

575

**Discussion**

577 This study investigated activity in the brain of a teacher when monitoring a student's responses, as

578 the student learnt from feedback provided by the teacher. In line with our hypothesis, activity in a

579 portion of the ACCg varied with prediction error (PE) values in a RL-based computational model.

580 Activity in insula cortex and in the VmPFC varied with the predicted value of the action according to

581 the student. These results suggest that the ACCg plays a specific role in signaling information about

582 how erroneous another's predictions about their actions are. In addition, we found that areas that

583 are monosynaptically interconnected with the ACCg also play important roles in the processing of

584 information about other people's learning.

585 Anatomical evidence supports the notion that the ACCg is sensitive to information that guides

586 reinforcement learning. The ACCg receives direct input from dopaminergic neurons in the Ventral

587 Tegmental Area (VTA) (Williams and Goldman-Rakic, 1998). It has been well established that the

588 firing properties of dopamine neurons in the VTA conform to the principles of RL. Specifically, they

589 show an increased spike frequency to unexpectedly positive outcomes, a decreased spike frequency

590 to unexpectedly negative outcomes and no activity change to predictable outcomes (Schultz and

591 Dickinson, 2000; Schultz, 2006). As such the VTA is believed to signal PEs in a manner that drives

592 one's own learning of rewarding outcomes. Interestingly, we found that the BOLD signal in the ACCg

593 showed similar response characteristics. However, whilst it is well known that dopamine neurons

594 signal this information for one's predictions about the outcomes of one's own decisions, we have

595 shown that the ACCg processes such PE signals when they pertain to others' predictions and the

596 outcomes of others' actions as well.

597 Anatomical evidence also supports the notion that the ACC processes social information. The portion

598 of the ACCg that was activated in this study (in the gyral, midcingulate cortex) has strong

599 connections to the posterior portions of the superior temporal sulcus (pSTS), the temporal poles

600 (TPs) (Markowitsch et al., 1985; Seltzer and Pandya, 1989; Barbas et al., 1999), and the paracingulate

601 cortex (Pandya et al., 1981; Vogt and Pandya, 1987; Petrides and Pandya, 2007). These three regions

602 are believed to form a core circuit that is engaged when processing information about the mental

603 states of others (Ramnani and Miall, 2004; Frith and Frith, 2006; Hampton et al., 2008). In addition,

604 the ACCg has monosynaptic connections to the portions of the insula and the VmPFC that were

605 found to covary with the student's prediction in this study (Mesulam and Mufson, 1982; Mufson and

606 Mesulam, 1982; Morecraft et al., 1992; Cavada et al., 2000). Previous studies have shown that

607 activity in the VmPFC, the insula, the pSTS, the paracingulate cortex and the TPs covaries with

608 parameters from RL-based computational models during other forms of social interactions (Ramnani

609  and Miall, 2004; Behrens et al., 2008; Hampton et al., 2008; Baumgartner et al., 2009; Klucharev et

610  al., 2009; Cooper et al., 2013; Gariépy et al., 2014). Thus, input from areas which appear to process

611  information in a manner that conforms to the principles of RL during social interactions and the

612  input from midbrain dopaminergic nuclei both highlight the ACCg as a candidate for processing PE

613  signals relating to the behaviour of others. Moreover, these results suggest that the ACCg may

614  process information in concert with the VmPFC and the insula in order to vicariously process

615  information about the predictions other people make when learning.

616  Functional evidence also supports the claim that an overarching functional property of the ACCg is

617  that it processes information about rewards during social interactions (Apps et al., 2013a). Lesions to

618  the ACCg in monkeys disrupt the processing of social stimuli (Hadland et al., 2003; Rudebeck et al.,

619  2006) by reducing the typical delay present when reaching for a rewarding stimulus in the presence

620  of another monkey. In addition, single-unit recording studies have shown that a large proportion of

621  neurons in the ACCg code for a reward that a conspecific will receive. Crucially, these neurons do not

622  change their firing rate when an identical reward is to be received by oneself (Chang et al., 2013).

