# Simultaneous representation of a spectrum of dynamically changing value estimates during decision making

3  David Meder[1,2], Nils Kolling[1,3], Lennart Verhagen[1], Marco K Wittmann[1,3],

4  Jacqueline Scholl[1], Kristoffer H Madsen[2], Oliver J Hulme[2], Timothy EJ

5  Behrens[3], Matthew FS Rushworth[1,3]

6

7  1    Department of Experimental Psychology, University of Oxford, South Parks

8       Road, Oxford OX1 3UD, UK

9  2    Danish Research Centre for Magnetic Resonance; Centre for Functional and

10      Diagnostic Imaging and Research, Copenhagen University Hospital

11      Hvidovre; Hvidovre, 2650; Denmark.

12  3    Oxford Centre for Functional MRI of the Brain, Nuffield Department of

13      Clinical Neurosciences, University of Oxford, John Radcliffe Hospital, Oxford

14      OX3 9DU, UK

15

16

17  **Corresponding Author:**

18  David Meder: davidm@drcmr.dk

## 19 Summary

20 Decisions are based on value expectations derived from experience. We show

21 that dorsal anterior cingulate cortex and three other brain regions hold multiple

22 representations of choice value based on different time-scales of experience

23 organized in terms of systematic gradients across the cortex. Some parts of each

24 area represent value estimates based on recent reward experience while others

25 represent value estimates based on experience over the longer term. The value

26 estimates within these four brain areas interact with one another according to

27 their temporal scaling. Some aspects of the representations change dynamically

28 as the environment changes. The spectrum of value estimates may act as a

29 flexible selection mechanism for combining experience-derived value

30 information with other aspects of value to allow flexible and adaptive decisions

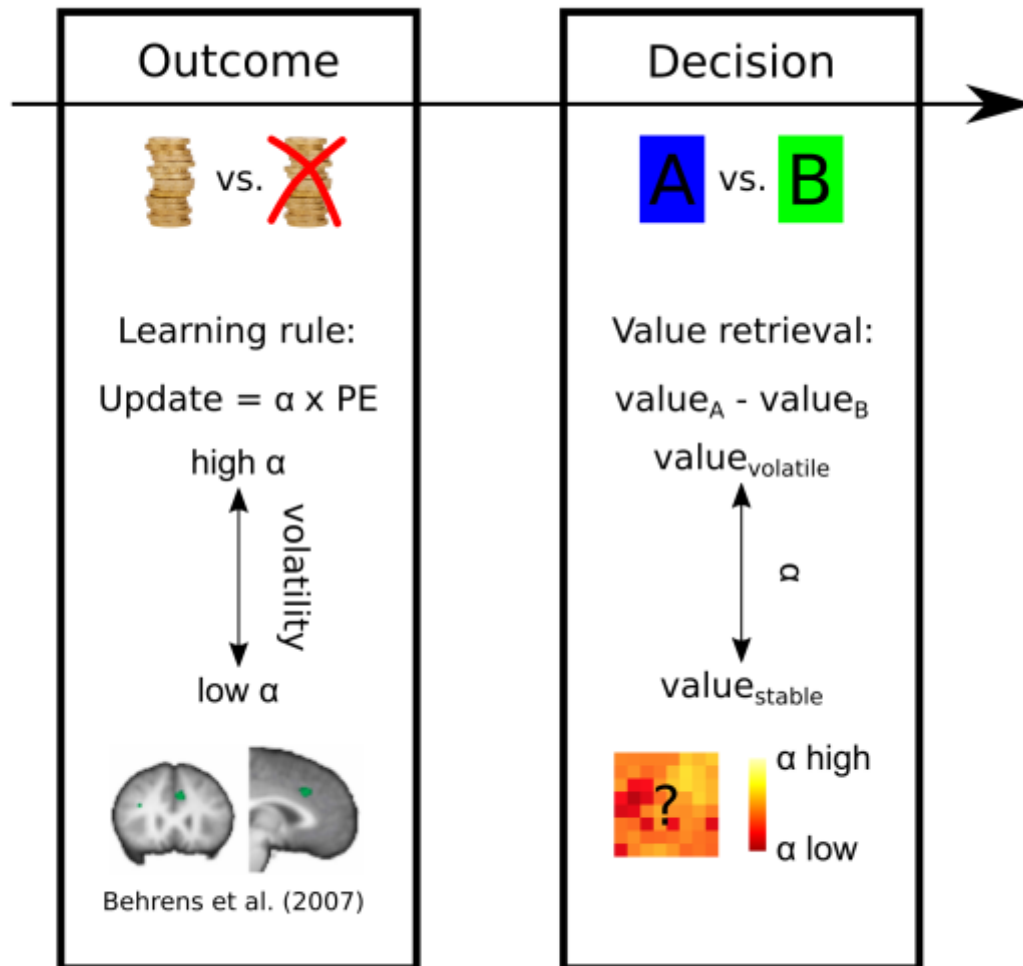31 in changing environments.

32

## Introduction

When an organism makes a decision, it is guided by expectations about the values of potential choices. Estimates of value are, in turn, often dependent on past experience. How past experience should be used when deriving value estimates to guide decisions is not, however, always clear. While it might seem ideal to use the most experience possible, from both the recent and more distant past, this is only true if the environment is stable. In a changing environment it may be better to rely only on most recent experience because earlier experience is no longer informative[1,2].

Previous studies have focused on value learning: how value estimates are updated after the choice is made and the choice outcome is witnessed[1,2]. These studies have emphasized that each outcome has a greater impact on value estimates when the environment is changeable or volatile; the learning rate (LR) is higher and so value estimates are updated more after each choice outcome. Similarly, each outcome has a greater effect on activity in brain areas such as dorsal anterior cingulate cortex (dACC) when the environment is volatile (Fig. 1).

However, while volatility affected dACC at the time of each decision-outcome, there was no evidence that it affected average dACC activity at the time of the next decision. It is therefore unclear how dACC activity might change as a function of the learning rate determining the choice value estimates that guide decision making at the point in time when decisions are actually made (Fig.1). This is this question that we address here. Rather than investigating dACC activity at the time of *decision outcomes* and in relation to learning we focus

56 instead on how dACC represents value estimates employed at the time of

57 *decision making.*

58



59

60 **Fig. 1. When outcomes of decisions are witnessed, average activity in dACC reflects**

61 **the environment's volatility. Under high volatility, the options' values are updated**

62 **with a high learning rate $a$. However, at the time of the actual decision on the next**

63 **trial, volatility no longer exerts a significant effect on average dACC activity.**

64 **However, the representation of choice value estimates necessary for decision-**

65 **making might be represented in some other way such as an anatomically**

66 **distributed pattern of activity.**

67

4

68      When making decisions, the brain might first attempt to determine the

69      best suited LR for the given environment and then calculate a single value

70      estimate based only on this LR.  If this is the case then there may be no overall

71      change in average dACC activity but variance in dACC might best be explained by

72      value estimates calculated at the best LR rather than other inappropriate LRs.

73      Alternatively dACC might hold simultaneous representations of value estimates

74      based on a broad spectrum of LRs. Although intuitively the former might seem

75      computationally simpler, there is evidence that neurons in macaque dACC reflect

76      recent reward experience with different time constants as might be expected if

77      they were each employing a different LR[3–5]. However, the role of such neurons in

78      behavior remains unclear. Here we sought evidence for the existence of value

79      estimates in dACC and elsewhere in the human brain, based on experience over

80      different time scales (and therefore employing different LRs), and examined how

81      such representations mediate decision making (Fig. 1).

82      We developed a new approach to analyse neural data going beyond the

83      typical use of computational models in investigation of brain behavior

84      relationships. Typically, the free parameters of a computational model (e.g. LR)

85      are fitted to the behavior of the subject from which trial-wise estimates of the

86      computed variables can be extracted (e.g. value estimates). However, here we

87      also test whether neuronal populations exist with responses that are better

88      characterised by parts of parameter space that are not overtly expressed in

89      current behavior. Identification of such representations is precluded by focusing

90      exclusively on the parameters currently expressed in behavior. Here we take the

91      approach of fitting LR values to each voxel independently, visualising those

92      parameters over anatomical space and computing their interactions. Instead of

93 investigating where in the brain clusters of voxels express similar neural activity

94 related to value estimates, here we examine the range of value estimates across

