

Pallidal neuromodulation of the explore/exploit trade-off in decision-making

Ana Luisa de A Marcelino^{1,2}, Owen Gray³, Bassam Al-Fatly¹, William Gilmour³, J Douglas Steele³, Andrea A Kühn^{1,2,4,5,6}, Tom Gilbertson^{3,7*}

¹Charité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin and Humboldt-Universität zu Berlin, Movement Disorder and Neuromodulation Unit, Department of Neurology, Charité Campus Mitte, Berlin, Germany; ²Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Core Facility Genomics, Berlin, Germany; ³Division of Imaging Science and Technology, Medical School, University of Dundee, Dundee, United Kingdom; ⁴Berlin School of Mind and Brain, Charité - University Medicine Berlin, Berlin, Germany; ⁵NeuroCure, Charité - University Medicine Berlin, Berlin, Germany; ⁶DZNE, German Centre for Degenerative Diseases, Berlin, Germany; ⁷Department of Neurology, Ninewells Hospital & Medical School, Dundee, United Kingdom

Abstract Every decision that we make involves a conflict between exploiting our current knowledge of an action's value or exploring alternative courses of action that might lead to a better, or worse outcome. The sub-cortical nuclei that make up the basal ganglia have been proposed as a neural circuit that may contribute to resolving this explore-exploit 'dilemma'. To test this hypothesis, we examined the effects of neuromodulating the basal ganglia's output nucleus, the globus pallidus interna, in patients who had undergone deep brain stimulation (DBS) for isolated dystonia. Neuromodulation enhanced the number of exploratory choices to the lower value option in a two-armed bandit probabilistic reversal-learning task. Enhanced exploration was explained by a reduction in the rate of evidence accumulation (drift rate) in a reinforcement learning drift diffusion model. We estimated the functional connectivity profile between the stimulating DBS electrode and the rest of the brain using a normative functional connectome derived from healthy controls. Variation in the extent of neuromodulation induced exploration between patients was associated with functional connectivity from the stimulation electrode site to a distributed brain functional network. We conclude that the basal ganglia's output nucleus, the globus pallidus interna, can adaptively modify decision choice when faced with the dilemma to explore or exploit.

*For correspondence: tgilbertson@dundee.ac.uk

Competing interest: [See page 18](#)

Funding: [See page 19](#)

Received: 21 April 2022

Preprinted: [22 April 2022](#)

Accepted: 01 February 2023

Published: 02 February 2023

Reviewing Editor: Birte U Forstmann, University of Amsterdam, Netherlands

© Copyright de A Marcelino et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

Editor's evaluation

This paper presents valuable data from 18 patients treated with GPi DBS for dystonia using a standard RL task. Their compelling main observation is that DBS reduced the impact of value on evidence accumulation leading to more exploratory choices which was supported by fitting a dynamic decision model to the data. This work will be interesting for scientists working in fundamental and clinical neurosciences.

Introduction

Flexible adaption of choice behaviour is crucial for decision-making when facing uncertainty. The balance between exploiting familiar options and exploring alternatives with potentially less favorable outcomes, represents a fundamental challenge of adaptive control ([Mehlhorn et al., 2015](#)). The

neural basis of how humans resolve this dilemma has been advanced significantly over the last decade (**Chakroun et al., 2020; Wilson et al., 2014**), but the specific contribution from sub-cortical nuclei that make up the basal ganglia (BG) remains poorly defined (**Cohen et al., 2007; Sheth et al., 2011**).

To investigate how humans solve this conflict in reinforcement learning, so-called 'bandit tasks' have been used, where subjects choose between several options associated with different reward probabilities to maximise reward payoff. In a 'reversal-learning' paradigm, the reward probabilities change (i.e. reverse), thus subjects need to constantly evaluate whether choice feedback represents a reversal false alarm (i.e. a lower than expected outcome from the option with the highest reward probability) or a true reversal in the reward probabilities. Even in reversal-learning tasks with low levels of payout uncertainty, human participants make a significant number of choices to the lower value option, rather than the choice with the higher expected value. This leads to the observation of probability matching (**Herrnstein, 1970**), where humans fail to maximise their payoff and adopt a sub-optimal policy which matches their response probabilities to that of the payoff probability. 'Non-greedy' choices, are by definition exploratory, as they are made to the option with the lower expected value (**Daw et al., 2006**). Conversely, 'greedy' choices are made to the option with the highest value (**Findling and Wyart, 2021; Costa et al., 2015; Izquierdo et al., 2017**) and can be considered to exploit current knowledge of the values of the actions (**Sutton and Barto, 2018**). One possibility is that exploration represents a cognitive strategy for active information seeking, to reaffirm existing beliefs in the relative value of the choices (**Wilson et al., 2014; Izquierdo et al., 2017; Sutton and Barto, 2018**). Alternatively, it may represent random choice perturbations caused by limitations in the brain's implementation of the decision process leading to decision noise (**Findling et al., 2019**). Either way, the strategies adopted by human participants performing bandit tasks, including reversal-learning (**Costa et al., 2015; Costa et al., 2019; Costa et al., 2014**), rely on exploiting stability whilst at the same time adapting to volatility by exploring alternative courses of action.

Several observations support a role of the BG circuit in explore-exploit decisions. Computational circuit models of the BG predict that its principal function is action selection – the process where by one action (or decision) is selected over other competing alternatives (**Gurney et al., 2001; Bogacz and Gurney, 2007**). An extension of this action selection function is the proposal that its precision (i.e. how selective the circuit is) can be influenced by the overall excitability of the BG circuit (**Suryanarayana et al., 2019**). In the context of the explore-exploit dilemma, this means the BG may implement a neural 'decision filter', with the bandwidth of this filter being adjusted by neuromodulators, such as dopamine, which influence BG excitability (**Gilbertson and Steele, 2021**). The same models also predict that the excitability of the basal ganglia's principle output nucleus, the globus pallidus interna (GPI) should be a read-out of whether the BG are supporting an exploratory or exploitative decision strategy (**Humphries et al., 2012; Chakravarthy et al., 2010**). Consistent with this prediction, the firing rate of GPI neurons in non-human primates encodes the transition from exploratory to exploitative decision making (**Sheth et al., 2011**). Transient decreases in GPI firing observed during learning have been proposed as a potential source of decision variability and could therefore represent a neurophysiological correlate of exploratory (non-greedy) choices observed in bandit tasks.

At a cortical level, functional imaging (fMRI) supports a role for distinct pre-frontal cortical regions in reversal-learning including the orbito-frontal cortex (**Hampshire et al., 2012; Remijnse et al., 2005; Cools et al., 2002; Ghahremani et al., 2010**), ventromedial prefrontal (vmPFC), dorsolateral prefrontal (DLPFC) and anterior cingulate (ACC) cortices. The same pre-frontal regions are also implicated in explore-exploit choices using (so-called 'restless') bandit tasks with high between-trial outcome volatility and associated choice uncertainty (**Chakroun et al., 2020; Daw et al., 2006**). Analysis of brain activation patterns in these tasks support distinct cortical regions are involved in the different types of decisions. Consistent with their role within a valuation network, ventromedial prefrontal cortex (vmPFC) and orbitofrontal cortex (OFC) (**Chakroun et al., 2020; Daw et al., 2006; Findling et al., 2019; Badre et al., 2012; Tomov et al., 2020**) are activated during exploitative choices. In contrast, exploratory choices activate frontopolar cortex (FPC), dorsolateral prefrontal cortex (DLPFC), anterior insula and anterior cingulate cortex, commensurate with their proposed functions in the encoding of uncertainty, behavioural switching and cognitive control (**Rushworth and Behrens, 2008; Boorman et al., 2009; Bartolo and Averbeck, 2020; Hayden et al., 2011; Shenhav et al., 2016**). The striatum, as the BG input nucleus, receives afferents from these same pre-frontal cortical regions forming segregated cortico-striatal-thalamo-cortical circuits which respect the functional anatomy of distinct cortical

circuits for explore-exploit decision making, thus serving as a point of convergence (*Draganski et al., 2008*).