623  Imaging studies have also shown that the ACCg signals the net-value of rewards that others will

624  receive (Apps and Ramnani, 2014), signals the unpredictibility of the relationship between another's

625  advice and the outcomes of another's choices (Behrens et al., 2008), and signals when the outcomes

626  of another's actions are unexpected (Apps et al., 2013b). These results all support the view that the

627  ACCg signals information relating to reward-based decisions during social interactions. However, the

628  new contribution that our study makes is to show that the ACCg processes information at the time

629  of others' actions and does so when a subject's behaviour is aimed at guiding another's learning.

630  It has been argued that there are two major social frames of reference within which brain areas

631  process social information. Whilst some areas process information when inferring the intentions and

632  mental states of other people ('other' reference frame), other regions process information when

633  updating one's own behaviour based on other's intentions or behaviour ('self' reference frame)

634  (Hunt and Behrens, 2011; Baez-Mendoza et al., 2013; Baez-Mendoza and Schultz, 2013; Chang,

635  2013; Chang et al., 2013). Understanding the reference frames present in a task is therefore

636  important for understanding the frame of reference within which a region, in this case the ACCg,

637  processes social information. In this task, subjects were monitoring the learning of others in order to

638  provide them with feedback. Importantly, the design of the task ensured that participants were not

639  processing information about the relationship between their own actions and the reward they

640  would receive themselves. Rather, they were processing information about the erroneous

641  predictions of another. Interestingly, this supports recent claims that the ACCg (areas 24a'/24b') may

642    in fact act as a nexus between these two frames of reference (Hunt and Behrens, 2011; Apps et al.,

643    2013a). Specifically, it has been claimed that the area is engaged when processing information about

644    (i) the rewards that others will receive, based on one's own or others' actions, and (ii) others'

645    predictions about rewards, when others' predictions can be used to guide one's own behaviour

646    (Apps et al., 2013a). Our results support this claim by showing that the ACCg processes the

647    erroneous predictions of others (i.e. inferring information about others), in order that a subject can

648    provide them with feedback (i.e. updating one's own behaviour based on another's intentions).

649    Thus, the ACCg appears to process information in a way that acts as a nexus between the two major

650    social reference frames.

651    The functional and computational properties of the whole ACC are still under considerable debate,

652    however, one common feature of several recent accounts of the ACC is that they are underpinned

653    by similar computational principles to those of RL theory (Silvetti et al., in press; Yeung and

654    Nieuwenhuis, 2009). Several theories of ACC function have recently been developed that account for

655    a diverse range of single-unit recording, EEG and fMRI data. Silvetti et al.'s (in press) reward-value

656    and prediction model (RVPM) and Alexander and Brown's (2011) Predicted-Response Outcome

657    (PRO) model both argue that the ACC acts as a 'critic', learning the value of stimuli or actions

658    through PE signals. Similarly, Shenhav et al.'s (2013) Expected Value of Control (EVC) model is based

659    around the notion that the ACC signals the value of the amount of cognitive control that will be

660    required and updates this valuation when an outcome suggests this is required. Each of these

661    models relies upon PE signals updating predictions. These models are largely supported by empirical

662    evidence reporting from activity in areas 24c'/32', which lie in the sulcus of the ACC - a different

663    region of the ACC from that found of this study. The area we identified was in the ACCg in areas

664    24a'/24b'. Thus, in line with other recent studies (Boorman et al., 2013; Apps et al., 2013b),  our

665    research has shown that this region may also process PEs, a key component of R-L based models and

666    also of computational accounts of other ACC regions. Whether this PE is signalled by neurons that

667    also signal fictive PEs – PEs for the outcomes of unchosen actions –  that have been found in the ACC

668    (Hayden et al., 2009) is yet to be determined. However,  our results suggest that whilst the ACCg

669    may have a degree of specialization for social information processing, the computational principles

670    that govern its operation are similar to those of other regions of the ACC.

671    In summary, this study provided the first characterisation of the neural and computational processes

672    that may operate in the brain of a teacher as they deliver reinforcement to a student. Our findings

673    have highlighted a novel PE processed in the ACCg of a teacher that may play a key role in signalling

674    how erroneous students' predictions are. Furthermore, our findings suggest that areas previously

675    implicated in RL for oneself may also be important for vicariously processing and understanding the

676    learning of others.