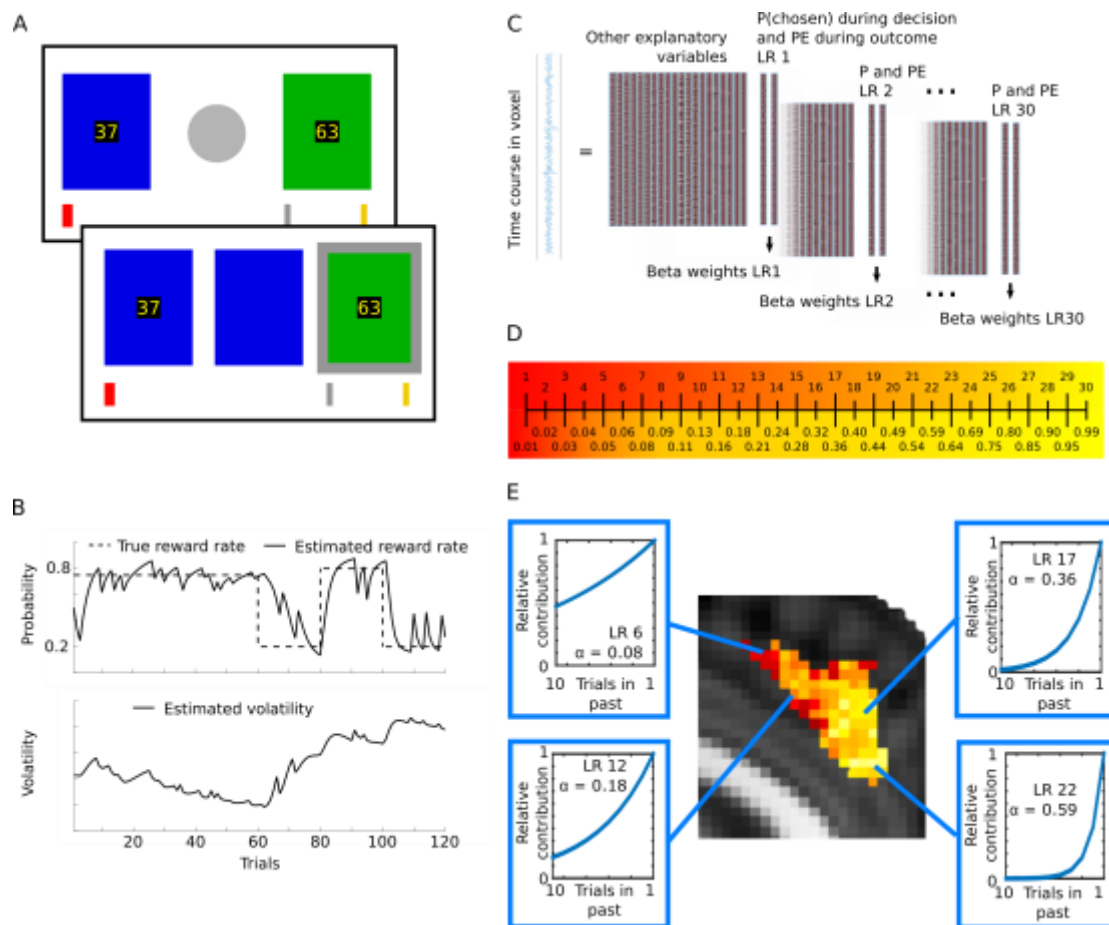95 voxels. We also examine changes to this pattern as a function of volatility.

## Results

### Experimental Strategy

98 We used fMRI data from 17 subjects acquired during a probabilistic reversal

99 learning task[1]. Subjects repeatedly chose between two stimuli with visible

100 reward magnitudes and hidden reward probabilities that had to be learned

101 through feedback (Fig. 2A). Thus in this experiment subjects had to use past

102 experience to estimate reward probabilities for each choice. Accordingly, reward

103 magnitude estimates should be based on the stimuli displayed on each trial but

104 the reward probability estimates should depend on recent experience over

105 several trials. The reward probability might be estimated with different LRs

106 depending on how quickly the environment is changing[1]. Each choice's value can

107 then be derived by combining the explicit reward magnitude with the estimated

108 probability of receiving the reward. Each session comprised two sub-sessions

109 (order counterbalanced across subjects): one where reward probabilities

110 remained stable and another sub-session where reward probabilities were

111 volatile (Fig. 2B). The transition between the two sub-sessions was not

112 announced to the subject.

113  In order to investigate whether the human brain represents multiple

114 reward probability estimates that are based on a spectrum of LRs, we used a

115 novel approach to analyse fMRI data. In addition to other regressors modelling

6

116  standard variables of interest (such as the reward magnitudes displayed to

117  subjects on the screen, the reward received, etc) and physiological noise, we

118  added two regressors, one modelling the estimated reward probability of the

119  chosen option during the decision phase, another one modelling the prediction

120  error during the outcome phase. We repeated this entire analysis 30 times for

121  probability estimates and prediction errors based on 30 different LRs ranging

122  from 0.01 to 0.99 (slow to fast LRs), deriving the best-fitting LR for every voxel

123  (Fig. 2C, D, E). In other words, the 30 repetitions of the analysis make it possible

124  to derive 30 different estimates of the reward probability based on 30 different

125  LRs.  The 30 different LRs were chosen so as to sample the entire LR space

126  between 0.01 (almost no learning) and 0.99 (almost complete revision of value

127  estimates on each trial) and to be equally spaced in terms of their correlation to

128  the neighbouring regressors (Fig. 2D; Methods).  In the previous study Behrens

129  et al.[1] assumed one dynamic, but unitary LR generating value estimates across

130  the brain. However, assigning a best-fitting LR to each voxel based on its own

131  data reveals a pattern of diverse value estimates based on different time periods

132  of experience (different LRs). The best-fitting LR of a voxel corresponds to the

133  value regressor calculated with an LR that explained most of the variance in the

134  voxel's time-course, compared to the other LR regressors, regardless of how

135  much variance it actually explains.  While such an approach is unlikely to capture

136  the full range of factors affecting activity in a voxel it has the potential to identify

137  relationships between brain activity and choice value estimates that cannot be

138  captured with standard analysis techniques.

139



140

**Fig. 2. Methods and analysis. (A) Probabilistic reversal learning task. Subjects had to choose between a green and a blue stimulus with different reward magnitudes (displayed at the centre of each stimulus). In addition to the reward magnitude, which changed randomly from trial to trial, the value of each stimulus was determined by the probability of reward associated with each stimulus which drifted during the course of the experiment and had to be learned from feedback. After choice (here: green on second panel), the red bar moved from left to right if the chosen option was rewarded. Subjects tried to reach the silver bar for £10 and the gold bar for £20. (B) Example of reward probability schedule and estimated volatility of the reward probability from a Bayesian learner when the stable phase came first[1]. Each session had a stable phase of 60 trials where one stimulus was rewarded 75% of trials, the other 25%, and a volatile phase with reward probabilities of 80% vs. 20%, swapping every 20 trials. The order was counterbalanced between subjects. (C) Analysis. As in a conventional fMRI analysis, the blood-oxygen-level-dependent (BOLD) signal time course in every voxel was analysed in a GLM with a design matrix containing relevant regressors. Additionally, one of the regressors modelled a key component of choice value, the estimated reward probability of the chosen option during the decision phase, another one the prediction error during the outcome phase. The same LRs used when deriving the reward probability estimates were used also for the prediction error regressors (the reward probability and prediction error regressors are referred to collectively as LR regressors). This analysis was repeated 30 times, deriving the beta-values for probability estimates and prediction errors based on 30 different LRs. Thus 30 different estimates of the reward probability based on**
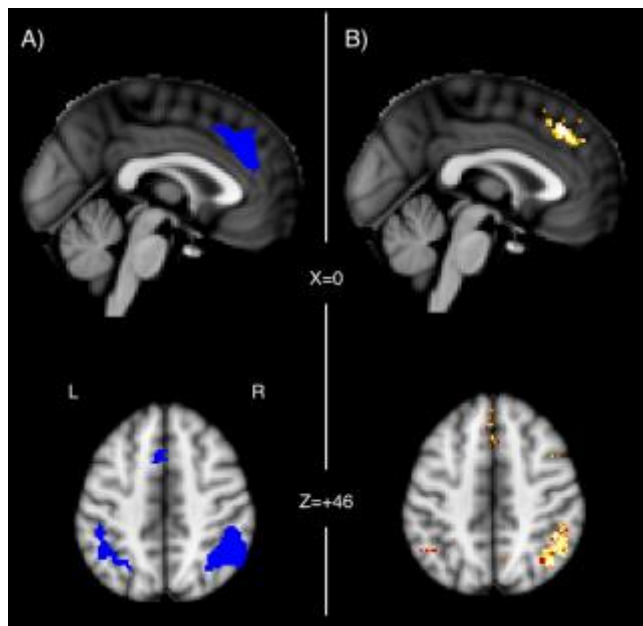
165 **30 different LRs were tested for their ability to explain BOLD signal variance. (D)**
166 **With equal distance separating LRs across the LR spectrum [0.01 to 0.99] the**
167 **regressors would be more strongly correlated at higher LRs, therefore we derived**
168 **30 LRs with larger intervals between higher LRs, resulting in uniform correlation**
169 **across the spectrum. (E) In an environment with high volatility, the stimulus-**
170 **reward history should be more steeply discounted (corresponding to a higher LR)**
171 **than in a stable environment because information from many trials ago is likely to**
172 **be outdated. The plots in the blue boxes show the relative contribution of the**
173 **previous trials' outcomes to the current reward probability estimation with**
174 **different LRs. We thus derived the best-fitting LR for every voxel in every subject,**
175 **averaging across the group. For example, within dACC the BOLD signal in some**
176 **voxels is best explained by a low LR (red) while in others it is best explained by a**
177 **high LR (yellow).**
178

179     We combined two approaches to define the brain areas that we

180 investigated in detail. First, *a priori* we anatomically defined two regions of

181 interest (ROIs) that are known to play important roles in decision-making:

182 dACC[1,6–13], and the inferior parietal lobule (IPL)[14–16]. The anatomical masks for

183 dACC and IPL were taken from connectivity-based parcellation atlases[17,18].

