1    **Behavioural variability and cortical electrophysiological signals depend on**

2    **recent outcomes during human reinforcement motor learning**

3

4

5    Patrick Wiegel*[1,2], Meaghan Elizabeth Spedden[1], Christina Ramsenthaler[3,4,5], Mikkel

6    Malling Beck[1] & Jesper Lundbye-Jensen[1,6]

7

8    [1]Department of Nutrition, Exercise and Sports, University of Copenhagen, Copenhagen,

9    Denmark

10   [2]Department of Sport Science, University of Freiburg, Freiburg, Germany

11   [3]Klinik für Palliativmedizin, Universitätsklinikum Freiburg, Freiburg, Germany

12   [4]Wolfson Palliative care Research Centre, Hull & York Medical School, University of Hull, Hull,

13   UK

14   [5]Cicely Saunders Institute, Department of Palliative care, Policy & Rehabilitation, King's

15   College London, London, UK

16   [6]Department of Neuroscience, University of Copenhagen, Copenhagen, Denmark

17

18

19

20

21   *Corresponding author:

22   E-mail: patrick.wiegel@sport.uni-freiburg.de

23   Phone: +49 761 203 4550

24

25

26

27

1

**Abstract**

The history of our actions and the outcomes of these represent important information, which can inform choices, and efficiently guide future behaviour. While unsuccessful (S-) outcomes are expected to lead to more explorative motor states and increased behavioural variability, successful (S+) outcomes lead to reinforcement of the previous action and thus exploitation. Here, we show that during reinforcement motor learning, humans attribute different values to previous actions when they experience S- vs. S+ outcomes. Behavioural variability after S- outcomes is influenced more by the previous outcomes compared to what is observed after S+ outcomes. Using electroencephalography, we show that neural oscillations of the prefrontal cortex encode the level of reinforcement (high beta frequencies) and reflect the detection of reward prediction errors (theta frequencies). The results suggest that S+ experiences 'overwrite' previous motor states to a greater extent than S- experiences and that modulations in neural oscillations in the prefrontal cortex play a potential role in encoding the (changes in) movement variability state during reinforcement motor learning.

53 **Keywords**

54 Reinforcement, motor learning, prefrontal cortex, variability, neural oscillations,

55 exploration and exploitation.

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

**Introduction**

A striking feature of the nervous system is its ability to plan motor actions, to monitor the consequences during and following motor actions and to integrate these in future actions. During repeated practice of a given task, unsuccessful (S-) outcomes of our actions require exploration of new task solutions, while successful (S+) actions should ideally be reproduced and the strategy should be exploited. This process is called reinforcement learning (RL) and it relies on the evaluation of performance outcomes (Sutton & Barto, 1998).

RL guides behavioural variability in subsequent actions and optimises future performance. The level of reinforcement can be estimated by measuring changes in trial-to-trial variability (TTV) of behavioural characteristics (e.g. kinematic parameters). Although TTV has long been considered and in some ways can be an unwanted by-product of a noisy sensorimotor system, it is also an important aspect of motor learning (Wu *et al.*, 2014). TTV may be separated into unintended variability (due to noise) and intended variability (due to exploration). Provided that the outcome is monitored and evaluated by the central nervous system, variability may indeed serve to guide learning processes towards those behavioural characteristics, that lead to desirable outcomes. In this way, changes in TTV across time are necessary components of motor learning. During motor learning, TTV is increased after failures but decreased after rewards and this points towards a reward-dependent modulation of TTV to maximise future rewards (Takikawa *et al.*, 2002; Galea *et al.*, 2013; Pekny *et al.*, 2015). However, the outcome of the current action is not the only information used to guide subsequent behaviour. Instead, findings in rodent experiments demonstrate that the past several outcomes can be integrated to control future movements in rodents (Dhawale *et al.*, 2019), but in

102  humans, knowledge on the regulation of TTV in response to previous outcomes during

103  reinforcement motor learning is largely limited.

104  To which degree past outcomes are relevant during reinforcement motor learning may

105  depend on the outcomes of – and thus the value attributed to - previously performed

106  actions. Past actions are only informative when they help the agent to perform better

107  in future trials. Previous S- actions help to delineate actions that should be avoided.

108  As a consequence, other solutions can be tested in future actions, and S- outcomes

109  can thus guide exploration, potentially leading to S+ outcomes. When an S+ outcome

110  is experienced, this should lead to exploitation. But are past motor actions also helpful

111  or merely disregarded when experiencing S+ outcomes during reinforcement motor

112  learning? In this case, information about earlier movements might be less valuable or

113  down-weighted since the agent naturally aim to reproduce the current S+ movement.

114  Nevertheless, history of previous actions may still inform future actions. Here, we

115  tested this assumption and designed a reinforcement motor learning task in which

116  human participants performed goal-directed wrist flexion movements. We measured

117  TTV in wrist angle at movement end point and investigated influences of different

118  previous outcomes as a measure of the level of reinforcement.

119  Previous studies hypothesised that reward-dependent learning is mediated by the

120  difference between expected and actual rewards, so-called reward prediction errors

121  (Schultz, 2017). These signals drive learning based on feedback on outcome and

122  serve as the basis for future behavioural adjustments. A variety of neural circuits in

123  subcortical and cortical systems have been implicated in this context (Watanabe, 1996;

124  Knutson *et al.*, 2000; Schultz, 2000; Schall *et al.*, 2002; Cohen *et al.*, 2007; Histed *et*

125  *al.*, 2009; Narayanan *et al.*, 2013; HajiHosseini & Holroyd, 2015; Levy *et al.*, 2020).

126  Lately, a growing body of evidence suggests that neural circuits of the prefrontal cortex

127    contribute to reward-based learning in monkeys (Watanabe, 1996; Kim & Shadlen,

128    1999; Barraclough *et al.*, 2004; Seo & Lee, 2008; Histed *et al.*, 2009) and humans

129    (Akitsuki *et al.*, 2003; Cohen *et al.*, 2007; Marco-Pallares *et al.*, 2008; HajiHosseini *et*

130    *al.*, 2012; HajiHosseini & Holroyd, 2015). Modulations of neural oscillations over frontal

131    cortical areas can be observed after outcome information in decision-making tasks

132    potentially reflecting outcome-guided learning (Luft, 2014). Specifically, increases in

133    oscillatory activity have been observed for S- outcomes in theta band frequencies (4 –

134    8 Hz) and for S+ outcomes in high beta/low gamma frequencies (25 – 35 Hz). Neural

135    oscillations have been suggested to subserve important functions for the regulation of

136    information transfer across the brain and for controlling synaptic plasticity during

137    learning (Buzsaki & Draguhn, 2004; Luft, 2014). Thus, modulations of oscillatory

138    activity during outcome-processing constitute a reasonable mechanism to drive

139    adjustments in behaviour during reinforcement motor learning. Although the role of

140    neural oscillations during feedback-based learning is well described in cognitive

141    (decision-making) tasks, modulations of neural oscillations during human

142    reinforcement motor learning are not well understood. To test whether different

143    behavioural outcome scenarios are associated with specific oscillatory reinforcement

144    signals in the prefrontal cortex during motor learning, we recorded

145    electroencephalography (EEG) while participants practiced the motor task.

146    Thus, in the present paper we tested the effects of different outcomes and history of

147    previous actions on TTV and oscillatory signals in the prefrontal cortex during human

148    reinforcement motor learning. We hypothesised that S+ actions would result in

149    behavioural reinforcement i.e. exploitation and reduced consideration of previous

150    outcomes. In contrast, we assumed that S- actions would lead to greater behavioural

151    exploration that is informed by and thus depends more on previous outcomes. In

152    addition, we hypothesised that S+ outcomes would lead to increased oscillatory activity

153    in high beta frequencies and S- outcomes to increased oscillatory activity in theta

154    frequencies in the prefrontal cortex.

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

**Results**

177

178 Twenty-six participants performed wrist flexion reaching movements to horizontally

179 move a computer cursor into a target area (Fig. 1a & 1b). Movements were guided by

180 a visual scene on the computer screen which informed participants about movement

181 outcomes (Fig. 1c). During the experiment, participants were confronted with different

182 target positions and they received no online feedback but only binary augmented

183 feedback on the outcome (i.e. S+ or S-) following each trial. Movement end points for

184 a representative participant during the main protocol are plotted in Fig. 1d. Participants

185 showed greater movement end point variability during the main protocol than during

186 performance of movements to a stationary target with visual online feedback (Fam1,

187 F-test, F = 319.5, P < 0.001) and without visual online feedback (Fam2, F-test, F =

188 239.7, P < 0.001, Fig. 1e; Fam2 vs. Fam1: F-test, F = 22.2, P < 0.001).

189 In general, cursor end point distance from the targets varied considerably. The grand

190 average data are illustrated in Fig. 1f. The data suggest that while participants

191 approximated the targets positions quite well during blocks when the target did not

192 move, the distance from the target became much greater when the target position

193 changed. Also, participants were generally better at approximating the target they were

194 familiarized with (grey target in Fam1 and Fam2) compared to the unknown targets

195 (grey vs. green target: t = -5.8, P < 0.001; grey vs. purple target: t = -6.9, P < 0.001;

196 green vs. purple target: t = -1.8, P = 0.080) (Fig. 1g).

197 To investigate the strength of the relationship between previous and future outcomes,

198 we performed PAC for the outcome time series of all movements of the main protocol

199 (Fig. 1h). The analysis showed that the association of previous outcomes (trial lags 1

200 – 20) with the current outcome (trial lag = 0) decays with increasing trial lags which is

201 in line with a previous study in rodents (Dhawale et al., 2019). It also confirms our

8

202 assumption that the past two outcomes have the greatest association with future

203 outcomes. In the remainder of the paper, we focused on the impact of the previous two

204 outcomes (trial$^{(n)}$ & trial$^{(n-1)}$) on TTV in the main protocol (Fig. 1i).
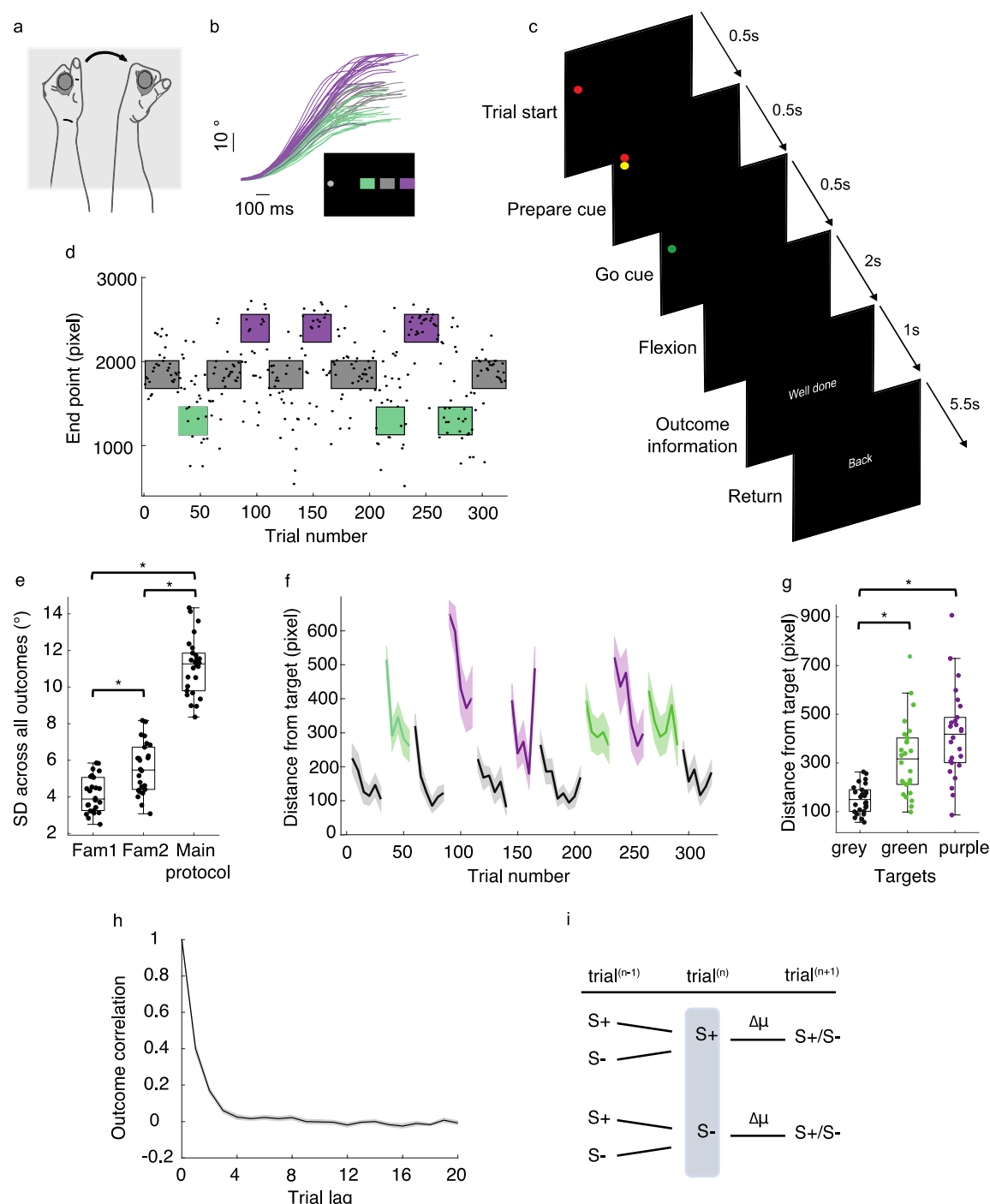


205

206 **Figure 1. Motor variability during reinforcement motor learning**

207 **a** Participants grabbed a handle to perform discrete wrist flexion movements with their left hand. The

208 sight of the hand and forearm was hidden by a custom-made box (grey shaded area) to remove visual

209 feedback of the moving arm. **b** The wrist angle was recorded by a goniometer that was integrated into

210 the handle. The panel shows 60 exemplary wrist angle traces from a participant aiming at the green,

211 grey and purple target box (20 movements each). Data were aligned to movement start and movement

212 end. The rectangle illustrates the computer cursor and the approximate positions of the different targets.