To date, no study has directly addressed the hypothesised role of the GPI in the human approaches to resolving explore-exploit decisions during reversal-learning. Deep Brain Stimulation (DBS) of the GPI is routinely performed as a treatment for isolated generalised dystonia and also in patients with focal and segmental dystonia if botulinum toxin treatment fails (*Volkman et al., 2014*). Dystonia is a movement disorder characterised by sustained or intermittent muscle contractions causing abnormal postures. In focal dystonia, it is restricted to one region such as the head or neck (*Bhatia et al., 2018*). Patients with dystonia have normal levels of striatal dopamine, intact reward prediction error (RPE) signalling (*Gilbertson et al., 2019*), and exhibit subtle cognitive biases (*Romano et al., 2014*). These include abnormalities of reinforcement learning (RL) such as increased risk taking and delayed flexibility to changes in reward contingency reversal (*Gilbertson et al., 2019; Arkadir et al., 2016*). With the caveat of these disease specific RL abnormalities in mind, dystonia patients with chronically implanted DBS-GPI electrodes represent a unique opportunity to test the influence of the GPI and, in turn, as its principle output nucleus, the basal ganglia's role in explore-exploit decision making in humans.

Our working hypothesis was that neuromodulation of the excitability of GPI neurons would modify explore-exploit decisions in a reversal-learning task. As DBS suppresses firing rates in human GPI neurons (*Cleary et al., 2013; LafreniereRoula et al., 2010*), we hypothesised that DBS-mediated reduced GPI excitability should drive greater proportion of exploratory choices consistent with previous findings by *Sheth et al., 2011*. To test this hypothesis, we examined choice behaviour in patients with chronically implanted DBS electrodes in both 'ON' and 'OFF' stimulation whilst they executed a two-armed bandit reversal-learning task. Given the more recognised role of the BG in habitual learning of stimulus-response associations (*Piron et al., 2016; Redgrave et al., 2010*), we purposely chose a low volatility probabilistic reversal-learning task to engage cortico-sub-cortical circuits in 'model-free' learning (*Daw et al., 2011*). Our choice of this task was also aimed at isolating the specific contribution of the GPI to random exploration. This form of exploration is more closely aligned with the sub-optimal, non-greedy choices identified during pauses in GPI firing in primates (*Sheth et al., 2011*). Exploration can also be 'directed' to the option which strategically aims to minimise uncertainty about a choice that a decision maker has least familiarity with (*Chakroun et al., 2020*). Accordingly, directed exploration relies upon brain regions that encode working memory (*Boorman et al., 2009*) and circuits including the external segment of the globus pallidus (GPe) which contribute to active information seeking to resolve uncertainty (*White et al., 2019*).

By fitting a model of the decision-making process (Reinforcement Learning Drift Diffusion Model: RLDDM), we aimed to test further mechanistic hypotheses regarding latent cognitive effects of DBS neuromodulation on GPI choice arbitration (*Pedersen et al., 2017*). The RLDDM incorporates both decision choices and decision time (DT) information, thus affording a richer interpretation of the mechanisms underlying choice compared to more traditional reinforcement learning models.

We further hypothesised that the influence of GPI-DBS neuromodulation in decision-making should be dependent on the connectivity of the stimulated site with cortical areas implicated in explore-exploit choice behaviour.

Results

Eighteen patients with isolated dystonia (for clinical details see *Supplementary file 1*) who had chronically implanted DBS electrodes (*Figure 1A*) targeted at the globus pallidus interna (GPI) performed a reversal-learning task (*Figure 1B&C*). The patients were instructed to try and win as many 'vouchers' as possible throughout the duration of the task by choosing one of two options presented to them on a computer screen. Patients were tested with their DBS device turned 'ON' or 'OFF' in a randomised order between the two task blocks. Each block consisted of three, 40-trial sessions where the probability associated with winning a 'voucher' switched from 80%:20% to 20%:80% midway through the second session (*Figure 1B*). Participants performed two versions of this task separated by an interval of 20 min after the initially randomised DBS condition (ON or OFF) had been inverted. The patient's performance was compared with a group of 18 age- and sex-matched healthy controls (HC) performing task block 1. The two blocks differed only by the fractals presented (*Figure 1B*). Participants were not informed about the contingency reversal in the task.