184 Subsequently, we checked that these regions were task-relevant by looking for

185 activity that was associated significantly with the reward magnitude of the

186 choice taken and constrained the ROIs by the conjunction of the anatomy and

187 task-relevant activity (Fig. 3A).

188     In order to confirm that the voxels in our ROIs reflected activity that was

189 related to probability estimates, we ran a singular value decomposition (SVD)

190 over the LR regressors (before HRF-convolution, normalisation and high-pass

191 filtering) to derive singular values capturing most of the variance associated with

192 the LR regressors. For every voxel we then derived the Akaike Information

193 Criterion (AIC) scores from our main GLM (in the absence of any LR regressors).

194 This reveals how well a model lacking multiple LRs accounts for activity

195 variation in every voxel in the brain. We also ran an identical GLM that contained

196 the same regressors but also the first three principle components from the SVD

197  (HRF-convolved, demeaned and high-pass filtered), and again computed the AIC

198  score. This reveals how well a model containing LR-based reward probability

199  estimates accounts for activity variation in every voxel in the brain. We then

200  compared the AIC scores of the two models of brain activity at every voxel using

201  random-effects Bayesian model comparison for group studies[19]. This procedure

202  returned protected exceedance probabilities for every voxel, revealing the

203  degree to which the model containing the singular values, reflecting value

204  estimates based on one or multiple LRs, was the more likely model of the neural

205  data (Fig. 3B). For voxels with a high exceedance probability we can state that

206  LRs have an impact on activity. Having established initial candidate areas of

207  interest in an unbiased way we then went on in subsequent analyses to establish

208  more specifically *how* reward probability estimates based on different LRs were

209  represented.

210



211
212  **Fig. 3. Regions of Interest. (A) dACC and IPL regions defined by conjunction of 1)**
213  **anatomical masks for dACC and IPL from the connectivity-based parcellation**
214  **atlases (http://www.rbmars.dds.nl/CBPatlases.htm)[17,18] and 2) significantly**
215  **decreasing activity (blue) associated with the magnitude of the chosen option**

216 **during decision (B) The dACC and IPL region showed high evidence for coding LRs**
217 **(posterior exceedance probability > 0.95).**
218

219     The relevance of the dACC and IPL regions that we had defined *a priori*

220 based on anatomy was confirmed: these ROIs showed high evidence of coding

221 reward probability estimates based on LRs. Accordingly, for subsequent analyses

222 we constrained the ROI masks to those voxels that fulfilled both the anatomical

223 and task-relevant exceedance probability criteria. We found two further clusters

224 with high evidence in the right frontal operculum (rFO) and bilateral lateral

225 frontopolar cortex (FPl) (Fig. S1A). We focus on reporting results for our primary

226 regions of interest, dACC and IPL, but in the supplemental information we show

227 related results for rFO and FPl. Using a different model, with an additional

228 regressor coding the outcome of the trial (win or loss), the evidence in favour of

229 an LR-based model in these regions was even stronger (Supplemental Material 2

230 and Fig. S1B). This finding is consistent with several other demonstrations that

231 value representations in dACC guide stay/switch or engage/explore decisions of

232 the sort that might be used to perform the current task in humans[9,20–24] and

233 other primates[25,26].

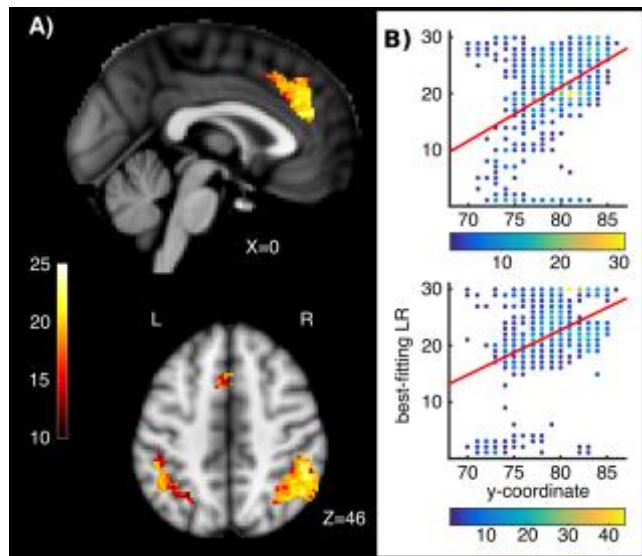234 ## Diversity and Topography of Value Representation

235 The high exceedance probabilities in dACC and IPL reveal that LRs have an

236 impact on activity in these regions, but not whether different voxels represent

237 probability estimates based on different LRs and whether there is any

238 topographic structure in such a representation. Using our multivariate mapping

239 approach, we found that in our ROIs, voxels did not homogeneously integrate the

240 reward history with the same LR, but that there was some degree of spatial

241    topographic organization of the diverse probability estimates (Fig. 4). In both IPL

242    and dACC, a significant amount of variability in the best-fitting LRs in voxels was

243    explained by the x, y, and z coordinates of the voxel when regression models

244    were fitted to each subject's data (t-test over the variance explained by every

245    subject's regression model ($r^2$) against the mean $r^2$ of 10,000 regression models

246    with randomly permuted coordinates. dACC: Mean $r^2$ true data = 0.101, mean $r^2$

247    permuted data = 0.002, $t_{16}$ = 5.071, p < 0.001, IPL right hemisphere: Mean $r^2$ true

248    data = 0.124, mean $r^2$ permuted data = 0.003, $t_{16}$ = 5.566, p < 0.001, IPL left

249    hemisphere: Mean $r^2$ true data = 0.182, mean $r^2$ permuted data = 0.006, $t_{16}$ =

250    5.040, p < 0.001).   The principle axis of anatomical organization in dACC in

251    humans and other primates is approximately rostrocaudally oriented[18,27].

252    Although this axis does not fully correspond to the cardinal axes in the standard

253    space for illustrating neuroimaging data (Montreal Neurological Institute [MNI]

254    space) we nevertheless examined whether LRs were also organized along the

255    MNI y-axis.   Consistently, across subjects, in the dACC, LRs showed a gradient

256    along the MNI y-axis with increasing LRs in the rostral direction (t-test of

257    subjects' regression coefficients of the y-coordinate regressor against 0, $t_{16}$ =

258    2.175, p = 0.045). No major direction of anatomical organization has been

259    reported for the IPL.

260        Previous studies have suggested that some brain regions may reflect a

261    particular time scale of experience or LR that is appropriate to its function[28] but

262    our analysis suggests dACC and IPL are, in addition, representing a spectrum of

263    different LRs. Other relatively abstract features, such as numerosity are known

264    to be represented topographically even though such representations do not map

265    onto sensory receptors or motor effectors in any simple manner[29]. The

12

distribution of LRs in dACC might approximately be related to the rostral-to-

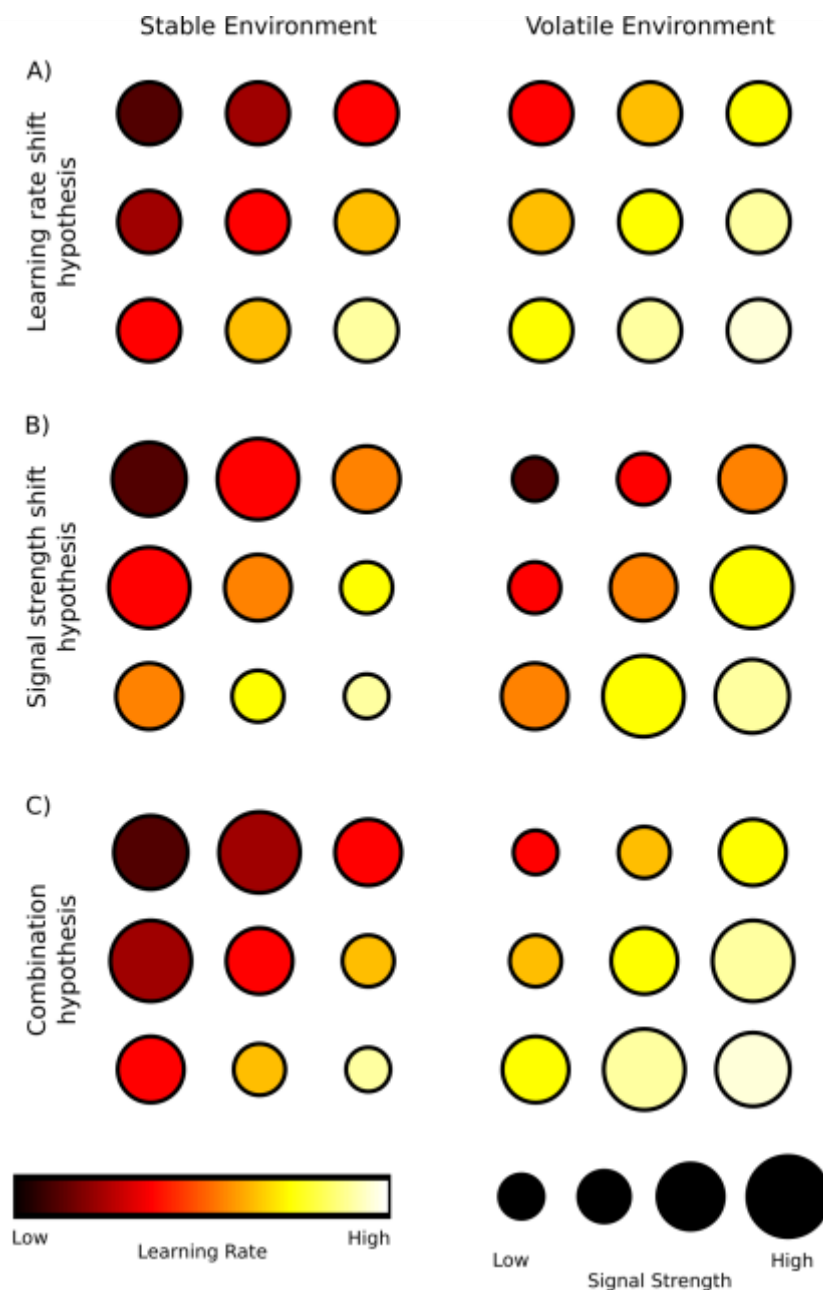caudal gradient in its connectivity with limbic versus motor areas[30].