677

678

679    **References**

680    Alexander WH, Brown JW (2011) Medial prefrontal cortex as an action-outcome predictor. Nat
681            Neurosci 14:1338-U1163.
682    Amiez C, Joseph JP, Procyk E (2005) Anterior cingulate error-related activity is modulated by
683            predicted reward. Eur J Neurosci 21:3447-3452.
684    Andersson JLR, Hutton C, Ashburner J, Turner R, Friston K (2001) Modeling geometric deformations
685            in EPI time series. Neuroimage 13:903-919.
686    Apps MA, Lockwood PL, Balsters JH (2013a) The role of the midcingulate cortex in monitoring others'
687            decisions. Frontiers in Neuroscience **7**.
688    Apps MAJ, Ramnani N (2014) The Anterior Cingulate Gyrus Signals the Net Value of Others' Rewards.
689            The Journal of Neuroscience 34:6190-6200.
690    Apps MAJ, Balsters JH, Ramnani N (2012) The anterior cingulate cortex: Monitoring the outcomes of
691            others' decisions. Social neuroscience 7:424-435.
692    Apps MAJ, Green R, Ramnani N (2013b) Reinforcement learning signals in the anterior cingulate
693            cortex code for others' false beliefs. Neuroimage 64:1-9.
694    Ashburner J, Friston KJ (2005) Unified segmentation. Neuroimage 26:839-851.
695    Baez-Mendoza R, Schultz W (2013) The role of the striatum in social behavior. Front Neurosci 7:233.
696    Baez-Mendoza R, Harris CJ, Schultz W (2013) Activity of striatal neurons reflects social action and
697            own reward. Proc Natl Acad Sci U S A 110:16634-16639.
698    Barbas H, Ghashghaei H, Dombrowski SM, Rempel-Clower NL (1999) Medial prefrontal cortices are
699            unified by common connections with superior temporal cortices and distinguished by input
700            from memory-related areas in the rhesus monkey. J Comp Neurol 410:343-367.
701    Baumgartner T, Fischbacher U, Feierabend A, Lutz K, Fehr E (2009) The Neural Circuitry of a Broken
702            Promise. Neuron 64:756-770.
703    Behrens TEJ, Hunt LT, Woolrich MW, Rushworth MFS (2008) Associative learning of social value.
704            Nature 456:245-U245.
705    Boorman ED, O'Doherty JP, Adolphs R, Rangel A (2013) The behavioral and neural mechanisms
706            underlying the tracking of expertise. Neuron 80:1558-1571.
707    Brovelli A, Laksiri N, Nazarian B, Meunier M, Boussaoud D (2008) Understanding the neural
708            computations of arbitrary visuomotor learning through fMRI and associative learning theory.
709            Cereb Cortex 18:1485-1495.
710    Burke CJ, Tobler PN, Baddeley M, Schultz W (2010) Neural mechanisms of observational
711            learning. Proc Natl Acad Sci U S A 107:14431-14436.
712    Bush G, Luu P, Posner MI (2000) Cognitive and emotional influences in anterior cingulate cortex.
713            Trends Cogn Sci 4:215-222.
714    Carter CS, Braver TS, Barch DM, Botvinick MM, Noll D, Cohen JD (1998) Anterior cingulate cortex,
715            error detection, and the online monitoring of performance. Science 280:747-749.
716    Cavada C, Company T, Tejedor J, Cruz-Rizzolo RJ, Reinoso-Suarez F (2000) The anatomical
717            connections of the macaque monkey orbitofrontal cortex. A review. Cerebral Cortex 10:220-
718            242.
719    Chang SW (2013) Coordinate transformation approach to social interactions. Front Neurosci 7:147.