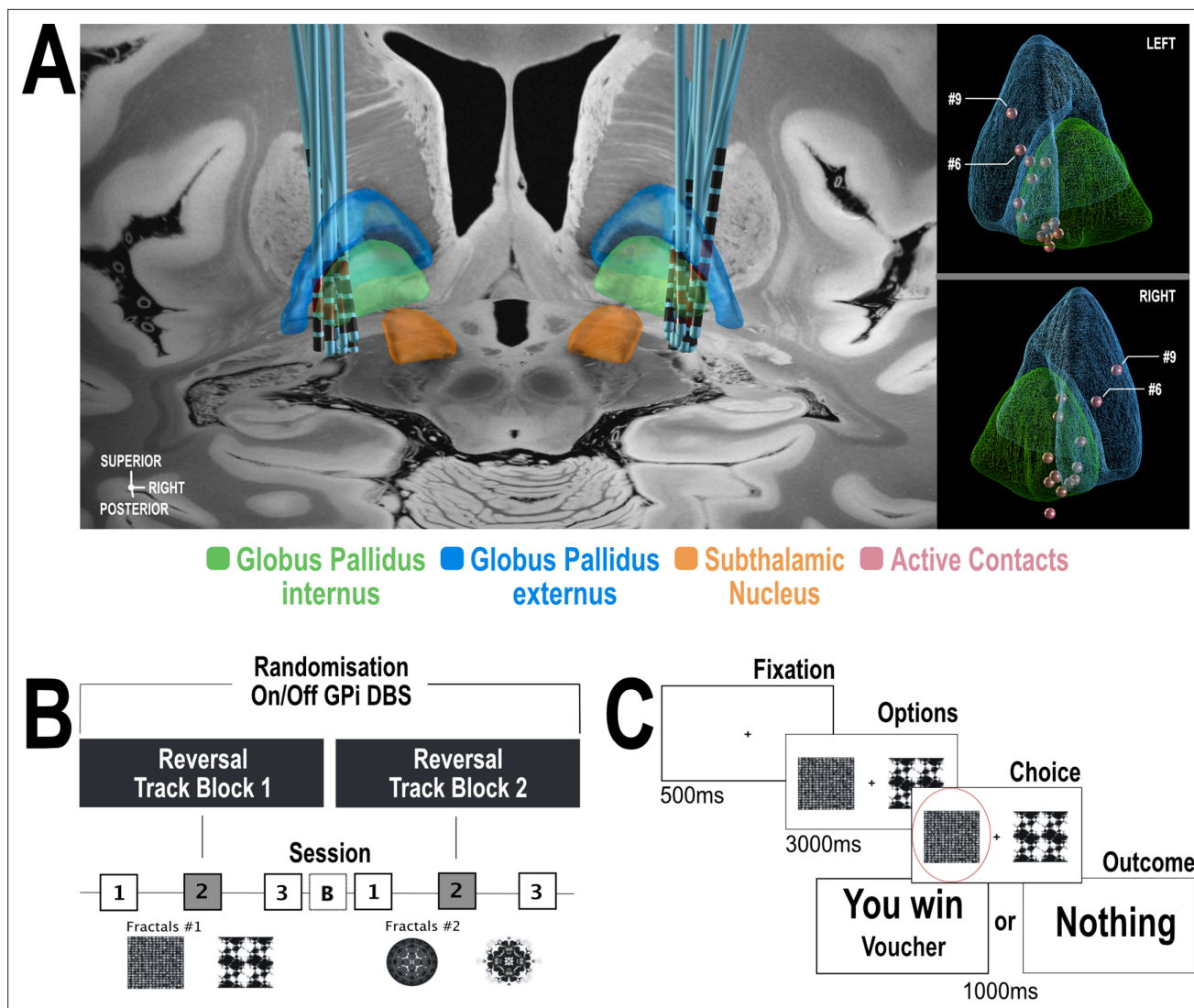


Figure 1. General overview of the study design and task. **(A)** Localisation of DBS electrodes was done with Lead-DBS (<https://www.lead-dbs.org/>) as previously described (*Horn et al., 2019*). The left panel depicts electrode localisation within the GPI projected to a 7T brain backdrop (*Edlow et al., 2019*). On the right panel, active contacts location is shown in red – only two patients (#6 and #9) show active contacts within the GPe (blue). **(B)** Each of the two reversal learning task blocks. The blocks differed only in the pair of fractals used as stimulus options. Each reversal learning task block consisted of three 40 trial sessions. Reversal of the probability of receiving a reward occurred half-way through session two. The task was performed once in the OFF-DBS state and once in the ON-DBS state in a counterbalanced manner with a 20 minute break ('B') in between blocks 1 and 2. During this, DBS stimulator was either switched ON or OFF. **(C)** Example of a single trial in the reversal learning task. On each trial, subjects chose either left or right fractal options, which were also counterbalanced, using their left or right hand to press the corresponding keyboard button. The selected cue was then shown surrounded by a red circle (in this example Task Block 1 the left-hand cue is chosen). Subjects were then presented with the outcome of their choice on the next screen, which could be either a reward ('You Win') or zero ('Nothing'). Outcome probabilities of receiving a reward on choosing either fractal were 80%:20%.

The online version of this article includes the following figure supplement(s) for figure 1:

Figure supplement 1. Spatial distribution of stimulation volume on a group level.

GPI DBS enhances exploratory choices without affecting task performance

A fixed effects analysis of performance (number of vouchers won) was conducted, using a two-way repeated measures ANOVA with stimulation state (ON-DBS versus OFF-DBS) and task session as fixed factors. Consistent with the detrimental effect of contingency reversal on performance in session 2, there was a main effect of session on performance ($F(2,34) = 6.26, p=0.002$, Average rewards:

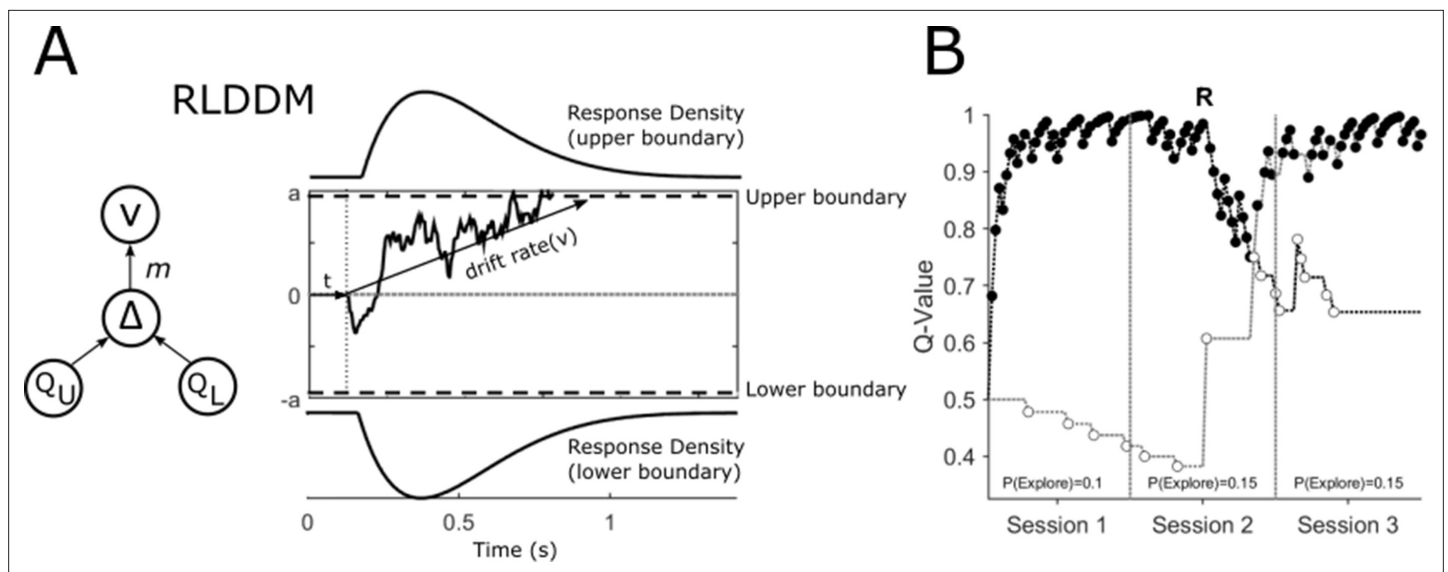


Figure 2. Reinforcement Learning Drift Diffusion Model (RLDDM). **(A)** The accumulation of evidence begins at a starting point at 0. The non-decision time is represented by t . Evidence accumulation is represented by a sample path with added Gaussian noise and is gathered until a decision boundary (a and $-a$) is reached and a response is initiated. The drift rate v determines the rate at which this evidence is accumulated. The extent to which the difference in the expected value of the options of the upper (Q_U) and lower response boundaries (Q_L) modifies the drift rate is determined by the drift rate scaling parameter, m . **(B)** An example of an individual patient's choices (OFF-DBS, Task Block 1) across the three sessions of the task with the expected value of choices represented by the upper Q_U and lower Q_L decision boundaries in black and grey respectively. Closed circles represent 'exploitative' choices where the choice with the highest expected value was chosen. Open circles are 'exploratory', representing the choice of the lower value of the two options. The change in value of the two choices halfway through session 2 of the task reflects the reversal, 'R', in outcome probabilities. The probability of exploring $P(\text{Explore})$ is the total number of choices made for the option with the lower expected value divided by the number of choices made in the session.

Sessions 1, 2 and 3; 27.5 ± 0.8 , 23.6 ± 0.6 , 24.8 ± 0.9) but no main effects of DBS state ($F(1,34) = 0.9$, $p=0.34$) or interaction between DBS state and performance in a session of the task ($F(2,34) = 1.03$, $p=0.30$) see Figure 3A&B. The same analysis applied to the decision time (defined as the interval between the choices being presented and keypress) confirmed a significant main effect of session ($F(2,34) = 5.21$, $p=0.005$, DTs: Sessions 1, 2 and 3; 1.04 ± 0.04 , 1.14 ± 0.05 , 1.18 ± 0.05 s) but no effect of DBS ($F(1,34) = 0.01$, $p=0.92$) or interaction between DBS and the decision times within a session of the task ($F(2,34) = 1.27$, $p=0.27$). This analysis confirmed that switching their DBS stimulator to the OFF state did not affect task performance by deteriorating the patient's symptom control. Applying the same analysis, with group (HC or patients) and task session as fixed factors, the number of rewards obtained in each session by the patients was no different compared with the HC group (HC average rewards: Sessions 1, 2 and 3; 24.7 ± 1.0 , 22.7 ± 0.8 , 28.1 ± 0.8 ; $F(1,70) = 0.01$, $p=0.91$) but the decision times in the HC group were on average faster compared to patients HC DTs: Sessions 1, 2 and 3 1.06 ± 0.02 , 1.03 ± 0.01 , 0.95 ± 0.01 s main effect of group ($F(1,70) = 5.36$, $p=0.02$).