**Fig. 4. Topographic maps of LRs. (A) A topography of diverse estimates of the reward probability based on different LRs exists in the ROIs. Bright yellow and white colors indicate voxels with high LRs while darker, redder voxels indicate voxels with lower LRs. The color bar on the left indicates the set of LRs (high LRs at top, low LRs at bottom) chosen in 30 steps to minimize correlation between regressors in LR space (see also figure 2d). (B) Spatial gradient along the rostro-caudal axis in dACC in two example subjects. Each voxel's best-fitting LR is plotted against its position on the y-coordinate. The color of the dots reflects the number of voxels having a given combination of values (see color bars beneath graph). Red lines: Regression of all voxels' best-fitting LR against their y-coordinate.**

## Mechanisms of Adaption to Changes in the Environment

As already explained, in a volatile environment, ideally decisions should be based

on probability estimates derived from voxels with higher LRs, while in a stable

environment, voxels with lower LRs might inform the decision. This suggests

that one of two changes to the representation might occur as volatility of the

reward environment changed. First, voxels might have dynamically changing

LRs, depending on the environment (Fig. 5A). Alternatively, each voxel might

retain its best-fitting LR regardless of volatility but the degree to which variance

13

289 in each voxel's activity was explained by reward probability estimates with the

290 best-fitting LR might get stronger in high LR voxels in volatile environment (or

291 stronger in low LR voxels in stable environments). In other words, the regressor

292 effect size (beta-weight) in high LR and low LR voxels might increase and

293 decrease in volatile and stable environments respectively (Fig. 5B). To probe

294 these hypotheses, we split the BOLD signal time course into stable and volatile

295 sub-sessions and again identified the best-fitting LR for every voxel in each of the

296 two sub-sessions. We then compared the best-fitting LR in each sub-session in
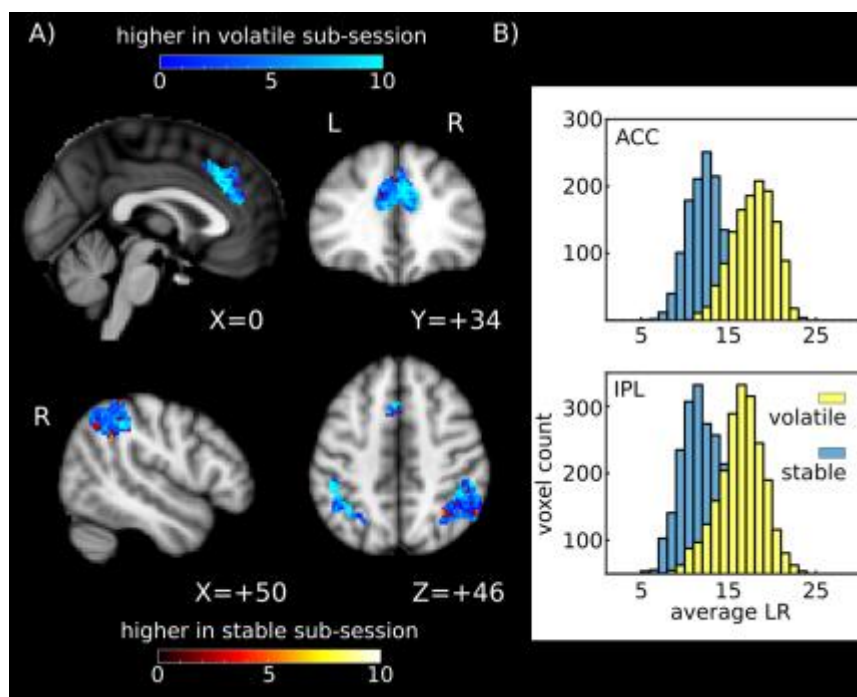
297 every voxel.

14

Fig. 5. Schematic figure depicting possible ways in which multiple value estimates, based on different periods of experience determined by different LRs, might be represented in the brain as indexed by fMRI. We consider how such representations might change as the environment's volatility changes. Each row shows the representation of value estimates in nine example voxels in a stable and in a volatile environment. A) According to the LR shift hypothesis, in a stable environment neurons in more voxels would compute value estimates based on lower LRs while they would shift towards higher LRs in a volatile environment. B) The signal strength shift hypothesis predicts that the value estimates computed by the neurons of each voxel remain constant in all environments, but that those voxels with value estimates that are currently most relevant for the environment (high LR voxels in volatile environments and low LR voxels in stable environments) increase their signal strength. C) The combination hypothesis suggests a combination of the two mechanisms in A) and B).

15

313

314        In the dACC and IPL, the LRs of the voxels' probability estimates were

315    approximately normally distributed (Lilliefors test: dACC p=0.363; IPL p=0.950)

316    but they had significantly higher LRs in the volatile compared to the stable sub-

317    session (average LR difference in dACC: 5.36 [details of LR scaling are shown in

318    Fig. 2D], t-test of each subject's mean change in LR's against 0: $t_{16}$ = 3.68,

319    p=0.002, average LR difference in IPL: 4.34, $t_{16}$ = 2.58, p=0.020) (Fig. 6). This

320    finding suggests an adaptation mechanism resembling the one outlined in the

321    shift-hypothesis (Fig. 5A). However, there might also be a change in how much of

322    the neural activity in a voxel can be explained by the best-fitting LR. This would

323    constitute a change in the effect size or beta-weight of the best fitting regressor
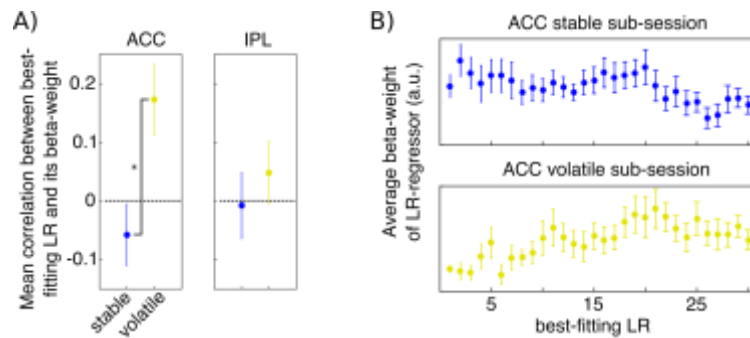
324    (Fig. 5B,C).

325



326

**Fig. 6. Dynamic changes in LR between stable and volatile sub-session. A) Change in LR in every voxel between stable and volatile sub-session. Values on the color bars show the change in LR rank. B) Distribution of number of voxels with best-fitting LRs in the two regions of interest.**

16

331

332        We therefore tested whether there was a dynamic change in the effect

333    sizes of the best-fitting LRs depending on which LRs were currently behaviorally

334    relevant. If such a boosting of relevant LR signals exists, then we would expect

335    voxels with lower best-fitting LRs to have higher beta-weights in the stable sub-

336    session (a negative correlation between best-fitting LR and beta-weight) and

337    voxels with higher best-fitting LRs to having the higher beta-values in the

338    volatile session (positive correlation between best-fitting LR and beta-weight).

339    We calculated the correlation between best-fitting LR and beta-weights for every

340    subject in the two sub-sessions and transformed the correlation coefficients to z-

341    scores (Fisher transformation). In the dACC, there was indeed such a dynamic

342    change in effect size (mean difference in z-scores stable minus volatile sub-

343    session -0.230, $t_{16}$ = -3.802, p = 0.002), while this was not the case for the IPL

344    (mean difference -0.056, $t_{16}$ = -0.818, p = 0.425.) (Fig. 7). This shows that in the

345    dACC, there is a combined adaptation of both the best-fitting LRs in voxels and a

346    change in the effect size of the best-fitting LR, depending on the behavioral

347    relevance of the best-fitting LR in a given environment (Fig. 5C). Thus, voxels

348    change so as to code LRs appropriate for the current environment and they

349    change so as to encode appropriate LRs more strongly than inappropriate LRs. In

350    the IPL, however, only the former adaptation to the environment seems to take

351    place (Fig. 5A).