213 Cursor movements were one-dimensional in the horizontal plane. **c** The motor task was guided by

214 different visual scenes on the computer screen. Trial start was indicated by a red light appearing on the

215 upper left side of the computer screen. After 0.5 s, a yellow light below the red light signalled that

216 participants had to prepare their movements. Finally, after another 0.5 s a green light indicated the Go

217 cue. Participants had a time window of 2 s to perform their movements. Both, the cursor (representing

218 the wrist angle) and the target area were invisible on the screen. Thus, after each trial participants

219 received binary feedback about the outcome of the movement. S+ trials were indicated by "Well done"

220 while S- trials were indicated by "Try again". In both cases, feedback was visible for 1 s. After this period,

221 participants were asked to move their wrist back to the starting position. A new trial started after 5.5 s.

222 **d** Movement end points (in pixel) from one exemplary participant. Note that we changed the horizontal

223 position of the target several times during the experiment (after 25, 30 or 40 trials) to stimulate exploring

224 motor behaviour of the participants. **e** Boxplots and individual data for standard deviations from Fam1,

225 Fam2 and the main protocol. * indicate significant differences between the conditions. Note that x-axis

226 values were jittered to more clearly present the data. **f** The grand average data for distance from target

227 (pixel) during the main protocol. Note that the figure shows averaged data in bins of 5 trials. Shading

228 represents the standard error of the mean. **g** Boxplots and individual data for the distance from the

229 different targets (in pixel) averaged across the main protocol. * indicate significant differences between

230 the conditions. **h** We performed partial autocorrelation of the outcome time series (S+ and S- outcomes).

231 The plot shows the grand averaged data. Shading represents the standard error of the mean. **i** Our

232 framework focused on the question how the past two trials (trial$^{(n)}$ and trial$^{(n-1)}$) influence regulations of

233 TTV ($\Delta\mu$) in movement endpoint and oscillatory reinforcement signals.

234

235 *TTV depends on the outcome of the previous movement*

236 In a first step, we calculated the number of S+$^{(n)}$ and S-$^{(n)}$ movements for each

237 participant. Table 1 summarizes the descriptive data of this analysis. There was no

10

238 significant difference in the number of S+ and S- movements (paired t-test, t = -0.7, P

239 = 0.499) suggesting that participants performed a similar number of S+ and S- trials.

240 We also tested the effect of the previous outcome$^{(n)}$ on motor performance in the

241 following trial$^{(n+1)}$. Participants had a greater proportion of S+$^{(n+1)}$ trials (i.e. hits) when

242 the preceding trial was S+$^{(n)}$ (71,3%) compared to when the preceding trial was S-$^{(n)}$

243 (28,1%) (Test of equal or given proportions, $X^2(1)$ = 1515.4, P < 0.001). This indicates

244 better motor performance after S+$^{(n)}$ trials compared to S-$^{(n)}$ trials.

245

246 **Table 1.** *Grand average descriptive data (n=26) of the number of S+$^{(n)}$ and S-$^{(n)}$ events and the number*

247 *of subsequent S+ events in trial$^{(n+1)}$. SD = standard deviation, Min = minimum, Max = maximum.*

| | Grand average number | Mean ± SD | Min | Max | Grand average number of S+ in trial$^{(n+1)}$ |
|---|---|---|---|---|---|
| S+$^{(n)}$ | 3951 | 151.9 ± 30.0 | 94 | 216 | 2818 (71.3%) |
| S-$^{(n)}$ | 4161 | 160.0 ± 30.0 | 96 | 218 | 1168 (28.1%) |

248

249 Next, we analysed TTV in movement endpoint after S+ (Δμ| S+$^{(n)}$) and S- (Δμ| S-$^{(n)}$)

250 trials. Fig. 2a and Fig. 2b show the grand average of the conditioned probability

251 distributions. Visual inspection of the data suggested greater kurtosis after S-$^{(n)}$ trials

252 than after S+$^{(n)}$ trials.

253 We tested this observation statistically and calculated the M (signed and unsigned

254 TTV) and SD (signed TTV). The analyses supported our assumption from the visual

255 inspection of the data and revealed significant differences in signed TTV (paired t-test,

256 t = -3.6, P = 0.001, Fig. 2c left panel), the SD of TTV (F-test, F = 232.3, P < 0.001, Fig.

257 2c middle panel) and the unsigned TTV (F-test, F = 223.5, P < 0.001, Fig. 2c right

258 panel). After S+$^{(n)}$ movements, TTV in movement endpoint (signed TTV: -0.5° ± 0.6°;

259 SD of TTV: 4.4° ± 1.2°; unsigned TTV: 3.4° ± 0.8°) was lower than after S- movements

11

260   (signed TTV: 0.3° ± 0.5°; SD of TTV: 13.5° ± 2.6°; unsigned TTV: 10.4° ± 2.1°). The

261   results indicate that variability of the movements and the absolute motor exploration

262   were greater after S-$^{(n)}$ movements than after S+$^{(n)}$ movements. The differences in

263   signed TTV suggest that participants tended to move the cursor slightly less after S+$^{(n)}$

264   movements than after S-$^{(n)}$ movements. Interestingly, differences in TTV were not only

265   constrained to movement end point, rather TTV in maximal movement speed and

266   movement time showed similar characteristics suggesting that participants

267   reinforcement processes generalize to outcome-irrelevant parameters

268   (Supplementary Fig. 1).

269   Moreover, SD of TTV in movement endpoint after S+$^{(n)}$ movements was a good

270   indicator of overall motor performance predicting the total individual score (i.e. total

271   number of hits) (linear regression: $F = 6.1$, $P = 0.021$, $R^2 = 0.20$, Fig. 2d top panel).

272   This was not the case for SD of TTV in movement endpoint after S-$^{(n)}$ movements

273   (linear regression: $F = 0.1$, $P = 0.720$, $R^2 = 0.01$, Fig. 2d bottom panel). Additional

274   analyses demonstrated that the unsigned TTV after S+$^{(n)}$ movements also predicted

275   motor performance but not the signed TTV (Supplementary Fig. 2).

276   These results suggest different mechanisms of trial-by-trial reinforcement motor

277   learning after S+$^{(n)}$ compared to S-$^{(n)}$ trials. S-$^{(n)}$ trials stimulate greater TTV while S+$^{(n)}$

278   trials lead to lower TTV. Moreover, participants who were able to "reproduce" S+

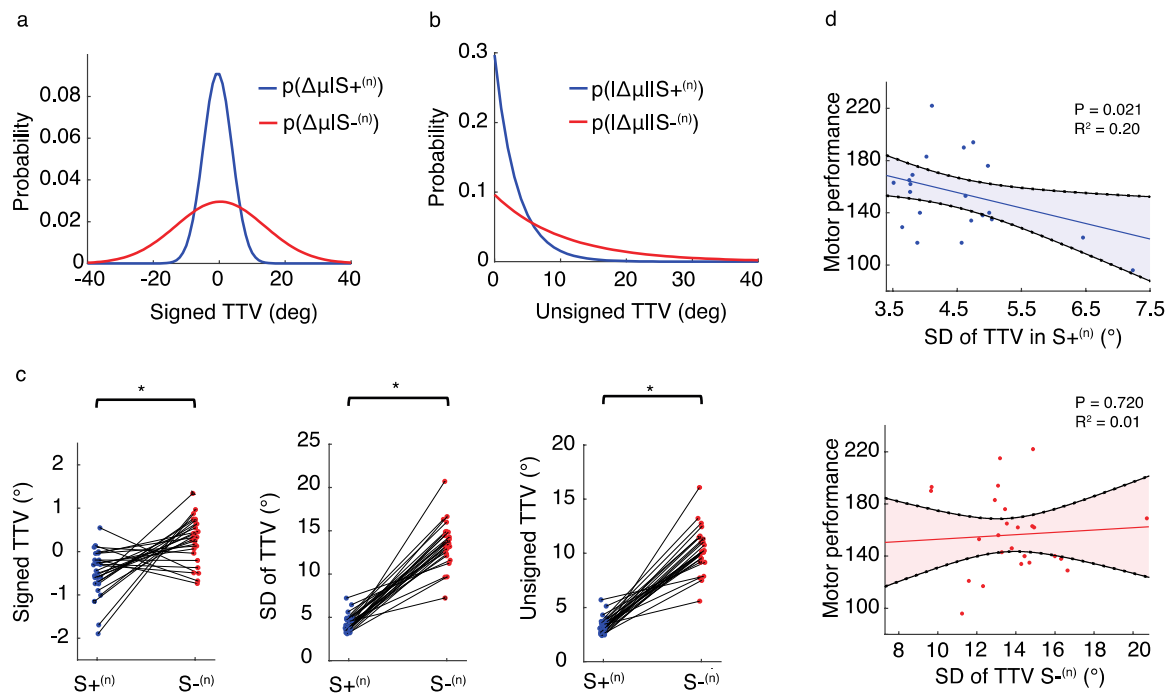279   movements more accurately also performed better overall.

**Figure 2. Behavioural variability depends on the previous outcome**

**a** The grand average normal distribution for signed changes in TTV after $S+^{(n)}$ and $S-^{(n)}$ trials. **b** The grand average exponential distribution for unsigned changes in TTV after $S+^{(n)}$ and $S-^{(n)}$ trials. **c** Individual differences in signed TTV (left panel), SD of TTV (middle panel) and unsigned TTV (right panel) after $S+^{(n)}$ and $S-^{(n)}$ trials. * indicate significant differences between the conditions (after correction for multiple comparisons). Note that x-axis values were jittered to more clearly present the data. **d** Scatter plot showing the association between motor performance and SD of TTV after $S+^{(n)}$ trials (top panel), and SD of TTV after $S-^{(n)}$ trials (bottom panel). Each dot represents an individual participant. The lines represent the fitted regression line. Shading is 95% confidence intervals of the regression.

*The influence of previous outcomes differs for S+ and S- motor actions*

It remains an open question whether the outcome of the second-to-last movement$^{(n-1)}$ changes the impact of the outcome of the previous movement$^{(n)}$. In other words, do outcomes prior to the current one have a different relevance when participants experience $S+^{(n)}$ and $S-^{(n)}$ movements? To investigate this, we extracted TTV in movement endpoint conditioned on the previous two trials. Thus, we analysed TTV for

297    four outcome scenarios: 1) $S+^{(n)}$ & $S+^{(n-1)}$, 2) $S+^{(n)}$ & $S-^{(n-1)}$, 3) $S-^{(n)}$ & $S+^{(n-1)}$ and 4) $S-^{(n)}$ & $S-^{(n-1)}$.

298

299    The four conditions contained different trial numbers suggesting that certain outcome

300    combinations were more likely than others. The rmANOVA with the number of trials as

301    the dependent variable yielded a significant effect of outcome history ($F_{[1.3, 31.4]} = 56.2$,

302    $P < 0.001$, $\eta2_{partial} = 0.69$). $S+^{(n)}$ movements were more often preceded by $S+^{(n-1)}$

303    movements than by $S-^{(n-1)}$ movements (paired t-test: $t = -9.2$, $P < 0.001$). In contrast,

304    $S-^{(n)}$ movements were more often preceded by $S-^{(n-1)}$ movements than by $S+^{(n-1)}$

305    movements (paired t-test: $t = 10.6$, $P < 0.001$). The descriptive data of this analysis are

306    shown in table 2.

307

308    **Table 2.** *Grand average descriptive data (n=26) of the number of* $S+^{(n)}$ & $S+^{(n-1)}$, $S+^{(n)}$ & $S-^{(n-1)}$, $S-^{(n)}$ &

309    $S+^{(n-1)}$ and $S-^{(n)}$ & $S-^{(n-1)}$ events *and the number of subsequent S+ events in trial$^{(n+1)}$. SD = standard*

310    *deviation, Min = minimum, Max = maximum.*

| | Grand average number | Mean ± SD | Min | Max | Grand average number of S+ in trial$^{(n+1)}$ |
|---|---|---|---|---|---|
| $S+^{(n)}$ & $S+^{(n-1)}$ | 2661 | 108.1 ± 31.4 | 48 | 177 | 1971 (75.1%) |
| $S+^{(n)}$ & $S-^{(n-1)}$ | 1136 | 45.5 ± 7.8 | 27 | 61 | 709 (62.4%) |
| $S-^{(n)}$ & $S+^{(n-1)}$ | 1145 | 45.2 ± 7.6 | 27 | 60 | 533 (46.6%) |
| $S-^{(n)}$ & $S-^{(n-1)}$ | 2910 | 116.2 ± 31.8 | 54 | 176 | 600 (20.1%) |

311

312    Next, we tested the effect of the different outcome histories on the proportion of S+

313    movements in trial$^{(n+1)}$. A test of equal or given proportions revealed significant

314    differences between the four outcome histories ($X^2(1) = 1691.9$, $P < 0.001$) suggesting

315    that motor performance depends on the past two outcomes. Indeed, the proportion of

316     S+ movements in trial$^{(n+1)}$ was highest in S+$^{(n)}$ & S+$^{(n-1)}$ (75.1%) and lowest in S-$^{(n)}$ & S-

317     $^{(n-1)}$ (20.1%).