Our hypothesis was that DBS induced modulation of GPI would enhance exploratory choices. We defined the probability of making an exploratory choice, 'P(Explore)', as the proportion of choices made in a session where the patient chose the low-value option as in an e-greedy or softmax choice rule (Sutton and Barto, 2018). This categorisation relies on estimates of the expected values of the available two options for low-value, exploratory choices to be defined. To identify these, we fitted the patients' choices and DT's using the RLDDM (Pedersen et al., 2017). This model represents a fusion of two traditionally separate models which combine the process of iterative updating the value of expectations (Q-value Figure 2A) using the delta learning rule (Rescorla, 1972) with a choice rule (substituted for the traditional softmax) based on the sequential sampling mechanism of the Drift Diffusion Model (Ratcliff, 1978).

The RLDDM model (Figure 2A) includes four main parameters of interest (free parameters). First, the drift rate, v , reflects the average speed with which the decision process approaches the response boundaries a (which take on positive and negative values to represent the two choices in the task).

Because the *drift rate* on any one trial is proportionate to the difference in the expected value of the two choices (**Figure 2A** $Q_U - Q_L$) the scaling parameter, m , determines the extent to which the diffusion process is weighed by the difference in values of the two options. It is analogous to, and closely related to the beta parameter (inverse temperature function) which governs the explore-exploit trade-off in a softmax choice rule (**Pedersen et al., 2017**). The expected values for the choices are in turn derived from the delta learning rule $\alpha \cdot (R - Q)$ with either a single learning rate, α , or a model variant with separate learning rates for positive, α_+ and negative α_- prediction errors. The nondecision time parameter, t , captures the time taken by stimulus encoding and motor processes. The model fit, as measured by the Deviance Information Criterion (DIC) (**Spiegelhalter et al., 2002**) indicated that the RLDDM with separate learning rates provided a better fit to the data compared to the model with a single learning rate (Patients: Dual learning rate DIC = 4344.98; single learning rate DIC = 4436.27; HC: Dual learning rate DIC = 3053.82; single learning rate DIC = 3056.00). Accordingly, all subsequent results reported are obtained using the dual learning rate model.

With the RLDDM model fitted to the choices separately for the ON and OFF-DBS states, we estimated the P(Explore) values for each session in the task (**Figure 2D**). A fixed effects analysis of P(Explore), was performed using a two-way repeated measures ANOVA with stimulation state (ON-DBS versus OFF-DBS) and task session as fixed factors. This demonstrated main effects of session ($F(2,34) = 3.14, p=0.04$) and DBS state ($F(1,34) = 4.64, p=0.03$) but no interaction between DBS state and the session of the task on the degree of exploration $F(2,34) = 0.77, p=0.46$. In the OFF-DBS state the average P(Explore) across the sessions was 0.13 ± 0.03 which increased to 0.2 ± 0.04 when the patients performed the task with the DBS stimulator switched ON (**Figure 3**). In **Figure 3**, the average P(Explore) values and reward performance for each session from 50 synthetic data sets is overlaid with the same values derived experimentally. These simulated choices were generated from RLDDM using the individual parameter estimates for each subject in the ON and OFF-DBS states from the tasks. This generated data was able to reproduce both the enhancing effect of DBS on P(Explore) and the preservation of reward learning performance in the ON and OFF-DBS states.

This analysis across the three experimental sessions in each task block may have averaged out more nuanced effects of DBS on exploration related to the reversal switch in contingencies, midway through session 2. To address this, we re-analysed the P(Explore) across twelve 10-trial bins (**Figure 3—figure supplement 1A**). A fixed effects analysis of P(Explore) was performed using a two-way repeated measures ANOVA with stimulation state (ON-DBS versus OFF-DBS) and task bin as fixed factors. This demonstrated main effect of bin ($F(11, 187) = 5.26, p < 0.001$), consistent with increased P(Explore) as an effect of contingency reversal and enhanced P(Explore) values in the ON-DBS state ($F(1,187) = 9.84, p = 0.001$). Despite this, there was no interaction between DBS state and a specific bin of the task on the degree of exploration ($F(11,187) = 0.51, p = 0.89$). In view of the marked increase in P(Explore) in the post-reversal bins 7–12 (**Figure 3—figure supplement 1A**), we also analysed the binned P(Explore) values with a fixed effect of reversal, allowing comparison of the P(Explore) values within pre-reversal bins (1–6) with those post-reversal (7–12). This demonstrated a main effect of reversal on the P(Explore) values ($F(1,187) = 27.4, p < 0.001$), a main effect of DBS state ($F(1,187) = 9.35, p = 0.002$), but no interaction between DBS state and reversal ($F(1,187) = 2.54, p = 0.11$). The mean P(Explore) value in each 10 trial bin in the ON-DBS state was 0.2 ± 0.04 and 0.14 ± 0.05 in the OFF-DBS condition.

We defined the probability of a reward, P(Reward), as the number of rewards obtained in a 10-trial bin divided by the number of responses, to assess the performance in the task on a finer time scale. This was used rather than the absolute number of rewards (as used in the analysis of session performance above) as it was a more accurate measure of performance in a small trial bin as missed trials were accounted for. Using P(Reward) for each subject at each bin, we found no main effect of DBS ($F(1,187) = 1.88, p = 0.17$), or interaction between the DBS and the task bin ($F(11,187) = 0.74, p = 0.7$). Again, this analysis confirmed that enhanced exploratory choices in the ON-DBS state were not associated with degradation in the ability of the patients to acquire rewards during the task (mean P(Reward) values ON-DBS were 0.63 ± 0.04 and 0.65 ± 0.04 OFF-DBS; **Figure 3—figure supplement 1B**).

To place the effects of DBS in our patients into the context of known abnormal reinforcement learning behaviour in isolated dystonia (**Gilbertson et al., 2019; Gilbertson et al., 2020**), we compared the binned P(Explore) values of the HC group, to those of the patients in the OFF-DBS condition (**Figure 3—figure supplement 2**). The probability of the HC group making an exploratory choice was greater than the patients in the OFF-DBS condition (mean HC P(Explore) = 0.19 ± 0.05),