720    Chang SWC, Gariepy J-F, Platt ML (2013) Neuronal reference frames for social decisions in primate
721        frontal cortex. Nature Neuroscience 16:243-250.
722    Cooper JC, Dunne S, Furey T, O'Doherty JP (2013) The Role of the Posterior Temporal and Medial
723        Prefrontal Cortices in Mediating Learning from Romantic Interest and Rejection. Cereb
724        Cortex.
725    Dayan P, Balleine BW (2002) Reward, motivation, and reinforcement learning. Neuron 36:285-298.
726    Dayan P, Daw ND (2008) Decision theory, reinforcement learning, and the brain. Cogn Affect Behav
727        Neurosci 8:429-453.
728    Deichmann R, Gottfried JA, Hutton C, Turner R (2003) Optimized EPI for fMRI studies of the
729        orbitofrontal cortex. Neuroimage 19:430-441.
730    Frith CD, Frith U (2006) The neural basis of mentalizing. Neuron 50:531-534.
731    Gariépy J-F, Watson KK, Du E, Xie DL, Erb J, Amasino D, Platt ML (2014) Social learning in humans and
732        other animals. Frontiers in Neuroscience 8.
733    Hadland KA, Rushworth MFS, Gaffan D, Passingham RE (2003) The effect of cingulate lesions on
734        social behaviour and emotion. Neuropsychologia 41:919-931.
735    Hampton AN, Bossaerts P, O'Doherty JP (2008) Neural correlates of mentalizing-related
736        computations during strategic interactions in humans. Proc Natl Acad Sci U S A 105:6741-
737        6746.
738    Hayden BY, Heilbronner SR, Pearson JM, Platt ML (2011) Surprise Signals in Anterior Cingulate
739        Cortex: Neuronal Encoding of Unsigned Reward Prediction Errors Driving Adjustment in
740        Behavior. J Neurosci 31:4178-4187.
741    Hayden, BY, Pearson, JM, Platt, ML (2009) Fictive reward signals in the anterior cingulate cortex.
742        Science 324: 948ding
743    Holroyd CB, Nieuwenhuis S, Yeung N, Nystrom L, Mars RB, Coles MGH, Cohen JD (2004) Dorsal
744        anterior cingulate cortex shows fMRI response to internal and external error signals. Nat
745        Neurosci 7:497-498.
746    Hoppitt WJE, Brown GR, Kendal R, Rendell L, Thornton A, Webster MM, Laland KN (2008) Lessons
747        from animal teaching. Trends Ecol Evol 23:486-493.
748    Hutton C, Bork A, Josephs O, Deichmann R, Ashburner J, Turner R (2002) Image distortion correction
749        in fMRI: A quantitative evaluation. Neuroimage 16:217-240.
750    Hunt LT, Behrens TEJ (2011) Frames of Reference in Human Social Decision Making. Neural Basis of
751        Motivational and Cognitive Control:409-424.
752    Jones RM, Somerville LH, Li J, Ruberry EJ, Libby V, Glover G, Voss HU, Ballon DJ, Casey BJ (2011)
753        Behavioral and neural properties of social reinforcement learning. J Neurosci 31:13039-
754        13045.
755    Kennerley SW, Behrens TEJ, Wallis JD (2011) Double dissociation of value computations in
756        orbitofrontal and anterior cingulate neurons. Nature Neuroscience 14:1581-U1119.
757    Klucharev V, Hytonen K, Rijpkema M, Smidts A, Fernandez G (2009) Reinforcement Learning Signal
758        Predicts Social Conformity. Neuron 61:140-151.
759    Markowitsch HJ, Emmans D, Irle E, Streicher M, Preilowski B (1985) cortical and subcortical afferent
760        connections of the primates temporal pole - a study of rhesus-monkeys, squirrel-monkeys,
761        and marmosets. J Comp Neurol 242:425-458.
762    Matsumoto M, Matsumoto K, Abe H, Tanaka K (2007) Medial prefrontal cell activity signaling
763        prediction errors of action values. Nature Neuroscience 10:647-656.
764    Mesulam MM, Mufson EJ (1982) insula of the old-world monkey .3. Efferent cortical output and
765        comments on function. Journal of Comparative Neurology 212:38-52.
766    Morecraft RJ, Geula C, Mesulam MM (1992) Cytoarchitecture and neural afferents of orbitofrontal
767        cortex in the brain of the monkey. Journal of Comparative Neurology 323:341-358.
768    Mufson EJ, Mesulam MM (1982) insula of the old-world monkey .2. Afferent cortical input and
769        comments on the claustrum. Journal of Comparative Neurology 212:23-37.