318     In a final step, we computed the TTV in movement endpoint for the different outcome

319     histories. Due to the differences in trial number between the outcome histories, we

320     performed bootstrapping with replacement (1,000 iterations) and matched the trial

321     numbers for each participant. For each participant and condition, normal and

322     exponential distributions were fitted from the bootstrapped datasets. The normal and

323     exponential distributions for TTV after S+ trials$^{(n)}$ conditioned on the trial$^{(n-1)}$ are shown

324     in Fig. 3a and Fig. 3b, respectively. These probability distributions appear to be very

325     similar. As can be seen in Fig. 3c and Fig. 3d, this was different when trial$^{(n)}$ was S-.

326     The probability distribution after S- trials$^{(n)}$ is broader when the preceding trial is also

327     S-$^{(n-1)}$ while it is narrower when the preceding trial is S+$^{(n-1)}$. Statistical analyses were

328     performed on the M and SD of the individual distributions. The results from the

329     rmANOVA are presented in table 3.

330

331     **Table 3.** rmANOVA results for the effect of the four outcome histories on the different dependent

332     variables

|  | *Signed change* | *SD of signed change* | *Unsigned change* |
|---|---|---|---|
| $F_{[DF,error]}$ | $15.1_{[1.8,44.7]}$ | $167.0_{[1.7, 42.3]}$ | $158.8_{[1.5, 36.9]}$ |
| P | <0.001 | <0.001 | <0.001 |
| $\eta2_{partial}$ | 0.38 | 0.87 | 0.86 |

333

334

15

335   Post hoc tests were consequently used. For this purpose, we compared the TTV in

336   movement endpoint between $S-^{(n)}$ & $S+^{(n-1)}$ and $S-^{(n)}$ & $S-^{(n-1)}$ and between $S+^{(n)}$ & $S+^{(n-1)}$

337   and $S+^{(n)}$ & $S-^{(n-1)}$ trials.

338   The signed TTV ($S+^{(n)}$ & $S+^{(n-1)}$: -0.3° $\pm$ 0.5°; $S+^{(n)}$ & $S-^{(n-1)}$: -0.6° $\pm$ 1.0°), SD of the TV

339   ($S+^{(n)}$ & $S+^{(n-1)}$: 4.4° $\pm$ 1.1°; $S+^{(n)}$ & $S-^{(n-1)}$: 5.0° $\pm$ 1.3°) and unsigned TTV ($S+^{(n)}$ & $S+^{(n-1)}$

340   : 3.3° $\pm$ 0.8°; $S+^{(n)}$ & $S-^{(n-1)}$: 3.9° $\pm$ 0.8°) were significantly different between $S+^{(n)}$

341   movements that were preceded by $S+^{(n-1)}$ and those that were preceded by $S-^{(n-1)}$

342   (paired t-test: signed TTV: t = -1.6, P = 0.112; F-test: SD of TTV: F = 13.7 , P = 0.001;

343   F-test : unsigned TTV: F = 36.3, P < 0.001, Fig. 3e). Likewise, there were significant

344   differences in signed TTV ($S-^{(n)}$ & $S+^{(n-1)}$: 1.3° $\pm$ 1.6°; $S-^{(n)}$ & $S-^{(n-1)}$: -0.1° $\pm$ 0.6°), SD

345   of TTV ($S-^{(n)}$ & $S+^{(n-1)}$: 8.9° $\pm$ 2.4°; $S-^{(n)}$ & $S-^{(n-1)}$: 14.9° $\pm$ 2.9°) and unsigned TTV ($S-^{(n)}$

346   & $S+^{(n-1)}$: 7.0° $\pm$ 1.7°; $S-^{(n)}$ & $S-^{(n-1)}$: 11.9° $\pm$ 2.7°) after S- movements $^{(n)}$ between trials

347   that were conditioned on $S+^{(n-1)}$ and $S-^{(n-1)}$ (paired t-test: signed TTV: t = -3.8, P =

348   0.001; F-test: SD of TTV: F = 96.1, P < 0.001; F-test: unsigned TTV: F = 84.6, P <

349   0.001, Fig. 3f). These results demonstrate that S+ and S- outcomes in trial$^{(n-1)}$ have

350   differential influences on changes in movement end point after S+ and S- outcomes in

351   trial$^{(n)}$.

352   We tested whether the differences in TTV conditioned on the past two trials were

353   different for $S+^{(n)}$ and $S-^{(n)}$ trials. Indeed, differences in TTV between $S+^{(n)}$ and $S-^{(n)}$

354   movements were greater when the previous outcome was $S-^{(n)}$ than when it was $S+^{(n)}$

355   (paired t-test: TTV: t = -2.4, P = 0.022; F-test: SD of TTV: F = 63.9, P < 0.001; F-test:

356   unsigned TTV: F = 64.9, P < 0.001). These results suggest different RL signals after

357   distinct outcome histories. When participants experienced $S+^{(n)}$ movements, the

358   outcome of the second-to-last trial$^{(n-1)}$ became less influential. Additionally, in case of

S-$^{(n)}$ movements, the experience of S- outcomes in trial$^{(n-1)}$ lead to greater motor exploration but to less motor exploration when trial$^{(n-1)}$ was S+.
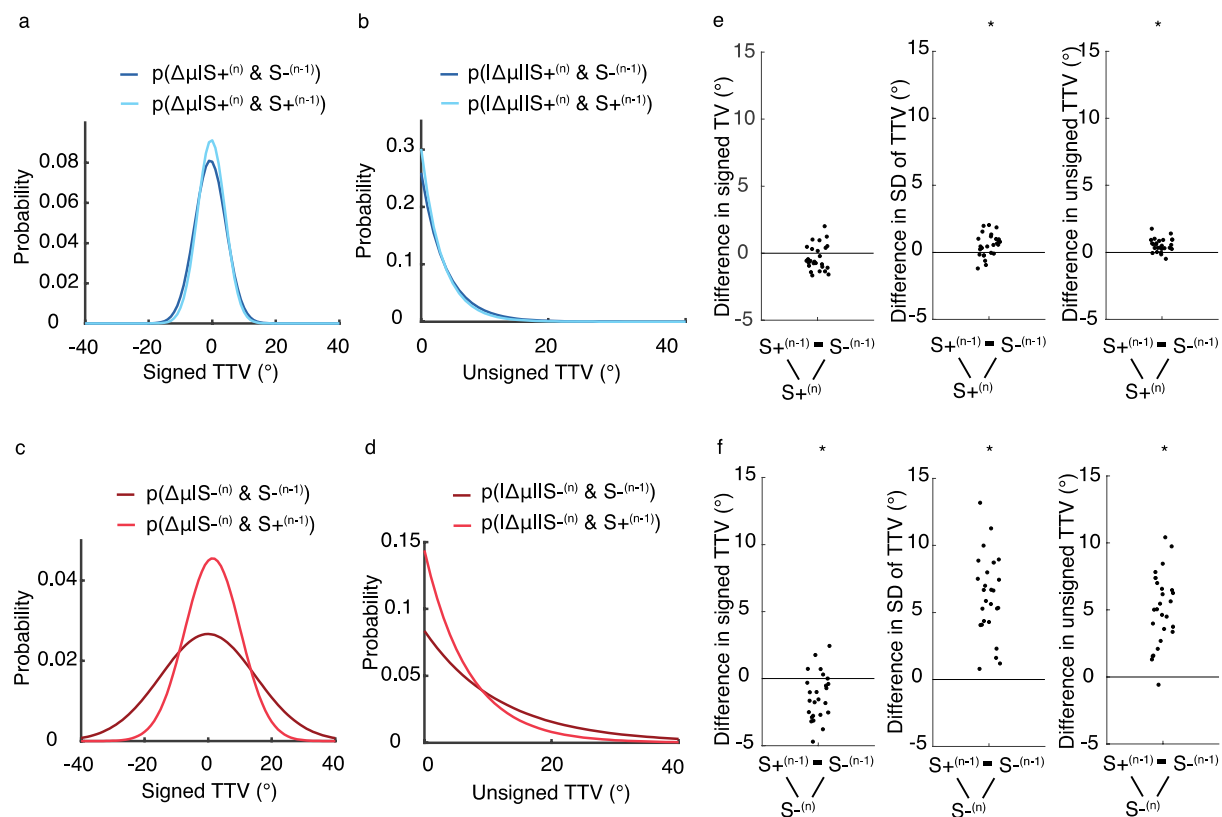


**Figure 3. The impact of previous outcomes differs for S+ and S- motor actions**

**a** The grand average normal distribution for signed TTV after S+$^{(n)}$ trials when the previous trial was S+$^{(n-1)}$ and S-$^{(n-1)}$. **b** The grand average exponential distribution for unsigned changes in TTV after S+$^{(n)}$ trials when the previous trial was S+$^{(n-1)}$ and S-$^{(n-1)}$. **c** The grand average normal distribution for signed TTV after S-$^{(n)}$ trials when the previous trial was S+$^{(n-1)}$ and S-$^{(n-1)}$. **d** The grand average exponential distribution for unsigned changes in TTV after S-$^{(n)}$ trials when the previous trial was S+$^{(n-1)}$ and S-$^{(n-1)}$. **e** Individual differences in signed TTV (left panel), SD of TTV (middle panel) and unsigned TTV (right panel) between S+$^{(n)}$ that were preceded by S+$^{(n-1)}$ and S-$^{(n-1)}$ movements. **f** Individual differences in signed TTV (left panel), SD of TTV (middle panel) and unsigned TTV (right panel) between S-$^{(n)}$ that were preceded by S+$^{(n-1)}$ and S-$^{(n-1)}$ movements. Each dot represents an individual. Note that x-axis values were jittered to more clearly present the data. * indicate significant differences between the different outcome histories (after correction for multiple comparisons).

376 *PFC oscillatory responses to outcomes during reinforcement motor learning*

377 Since TTV clearly depends on the different outcomes (as evident from the behavioural

378 results), it is indeed plausible that the brain generates different reinforcement signals

379 accordingly to regulate future behavioural adjustments during reinforcement motor

380 learning. Here, we tested this assumption and analysed modulations in neural

381 oscillations during outcome processing.

382 We concentrated on pre-selected frequency ranges (theta: 4 – 8 Hz, high beta: 25 –

383 35 Hz) and the time of outcome processing (250 ms – 550 ms after outcome

384 information). First, we plotted the power data in sensor space and observed the

385 greatest increases over frontal sensors for both frequencies (theta: Fig 4a, high beta:

386 4b). Subsequent source space analyses confirmed our initial assumption that the

387 greatest power increases (relative to pre-feedback) were observed over the prefrontal

388 cortex (theta: Fig 4c, high beta: 4d). In the remaining analyses, we focused on

389 oscillatory activity in the superior frontal gyrus (SFG) and rostral middle frontal gyrus

390 (RMFG), two areas that were previously shown to engage in RL (Garrison *et al.*, 2013)

391 (Supplementary Fig. 3). Finally, we extracted averaged time frequency data of the ROI

392 and found similar power time courses as in the sensor space (theta: Fig 4e, high beta:

393 4f). The data support the assumption that information about motor outcomes result in