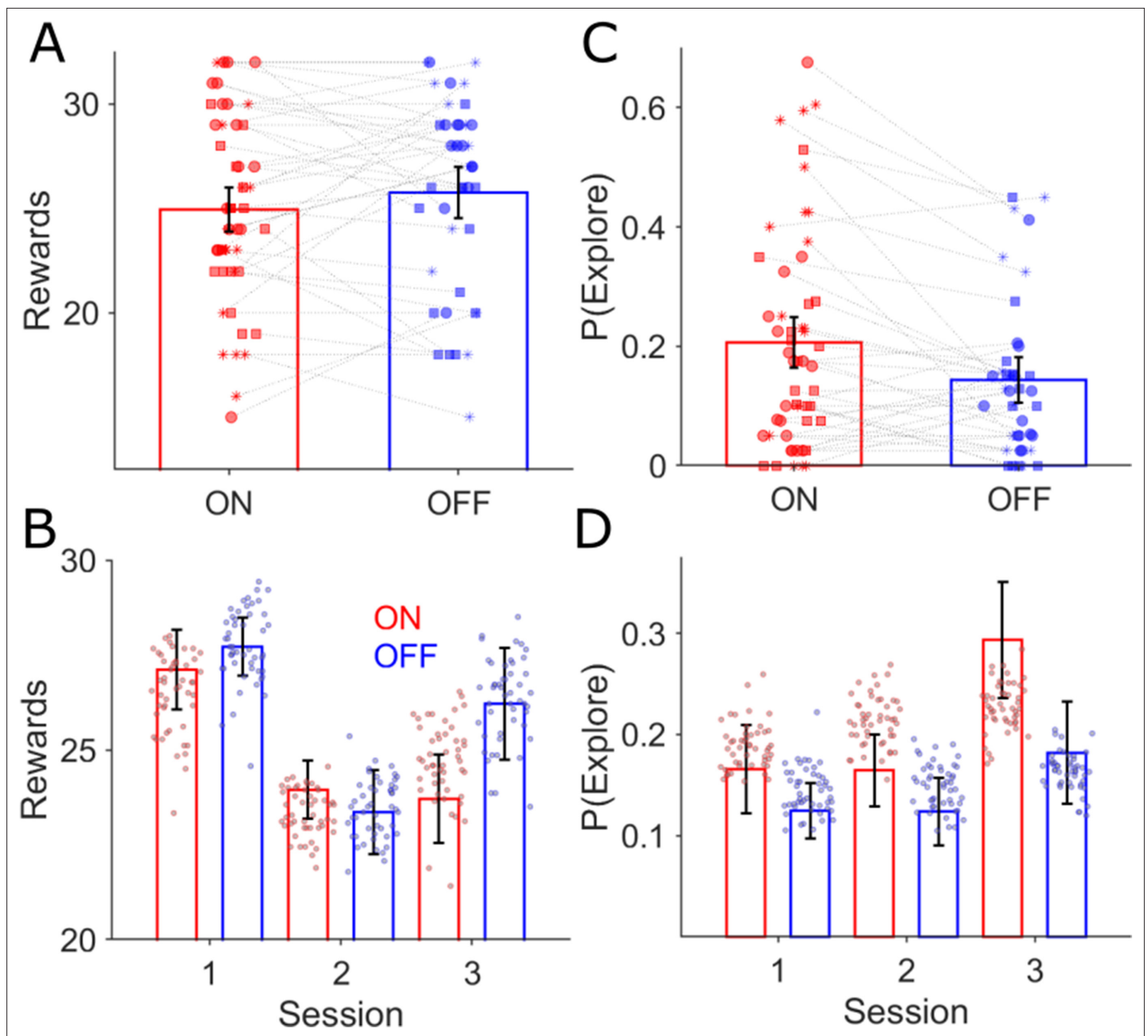


Figure 3. Behavioural effects of GPI DBS on task performance. There was no difference in the mean rewards (ANOVA, $p > 0.05$) won across the three task sessions for the ON (red bars, $n = 18$) and OFF-DBS (blue, $n = 14$) conditions (**A**) and (**B**) mean \pm s.e.m. Each symbol represents a subject's number of rewards won for that session, with sessions one, two and three represented by the circle, square and asterisk symbols. During ON-DBS testing the probability of exploring the lower value choice was significantly greater (ANOVA, $p < 0.05$) (**C**). Behavioural performance is plotted in (**B**) (number of rewards) and (**D**) (probability of exploring the lower value choice) for both DBS conditions with the superimposed scatter plots of the 50 simulated experiments generated using the RLDDM fitted to the ON and OFF experimental choices.

The online version of this article includes the following figure supplement(s) for figure 3:

Figure supplement 1. Behavioural effects of GPI DBS on task performance peri- contingency reversal.

Figure supplement 2. Behavioural performance in Healthy Controls (HC) compared to patients in the DBS-OFF condition.

confirmed by a main effect of group ($F(1,341) = 7.34$, $p = 0.005$). By contrast, there was no difference in the HC $P(\text{Explore})$ values when these were compared to the ON-DBS ($F(1,385) = 1.1$, $p = 0.3$) and no effect of dystonia on the overall performance of the task, as $P(\text{reward})$ values in the HC group were the similar (HC; mean $P(\text{Reward}) = 0.63 \pm 0.03$) to those of the patients in the OFF-DBS condition ($F(1,341)$

= 2.48, $\eta^2 = 0.2$). Overall, this analysis supported that GPI-DBS enhances exploratory choices in dystonia patients from a disease related baseline level of exploitation that is higher than age matched controls.

Finally, to ensure that the order of the two reversal tasks performed OFF and ON-DBS did not contribute in any way to this result, we ran a fixed effects analysis of P(Explore) across the three sessions with stimulation state (ON-DBS versus OFF-DBS) and task order as fixed factors. No main effect of task order on P(Explore) was evident from this analysis $F(1,34) = 0.02$, $p = 0.87$, therefore, the enhanced exploration ON-DBS was unrelated to whether the first of the two reversal tasks were performed with DBS stimulator switched ON or OFF.

Cortical functional connectivity correlates of the effect of DBS on exploration

When the DBS enhancing effects on exploratory choices on an individual subject level was investigated, this varied from a maximum within session increase of P(Explore) of 0.34, equating to 14 additional exploratory choices out of 40 in one subject, to no influence of DBS on exploration in any session in another. We leveraged the known influence of DBS on brain networks (*Horn et al., 2017*) to test the hypothesis that this variance could be explained by individual differences in connectivity of the DBS electrode with cortical regions which correlated strongly with exploratory decisions.

The connectivity analysis was performed based on established methodology previously introduced to assess brain networks affected by distributed lesions or stimulation effects in neurological and psychiatric conditions (*Boes et al., 2015; Corp et al., 2019; de Almeida Marcelino et al., 2019*).