770 Pandya DN, Vanhoesen GW, Mesulam MM (1981) efferent connections of the cingulate gyrus in the
771      rhesus-monkey. Exp Brain Res 42:319-330.
772 Petrides M, Pandya DN (2007) Efferent association pathways from the rostral prefrontal cortex in the
773      macaque monkey. J Neurosci 27:11573-11586.
774 Ramnani N, Miall RC (2004) A system in the human brain for predicting the actions of others. Nat
775      Neurosci 7:85-90.
776 Rescorla RA, Wagner AR (1972) Classical Conditioning II: Current Research andTheory. In, pp 64–99.
777      New York: Appleton-Century Crofts.

778 Rudebeck PH, Buckley MJ, Walton ME, Rushworth MFS (2006) A role for the macaque anterior
779      cingulate gyrus in social valuation. Science 313:1310-1312.
780 Ruff CC, Fehr E (2014) The neurobiology of rewards and values in social decision making. Nat Rev
781      Neurosci 15:549-562.
782 Rushworth MFS, Mars RB, Summerfield C (2009) General mechanisms for making decisions? Curr
783      Opin Neurobiol 19:75-83.
784 Schultz W (2006) Behavioral theories and the neurophysiology of reward. Annual Review of
785      Psychology 57:87-115.
786 Schultz W, Dickinson A (2000) Neuronal coding of prediction errors. Annual Review of Neuroscience
787      23:473-500.
788 Seltzer B, Pandya DN (1989) frontal-lobe connections of the superior temporal sulcus in the rhesus-
789      monkey. J Comp Neurol 281:97-113.
790 Shane MS, Stevens M, Harenski CL, Kiehl KA (2008) Neural correlates of the processing of another's
791      mistakes: A possible underpinning for social and observational learning. Neuroimage 42:450-
792      459.
793 Shenhav, A, Botvinick, MM, Cohen, JD (2013) The expected value of control: an integrative theory of
794      anterior cingulate cortex function. Neuron 79: 217
795 Silvetti M, Alexander W, Verguts T, Brown J (in press) From conflict management to reward-based
796      decision making: Actors and critics in primate medial frontal cortex. Neuroscience &
797      Biobehavioral Reviews.
798 Singer T, Seymour B, O'Doherty J, Kaube H, Dolan RJ, Frith CD (2004) Empathy for pain involves the
799      affective but not sensory components of pain. Science 303:1157-1162.
800 Somerville LH, Heatherton TF, Kelley WM (2006) Anterior cingulate cortex responds differentially to
801      expectancy violation and social rejection. Nat Neurosci 9:1007-1008.
802 Stanley DA, Adolphs R (2013) Toward a neural basis for social behavior. Neuron 80:816-826.
803 Sutton RS, Barto AG (1981) toward a modern theory of adaptive networks - expectation and
804      prediction. Psychol Rev 88:135-170.
805 Sutton RS, Barto AG (1998) Reinforcement learning: an introduction. Cambridge, Massachusetts: MIT
806      press.
807 Vogt BA, Pandya DN (1987) cingulate cortex of the rhesus-monkey .2. Cortical afferents. J Comp
808      Neurol 262:271-289.
809 Vogt BA, Nimchinsky EA, Vogt LJ, Hof PR (1995) human cingulate cortex - surface-features, flat maps,
810      and cytoarchitecture. J Comp Neurol 359:490-506.
811 Williams SM, Goldman-Rakic PS (1998) Widespread origin of the primate mesofrontal dopamine
812      system. Cereb Cortex 8:321-345.
813 Yeung N, Nieuwenhuis S (2009) Dissociating Response Conflict and Error Likelihood in Anterior
814      Cingulate Cortex. Journal of Neuroscience 29:14506-14510.
815 Yoshida K, Saito N, Iriki A, Isoda M (2012) Social error monitoring in macaque frontal cortex. Nat
816      Neurosci 15:1307-U1180.
817 Zhu L, Mathewson KE, Hsu M (2012) Dissociable neural representations of reinforcement and belief
818      prediction errors underlie strategic learning. Proc Natl Acad Sci U S A 109:1419-1424.