352

353



354 **Fig. 7. Change in the correlation between beta-weights of the best-fitting LR**
355 **regressors and the best-fitting LR between sub-sessions. A) In the dACC, the**
356 **correlation was significantly positive for the volatile sub-session and significantly**
357 **different from the negative correlation seen in the stable phase. B) Average beta-**
358 **weights across the whole spectrum of LRs in stable and volatile sub-session in the**
359 **dACC.**

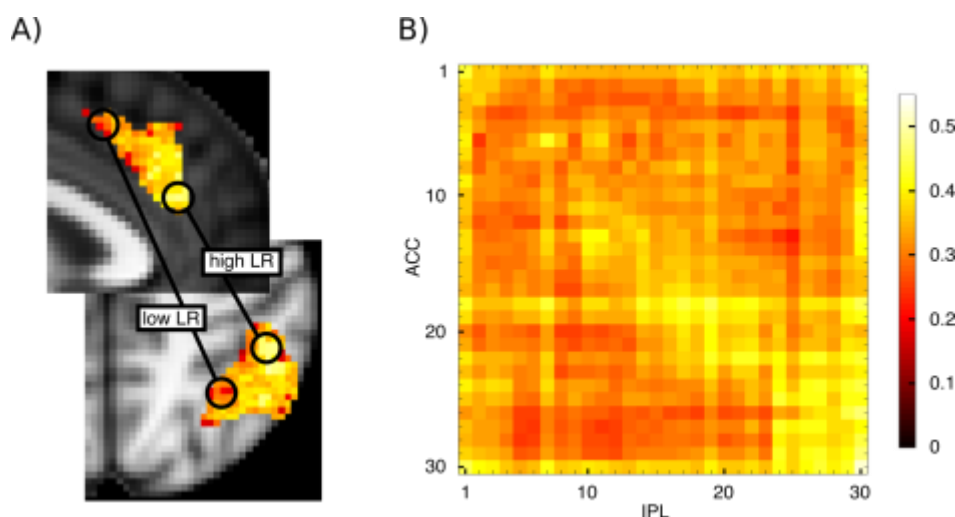360 ## LRs as Organizational Principle of Interregional Interaction

361 So far we have seen that four brain regions carry multiple estimates of the value

362 of choices that are based on different time constants of experience

363 corresponding to different LRs. Thus, multiple LRs constitute an organizing

364 principle determining distribution of activity patterns within these areas. We

365 therefore next asked whether multiple LRs exerted a similar influence over the

366 manner in which the areas interacted with one another. In other words, do

367 voxels that code recent reward probability experience with a small time constant

368 (high LR) in one brain region (e.g. dACC) interact preferentially with voxels with

369 high LRs elsewhere? Similarly, are low LR voxels in different brain areas

370 preferentially interacting with one another?

371 For every subject, we extracted the mean residual BOLD time course for

372 all voxels after regressing out all the information contained in our original design

373 matrix (coding, for example, for the various task events) and additionally all 30

374 LR regressors indexing the estimated reward probability in the decision phase

375 and all 30 LR regressors indexing prediction error in the outcome phase. Thus,

376    the residual time course no longer contained any LR related information. We

377    then created a mean residual time course for all voxels originally identified as

378    being of the same LR within each ROI and correlated these 30 mean residual

379    time courses with the 30 mean residual time courses of another region. We

380    found that the more similar the best-fitting LRs, the higher was the correlation of

381    these voxels' residual time courses between the dACC and the IPL, as reflected in

382    higher average correlation values along the diagonal (Fig. 8). For example, voxels

383    with high LRs in the dACC were more correlated with high-LR voxels compared

384    to low-LR voxels in the IPL (Fig. 8; bright yellow diagonal line running from top

385    left to bottom right).

386          The statistical test for demonstrating the significance of the effect is best

387    understood with reference to figure 8. It is to examine whether the subjects' z-

388    transformed correlation coefficients are correlated positively with their

389    closeness to the diagonal; this was indeed the case (negative Euclidian distance,

390    one-tailed t-test of z-transformed correlation values $t_{16}$ = -2.944, p = 0.005); the

391    correlation between the brain areas' signals became greater the more that the

392    signals were drawn from voxels with similar LRs.

393

394

395 **Fig. 8. LR topography as an organizing principle for interaction between regions.**

396 **A) We investigated whether voxels that represent choice values with similar LRs**

397 **also show stronger connectivity between regions. B) Correlation plot depicting**

398 **the correlation of the residual BOLD time course averaged over all voxels with the**

399 **same best-fitting LR within dACC with the residual BOLD time course over all**

400 **voxels with the same best-fitting LR within IPL, averaged over all subjects. The**

401 **subjects' z-transformed correlation coefficients were correlated positively with**

402 **their closeness to the diagonal.**

403

404 In summary, even after removing all linear task-related information

405 (activity linearly related to task variables and value estimates), voxels with the

406 same best-fitting LR shared a more similar pattern of activity in dACC and IPL.

407 Thus, LRs are not just an organizational feature of individual brain regions but

408 also an organizing principle determining how these regions interact with one

409 another. This feature of interactions between areas was also apparent in all

410 combinations of interactions between all the four regions that showed high

411 evidence for the coding of reward probabilities based on multiple LRs (ACC, IPL,

412 FPl and rFO; Fig. S6, Table S1).

413 ## Ubiquity or Localization of Dynamic Topographic Value Representations

414 We have presented evidence for topographic organization of value estimates as a

415 function of different LRs and shown LRs are an organizational principle of

416 connectivity between regions such as dACC and IPL. We next asked whether such

417 representations and interaction patterns are ubiquitous in all brain areas

418 signalling value. We therefore performed the same analyses in another brain

419 region that has repeatedly been linked to value and decision making, the

20

420     ventromedial prefrontal cortex (vmPFC)[9,14,31–37]. In most studies, the strongest

421     value-related activation was found in the anterior part of the vmPFC. We

422     examined two vmPFC regions: anterior vmPFC and posterior vmPFC

423     (Supplemental Materials 3). We found some, albeit weak, evidence for LR related

424     activity in anterior vmPFC (Fig. S1C). Unlike in dACC and IPL, in vmPFC the

425     amount of BOLD variance explained by SVD-derived singular values reflecting

426     the LR regressors was not significantly greater than the amount of variance

427     explained by a model lacking LR information.  In fact, when the same statistical

428     approaches were used as in our investigation of dACC and IPL we found that

429     activity in many voxels in vmPFC was better explained by a model lacking the LR

430     regressors.  Value estimates with different LRs could be fit to voxels in vmPFC

431     (Fig. S2) but there was no shift in the distribution of LRs depending on the

432     volatility of the environment (Fig. S3, compare to Fig. 6) and there was no change

433     in the correlation between the best-fitting LR and its beta-weight as seen in the

434     dACC (Fig. S3, compare to Fig. 7) in either vmPFC region. Additionally, unlike

435     dACC, IPL, rFO, and FPl, there was no evidence that voxels in either vmPFC

436     region preferentially interacted with voxels with similar LRs in other brain

437     regions (i.e., no diagonal with high correlation values; Supplemental Materials 5;

438     Fig. S5, Table S1, compare to Fig. 8). In general, the average correlation over all

439     voxels between two regions was significantly higher for dACC, IPL, rFO, and FPl

440     than between any of these areas and either vmPFC subdivision (Table S2).

441        In summary, there is only comparatively weak evidence for the vmPFC

442     holding value related information that reflects recent experience of reward

443     probability and the value estimates it held were not as sensitive to

444     environmental volatility. Thus, the neuroanatomical gradients of probability

21

445     estimates calculated with different LRs in dACC and IPL, their sensitivity to

446     environmental volatility, and their inter-regional LR-specific connectivity are not

447     ubiquitous features of all value encoding brain regions. This supports the notion

448     that the spectrum of value estimates based on multiple LRs that we find in some

449     brain regions cannot be attributed to noise over subjects, time, or voxels.

450

451     **LR-based representation at decision outcome**

452     Finally, while the current investigation is focussed on the decision-making

453     process, rather than the outcome monitoring phase of the task, we wanted to

454     know whether we could observe comparable dynamic adaptations to

455     environmental volatility during the outcome phase. We therefore investigated

456     whether prediction error coding in ventral striatum (VS) would also reflect

457     adaptations of which LRs should be expressed as a function of volatility. A model

458     containing the first three singular values from an SVD over the prediction error

459     regressors provided a good model of right VS activity during the outcome phase

460     of the trials (Fig. S6A). However, using a bilateral anatomical mask of the VS

461     (Automated Anatomical Labeling (AAL) atlas[38]), the distributions of the LRs

462     generating the prediction error were stable and did not change between the

463     stable and volatile sub-sessions (Supplemental Materials 6; Fig. S6B). While

464     Behrens et al.[1] found an overall change in dACC activity during outcome, there

465     was no evidence in the current study for a prediction error signal in dACC, using

466     either standard analysis procedure similar to those used before[1] nor based on

467     Bayesian group model comparisons such as those employed here.

## Discussion

A number of cortical regions have been implicated in reward-guided decision making and it is possible that they operate partly in parallel[12,31]. For example, some aspects of decision making behavior are predicted by activity in vmPFC while others, even in the same task and at the same time, are better predicted by activity in the intraparietal sulcus[31].