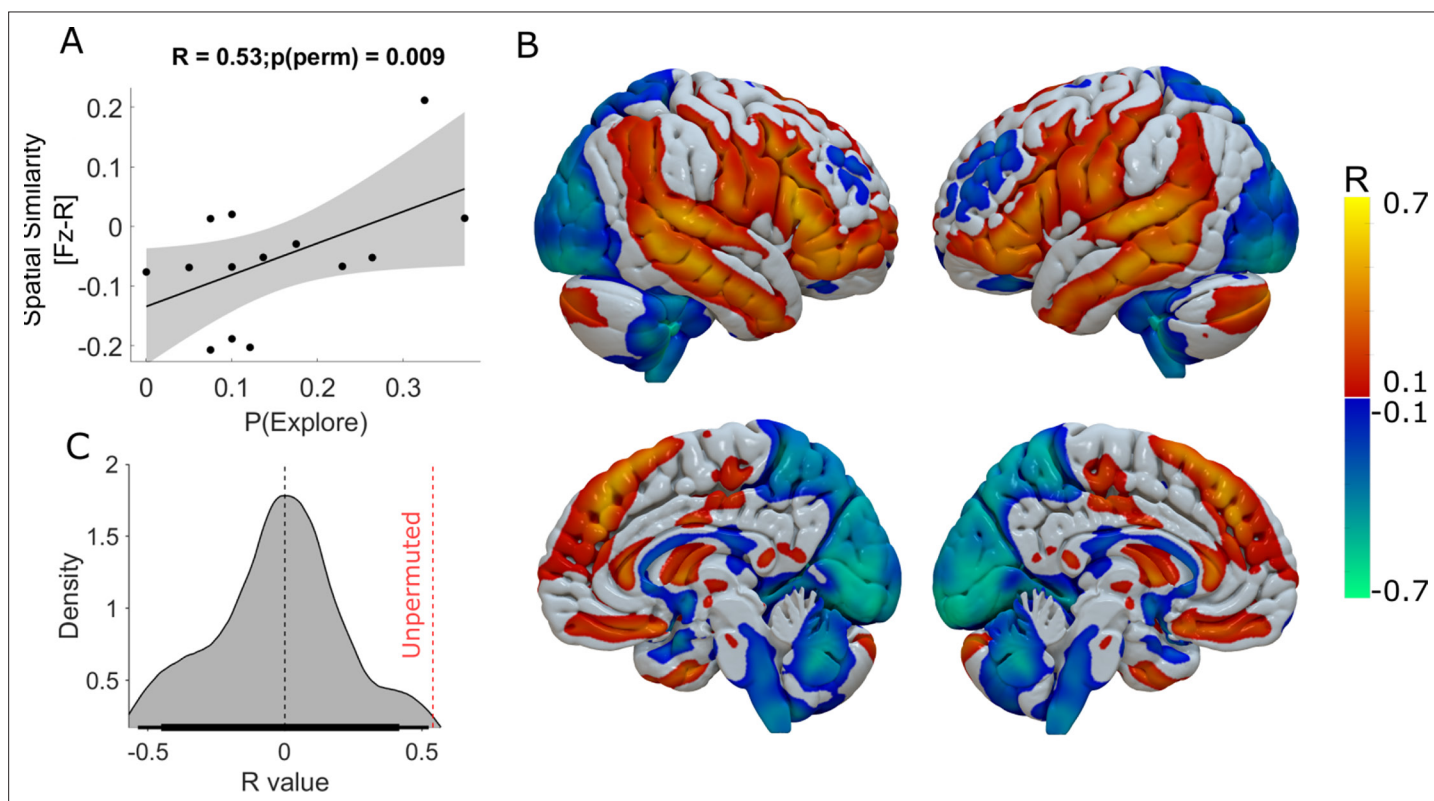


Figure 4. Functional connectivity of DBS-Induced exploration. The whole-brain voxel-wise R-Map demonstrates the optimal functional connectivity profile for DBS-induced enhancement of exploration (B). In this analysis, the maximal increase in DBS induced exploration (ON-OFF DBS) in one of the three experimental sessions was used as the regressor. Warm colours show voxels where functional connectivity to the DBS stimulation volumes was associated with greater exploration. Cool colours indicate voxels where functional connectivity to the DBS stimulation volumes was associated with lesser exploration. The more the individual functional connectivity profile matched the 'optimal' R-Map, the greater was the DBS-induced exploration (A) ($R^2 = 0.28$, $p = 0.04$). In (C) we plot the R value distribution derived from 1000 repermuted correlations between the enhancing effect DBS and the R-map. The probability of seeing the same correlation by chance was $p < 0.01$ (represented by the red dashed vertical line).

The online version of this article includes the following figure supplement(s) for figure 4:

Figure supplement 1. Functional connectivity of DBS-induced exploration.

Using the stimulation volumes of the DBS electrodes as seed regions to a group connectome derived from resting state functional connectivity of healthy volunteers (Yeo *et al.*, 2011), we derived whole-brain functional connectivity maps for each patient. Then, the correlation between functional connectivity and DBS-induced change in exploratory behaviour was calculated voxel-wise. The R-Map in **Figure 4B** depicts thresholded correlations between functional connectivity from DBS stimulation volumes to all brain voxels and DBS-induced exploration.

We used the difference in P(Explore) values ON minus OFF-DBS as the behavioural parameter for this analysis to normalise for between subject variation in the baseline level of exploration and to identify the super-added effect of DBS on this behaviour. For each subject, we chose the experimental session with the largest increase in DBS induced exploration. This aimed to minimise the effect of averaging this behaviour across the task as a whole, whilst allowing for individual differences in the effect of DBS across different sessions of the task not evident at a group level from the analysis of P(Explore). Spatial correlation of individual patients' DBS connectivity profiles with this R-Map significantly explained variance in DBS-induced exploration (**Figure 4A**) $n=14$ $R^2=0.28$ $p=0.04$, Permutation test: $R=0.53$, $p=0.009$. **Supplementary file 2** includes the details of the peak R-map values and voxel co-ordinates of the cortical regions whose connectivity with the DBS electrode predicted increasing exploration. To ensure that choosing the maximum within session difference in P(Explore) for this R-map did not arbitrarily identify network effects of pallidal neuromodulation, we performed an additional R-map analysis, averaging across the whole task and subtracting the difference in P(Explore) values in the ON minus OFF-DBS conditions (**Figure 4—figure supplement 1**). This revealed cortical connectivity patterns with similar topography. This R-map had a poorer linear predictive value of P(Explore), $R^2=0.17$, $P=0.12$ but met the permutation threshold (Permutation test: $R=0.41$, $p=0.004$) suggesting this was unlikely to occur by chance.

Effects of DBS on decision choices revealed by influence on RLDDM parameters: model validation and parameter recovery

We performed posterior predictive checks to test whether the dual learning rate RLDDM was a good model of the experimental data. First, we compared the observed data, in our case the performance in the reversal task and the decision time distributions, to simulated data generated by the model. The learning curves (which illustrate the probability of choosing the highest value choice) are illustrated in **Figure 5A**. These represent the model's generative choices for fifty simulated 'experiments' with the RLDDM parameters estimated for each individual subject in both the ON and OFF-DBS conditions. The good model fit is indicated by the overlap between the generated data performance (shaded) with the observed choices (solid line) for both conditions. Furthermore, because the RLDDM includes a prediction of the decision time (DT) at which the choice was made, we also overlaid the observed DT histograms with the corresponding density for 50 generated simulations (**Figure 5C**). The close overlap between the observed and generated data supports that the RLDDM captured the salient features of the patients' choices and decision times whilst observing the heightened exploratory choice tendency in the ON-DBS state, as the same generative data was used for the overlaid P(Explore) simulated data in **Figure 3D**. As an additional check, we performed parameter recovery by re-fitting the first five generated data sets, plotting the model parameter estimates for the observed experimental data against the estimates derived from generated choices (**Figure 6D**). For all parameters, we found significant correlations between the parameters used to generate the behaviour and the recovered parameters obtained from re-fitting the model to the generated behaviour.

We estimated group and individual parameter values for the two learning rates (α_+ , α_-), the boundary separation parameter (a), and the drift rate scaling parameter, (m) as dependent variables of DBS state (ON or OFF). The non-decision time, t_0 (seconds), was estimated for each subject but was assumed not to be influenced by DBS state (as DBS did not appear to affect the absolute DT in the behavioural analysis above). Consistent with the finding of slower DTs in the patients compared to HC, the estimate for t_0 in the patients was correspondingly longer (Patients = 0.57) [95% HDI 0.46, 0.69], HC = 0.42 [0.36 0.50] a difference in the posterior difference of the means which was statistically significant $M_{diff} = 0.13$ [0.001, 0.27].

Group level parameters of the within-subject effects of DBS were used to assess how pallidal neuromodulation influenced exploratory choices. We found no statistically significant differences between the two DBS states when the dual learning rate model was fitted to all trials across the