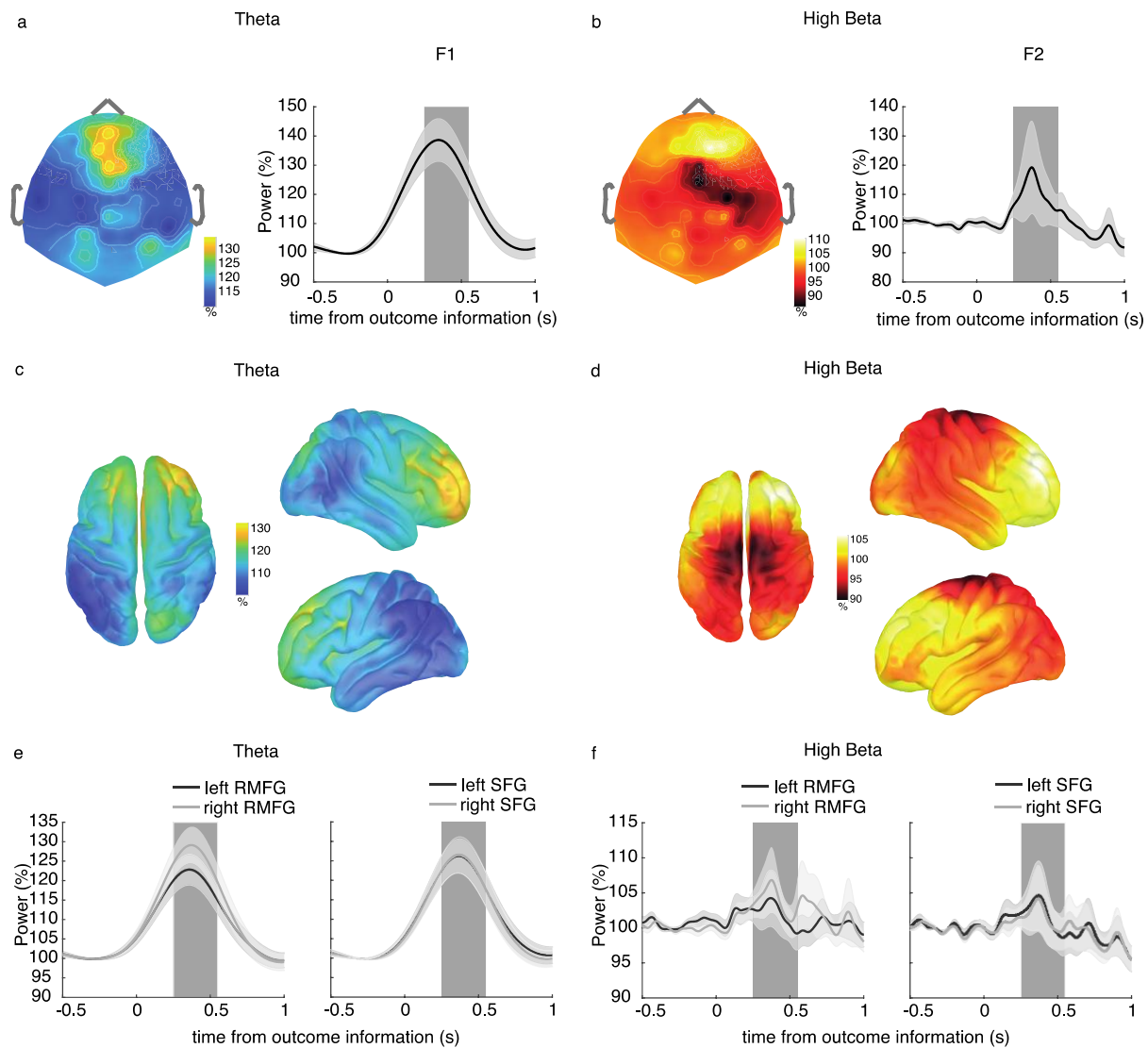394 changes neural oscillations in the PFC.

**Figure 4. Oscillatory responses in PFC to outcome presentation during reinforcement motor learning**

**a & b** Topographical distribution of theta (a) and high beta (b) power during outcome processing (250 – 550 ms after outcome information) (left panel). Grand average power time courses for theta (a) and high beta (b) frequencies during outcome processing (-0.5 s – 1 s relative to outcome information) (right panel). The data show the mean of all trials independent of the prior outcome. Data were plotted for the channel displaying the greatest power increase relative to pre-feedback (theta: F1, high beta: F2). Note that we used different colours to plot the topographical power distribution for theta and high beta frequencies since the data have different scales. **c & d** Source localisation results of theta power (c) and high beta power (d) during outcome processing (average from 250 – 550 ms after outcome information). Data were interpolated on the MRI template. Data are shown from the left, right and front.

19

408 **e & f** Power time courses of the ROI during outcome processing (-0.5 s – 1 s relative to outcome

409 information). Theta power (e) time courses of bilateral RMFG are shown in left panel and of the bilateral

410 SFG in the right panel. High beta (f) power time courses of bilateral RMFG are shown in the left panel

411 and of the bilateral SFG in the right panel. The shaded rectangle highlights the time window of interest

412 250 ms – 550 ms. Shading around the mean represent the standard error of the mean.

413

414 _Different oscillatory reinforcement signals following S+$^{(n)}$ and S-$^{(n)}$ movements_

415 If oscillatory responses in SFG and RMFG in response to the outcome represent

416 different reinforcement signals, then we would expect different oscillatory responses

417 to S+ and S- outcomes in trial$^{(n)}$. The time courses of theta and beta power changes

418 after both outcomes are shown in Fig 5. Presentation of movement outcomes had

419 differential effects on theta and high beta frequencies depending on the outcome of

420 the movement.

421 In comparison to S+$^{(n)}$ movements, S-$^{(n)}$ movements resulted in greater theta

422 oscillations. While there were no significant differences in theta power over the left

423 RMFG (Fig. 5a, left panel) and left SFG (Fig. 5b, left panel) (all $P > 0.05$), significant

424 differences in power between S+$^{(n)}$ and S-$^{(n)}$ movements were revealed between 250

425 ms and 550 ms in the right RMFG (critical P-value: 0.037, Fig. 5a, right panel) and

426 between 300 ms and 500 ms in the right SFG (critical P-value: 0.034, Fig. 5b, right

427 panel). In contrast, S+ movements$^{(n)}$ resulted in greater high beta oscillatory responses

428 compared to S-$^{(n)}$ movements. This was the case in the left RMFG (critical P-value:

429 0.009, Fig. 5c, left panel) and left SFG from 450 ms and 550 ms (critical P-value: 0.015,

430 Fig. 5d, left panel) and in the right SFG at 450 ms (critical P-value: 0.002, Fig. 5d, right

431 panel). No significant differences in high beta oscillatory response between S+$^{(n)}$ and

432 S-$^{(n)}$ movements were observed in the right RMFG ($P > 0.05$) (Fig. 5c, right panel).

433 Similar results were obtained for sensor-space data (Supplementary Fig. 4). The data

20

434    suggests different oscillatory responses to $S+^{(n)}$ compared to $S-^{(n)}$ outcomes in the

435    prefrontal cortex. While $S+^{(n)}$ outcomes lead to greater high beta band power, $S-^{(n)}$

436    outcomes lead to greater theta band power.
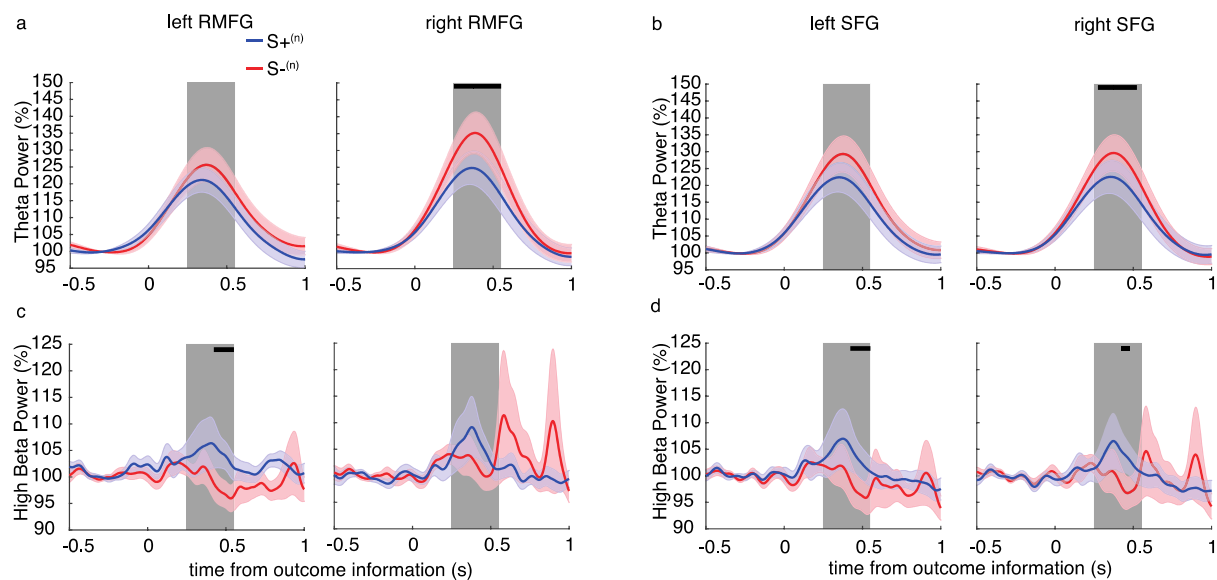
437



438

**Figure 5. $S+(n)$ and $S-(n)$ outcomes evoke different oscillatory signals in PFC**

440    **a** Theta power time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome

441    processing for $S+^{(n)}$ and $S-^{(n)}$ movements. **b** Theta power time courses of the left SFG (left panel) and

442    right SFG (right panel) during outcome processing for $S+^{(n)}$ and $S-^{(n)}$ movements. **c** High beta power

443    time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome processing for

444    $S+^{(n)}$ and $S-^{(n)}$ movements. **d** High beta power time courses of the left SFG (left panel) and right SFG

445    (right panel) during outcome processing for $S+^{(n)}$ and $S-^{(n)}$ movements. In all plots. the shaded rectangle

446    highlights the time window of interest 250 ms – 550 ms. Shading around the mean represent the

447    standard error of the mean. Significant differences between the outcomes is indicated by the horizontal

448    lines.

449

450    *Oscillatory reinforcement signals in PFC depend on the outcomes of the past two*

451    *movements*

452    Next, we asked whether neural oscillations during outcome processing depend on the

453    outcome of the past two trials. Due to differences in the number of trials per each

454    outcome scenario and potential differences in signal-to-noise ratios, we performed a

455    bootstrapping with replacement for each individual EEG dataset. The bootstrapped

456    grand average time frequency responses of our ROI are shown in Fig. 6.

457    Theta power was significantly greater between 250 ms and 550 ms in the left RMFG

458    (critical P-value: 0.005, Fig. 6a, left panel), the right RMFG (critical P-value: 0.001, Fig.

459    6a, right panel), the left SFG (critical P-value: 0.003, Fig. 6b, left panel) and the right

460    SFG (critical P-value: 0.001, Fig. 6b, right panel) in $S+^{(n)}$ movements that were

461    preceded by $S-^{(n-1)}$ movements than those that were preceded by $S+^{(n-1)}$ movements.

462    Likewise, high beta power was significantly greater at 500 ms in the left RMFG (critical

463    P-value: 0.005, Fig. 6c, left panel) in $S+^{(n)}$ movements that were preceded by $S-^{(n-1)}$

464    movements than those that were preceded by $S+^{(n-1)}$ movements. Though there were

465    similar tendencies in the other regions, the results did not reach statistical significance

466    ($P > 0.05$) in the right RMFG (Fig. 6c, right panel), the left SFG (Fig. 6d, left panel) and

467    the right SFG (Fig. 6d, right panel).

468    Finally, the time course of theta power was not significantly different between $S-^{(n)}$

469    movements that were preceded by $S-^{(n-1)}$ movements and those that were preceded by

470    $S+^{(n-1)}$ movements in the left RMFG (Fig. 6e, left panel), the right RMFG (Fig. 6e, right

471    panel) and the left SFG (Fig. 6f, left panel) (all $P > 0.05$). Theta power between 400

472    ms and 550 ms was significantly greater in $S-^{(n)}$ movements that were preceded by

473    $S+^{(n-1)}$ movements compared to those that were preceded by $S+^{(n-1)}$ movements in the

474    right SFG (critical P-value: 0.028, Fig. 6f, right panel). No significant differences were

475    observed for high beta power in the left RMFG (Fig. 6g, left panel), the right RMFG

476    (Fig. 6g, right panel), the left SFG (Fig. 6h, left panel) and the right SFG (Fig. 6h, right

477    panel) between $S-^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and those

478    that were preceded by $S+^{(n-1)}$ movements (all $P > 0.05$). Additional information on

479    sensor-space data can be found in Supplementary Fig. 5. The results suggest that high

480    beta band power increases especially in cases were $S+^{(n)}$ trials are preceded by $S-^{(n-1)}$

481    trials. On the contrary, theta band power increases were most prominent when the