DACC may be particularly important when deciding whether to switch and change between choices and behavioral strategies[9,10,12,20–26]. A flexible behavioral repertoire would be promoted by having multiple experience dependent value estimates, estimated over different time scales: representations of how well things have been recently and, simultaneously, how well they have been over the longer term. By contrasting the strength of such representations a decision-maker would be able to know whether the value of their environment is stable or improving or whether it is declining and that it might be time to explore elsewhere[24].

In the present study we have found evidence that indeed multiple value representations, with different time constants, are especially prominent in dACC and IPL. A diversity of value estimates based on a spectrum of LRs could either reflect features of the neural representation guiding decision making, or it might simply be a reflection of natural variability over samples, trials, and voxels. Several aspects of our findings suggest that they reflect features of neural activity rather than noise. First, multiple LR-based representations were not ubiquitous; they were prominent in only a subset of regions implicated in value representation and decision making (Supplemental Materials 3-5; Figs. S1-S5).

492 Second, the multiple LR representations were structured; they were

493 topographically organized within areas (Fig. 4) and they were an organizing

494 feature of interaction patterns between areas (Fig. 8). The conclusion that there

495 are multiple LR-based value estimates is derived from averaging data over trials;

496 in the future it might be interesting to examine the nature of these

497 representations on a trial-by-trial basis.

498     While the parallel information processing entailed by such a

499 representation might appear an unnecessary waste of computational resources,

500 it may be advantageous when the volatility of the environment is changing and

501 other LRs generate better value estimates than the one currently employed to

502 guide behaviour. Imagine a decision-maker that has estimated that the current

503 environment is volatile and estimates choice values only on the basis of recent

504 experience (high LR). If the decision-maker realises that actually the

505 environment is more stable than suspected, then it needs to retrieve the

506 outcomes of earlier decisions and reweigh each of them according to the LR that

507 is now optimal for estimating choice values. Our evidence suggests that the brain

508 may compute many values estimates in parallel over different time scales and

509 that such longer term time scale estimates (lower LR estimates) are immediately

510 available for the decision-maker to switch to on realising the true level of

511 environmental volatility. Since these value estimates are derived in a Markov

512 decision process, only the most recent value estimate has to be remembered and

513 updated so that it is not necessary to remember preceding outcomes.

514     The co-existence of multiple experience dependent value estimates guiding

515 decisions is also consistent with the results of single unit recordings made in

516 macaques[3] in a dACC region homologous with the one we investigated here[18].

517 Neurons that varied in the degree to which their activity reflected just recent

518 outcomes or also outcomes in the more distant past were also reported in the

519 intraparietal sulcus and dorsolateral prefrontal cortex[3]. In the present study we

520 also found evidence for such response patterns in fMRI activity in an adjacent

521 part of the parietal cortex (IPL), a very rostral part of prefrontal cortex (FPl), and

522 in FO. By recording activity in individual neurons it is possible to demonstrate

523 precisely how different neurons, even closely situated ones, can code both recent

524 and more distant rewards with different weights. In our study, however, by

525 manipulating the reward environment that subjects experienced in volatile and

526 stable sub-sessions, it was possible to show how such experience dependent

527 reward representations changed with environment and behavior.

528 The evidence for value learning using multiple LRs in several cortical areas

529 fits well with the idea that there exists a hierarchy of information accumulation

530 from short time scales in sensory areas to long time scales in prefrontal, dACC,

531 and parietal association areas[39–43]. In reinforcement learning, information

532 obtained many trials ago in the past can still influence probability estimates

533 when LRs are low. In our task, with an average trial duration of 20s[1], information

534 from several minutes ago has to be remembered. However, we can also show

535 that even within a single area, there are gradients of time scale representation

536 and that these representations are not fixed, but dynamically responding to the

537 environment.

538 In situations in which dACC value representations guide behavior there are

539 often also value-related activations in FPl and IPL[10,11,14,44,45]. Typically, these

540 areas differ from others such as vmPFC in that they encode the value of

541 behavioral change and exploration. In addition, in the present experiment we

542 were able to show that there are links between the value representations in

543 dACC and other brain regions. This suggests that multiple value representations

544 of recent experience constitute an organizing feature of inter-areal interaction. It

545 is not just that average activity throughout one region is related to the average

546 activity of another. Instead parts of dACC employing the fastest and slowest LRs

547 are interacting with corresponding subdivisions of FPl, IPL, and rOP. The pattern

548 of results is suggestive of a distributed representation across multiple brain

549 regions in which the value of initiating and changing behavior is evaluated over

550 multiple time scales simultaneously[46].

551    In a longer behavioral testing session (without fMRI acquisition) it was

552 shown that subjects do adapt their LR in response to changes in the volatility of

553 the environment[1]. The change in best-fitting LRs that we observe between the

554 stable and the volatile sub-session is in accordance with just such a shift in

555 behavior. The exact mechanism by which the broad spectrum of LR parameters

556 present in dACC, concerning many possible choice values estimated at different

557 time scales, is integrated into one eventual decision needs further elucidation.

558    In conclusion, there are multiple experience dependent value estimates with

559 coarse but systematic topographies in dACC and three other regions. Interactions

560 between these regions occur in relation to this pattern of specific time scales.

561 The distributions of value estimates are dynamically adjusted when there are

562 changes in the environment's volatility. Dynamic adjustment based on

563 environmental statistics might be critical for adjusting behavior to a particular

564 LR and for selecting a particular choice on a given trial.

## Experimental Procedures

The behavioral task and scanning procedures have been described in detail before[1]. In the task, subjects were presented with two choice options, a green and a blue rectangle (Fig. 2A). The potential reward magnitudes were presented in the centre of each stimulus while the reward probabilities had to be learned by the subjects. Reward probabilities were changing throughout the experiment. There was a stable sub-session of 60 trials where one of the stimuli was rewarded 75% of trials and the other one 25% and a volatile sub-session where reward probabilities for the stimuli were 80% and 20%, changing every 20 trials. The order of the sub-sessions was counterbalanced between subjects. Reward information was coupled between the stimuli, i.e. the feedback that the chosen stimulus was rewarded also implied that the choice of the other stimulus would not have led to a reward, and *vice versa*. If the chosen stimulus was rewarded, the presented reward magnitude was added to the subjects accumulating points and a red bar at the bottom of the screen increased in proportion to the points acquired. When the red bar reached a vertical silver bar, subjects received £10, if it reached a golden bar, they receive £20 at the end of the experiment. Subjects were presented with the two options for 4-8 s (jittered). When a question-mark appeared, they could signal their choice with a button press. As soon as the button press was registered, subjects had to wait for 4-8 s (jittered) until the rewarded stimulus was presented in the middle. After a jittered inter-trial-interval of 3-7 s, the next trial began. EPI images were acquired at 3 mm$^3$ voxel resolution with a repetition time (TR) of 3.0 s and an echo time (TE) of 30 ms, a flip angle of 87°. The slice angle was set to 15° and a local z-shim was applied

589    around the orbitofrontal cortex in order to reduce signal drop-out[1]. Since the

590    response was self-timed, the experiment's duration was variable. On average,

591    830 volumes (41.5 min) were acquired. A T1 structural image was acquired with

592    an MPRAGE sequence with 1mm$^3$ voxel resolution, a TE of 4.53 ms, an inversion

593    time(TI) of 900 ms and a TR of 2.2 s[1].

594          We used FMRIB's Software Library (FSL)[47] for image pre-processing

595    and the first level data analysis (see Supplemental Materials 1). Subsequent

596    analysis steps relating to the LR regressors were performed with MATLAB

597    (R2015a 8.5.0.197613).

598          The preprocessing was performed on the functional images of the entire

599    session (for the initial analysis), and of the stable and the volatile sub-sessions

600    (for subsequent analyses). In order to analyse the sub-sessions, we split the time

601    series of BOLD data into those portions that were collected when the reward

602    environment was in a stable or volatile sub-session. The data assigned to the first

603    sub-session encompassed all MRI volumes collected up to and including the

604    onset of the last outcome of that sub-session of the experiment plus two

605    additional volumes to account for the delay of the hemodynamic response

606    function.

607          The data were pre-whitened before analysis to account for temporal

608    autocorrelation[48]. For the subsequent mapping of LRs, we ran three GLM's for

609    the whole session, and separately for the stable and the volatile sub-sessions, at

610    the first level for each participant with the following regressors:

611    1) Decision phase main effect (duration: stimuli onset until response)

612    2) Predict phase main effect (duration: response until outcome)

613    3) Outcome monitor phase main effect (duration: 3s)

28

614     4) Parametric modulation of decision phase with reward magnitude of

615        chosen stimulus

616     5) Parametric modulation of decision phase with log of reaction time

617     6) Parametric modulation of decision phase with stay (0) or switch (1)

618        decision

619     7) Parametric modulation of outcome monitor phase with the reward

620        magnitude of the chosen stimulus

621 We also added the temporal derivative of each regressor to the design matrix in

622 order to explain variance related to possible differences in the timing between

623 the assumed and the actual hemodynamic response function (HRF).