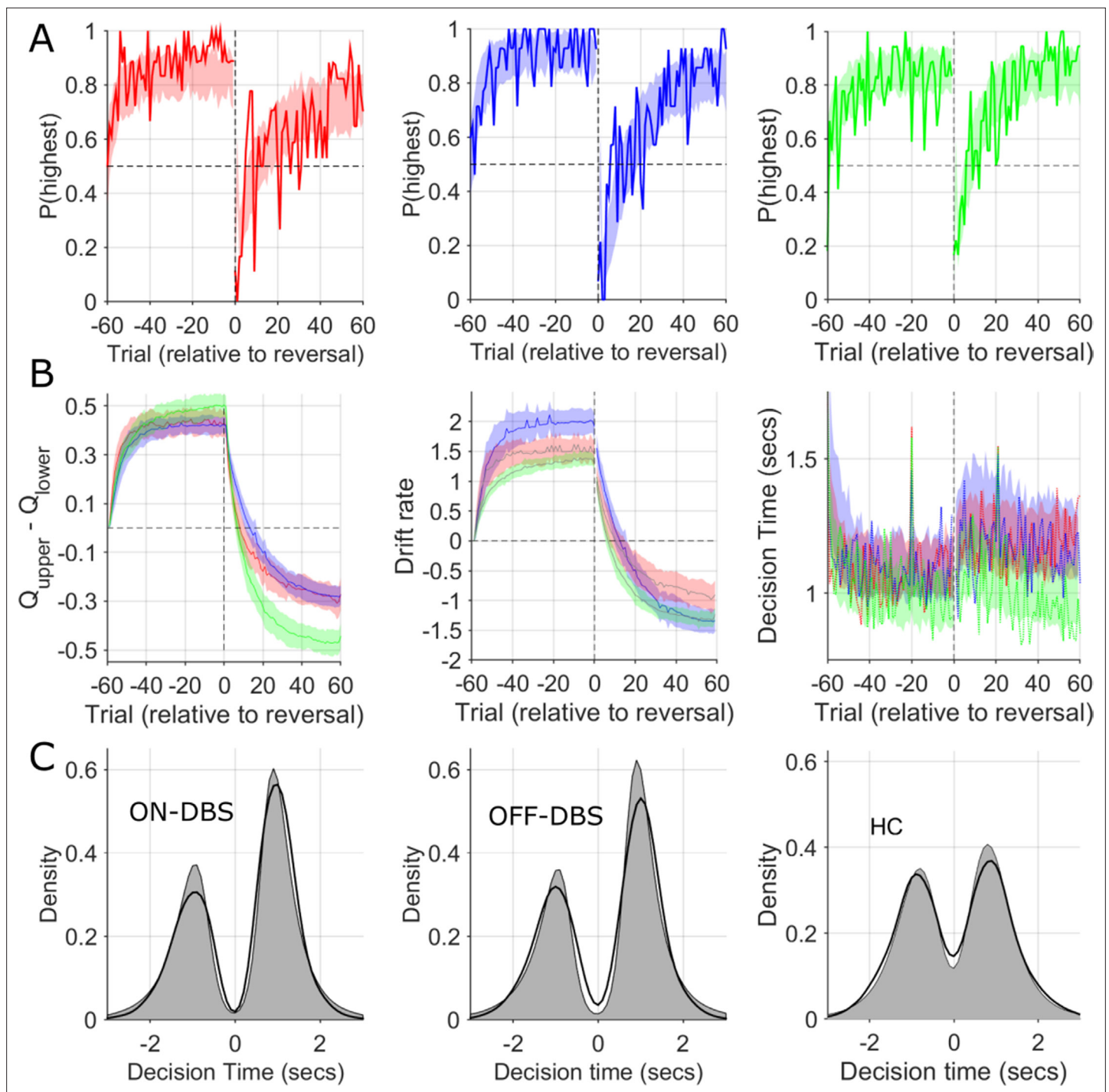


Figure 5. Observed and simulated RLDDM data. **(A)** The probability of choosing the highest value stimulus in the ON-DBS state is represented by the solid red line. Overlaid shaded area represents 95% confidence limits of simulated choices of the RLDDM (averaged across 50 simulated experiments). The same analysis for the OFF-DBS state is plotted in blue in the middle panel and for HC in green. This supported the interpretation of the model being a good fit to the experimental data due to the strong overlap between the synthetic data choices and those observed for both DBS conditions and in both HC and patient groups. **(B)** The difference in the expected value (Q) of the two choices in the RLDDM (mean \pm 95% confidence limits) are represented by the solid blue (OFF-DBS), red (ON-DBS) and green (HC) lines with the shaded overlay showing the confidence limits estimated from the synthetic (simulated) data. The middle panel in **(B)** illustrates the trial-to-trial variation in the mean drift rate, v , across the simulations in all conditions and groups. This demonstrates a reduction in the ON-DBS drift rate consistent with the group level effect of DBS on reducing the drift rate scaling parameter, m , (see **Figure 6—figure supplement 1** and **Supplementary file 2**) and the lower drift rate in the HC group. The mean experimental trial-to-trial variation in DT in each DBS condition and the HC group are plotted as a solid line with the simulated model DT's overlaid. The model captures

Figure 5 continued on next page

Figure 5 continued

both the within task cost of the contingency reversal on DT and the faster DT in the HC group through the task. In (C), thicker black lines display the observed decision time distribution across all patients and sessions in the ON-DBS (red) and OFF-DBS (blue) conditions and HC groups, with the simulated DT from the RLDDM represented by the grey shaded regions. The decision time (DT) for choices which have an initially lower value are shown as negative. The reliability of the RLDDM in capturing the decision mechanisms in the task are supported by the overlap in observed and simulated DT distributions.

reversal learning task (Figure 6 & Table 1). This model was therefore able to adequately capture the patient's exploratory decisions in the task and the effect of DBS (Figure 3B&D, Figure 5, Figure 6D) but provided no mechanistic insight into how pallidal neuromodulation drove an increase in these choices. Given the strong influence of the contingency reversal mid-way through session 2 of the task on exploration, (Figure 3—figure supplement 1), we refitted the model separately to the pre-(1-60) and post-reversal trials (61-120). The posterior distributions of the four parameters and their posterior differences (ON minus OFF-DBS) are illustrated in Table 1—source data 1 and Figure 6—figure supplement 1. We found a similar pattern of posterior directional effects of DBS when fitting to the whole task or separately to pre- and post-reversal. DBS led to increases in both learning rates, decreases in the boundary separation parameter (a) and drift rate scaling parameter (m). However, analysing the pre- and post-reversal trials independently, we detected the only statistically significant group level difference in the drift rate scaling parameter, m , in the post-reversal phase of the task; $M_{\text{diff}} = -1.71$ (95% HDI [-3.5, -0.18]). This reduction in the drift rate scaling parameter m ON-DBS is consistent with enhanced exploration, as this parameter amplifies the influence of the difference in the expected value of the two choices ($Q_u - Q_l$), by driving the diffusion process to the decision boundary of the higher value option (Figure 2A).

Comparing the RLDDM parameter estimates between the HC and patients in the OFF-DBS condition, the estimated value of, m , was also reduced in the HC group relative to the patients in the DBS-OFF state Figure 6—figure supplement 2, Table 1—source data 2; $M_{\text{diff}} = -1.71$ (95% HDI [-3.5, -0.18]). This explains the higher proportion of exploratory choices by HC's in the task for the same reason that in the ON-DBS state exploration is heightened – this parameter proportionately scales the influence of the difference in values between the two options on the eventual choice of decision. Patients with dystonia have previously been shown to have impaired learning from negative feedback (; Gilbertson et al., 2020). Consistent with these previous studies we also found statistically significant difference in the negative learning parameter, α_- , with a higher value HC group compared to the patients in the OFF-DBS condition $M_{\text{diff}} = 1.32$ (95% HDI [0.26, 2.5]).

This comparison of RLDDM parameter estimates between the HC and OFF-DBS groups confirmed that greater exploratory choice tendency in HC's could be explained by a lower influence of the difference in expected value of two choices. In both analyses, lower levels of exploration in HC's compared to patients with dystonia and enhanced exploration by GPI neuromodulation in the patients ON-DBS, were most likely explained in the RLDDM by relative reductions in the drift rate scaling parameter (m). In turn, this meant that in both the HC group and in the patients in the ON-DBS condition, the decision to choose one of the two options in each trial was much less influenced by encoding of their value. Accordingly, choices were influenced proportionately more by noise intrinsic to the decision process, resulting in greater proportion of random, exploratory, (and accordingly less greedy, exploitative) choices. This was in turn reflected by a higher level of P(Explore) values both ON-DBS and in HC's relative to the OFF-DBS condition in our behavioural analysis above.