482    outcomes in two subsequent trials changed (from $S+^{(n-1)}$ to $S-^{(n)}$ or from $S-^{(n-1)}$ to $S+^{(n)}$).
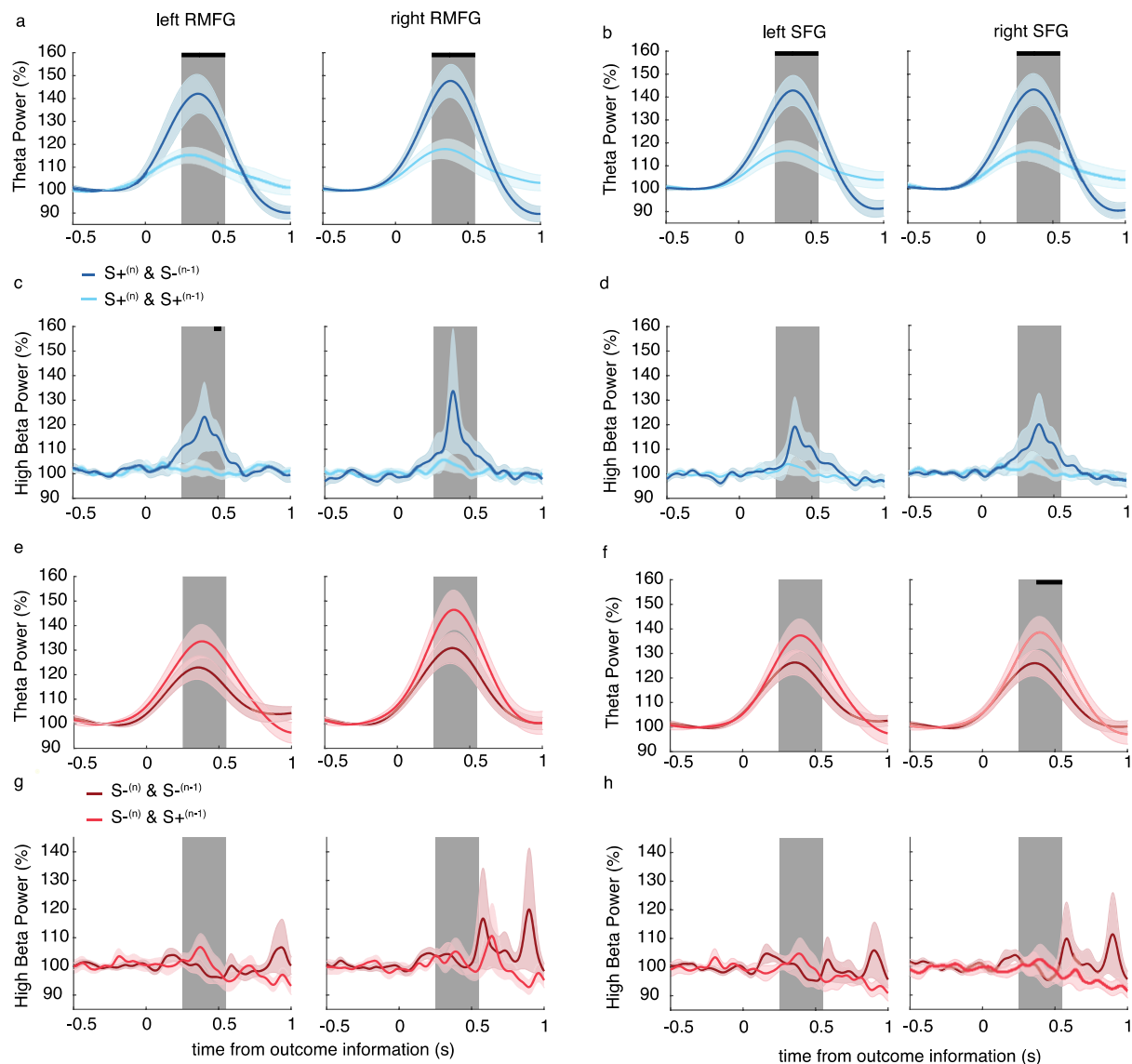


483

**Figure 6. Oscillatory reinforcement signals in PFC depend on the outcome history**

485    **a** Theta power time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome

486    processing for $S+^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **b**

487    Theta power time courses of the left SFG (left panel) and right SFG (right panel) during outcome

488    processing for $S+^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **c** High

489    beta power time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome

490    processing for $S+^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **d** High

491    beta power time courses of the left SFG (left panel) and right SFG (right panel) during outcome

492    processing for $S+^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **e**

493    Theta power time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome

494    processing for $S-^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **f** Theta

495    power time courses of the left SFG (left panel) and right SFG (right panel) during outcome processing

496    for $S-^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **g** High beta power

497    time courses of the left RMFG (left panel) and right RMFG (right panel) during outcome processing for

498    $S-^{(n)}$ movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. **h** High beta power

499    time courses of the left SFG (left panel) and right SFG (right panel) during outcome processing for $S-^{(n)}$

500    movements that were preceded by $S+^{(n-1)}$ movements and $S-^{(n-1)}$ movements. In all plots. the shaded

501    rectangle highlights the time window of interest (250 ms – 550 ms after outcome information). Shading

502    around the mean represent the standard error of the mean. Significant differences between the

503    outcomes is indicated by the horizontal lines.

504

505

506

507

508

509

510

511

512

513

514

515

516

517

**Discussion**

The present study revealed outcome-specific adjustments of TTV in movement endpoint during reinforcement motor learning. We found that $S+^{(n)}$ movements were followed by relatively little TTV and a greater proportion of S+ movements in the subsequent trial$^{(n+1)}$. In contrast, $S-^{(n)}$ movements caused increased TTV and had a lower proportion of S+ movements in the subsequent trial$^{(n+1)}$. We investigated the regulation of TTV in further detail and found that TTV following $S+^{(n)}$ movements was less influenced by the outcome of the previous trial$^{(n-1)}$. In contrast, TTV was greater after $S-^{(n)}$ movements when they were preceded by $S-^{(n-1)}$ trials compared to when they were preceded by $S+^{(n-1)}$ trials. These results suggest a change in values of previous outcomes when human participants experience $S+^{(n)}$ and $S-^{(n)}$ movements during a reinforcement-based motor learning task. We hypothesised that these behavioural effects involve different cortical reinforcement signals and analysed neural oscillations from the prefrontal cortex (bilateral RMFG and SFG). In general, $S+^{(n)}$ movements led to increased high beta oscillatory activity while $S-^{(n)}$ movements caused greater theta oscillatory activity. However, when the power data were conditioned on the past two trials we found that these oscillatory differences are driven by distinct outcome histories. Beta power following feedback presentation was greatest in $S+^{(n)}$ movements that were preceded by $S-^{(n)}$ movements providing a potential mechanism to reinforce the current $S+^{(n)}$ movement and disregard the behaviour from the previous $S-^{(n-1)}$ trial. Theta oscillatory activity was greatest after $S-^{(n)}$ movements when the preceding movement was $S+^{(n-1)}$ and after $S+^{(n)}$ movements when the preceding movement was $S-^{(n-1)}$ suggesting that increased theta oscillations reflect error detection, i.e. large discrepancies between expected and actual outcomes (reward prediction errors) and thereby a change in the variability state. Therefore, our results provide evidence for an

543  outcome-specific regulation of TTV and oscillatory reinforcement signals in the

544  prefrontal cortex during human reinforcement motor learning.

545  Learning from the outcome of recent trials is an important component of motor learning

546  and helps to adjust future motor behaviour. For instance, reinforcement of S+

547  movements through rewards causes stable performance gains (Pekny *et al.*, 2011)

548  and results in greater motor retention than learning from punishment (Galea *et al.*,

549  2015). Negative (punishment) feedback helps us to avoid S- movements and

550  accelerates online improvements in motor performance (Galea *et al.*, 2015). These

551  results suggest that learning from S+ and S- outcomes involves distinct processes. In

552  the present study, we showed that S+ and S- motor outcomes lead to different

553  adjustments in TTV. In agreement with an earlier report (Pekny *et al.*, 2015), TTV was

554  decreased after S+ movements and increased after S- movements. The outcome-

555  dependent modulation of TTV was additionally an important aspect of efficient

556  reinforcement motor learning as individual differences in TTV (after S+$^{(n)}$ trials but not

557  after S-$^{(n)}$ trials) predicted overall motor performance.

558  However, it is unlikely that updating of motor behaviour uses only the outcome of the

559  current action. Rather, motor variability is causally regulated by recent trial outcomes

560  in rats, that is the integrated outcomes of the past ~10 trials (Dhawale *et al.*, 2019).

561  The relevance of trial outcome on motor variability decays in an exponential weighted

562  manner, i.e. very recent outcomes have a greater effect compared to "older" outcomes.

563  In the present study, we focused on the effects of the past two trial outcomes.

564  In cognitive decision-making tasks, participants are often presented a finite number of

565  available options, e.g. in gambling games (Cohen *et al.*, 2007; HajiHosseini & Holroyd,

566  2015). Since there is no direct association between the values of the different options,

567  choosing an option does not inform about the value of the other options in these tasks.

568    This is different in the motor domain as the agents act in a continuous motor space

569    (Dhawale *et al.*, 2019). Here, a S+ movement end point implies that neighbouring

570    movement end points likely have a similar positive value. In contrast, an S- movement

571    suggests that S+ movement end points are located at more distant positions. However,

572    the value of previous outcomes may depend on the acutely experienced outcomes.

573    Here, we showed that TTV is less influenced by the outcome of the second-to-last trial

574    when the current trial was S+. The findings suggest that previous outcomes are

575    assigned less value as soon as humans experience positive motor outcomes, a

576    process that might promote acute improvements in short-term motor performance. In

577    contrast, TTV was increased after S- movements, especially when the second-to-last

578    trial was also S-. S- trials that were preceded by S+ trials showed less TTV. This implies

579    that participants consider previous S+ outcomes for adjustments in TTV, even when

580    they acutely experience S- outcomes. We speculate that motor outcomes prior to

581    positive outcomes are disregarded since they contain only little relevant information on

582    potential future behavioural adjustments in the task.

583    The efficacy of this mechanisms may depend on the certainty of the task and

584    environment (Dhawale *et al.*, 2019). In task situations where the conditions are

585    stationary, a greater reliance on previous movements might be a reasonable strategy

586    to reproduce the positive outcomes. In non-stationary conditions, too much reliance on

587    previous positive outcomes may be ineffective since the task conditions continuously

588    change and new task solution have to be explored.

589    Our results imply that S+ outcomes reinforce recent behaviour and inform an

590    exploitation strategy (low variability state) while S- outcomes will stimulate and inform

591    an exploration strategy (high variability state). These distinct motor states might be

592    differently encoded at the level of neural oscillations. Here, we demonstrate that

593    prefrontal neural oscillations in theta and high beta frequencies respond selectively to

594    different outcome histories during reinforcement motor learning, a mechanism that

595    potentially gates future adjustments in motor output. Activity in the prefrontal cortex is

596    sensitive to distinct motor states, characterized by marked differences in variability.

597    For instance, activity in regions of the PFC is modulated when participants switch

598    between exploratory and exploitative behavioural modes during decision-making tasks

599    (Daw *et al.*, 2006). Moreover, neural activity in the prefrontal cortex is sensitive to the

600    recent history of rewards, which might function as an update on estimates of predicted

601    rewards (Kim & Shadlen, 1999; Barraclough *et al.*, 2004; Padoa-Schioppa & Assad,

602    2006; Rushworth & Behrens, 2008; Seo & Lee, 2008; Histed *et al.*, 2009). The cerebral

603    regions that are involved in RL include prefrontal cortical areas such as SFG, RMFG

604    and cingulate cortex (Garrison *et al.*, 2013). Recently, it has been suggested that the

605    prefrontal cortex may serve as a meta-reinforcement learning system, which is

606    controlled by the midbrain dopamine system but acts as an independent learning

607    system (Wang *et al.*, 2018). Previous human decision-making studies (involving e.g.

608    gambling tasks) have shown that neural oscillations in prefrontal cortex respond to

609    reward-related feedback stimuli (Luft, 2014). Neural oscillations reflect the

610    synchronous activity of neural assemblies and have been considered important to

611    integrate large-scale networks (Fries, 2005), to facilitate synaptic plasticity during

612    learning (Buzsaki & Draguhn, 2004) and to engage in the control of top-down

613    information flow (Engel *et al.*, 2001). These studies suggest that the PFC might also

614    be Involved in processing the outcome of previous motor actions that can be used to

615    guide subsequent movements during motor tasks.

616    The results of the present study support this assumption by showing that neural

617    oscillations in the SFG and RMFG respond selectively to positive and negative motor

618    experiences. We showed that theta oscillations were most prominent after S-

619    movements when they were preceded by S+ outcomes and after S+ movements when

620    they were preceded by S- movements. Theta oscillations are present during cognitive

621    tasks and modulated by attentional and working memory demands (Kubota *et al.*,

622    2001; Onton *et al.*, 2005). The increase could reflect an increased cognitive load that

623    is caused by the detection of conflicts and/or errors, e.g. a mismatch between predicted

624    and actual rewards and thereby induce a change in the variability state (from low TTV

625    to high TTV or vice versa).

626    In contrast, high beta oscillations in prefrontal cortex were most prominent after S+

627    movements; and especially so when they were preceded by an S- movement. This

628    indicates that the reinforcement signal is indeed more distinct when the learner

629    experiences unexpected (positive) outcomes (Akitsuki *et al.*, 2003; HajiHosseini *et al.*,

630    2012). The increased beta activity could constitute a potential mechanism that ensures

631    that the current S+ action is reinforced while previous negative actions are assigned

632    less value and, in this way, beta activity could inform the model of future motor

633    behaviour. In other words, greater beta oscillations in the prefrontal cortex after S+

634    trials reinforce the most recent motor behaviour and inherently also signal less weight

635    on previous outcomes in planning of the next movement.

636    It is likely that oscillations in different frequency ranges and neuronal circuits interact

637    to support a common goal. For example, the interaction of theta and beta band

638    oscillations has been proposed to process and store short-term memories (Lisman &

639    Idiart, 1995; Axmacher *et al.*, 2010). Likewise, increases in the functional connectivity

640    between the striatum and the prefrontal cortex have been observed during categorical

641    learning (Antzoulatos & Miller, 2014). A role for the dopaminergic system is supported

642    by the fact that activity of dopaminergic neurons is modulated by reward stimuli in

643    primates (Mirenowicz & Schultz, 1994). Future studies have to clarify how interactions

644    between different neural circuits engage in RL.