624       Since reward magnitudes are changing unpredictably, participants

625 estimate reward probabilities and not action values. Thus, for each subject, we

626 then calculated the probability estimates for each stimulus from a simple

627 reinforcement learning model[49], based on all 99 LRs ($\alpha$) between 0.01 and 0.99.

628 The model estimates the probability of one of the stimuli leading to a reward by

629 updating the stimulus-reward probability $p(a)$ with LR $\alpha$, where R = 1 when the

630 stimulus was rewarded and R = 0 if not:

631

632 $$p(a_i) = p(a_{i-1}) + \alpha[R - p(a_{i-1})]$$

633

634 The probability estimate of the other stimulus $p(B)$ is $1 - p(A)$. From these

635 values, we also calculated the prediction error (PE) corresponding to the

636 outcome of that trial by subtracting the probability estimate of the chosen

637 stimulus from the outcome (1 for rewarded trials, 0 for non-rewarded trials).

638 Thus, the PE is a "probability PE" that is not weighted with the magnitude of the

29

639　(foregone) reward. After normalising the probability estimates for all LRs for

640　both stimuli, we derived the probability estimate of the chosen stimulus

641　p(chosen). These p(chosen)-regressors (hereafter "LR regressors") and the PE

642　regressors were convolved with the HRF, normalised and high-pass filtered in

643　the same way (in the same manner as in FSL). We calculated a correlation matrix

644　for the 99 resulting LR regressors for every subject and for the whole session as

645　well as the two sub-sessions. Since the correlation between regressors is not the

646　same for all levels of LR, we chose 30 regressors that were equally spaced in

647　terms of their correlation to the neighbouring regressors. We did so by averaging

648　the 30 LR regressors with equal correlation for every subject in all three sessions

649　and subsequently rounding them to two decimals. This procedure resulted in 30

650　LR regressors corresponding to the following LRs (see also Fig. 2):

651　[0.01 0.02 0.03 0.04 0.05 0.06 0.07 0.08 0.09 0.11 0.12 0.14 0.15 0.17 0.20 0.22

652　0.25 0.28 0.32 0.36 0.40 0.46 0.51 0.57 0.64 0.71 0.78 0.85 0.93 0.99].

653　　We used the BET procedure[50] on the high-pass filtered and motion corrected

654　functional MRI data to separate brain matter from non-brain matter. For each of

655　the (sub-)sessions in every subject, we explained activity in the filtered fMRI

656　data with 30 separate GLM's, each with the design matrix described above

657　together with one of the 30 LR regressors (onset during the decision phase) and

658　the corresponding PE regressor (onset during outcome monitoring phase).

659　　In each GLM, we retrieved the parameter estimate for the LR regressor and

660　we mapped the following three measures to every voxel in the brain:

661　　1) best-fitting LR: the regressor with the highest beta-value (regression

662　　　　weights indicative of the relationship between the regressor and the

663　　　　BOLD signal) in the GLM. For example, if regressor 20 had the highest

664        beta-values amongst the 30 LR regressors, that voxel would be assigned a

665        LR of 20.

666    2) the change in the best-fitting LR between the stable and the volatile sub-

667        sessions (measured as best-fitting LR in the stable sub-session minus the

668        best-fitting LR in the volatile sub-session).

669    3) the beta-weight of the best-fitting LR regressor in the entire session and

670        in the stable and the volatile sub-sessions

671    The resulting images were registered to MNI-space using the non-linear

672    warping field using nearest-neighbour interpolation. Subsequently, the single-

673    subject images were averaged across all subjects to create group-average

674    images.

675    We also used a standard FSL analysis with a GLM similar to the one above but

676    with two additional regressors corresponding to the probability of the chosen

677    stimulus during the decision phase and during the outcome monitoring phase as

678    derived from a Bayesian learner model[1] as well as a regressor coding the

679    outcome of the trial (won or lost). This analysis was used for retrieving the beta-

680    weight of the magnitude of the chosen option's potential reward of each voxel for

681    the correlation analysis with the best-fitting LR regressor's beta-weight.

682    The magnitude regressor was also used for generating regions of interest

683    (ROIs; Fig. 3). We defined our ROIs by the overlap of the contrast over this

684    regressor (cluster-corrected results with the standard threshold of z=2.3,

685    corrected significance level p=0.05) and anatomical masks derived from the

686    connectivity-based                    parcellation                    atlases[17,18]

687    (http://www.rbmars.dds.nl/CBPatlases.htm) (Fig. 3). For dACC, this included

688    bilateral areas 24a/b, d32 as well as the bilateral anterior rostral zones of the

689    cingulate motor areas. For posterior vmPFC, this included bilateral area 14m and

690    for anterior vmPFC it included 11m[18]. For IPL, this included inferior parietal

691    lobule areas c and d as defined by Mars and colleagues[17]. The atlas only contains

692    IPL regions for the right hemisphere, we therefore mirrored the regions along

693    the midline to create masks for the left hemisphere. Since the anatomical masks

694    are defined by white matter connectivity, they do not cover the entire cortical

695    area. Therefore, the dACC and vmPFC masks were extended with 2 voxels

696    medially, while the IPL masks were extended laterally and caudally to enssure

697    that all grey matter voxels were covered by the masks.

698    **Evidence for variability in voxels' activity related to reinforcement learning**

699    In order to confirm that the voxels in our ROIs actually reflected activity that was

700    related to probability estimates, we ran a singular value decomposition (SVD)

701    over the 99 LR regressors (before HRF-convolution, normalisation and high-pass

702    filtering) to derive singular values capturing most of the variance associated with

703    the variability in the 99 LR regressors. For every voxel we then derived the

704    Akaike Information Criterion (AIC) scores from our main GLM (not containing

705    any LR regressors) as well as from a GLM that contained the first three singular

706    values from the SVD (HRF-convolved, demeaned and high-pass filtered). We then

707    used random-effects Bayesian model comparison for group studies[19] by passing

708    each subject's AIC scores for the two models to the spm_bms matlab function

709    from    SPM12    (http://www.fil.ion.ucl.ac.uk/spm/software/spm12/).    This

710    procedure returned protected exceedance probabilities for every voxel, showing

711    the probability that the model containing the singular values was a more likely

712    model of the data than the model without those components.

## Author Contributions

714    D.M. and L.V. analysed data; T.E.J.B. acquired the data; D.M., K.H.M., O.J.H., N.K.,

715    L.V., M.K.W. and M.F.S.R developed the analysis approach; D.M., N.K., L.V., M.K.W.,

716    K.H.M, O.J.H., T.E.J.B. and M.F.S.R. discussed the results and wrote the manuscript.

## Acknowledgments

721

## References

1. Behrens, T. E. J., Woolrich, M., Walton, M. E. & Rushworth, M. F. S. Learning the value of information in an uncertain world. *Nat. Neurosci.* **10,** 1214–1221 (2007).

2. Nassar, M. R., Wilson, R. C., Heasly, B. & Gold, J. I. An Approximately Bayesian Delta-Rule Model Explains the Dynamics of Belief Updating in a Changing Environment. *J. Neurosci.* **30,** 12366–12378 (2010).

3. Bernacchia, A., Seo, H., Lee, D. & Wang, X.-J. A reservoir of time constants for memory traces in cortical neurons. *Nat. Neurosci.* **14,** 366–372 (2011).

4. Seo, H. & Lee, D. Cortical mechanisms for reinforcement learning in competitive games. *Philos. Trans. R. Soc. B Biol. Sci.* **363,** 3845–3857 (2008).

5. Seo, H. & Lee, D. Temporal Filtering of Reward Signals in the Dorsal Anterior Cingulate Cortex during a Mixed-Strategy Game. *J. Neurosci.* **27,** 8366–8377 (2007).

6. Walton, M. E., Devlin, J. & Rushworth, M. F. S. Interactions between decision making and performance monitoring within prefrontal cortex. *Nat. Neurosci.* **7,** 1259–1265 (2004).

7. Kennerley, S. W., Walton, M. E., Behrens, T. E. J., Buckley, M. J. & Rushworth, M. F. S. Optimal decision making and the anterior cingulate cortex. *Nat. Neurosci.* **9,** 940–947 (2006).

8. Kennerley, S. W., Dahmubed, A. F., Lara, A. H. & Wallis, J. D. Neurons in the Frontal Lobe Encode the Value of Multiple Decision Variables. *J. Cogn. Neurosci.* **21,** 1162–1178 (2009).

34

745   9.  Kolling, N., Behrens, T. E. J., Mars, R. B. & Rushworth, M. F. S. Neural

746       Mechanisms of Foraging. *Science* **336,** 95–98 (2012).

747   10. Kolling, N., Wittmann, M. & Rushworth, M. F. S. Multiple Neural Mechanisms

748       of Decision Making and Their Competition under Changing Risk Pressure.