819

**Figure Legends**

821 **Figure 1. (A) Trial Structure.** Participants performed trials as a teacher, guiding the associative
822 learning of a student. Each trial began a with a green instruction cue (one of ten that the teacher had
823 learnt the associations for during training), followed by the association cue informing the teacher of
824 the correct response for the stimulus. This was displayed in the corner of the teacher's screen. The
825 corresponding corner of the student's screen outside the scanner was covered, such that this cue
826 was shown only to the teacher inside the scanner. Following this, the teacher saw the student's
827 response. They were required to indicate to the student whether this response was correct or
828 incorrect. The teacher's indicated their response on a keypad at the time of a screen where a pound
829 coin (correct) or a crossed out pound coin (incorrect) were presented. Participants had to select the
830 corresponding stimulus to deliver to the student. This stimulus was also presented in the corner of
831 the screen, ensuring that the student could not see the teacher's decision at that time. The chosen
832 feedback was delivered to the student at the time of the outcome stimulus. **(B) Example model**
833 **data.** Plot of the data of the example output from the R-W model. In this example the learning rate
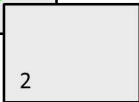834 was set to 1 for clarity.

835 **Figure 2. Student prediction errors**. (A) Activity shown in the ACC time-locked to the student's
836 response in which activity covaried with the prediction error parameter from the R-W model on the
837 mean anatomical image. (B) Parameter estimates in the peak ACC voxel. Activity in this region
838 correlated only with the prediction error parameter and not with the student's prediction or the
839 actual value of the outcome. Activity in this region also did not significantly covary with the unsigned
840 prediction error parameter or a parameter that simply coded for student erroneous responses. Error
841 bars depict standard error of the mean. (C) Peristimulus time histogram (PSTH) of activity time-
842 locked to the student's action in the brain of the teacher. Activity plotted for when the student's
843 prediction was erroneously positive (light green triangles) or erroneously negative (dark green
844 circles). The values of the prediction error were taken from the R-W computational model. Error bars
845 depict standard error of the mean.

846

847 **Figure 3. Simulating the student prediction.** Activity shown in the ventromedial prefrontal cortex (A)
848 and the right short insula gyrus (B) covarying with the predicted value according to the student,
849 taken from the R-W model. Plots of the parameter estimates from the peak voxel in the VmPFC (C)
850 and the insula (D) for the prediction error, the student predicted value and the actual value of the
851 outcome known by the teacher. Parameter estimates for the predicted value parameter are for the
852 unique variance explained by the regressor once orthogonalised with respect to the actual outcome
853 parameter. Parameter estimates for the prediction error parameter and the actual outcome
854 parameter are from regressors which have not been orthogonalised. Error bars depict standard error
855 of the mean. PSTH plots from the VmPFC (E) and the Insula (F) time-locked to the student's
856 prediction. Activity in these regions is broken down into low (<0.5) predicted value (light red
857 triangles) vs high (>0.5) predicted value (dark red circles) according to the model. Error bars depict
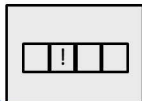858 standard error of the mean.

859

**A**

Instruction cue (750ms) – variable onset over 1st and 2nd scan of each trial (0 – 4500ms after first scan onset)

Association cue (750ms) – The correct action, visible only to the teacher in the scanner. Onset immediately after the instruction cue

Student response (1500ms) – The choice made by the student. Onset varied over the 3rd, 4th and 5th scans of each trial (0-7500ms after 3rd scan onset)

Teacher response (1500ms) – Teacher indicates whether student response is correct or incorrect., visible only to them. Onset varied over 6th and 7th scans (0-4500ms after 6th scan onset)

Outcome (750ms) – one pound, no pound or "no feedback" if the teacher's response was incorrect Onset varied over 8th, 9th or 10th scans (0-7500ms after 8th scan onset)

**B**

**A**

**B**

**C**