645    The present study contains a number of limitations that should be considered. Due to

646    the experimental requirements when combining reinforcement motor learning with

647    EEG measurements in the present study, we were not able to sample more than 320

648    movements per participant in the main protocol. Based on this, it was not feasible to

649    analyse further outcome sequences including "older" outcomes since these would

650    have contained a considerably smaller number of trials. However, on the basis of a

651    previous study in rats (Dhawale *et al.*, 2019), it can be assumed that "older" outcomes

652    also have an impact on the behavioural trial-to-trial adjustments during reinforcement

653    motor learning in humans.

654    Based on the behavioural measures, we are not able to discern to which extent

655    changes in TTV are the consequence of intended motor variability or simply a by-

656    product of sensorimotor noise or error in motor acuity i.e. unintended motor variability.

657    While changes in TTV after S- trials can be due to sensorimotor noise and active

658    exploration, changes in TTV after S+ trials should be mainly caused by sensorimotor

659    noise (van Mastrigt *et al.*, 2020). Thus, motor variability is influenced by a number of

660    factors including intended and unintended variability regulatory mechanisms (Pekny *et*

661    *al.*, 2015; Therrien *et al.*, 2018; Dhawale *et al.*, 2019; van Mastrigt *et al.*, 2020).

662    Nevertheless, the results support the notion that negative outcomes inform an

663    explorative strategy leading to higher TTV whereas positive outcomes lead to

664    exploitation, which is characterized by low TTV.

665    We focused our analyses of neural oscillations to prefrontal cortical regions since we

666    had a strong a-priori hypothesis that a prefrontal cortical network is involved in RL. It

667    is however very likely that the cortical activity is influenced also by other circuits

668    involving the basal ganglia, which also play an important role in reward-based learning.

669    For instance, regulation of motor variability is impaired after inactivation of a cortico-

670    basal ganglia circuit of songbirds (Olveczky *et al.*, 2005) and in human Parkinson

671    patients (Pekny *et al.*, 2015). Future studies should investigate the interactions of the

672    different network nodes and their influence on motor output. While the present results

673    do not provide a causal link between the behavioural observations and the modulations

674    in neural oscillations, previous studies do however suggest a clear association

675    between modulations in neural oscillations and behavioural adjustments (Luft, 2014).

676    In conclusion, we provide evidence that positive outcomes "overwrite" previous motor

677    states to a greater extent than negative outcomes and that changes in high beta and

678    theta oscillatory activity in the prefrontal cortex potentially reflect changes in the

679    movement variability state during reinforcement motor learning. Therefore, the results

680    provide novel insights on the behavioural and neural mechanisms that underline

681    learning from previous outcomes.

682

683

684

685

686

687

688

689

690

691

692

693  **Methods**

694  *Participants and study design*

695  Twenty-six participants were recruited for the study by convenience sampling using

696  flyers and social media. All participants were healthy young adults (mean age: 25.4 $\pm$

697  2.7 years old, 16 females). Participants had to be free of any neurological or psychiatric

698  diseases and have normal or corrected-to normal vision to be included in the study.

699  According to the Edinburgh Handedness questionnaire (Oldfield, 1971), participants

700  were right-handed (laterality index of 82.2 $\pm$ 42.4). The study was approved by the

701  regional ethical committee (H-17019671). The study conformed to the standards set

702  by the Declaration of Helsinki (latest revision in Fortaleza, Brazil). All participants gave

703  their written informed consent to the procedures of the study prior to participation.

704  In summary, each participant performed a reinforcement motor task and experienced

705  different outcomes (target hit or miss) dependent on their motor performance. We

706  compared the effects of the outcome on behavioural and oscillatory responses in an

707  experimental within-participant design with random effects.

708

709  *Experimental setup*

710  Participants sat in a laboratory chair approximately 50 cm in front of a computer screen

711  (27" monitor with 60 Hz frame rate and 2,560 x 1,440-pixel resolution). The height of

712  the computer screen was individually adjusted. Participants positioned their left

713  forearm in a neutral position in a splint that was placed on a table next to them. The

714  left hand grabbed a handle with a built-in goniometer and could be moved by

715  performing wrist flexion/extension movements. The forearm and hand position were

716  stabilized and supported with Velcro® straps to avoid changes in elbow and shoulder

717  joint angles. The view of the forearm and hand was hidden by a custom-built box to

718    prevent visual feedback of the moving hand. Goniometer data were recorded at a

719    sampling frequency of 2048 Hz (CED 1401+ with Signal 3.09 software, Cambridge

720    Electronic Design Ltd., UK) and stored offline for further analyses.

721

722    *Motor task*

723    All participants performed wrist flexion movements with their left wrist (Fig. 1a). The

724    goal of the task was to move a circular cursor (radius of 15 pixels) into a target area

725    on a computer screen. The position of the target area was changed several times

726    during the experiment. There were three different target positions (see experimental

727    protocol) and each target had a horizontal size of 330 pixels (green target: 1130 to

728    1460 pixels, grey target: 1680 to 2010 pixels, purple target: 2230 to 2560 pixels) (Fig.

729    1b).

730    In each trial, participants performed one wrist flexion movement. Visual traffic lights

731    and text on the computer screen marked the beginning of a trial (Fig. 1c). The computer

732    cursor always started on the left side of the screen. The start of a trial was signalled

733    by the appearance of a red dot that was positioned on the left side of the screen. After

734    500 ms a yellow dot appeared underneath the red dot and instructed participants to

735    prepare the movements. After another 500 ms a green dot appeared as the final GO

736    cue. Next, participants could move the cursor to the right by performing wrist flexion

737    movements.

738    Participants were instructed to not adjust the movement end position once they

739    finished their movement. To avoid post-movement corrections, cursor motion was

740    insensitive to wrist movements once the differences in cursor position fell below five

741    pixels in five consecutive samples. Since the goal of the motor task was an accurate

742    movement end position, participants were informed prior to the experiment that the

743　task was not a reaction time task. Binary feedback about motor performance (S+ or S-

744　) was given 2500 ms after the GO signal for 1000 ms. Participants were explicitly

745　instructed to relax their arm during that period. After information of the outcome,

746　participants were instructed to move back to the starting position and wait for the start

747　of the next trial. If flexion movement time was longer than 800 ms in two consecutive

748　trials, we asked participants to move faster in the next trial. The next trial was started

749　after 5500 ms. Thus, each trial lasted 10 s.

750　A movement was considered S+ if the centre of the cursor was within the target area.

751　Success was indicated by a text box indicating "Good job". After S- trials, participants

752　were prompted with a text box stating "Try again". When the trial was S+, a point was

753　added to the participant's score. The participant's score was displayed continuously

754　throughout a training block (40 trials) with the aim of motivating the participants in each

755　block. All participants were compensated equally for the time they had spent in the

756　laboratory and compensation was not based on their motor performance. The positions

757　and size of the target areas were determined in pilot experiments (n = 8, results are

758　not reported in this paper) such that participants yielded on average approximately an

759　equal number of S+ and S- trials. The task was created using MATLAB R2019a (The

760　MathWorks, Natick, Massachusetts, United States, R2019a) and the Psychophysics

761　Toolbox extensions (Brainard, 1997). Cursor movements (in pixels), distance from the

762　targets (in pixels) and information on movement outcome for each trial were logged

763　online and saved for further offline analyses.

764

765　*Experimental protocol*

766　Prior to the study, all participants were informed about the experimental protocol and

767　the motor task. Subsequently, participants were introduced to the task by performing

768 ten wrist flexions with no targets displayed on the screen. After this short introduction,

769 participants performed 40 movements with online visual feedback (Fam1) and 40

770 without online visual feedback (Fam2) of cursor motion and target position. In both

771 familiarization blocks, participants had to reach for the grey middle target and they

772 received offline binary feedback about their performance at the end of the trial

773 (knowledge of result). The target position did not change during these blocks. These

774 initial blocks were performed to ensure that participants could adhere to the time

775 course of a trial, to familiarize them with the handle's sensitivity and to quantify baseline

776 movement variability.

777 In the main protocol, participants performed 320 wrist flexion movements in blocks of

778 40 movements. Between blocks, we incorporated breaks of 2 minutes to ensure that

779 participants focused on the motor task the whole experiment. Pilot experiments

780 demonstrated that for most participants 40 consecutive movements (approx. 7

781 minutes) was feasible and could be performed with sufficient and sustained attention.

782 Moreover, the target position was changed several times during the experiment. The

783 green and purple targets were present in 80 trials, respectively, (i.e. 160 trials) and the

784 grey target in the remaining 160 trials. Thus, in total, participants performed 160

785 movements to an already familiarized target position (grey target) and 160 movements

786 to an unfamiliar target position (green or purple target). Targets were changed every

787 25, 30 or 40 trials such that there was no systematic order of the different conditions.

788 The size of the target area remained constant throughout the experiment.

789 In the main protocol, participants did receive knowledge of result but not knowledge of

790 performance. All participants were informed that the target position could change at

791 any time during the experiment. Giving this information, we wanted to avoid too much

792    frustration (which could influence participant's motivation) but also stimulate the

793    exploration of different motor actions following a history of unsuccessful trials.

794

795    *Pre-processing of behavioural data*

796    In the offline analysis, we analysed movement kinematics from the goniometer signal

797    of the wrist. In brief, continuous data was epoched from -1000 ms to 3000 ms relative

798    to the GO cue to capture start and end points of all movements. For these epochs, we

799    calculated the first derivate of the goniometer signal for all single trials. Subsequently,

800    we smoothed the derivates by calculating the moving average of 400 samples and

801    applying a 3$^{rd}$ order Butterworth lowpass filter (20 Hz). We used the smoothed signals

802    to determine the movement start and end points by calculating the maximum

803    movement speed within each epoch. Movement start was defined as the sample for

804    which signals exceeded 15% of its respective speed maximum. Movement end was

805    defined as the sample for which the signals dropped below 15% of its respective speed

806    maximum. The movement end point angle was calculated by subtracting movement

807    end angles from movement start angles (in °). During this process, all data were

808    continuously visually inspected and checked. In addition to the kinematic analysis, we

809    also extracted the information on trial outcomes (hits or misses). Trials in which the

810    participants did not move at all were discarded (movement angle < 2°).

811

812    *Behavioural data analyses*

813    First, we were interested in the association of the outcome time series (S+ and S- trials)

814    with lagged versions of the same outcome time series. Thus, we performed partial

815    autocorrelation (PAC) for each participant (Matlab function: *parcorr*). This analysis was

816    performed to obtain an initial estimate of the impact of previous outcomes on future

817    outcomes and to test the assumption that the past two outcomes have the greatest

818    association with future outcomes.

819    In the remainder of our analysis, we focused on changes in movement endpoint ($\Delta\mu$)

820    as a function of outcomes from the previous two trials. There were two reasons for

821    restricting the analyses to the last two outcomes. First, the PAC analysis confirmed

822    that the past two outcomes have the greatest correlation with future outcomes.

823    Second, outcome histories considering even "older" actions would have contained too

824    few trials (<10) to perform valid comparisons.

825    In the following, trial$^{(n)}$ refers to the most recent trial and trial$^{(n+1)}$ to the subsequent trial.

826    We measured the differences in movement endpoint ($\Delta\mu$) between trial$^{(n+1)}$ and trial$^{(n)}$

827    as a measure of TTV using the following formula:

828

829
$$\Delta\mu \ = \ \mu^{(n+1)} - \mu^{(n)}$$

830

831    For each participant, we extracted the signed TTV in movement endpoint ($\Delta\mu$)

832    conditioned on the outcome of the previous trial$^{(n)}$. The signed TTV gives an estimate

833    on the directional changes in movement endpoint and potentially reveal whether

834    changes in movement endpoint are biased in one or the other direction. To quantify

835    unsigned TTV in movement endpoint, we also calculated the absolute TTV in

836    movement endpoint ($|\Delta\mu|$). The unsigned TTV is a measure of the absolute motor

837    exploration independent of the direction of the change.

838

839        *Signed: $\Delta\mu|S+^{(n)}$*              *Unsigned: $|\Delta\mu||S+^{(n)}$*

840        *Signed: $\Delta\mu|S-^{(n)}$*              *Unsigned: $|\Delta\mu||S-^{(n)}$*

841

842   The first movement of every block were discarded, as these had no previous trial.

843   Subsequently, individual data histograms from the different conditions were visually

844   inspected. In accordance with a previous study (Pekny *et al.*, 2015), signed data

845   showed a normal distribution and unsigned data had a negative exponential

846   distribution. Thus, we fitted normal distributions as the conditional probability

847   distribution for the signed data $p(\Delta\mu|S+^{(n)})$ and $p(\Delta\mu|S-^{(n)})$. For the unsigned data, we

848   fitted exponential distributions as the conditional probability distribution $p(|\Delta\mu||S+^{(n)})$

849   and $p(|\Delta\mu||S-^{(n)})$. The mean (M) and standard deviation (SD) were calculated from the

850   individual normal distributions as a measure of the signed TTV and the variations in

851   the TTV, respectively. The M was calculated from the individual exponential distribution

852   as a measure of the absolute TTV.