749       *Neuron* **81,** 1190–1202 (2014).

750   11. Scholl, J. *et al.* The Good, the Bad, and the Irrelevant: Neural Mechanisms of

751       Learning Real and Hypothetical Rewards and Effort. *J. Neurosci.* **35,** 11233–

752       11251 (2015).

753   12. Rushworth, M. F., Kolling, N., Sallet, J. & Mars, R. B. Valuation and decision-

754       making in frontal cortex: one or many serial or parallel systems? *Curr. Opin.*

755       *Neurobiol.* **22,** 946–955 (2012).

756   13. Hunt, L. T., Behrens, T. E., Hosokawa, T., Wallis, J. D. & Kennerley, S. W.

757       Capturing the temporal evolution of choice across prefrontal cortex. *eLife* **4,**

758       e11945 (2015).

759   14. Boorman, E. D., Behrens, T. E. J., Woolrich, M. W. & Rushworth, M. F. S. How

760       Green Is the Grass on the Other Side? Frontopolar Cortex and the Evidence in

761       Favor of Alternative Courses of Action. *Neuron* **62,** 733–743 (2009).

762   15. Waskom, M. L., Kumaran, D., Gordon, A. M., Rissman, J. & Wagner, A. D.

763       Frontoparietal Representations of Task Context Support the Flexible Control

764       of Goal-Directed Cognition. *J. Neurosci.* **34,** 10743–10755 (2014).

765   16. Medic, N. *et al.* Dopamine Modulates the Neural Representation of Subjective

766       Value of Food in Hungry Subjects. *J. Neurosci.* **34,** 16856–16864 (2014).

767   17. Mars, R. B. *et al.* Diffusion-Weighted Imaging Tractography-Based

768       Parcellation of the Human Parietal Cortex and Comparison with Human and

769    Macaque Resting-State Functional Connectivity. *J. Neurosci.* **31,** 4087–4100

770    (2011).

771    18. Neubert, F.-X., Mars, R. B., Sallet, J. & Rushworth, M. F. S. Connectivity reveals

772    relationship of brain areas for reward-guided learning and decision making

773    in human and monkey frontal cortex. *Proc. Natl. Acad. Sci.* **112,** E2695–E2704

774    (2015).

775    19. Rigoux, L., Stephan, K. E., Friston, K. J. & Daunizeau, J. Bayesian model

776    selection for group studies — Revisited. *NeuroImage* **84,** 971–985 (2014).

777    20. Kolling, N., Behrens, T., Wittmann, M. & Rushworth, M. Multiple signals in

778    anterior cingulate cortex. *Curr. Opin. Neurobiol.* **37,** 36–43 (2016).

779    21. Meder, D. *et al.* Tuning the Brake While Raising the Stake: Network Dynamics

780    during Sequential Decision-Making. *J. Neurosci.* **36,** 5417–5426 (2016).

781    22. Rudebeck, P. H. *et al.* Frontal Cortex Subregions Play Distinct Roles in Choices

782    between Actions and Stimuli. *J. Neurosci.* **28,** 13775–13785 (2008).

783    23. Kolling, N. *et al.* Value, search, persistence and model updating in anterior

784    cingulate cortex. *Nat. Neurosci.* **19,** 1280–1285 (2016).

785    24. Wittmann, M. K. *et al.* Predictive decision making driven by multiple time-

786    linked reward representations in the anterior cingulate cortex. *Nat. Commun.*

787    **7,** 12327 (2016).

788    25. Quilodran, R., Rothé, M. & Procyk, E. Behavioral Shifts and Action Valuation in

789    the Anterior Cingulate Cortex. *Neuron* **57,** 314–325 (2008).

790    26. Stoll, F. M., Fontanier, V. & Procyk, E. Specific frontal neural dynamics

791    contribute to decisions to check. *Nat. Commun.* **7,** 11990 (2016).

792    27. Procyk, E. *et al.* Midcingulate Motor Map and Feedback Detection: Converging

793    Data from Humans and Monkeys. *Cereb. Cortex* **26,** 467–476 (2016).

794    28. Gläscher, J. & Büchel, C. Formal Learning Theory Dissociates Brain Regions

795        with Different Temporal Integration. *Neuron* **47,** 295–306 (2005).

796    29. Harvey, B. M., Klein, B. P., Petridou, N. & Dumoulin, S. O. Topographic

797        Representation of Numerosity in the Human Parietal Cortex. *Science* **341,**

798        1123–1126 (2013).

799    30. Kunishio, K. & Haber, S. N. Primate cingulostriatal projection: Limbic striatal

800        versus sensorimotor striatal input. *J. Comp. Neurol.* **350,** 337–356 (1994).

801    31. Chau, B. K. H., Kolling, N., Hunt, L. T., Walton, M. E. & Rushworth, M. F. S. A

802        neural mechanism underlying failure of optimal choice with multiple

803        alternatives. *Nat. Neurosci.* **17,** 463–470 (2014).

804    32. Economides, M., Guitart-Masip, M., Kurth-Nelson, Z. & Dolan, R. J. Anterior

805        Cingulate Cortex Instigates Adaptive Switches in Choice by Integrating

806        Immediate and Delayed Components of Value in Ventromedial Prefrontal

807        Cortex. *J. Neurosci.* **34,** 3340–3349 (2014).

808    33. Hunt, L. T. *et al.* Mechanisms underlying cortical activity during value-guided

809        choice. *Nat. Neurosci.* **15,** 470–476 (2012).

810    34. Jocham, G. *et al.* Dissociable contributions of ventromedial prefrontal and

811        posterior parietal cortex to value-guided choice. *NeuroImage* **100,** 498–506

812        (2014).

813    35. Jocham, G., Hunt, L. T., Near, J. & Behrens, T. E. A mechanism for value-guided

814        choice based on the excitation-inhibition balance in prefrontal cortex. *Nat.*

815        *Neurosci.* **15,** 960–961 (2012).

816    36. Noonan, M. P. *et al.* Separate value comparison and learning mechanisms in

817        macaque medial and lateral orbitofrontal cortex. *Proc. Natl. Acad. Sci.* **107,**

818        20547–20552 (2010).

819     37. Rushworth, M. F. S., Noonan, M. P., Boorman, E. D., Walton, M. E. & Behrens, T.

820         E. Frontal Cortex and Reward-Guided Learning and Decision-Making. *Neuron*

821         **70,** 1054–1069 (2011).

822     38. Tzourio-Mazoyer, N. *et al.* Automated Anatomical Labeling of Activations in

823         SPM Using a Macroscopic Anatomical Parcellation of the MNI MRI Single-

824         Subject Brain. *NeuroImage* **15,** 273–289 (2002).

825     39. Chaudhuri, R., Knoblauch, K., Gariel, M.-A., Kennedy, H. & Wang, X.-J. A Large-

826         Scale Circuit Mechanism for Hierarchical Dynamical Processing in the

827         Primate Cortex. *Neuron* **88,** 419–431 (2015).

828     40. Hasson, U., Chen, J. & Honey, C. J. Hierarchical process memory: memory as

829         an integral component of information processing. *Trends Cogn. Sci.* **19,** 304–

830         313 (2015).

831     41. Kiebel, S. J., Daunizeau, J. & Friston, K. J. A Hierarchy of Time-Scales and the

832         Brain. *PLOS Comput Biol* **4,** e1000209 (2008).

833     42. Murray, J. D. *et al.* A hierarchy of intrinsic timescales across primate cortex.

834         *Nat. Neurosci.* **17,** 1661–1663 (2014).

835     43. Wang, X.-J. & Kennedy, H. Brain structure and dynamics across scales: in

836         search of rules. *Curr. Opin. Neurobiol.* **37,** 92–98 (2016).

837     44. Boorman, E. D., Rushworth, M. F. & Behrens, T. E. Ventromedial Prefrontal

838         and Anterior Cingulate Cortex Adopt Choice and Default Reference Frames

839         during Sequential Multi-Alternative Choice. *J. Neurosci.* **33,** 2242–2253

840         (2013).

841     45. Boorman, E. D., Behrens, T. E. & Rushworth, M. F. Counterfactual Choice and

842         Learning in a Neural Network Centered on Human Lateral Frontopolar

843         Cortex. *PLOS Biol* **9,** e1001093 (2011).

844    46. Neta, M. *et al.* Spatial and Temporal Characteristics of Error-Related Activity

845        in the Human Brain. *J. Neurosci.* **35,** 253–266 (2015).

846    47. Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W. & Smith, S. M.

847        FSL. *NeuroImage* **62,** 782–790 (2012).

848    48. Woolrich, M. W., Ripley, B. D., Brady, M. & Smith, S. M. Temporal

849        Autocorrelation in Univariate Linear Modeling of FMRI Data. *NeuroImage* **14,**

850        1370–1386 (2001).

851    49. Watkins, C. J. C. H. & Dayan, P. Q -learning. *Mach. Learn.* **8,** 279–292 (1992).

852    50. Smith, S. M. Fast robust automated brain extraction. *Hum. Brain Mapp.* **17,**

853        143–155 (2002).

854