853   In the next step, we asked whether TTV is dependent on the previous two outcomes,

854   i.e. the most recent trial$^{(n)}$ and the second-to-last trial$^{(n-1)}$. Thus, we extracted TTV from

855   four different outcome histories:

856

857   *Signed: $\Delta\mu|S+^{(n)}$ & $S+^{(n-1)}$*       *Unsigned: $|\Delta\mu||S+^{(n)}$ & $S+^{(n-1)}$*

858   *Signed: $\Delta\mu|S+^{(n)}$ & $S-^{(n-1)}$*       *Unsigned: $|\Delta\mu||S+^{(n)}$ & $S+^{(n-1)}$*

859   *Signed: $\Delta\mu|S-^{(n)}$ & $S+^{(n-1)}$*       *Unsigned: $|\Delta\mu||S-^{(n)}$ & $S+^{(n-1)}$*

860   *Signed: $\Delta\mu|S-^{(n)}$ & $S-^{(n-1)}$*       *Unsigned: $|\Delta\mu||S-^{(n)}$ & $S-^{(n-1)}$*

861

862   Since these four conditions contained a significantly different number of trials, we

863   performed bootstrapping with replacement for each individual (Matlab function:

864   *datasample*). As suggested previously (Hesterberg, 2011), we sampled 1,000 new

865   datasets per participant and condition, each of them containing 100 randomly chosen

866   trials from the original data. We chose to use this bootstrapping approach since

867    differences in trial numbers can significantly bias measures of variability in the distinct

868    outcome histories (e.g. the SD). From the individual bootstrapped probability

869    distributions, we computed the M and SD as described above to derive robust

870    confidence intervals for estimating standard errors and for hypothesis testing. Again,

871    we fitted normal ($\Delta\mu$) and exponential ($|\Delta\mu|$) distributions for changes in movement

872    endpoint after $S+^{(n)}$ and $S-^{(n)}$ movements that were either preceded by $S+^{(n-1)}$ or $S-^{(n-1)}$

873    movements, respectively. Probability distributions were fitted using the MATLAB

874    functions *fitdist* and *makedist.* All described analyses were performed offline using

875    MATLAB R2019a (The MathWorks, Natick, Massachusetts, United States, R2019a).

876

877    *Statistical analyses of TTV*

878    Statistical analyses were performed using SPSS software 27 (SPSS®, Chicago, IL,

879    USA) and RStudio (RStudio Team, 2018). All dependent variables were tested for

880    normality by the Kolmogorov-Smirnov test. Homogeneity of variances was tested using

881    the Fisher's F-test. Repeated-measures analyses of variance (rmANOVA) with random

882    effects were calculated to test the effect of outcomes on the signed and unsigned TTV

883    and the SD of the signed TTV. The rmANOVA design had one within- participant factor

884    with 4 levels defining the outcome history ($S+^{(n)}$ & $S+^{(n-1)}$, $S+^{(n)}$ & $S-^{(n-1)}$, $S-^{(n)}$ & $S+^{(n-1)}$,

885    $S-^{(n)}$ & $S-^{(n-1)}$). The Greenhouse–Geisser correction was used for rmANOVAs if the

886    assumption of sphericity was violated (Mauchly's test). Effect sizes were estimated

887    using Eta-squared ($\eta2$ partial). Within-participant comparisons between the different

888    outcomes were performed with appropriate paired t-tests or F-tests. Paired t-tests were

889    calculated for comparisons of the distance from target (in pixels), the number of S+

890    and S- movements and the M values from the normal distribution of the signed

891    changes in movement endpoint. F-tests were calculated for of the SD values and the

892   M values from the exponential distribution. The test of equal or given proportions was

893   used to compare the proportion of S+$^{(n+1)}$ movements after the different outcomes. A

894   P-value correction using the Benjamini-Hochberg procedure was applied to control the

895   False discovery rate (FDR). Uncorrected p-values are presented in the Results section.

896   A significance level of $P < 0.05$ was assumed. All data are reported as M $\pm$ SD if not

897   stated otherwise.

898

899   *EEG recordings*

900   EEG recordings were performed with a 64 channel Biosemi system (BioSemi,

901   Amsterdam, The Netherlands) using the software ActiView (version 8.06). Electrodes

902   were embedded into an electrode cap and positioned according to the extended 10-

903   20 system. Data recording was performed at 2048 Hz and the online reference was

904   the Common Mode Sense (CMS)/Drive Right Leg (DRL). It was ensured that the

905   electrode offsets (voltage differences between single electrodes and CMS) were less

906   than $\pm20$ µV. EEG signals were recorded continuously throughout the experiment.

907   During the motor task, information of task timing and movement outcome was inputted

908   (via a trigger signals) from Matlab to synchronize the task and the EEG recording and

909   allow time-locked offline analyses of the EEG data. All participants were asked to relax

910   their face and neck muscles during the experiment to minimize signal artefacts.

911

912   *EEG pre-processing*

913   EEG data was pre-processed using EEGlab (Delorme & Makeig, 2004). Data were

914   imported and the DC offset was corrected by subtracting the mean from each channel.

915   Noisy channels were identified by visual analysis of the time-series (high amplitude

916   time series, e.g. from tonic muscle activity) and frequency plots of the entire recording

917   period (strong power deviations and/or unusual spikes in the frequency range 0.5 Hz

918   – 48 Hz). Subsequently, we removed 4.48 ± 1.96 channels per participant which were

919   mostly located temporally and/or occipitally. EEG signals were i) re-referenced to the

920   average of all channels, ii) bandpass-filtered between 0.5 Hz and 48 Hz using a FIR

921   filter and iii) down-sampled to 256 Hz. The pre-processed EEG data were segmented

922   into epochs ranging from -3500 ms to +3500 ms relative to the beginning of information

923   of the outcome. Using the segmented data, we performed independent component

924   analyses (ICA) with the runica algorithm. The ICA was performed on segmented data

925   as our continuous recordings also encompassed a great amount of irrelevant data (e.g.

926   breaks). The aim of this step was to identify and remove components indicating

927   horizontal and vertical eye movements and saccades based on the topological and

928   spectral characteristics of the components (Chaumon *et al.*, 2015). On average, we

929   selected and removed 2.12 ± 0.33 components per participant. Finally, removed

930   channels were interpolated with the standard spherical method.

931   All epochs were visually inspected from each participant to identify and remove trials

932   that were contaminated by signal artefacts (e.g. high amplitude signal deviations from

933   strong muscle artefacts in most channels). Following this assessment, we had to

934   remove the data from one participant since more than 25% of all trials had to be

935   excluded. For the remaining 25 participants, we removed on average 4.26 ± 4.4% trials

936   (out of 320) per participant.

937   For each participant, we created datasets encompassing i) all trials independent of the

938   outcome, ii) $S+^{(n)}$ and $S-^{(n)}$ trials and ii) $S+^{(n)}$ & $S+^{(n-1)}$ trials, $S+^{(n)}$ & $S-^{(n-1)}$ trials, $S-^{(n)}$ &

939   $S+^{(n-1)}$ trials and $S-^{(n)}$ & $S-^{(n-1)}$ trials. All datasets were saved offline for further analyses.

940   The subsequent EEG analyses comprised the following steps. First, we tested the

941   assumption that the information of movement outcome leads to changes in oscillatory

942    activity (normalised to a "pre-feedback" period) in frontal sensors and in prefrontal

943    cortex in source space. Subsequently, we constructed power time courses of two

944    prefrontal cortical regions previously shown to engage in RL, i.e. SFG and rostral

945    middle frontal gyrus RMFG (Garrison *et al.*, 2013), to compare power time series

946    between different outcomes. The regions of interest (ROI) are highlighted in

947    Supplementary Fig. 3. In our analyses, we focussed on pre-specified frequency ranges

948    and time windows of interest (see Time-frequency analyses in sensor space).

949

950    *Time-frequency analyses in sensor space*

951    Further data analyses were performed on the pre-processed segmented data using

952    Brainstorm (Tadel *et al.*, 2019), which is documented and freely available for download

953    online under the GNU general public license (http://neuroimage.usc.edu/brainstorm).

954    All datasets were imported to Brainstorm in the time range -2000 ms to +2000 ms

955    relative to information of the outcome. The reason for this time window was a

956    compromise between avoiding filtering related edge artefacts of time-frequency

957    analyses and minimising computational demands.

958    All single trials were transformed to the time-frequency domain by convolving the signal

959    with a set of complex Morlet wavelets, which are defined as complex sine waves

960    tapered by a Gaussian. The full-width at half-maximum was 3000 ms and the sine

961    waves were created at 1 Hz (corresponding to 7 cycles). The analyses were

962    constrained to the theta band (4 – 8 Hz) and the higher beta band (25 – 35 Hz). The

963    rationale of this choice was that there is ample evidence that oscillatory modulations

964    in both frequency bands represent different reinforcement signals during learning from

965    feedback (Luft, 2014). For simplicity, the frequency range 25 – 35 Hz will be referred

966    to as high beta although it also includes low gamma frequencies per some definitions.

967    Single time-frequency series were averaged for each participant and frequency band

968    (theta and high beta) at each electrode. In order to calculate event-related changes in

969    power, the averaged data were scaled with the mean of a "pre-feedback" period prior

970    to outcome information (-400 ms to -200 ms). This period covers the time where

971    participants had already finished the movement and were awaiting outcome

972    information. This normalisation procedure ensures that the data across all time points,

973    sensor/source points, conditions and subjects are in the same scale and hence

974    comparable. Changes in power were calculated in %, that is $(x/u) *100$ where x is the

975    data and u is the mean over the "pre-feedback" period. This analysis was run on the

976    dataset encompassing all trials independent of previous outcomes. Data were visually

977    inspected by plotting selected sensors and topographical plots on pre-specified time

978    windows (250 – 550 ms) on the basis of previous studies (Cohen *et al.*, 2007; Marco-

979    Pallares *et al.*, 2008; HajiHosseini *et al.*, 2012; Luft, 2014; HajiHosseini & Holroyd,

980    2015; Marco-Pallares *et al.*, 2015).

981

982    *Source modelling*

983    Motivated by the results in sensor-space, we performed source analyses on the pre-

984    processed sensor time series to localize and reconstruct the cortical regions

985    contributing to oscillatory modulations during RL of motor skills. In a first step, we

986    created our forward model describing how neuronal activity propagates from each

987    cortical position to the EEG sensors, also called the lead field matrices. Since individual

988    MRIs were not available, we used the MNI International Consortium of Brain

989    Mapping152 brain template which is a non-linear average of 152 participants (Fonov

990    *et al.*, 2009). The forward model was constructed using the symmetric boundary

991    element method (BEM) from the open-source software OpenMEEG (Gramfort *et al.*,

992  2010). The BEM uses three realistic layers (head, outer skull and inner skull; 1922

993  vertices per layer) to calculate the volume conduction model. The relative conductivity

994  of the layers was [1 0.0125 1] which describes the relative conductivities of each layer

995  and is the default in Brainstorm. Standard BioSemi sensor positions were aligned to

996  the template head space. The source space contained 15002 elements constrained to

997  the cortical sheet. The number of vertices has been suggested to be sufficient to

998  sample the folded details of the cortex (Tadel *et al.*, 2019).

999  The inverse solution was calculated using the weighted minimum norm estimation

1000  method and the measure standardized low resolution brain electromagnetic

1001  tomography as implemented in the Brainstorm software (Hamalainen & Ilmoniemi,

1002  1994; Baillet *et al.*, 2001; Pascual-Marqui, 2002). We chose to use a distributed source

1003  imaging model rather than a single dipole model since we expected multiple cortical

1004  regions to be modulated during feedback processing. Source activity (in the time

1005  domain) was estimated for sources with unconstrained orientation, that is source

1006  activity was calculated for three dipole orientations at each cortical location. The

1007  estimation of source activity via minimum norm estimators takes into account the level

1008  of noise in the sensors and hence requires an estimation of the noise in the recordings

1009  (Hauk, 2004). Thus, noise statistics (noise covariance across all sensors) were

1010  calculated across all trials from a pre-stimulus time, i.e. prior to information of the

1011  outcome (-2000 ms to 0 ms). Finally, we obtained source space time series at each

1012  cortical position.

1013

1014  *Time-frequency analyses in source space*

1015  Since cortical sources cannot be directly inferred from scalp time-frequency data, we

1016  computed time-frequency decomposition on the source time series using Morlet-

1017    wavelet analyses. Wavelets had a full-width at half-maximum of 3000 ms and a

1018    frequency of 1 Hz (corresponding to 7 cycles). Again, the analyses were constrained

1019    to the averaged theta band (4 – 8 Hz) and the averaged higher beta band (25 – 35

1020    Hz).

1021    In the first step of the analyses, these analyses were run on the full cortical map (entire

1022    cortical sheet) for all trials independent of the outcome to test our assumption that

1023    regions of the prefrontal cortex respond to outcome information in general. In support

1024    of this view, we observed strong oscillatory power over the SFG and RMFG suggesting

1025    that both regions engage in processing motor outcomes. Thus, in the following steps

1026    we constrained the analyses to SFG and RMFG. These were defined from pre-

1027    specified cortical areas of the Desikan-Killiany parcellation scheme which subdivides

1028    the cortex into gyral based ROI (Desikan *et al.*, 2006).

1029    Subsequently, time-frequency data of the ROI were computed for $S+^{(n)}$ and $S-^{(n)}$ trials.

1030    In our second analyses, we again faced the problem of differences in trial numbers

1031    between the different outcomes for every participant. This could result in meaningful

1032    differences in signal-to-noise ratios and bias the comparisons between the conditions

1033    (see also behavioural data analysis). Thus, we performed bootstrapping by sampling

1034    multiple new datasets with replacement. This resulted in 100x1,000 random trials per

1035    participant and condition. As for the behavioural analysis, trials were grouped

1036    according to the outcomes of the previous two trials in $S+^{(n)}$ & $S+^{(n-1)}$, $S+^{(n)}$ & $S-^{(n-1)}$, $S-^{(n)}$

1037    $^{(n)}$ & $S+^{(n-1)}$ and $S-^{(n)}$ & $S-^{(n-1)}$. The averaged and normalised data were used for

1038    statistical analyses.

1039

1040

1041

*Statistics on time-frequency data in source space*

Statistical analyses were performed on time-frequency source time series in the pre-selected time window of interest (250 ms – 550 ms after outcome information). To reduce the number of tests, statistical analyses was performed on specified data points (in steps of approximately 50 ms, i.e. at 250 ms, 300 ms, … 550 ms, according to the sampling frequency). Further, we constrained our analyses to selected comparisons of interest. In the first step, we compared theta and high beta power time courses for all ROI between $S+^{(n)}$ and $S-^{(n)}$ trials. In the second step, we compared theta and high beta time courses for all ROI between outcomes i) $S+^{(n)}$ & $S+^{(n-1)}$ *and* $S+^{(n)}$ & $S-^{(n-1)}$ and ii) $S-^{(n)}$ & $S+^{(n-1)}$ and $S-^{(n)}$ & $S-^{(n-1)}$. In all cases, we tested whether the data fulfil the criteria for statistics on normally distributed data. For this purpose, we plotted histograms and performed Kolmogorov-Smirnov tests. Given that the data were not normally distributed, we performed Wilcoxon signed rank tests between all comparisons of interest. The FDR method was used to correct for multiple comparisons. Critical P-values are presented in the Results section, i.e. the adjusted significance threshold after correcting for the FDR. MATLAB R2019a (The MathWorks, Natick, Massachusetts, United States, R2019a) was used to compute all statistical analyses of the EEG data.

**Author contributions**

1073

1074　P.W., M.E.S., C.R., M.M.B. and J.L.J. conceived and designed the research. P.W.

1075　performed the experiments. P.W. and C.R. analysed the data. P.W., M.E.S., C.R.,

1076　M.M.B. and J.L.J interpreted the results of experiments. P.W. drafted the paper. P.W.,

1077　M.E.S., C.R., M.M.B. and J.L.J. edited and revised the manuscript. All authors

1078　approved the final version of manuscript.

1079

**Competing interests**

1080

1081　The authors declare no conflicts of interest.

1082

**Data and code availability**

1083

1084　Data (Behaviour) and code (Behaviour and EEG) are available at

1085　https://data.mendeley.com (DOI: 10.17632/cw73pv9ct4.1). EEG data were not

1086　uploaded due to their large size but are available from the corresponding author on

1087　request.

1088

1089

1090

1091

## References

Akitsuki, Y., et al. (2003). Context-dependent cortical activation in response to financial reward and penalty: an event-related fMRI study. *Neuroimage* **19,** 1674-1685.

Antzoulatos, E. G. & Miller, E. K. (2014). Increases in functional connectivity between prefrontal cortex and striatum during category learning. *Neuron* **83,** 216-225.

Axmacher, N., et al. (2010). Cross-frequency coupling supports multi-item working memory in the human hippocampus. *Proc Natl Acad Sci U S A* **107,** 3228-3233.

Baillet, S., Mosher, J. C. & Leahy, R. M. (2001). Electromagnetic brain mapping. *IEEE Signal Processing Magazine* **18, no. 6, 14-30**.

Barraclough, D. J., Conroy, M. L. & Lee, D. (2004). Prefrontal cortex and decision making in a mixed-strategy game. *Nat Neurosci* **7,** 404-410.

Brainard, D. H. (1997). The Psychophysics Toolbox. *Spat Vis* **10,** 433-436.

Buzsaki, G. & Draguhn, A. (2004). Neuronal oscillations in cortical networks. *Science* **304,** 1926-1929.

Chaumon, M., Bishop, D. V. & Busch, N. A. (2015). A practical guide to the selection of independent components of the electroencephalogram for artifact correction. *J Neurosci Methods* **250,** 47-63.

Cohen, M. X., Elger, C. E. & Ranganath, C. (2007). Reward expectation modulates feedback-related negativity and EEG spectra. *Neuroimage* **35,** 968-978.

Daw, N. D., et al. (2006). Cortical substrates for exploratory decisions in humans. *Nature* **441,** 876-879.

Delorme, A. & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J Neurosci Methods* **134,** 9-21.

Desikan, R. S., et al. (2006). An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage* **31,** 968-980.

Dhawale, A. K., Miyamoto, Y. R., Smith, M. A. & Olveczky, B. P. (2019). Adaptive Regulation of Motor Variability. *Curr Biol* **29,** 3551-3562 e3557.

Engel, A. K., Fries, P. & Singer, W. (2001). Dynamic predictions: oscillations and synchrony in top-down processing. *Nat Rev Neurosci* **2,** 704-716.

1136 Fonov, V. S., et al. (2009). Unbiased nonlinear average age-appropriate brain templates from
1137     birth to adulthood. *NeuroImage 47(Suppl 1):S102 doi: 101016/S1053-8119(09)70884-*
1138     *5.*

1140 Fries, P. (2005). A mechanism for cognitive dynamics: neuronal communication through
1141     neuronal coherence. *Trends Cogn Sci* **9,** 474-480.

1143 Galea, J. M., Mallia, E., Rothwell, J. & Diedrichsen, J. (2015). The dissociable effects of
1144     punishment and reward on motor learning. *Nat Neurosci* **18,** 597-602.

1146 Galea, J. M., et al. (2013). Punishment-induced behavioral and neurophysiological variability
1147     reveals dopamine-dependent selection of kinematic movement parameters. *J*
1148     *Neurosci* **33,** 3981-3988.

1150 Garrison, J., Erdeniz, B. & Done, J. (2013). Prediction error in reinforcement learning: a meta-
1151     analysis of neuroimaging studies. *Neurosci Biobehav Rev* **37,** 1297-1310.

1153 Gramfort, A., Papadopoulo, T., Olivi, E. & Clerc, M. (2010). OpenMEEG: opensource software
1154     for quasistatic bioelectromagnetics. *Biomed Eng Online* **9,** 45.

1156 HajiHosseini, A. & Holroyd, C. B. (2015). Reward feedback stimuli elicit high-beta EEG
1157     oscillations in human dorsolateral prefrontal cortex. *Sci Rep* **5,** 13021.

1159 HajiHosseini, A., Rodriguez-Fornells, A. & Marco-Pallares, J. (2012). The role of beta-gamma
1160     oscillations in unexpected rewards processing. *Neuroimage* **60,** 1678-1685.

1162 Hamalainen, M. S. & Ilmoniemi, R. J. (1994). Interpreting magnetic fields of the brain:
1163     minimum norm estimates. *Med Biol Eng Comput* **32,** 35-42.

1165 Hauk, O. (2004). Keep it simple: a case for using classical minimum norm estimation in the
1166     analysis of EEG and MEG data. *Neuroimage* **21,** 1612-1621.

1168 Hesterberg, T. (2011). Bootstrap. *WIREs Computational Statistics* **3,** 497-526.

1170 Histed, M. H., Pasupathy, A. & Miller, E. K. (2009). Learning substrates in the primate
1171     prefrontal cortex and striatum: sustained activity related to successful actions. *Neuron*
1172     **63,** 244-253.

1174 Kim, J. N. & Shadlen, M. N. (1999). Neural correlates of a decision in the dorsolateral prefrontal
1175     cortex of the macaque. *Nat Neurosci* **2,** 176-185.

1177 Knutson, B., Westdorp, A., Kaiser, E. & Hommer, D. (2000). FMRI visualization of brain activity
1178     during a monetary incentive delay task. *Neuroimage* **12,** 20-27.

1180 Kubota, Y., et al. (2001). Frontal midline theta rhythm is correlated with cardiac autonomic
1181     activities during the performance of an attention demanding meditation procedure.
1182     *Brain Res Cogn Brain Res* **11,** 281-287.

1183
1184    Levy, S., et al. (2020). Cell-Type-Specific Outcome Representation in the Primary Motor Cortex.
1185        *Neuron* **107,** 954-971 e959.
1186
1187    Lisman, J. E. & Idiart, M. A. (1995). Storage of 7 +/- 2 short-term memories in oscillatory
1188        subcycles. *Science* **267,** 1512-1515.
1189
1190    Luft, C. D. (2014). Learning from feedback: the neural mechanisms of feedback processing
1191        facilitating better performance. *Behav Brain Res* **261,** 356-368.
1192
1193    Marco-Pallares, J., et al. (2008). Human oscillatory activity associated to reward processing in
1194        a gambling task. *Neuropsychologia* **46,** 241-248.
1195
1196    Marco-Pallares, J., Munte, T. F. & Rodriguez-Fornells, A. (2015). The role of high-frequency
1197        oscillatory activity in reward processing and learning. *Neurosci Biobehav Rev* **49,** 1-7.
1198
1199    Mirenowicz, J. & Schultz, W. (1994). Importance of unpredictability for reward responses in
1200        primate dopamine neurons. *J Neurophysiol* **72,** 1024-1027.
1201
1202    Narayanan, N. S., Cavanagh, J. F., Frank, M. J. & Laubach, M. (2013). Common medial frontal
1203        mechanisms of adaptive control in humans and rodents. *Nat Neurosci* **16,** 1888-1895.
1204
1205    Oldfield, R. C. (1971). The assessment and analysis of handedness: the Edinburgh inventory.
1206        *Neuropsychologia* **9,** 97-113.
1207
1208    Olveczky, B. P., Andalman, A. S. & Fee, M. S. (2005). Vocal experimentation in the juvenile
1209        songbird requires a basal ganglia circuit. *PLoS Biol* **3,** e153.
1210
1211    Onton, J., Delorme, A. & Makeig, S. (2005). Frontal midline EEG dynamics during working
1212        memory. *Neuroimage* **27,** 341-356.
1213
1214    Padoa-Schioppa, C. & Assad, J. A. (2006). Neurons in the orbitofrontal cortex encode economic
1215        value. *Nature* **441,** 223-226.
1216
1217    Pascual-Marqui, R. D. (2002). Standardized low-resolution brain electromagnetic tomography
1218        (sLORETA): technical details. *Methods Find Exp Clin Pharmacol* **24 Suppl D,** 5-12.
1219
1220    Pekny, S. E., Criscimagna-Hemminger, S. E. & Shadmehr, R. (2011). Protection and expression
1221        of human motor memories. *J Neurosci* **31,** 13829-13839.
1222
1223    Pekny, S. E., Izawa, J. & Shadmehr, R. (2015). Reward-dependent modulation of movement
1224        variability. *J Neurosci* **35,** 4015-4024.
1225
1226    RStudio Team. (2018). RStudio: Integrated Development for R. RStudio Team, Inc., Boston, MA
1227        URL http://www.rstudio.com/.
1228

1229 Rushworth, M. F. & Behrens, T. E. (2008). Choice, uncertainty and value in prefrontal and
1230     cingulate cortex. *Nat Neurosci* **11,** 389-397.

1232 Schall, J. D., Stuphorn, V. & Brown, J. W. (2002). Monitoring and control of action by the frontal
1233     lobes. *Neuron* **36,** 309-322.

1235 Schultz, W. (2000). Multiple reward signals in the brain. *Nat Rev Neurosci* **1,** 199-207.

1237 Schultz, W. (2017). Reward prediction error. *Curr Biol* **27,** R369-R371.

1239 Seo, H. & Lee, D. (2008). Cortical mechanisms for reinforcement learning in competitive
1240     games. *Philos Trans R Soc Lond B Biol Sci* **363,** 3845-3857.

1242 Sutton, R. S. & Barto, A. G. (1998). *Introduction to Reinforcement Learning*. MIT Press.

1244 Tadel, F., et al. (2019). MEG/EEG Group Analysis With Brainstorm. *Front Neurosci* **13,** 76.

1246 Takikawa, Y., et al. (2002). Modulation of saccadic eye movements by predicted reward
1247     outcome. *Exp Brain Res* **142,** 284-291.

1249 Therrien, A. S., Wolpert, D. M. & Bastian, A. J. (2018). Increasing Motor Noise Impairs
1250     Reinforcement Learning in Healthy Individuals. *eNeuro* **5**.

1252 van Mastrigt, N. M., Smeets, J. B. J. & van der Kooij, K. (2020). Quantifying exploration in
1253     reward-based motor learning. *PLoS One* **15,** e0226789.

1255 Wang, J. X., et al. (2018). Prefrontal cortex as a meta-reinforcement learning system. *Nature
1256     Neuroscience* **21,** 860-+.

1258 Watanabe, M. (1996). Reward expectancy in primate prefrontal neurons. *Nature* **382,** 629-
1259     632.

1261 Wu, H. G., et al. (2014). Temporal structure of motor variability is dynamically regulated and
1262     predicts motor learning ability. *Nat Neurosci* **17,** 312-321.