Discussion

In this study, we examined the role of the human pallidum (Globus Pallidus Interna) in mediating the trade-off in explore-exploit decision-making. Our two predictions were confirmed: (1) Neuromodulation by deep brain stimulation of the GPI increased the likelihood of patients exploring the lower value choice in a two-armed probabilistic reversal learning task; (2) DBS-induced enhanced exploration correlated with the functional connectivity of the stimulation volume in the GPI to a distributed brain network including frontal cortical regions identified previously in functional imaging studies of explore-exploit decision-making. Furthermore, a recently proposed reinforcement learning model successfully predicted the behaviour, enabling a more mechanistic interpretation of experimental results.

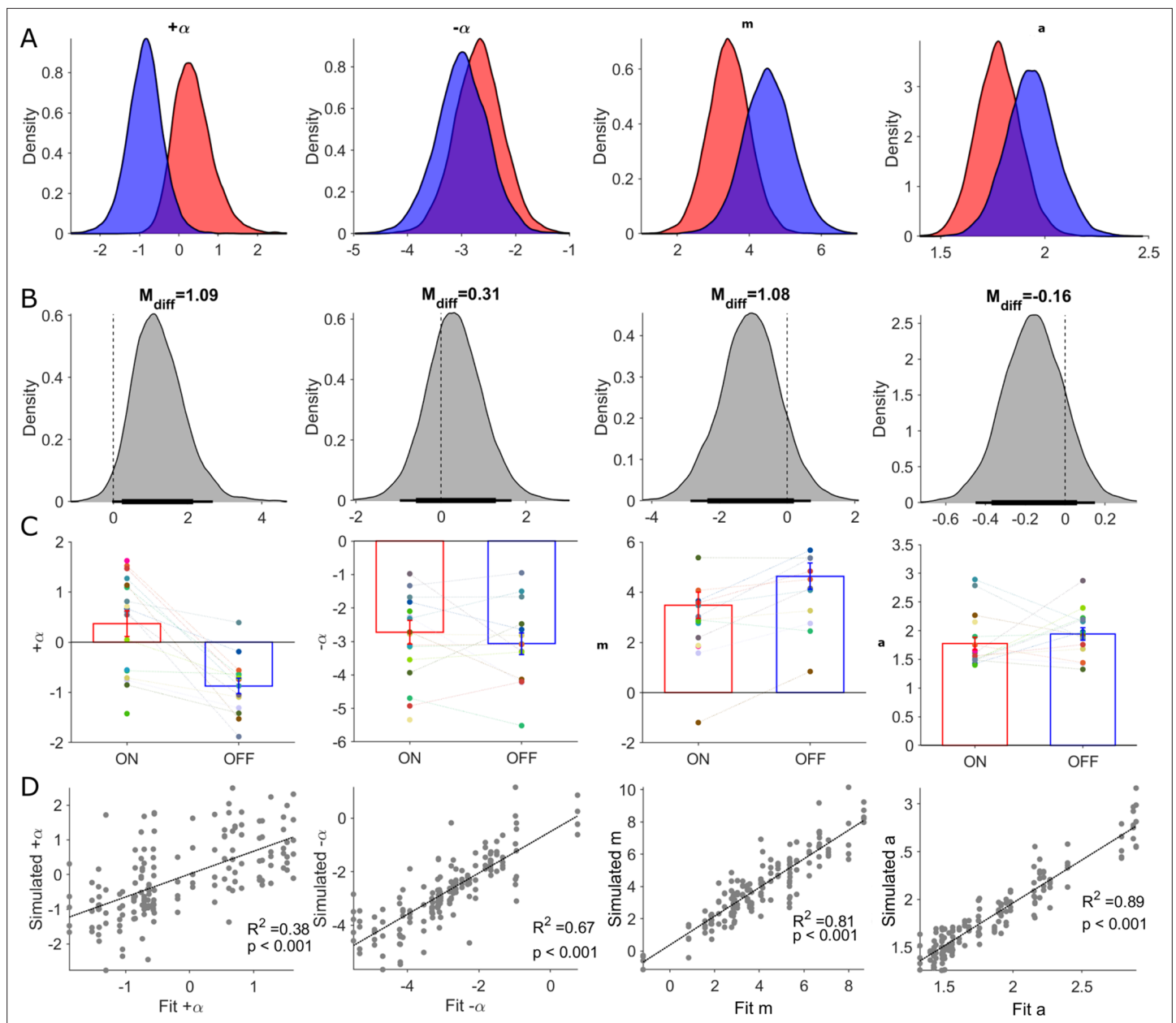


Figure 6. Posterior distributions of RLDDM parameters and parameter recovery. Columns in this figure display results for each of the estimated parameters in the RLDDM from left to right, including positive learning rate (α_+), negative learning rate (α_-), drift rate scaling parameter (m), boundary separation parameter (a). Rows, from top to bottom, correspond to **(A)** posterior distributions ON-DBS (red), OFF-DBS (blue). **(B)** Posterior distributions of differences ON-DBS versus OFF-DBS. Thick and thin horizontal bars below the distributions represent the 85% and 95% highest density intervals, respectively. **(C)** Mean \pm S.E.M. parameter estimate ON-DBS (red) and OFF-DBS (blue). Each individual's parameter estimate represented by a different colour in the scatter plot for each parameter. **(D)** Results of parameter recovery analysis, plotting the estimated parameter values for the observed data against the parameter values re-estimated from simulated data. Significant ($p < 0.001$) linear correlations were observed for all four parameters supporting successful recovery and validation of applying the RLDDM model of decision-making.

The online version of this article includes the following figure supplement(s) for figure 6:

Figure supplement 1. Posterior distributions of RLDDM parameters and parameter recovery – separately fitted to pre- and post-reversal trials.

Figure supplement 2. Posterior distributions of RLDDM parameters from Healthy Controls (HC) and parameter recovery.

Table 1. Summary of posterior distributions (RLDDM model fit to whole task).

Parameter	ON-DBS		OFF-DBS			Contrast			
	mean	HDI	mean	HDI	mean	HDI	mean	HDI	
Boundary separation	1.78	1.56	1.98	1.96	1.7	2.18	-0.16	-0.45	0.15
Drift rate scaling	3.4	2.26	4.54	4.51	3.20	5.88	-1.08	-2.86	0.71
Learning rate +	0.22	-0.55	1.37	-0.85	-1.79	0.01	1.09	-0.03	2.68
Learning rate -	-2.65	-3.58	-1.8	-2.97	-3.97	-2.05	0.31	-0.97	1.66

Estimated means (m) for ON and OFF-DBS conditions as well as contrast for ON-OFF. HDI values are for 95% highest density Interval.

The online version of this article includes the following source data for table 1:

Source data 1. Summary of Posterior distributions (RLDDM model fitted separately to pre- and post-reversal trials).

Source data 2. Summary of Posterior distributions (RLDDM model fitted to whole task).

DBS effects on GPI and their contribution to explore-exploit choice behaviour

The exploration – exploitation ‘dilemma’ is considered one of the fundamental challenges of adaptive control and behaviour (Mehlhorn *et al.*, 2015; Cohen *et al.*, 2007). We took advantage of the opportunity leveraged by DBS to address the paucity of experimental evidence (Sheth *et al.*, 2011) relative to theoretical support for the BG role in explore-exploit choice arbitration (Gilbertson and Steele, 2021; Humphries *et al.*, 2012; Chakravarthy *et al.*, 2010).

Patients and non-human primates with lesions of the GPI exhibit an inability to learn novel stimulus response contingencies and random choice behaviour which prevents the acquisition of choice preference or implicit learning of sequences (Piron *et al.*, 2016; Obeso *et al.*, 2009). We did not observe levels of performance degradation in our patients that would be comparable to complete lesioning of pallidal output.

Neurophysiological data support a net inhibitory effect of DBS on GPI neuron firing rates without complete firing suppression at stimulation intensities comparable to those used in our patients (Bar-Gad *et al.*, 2004). We suggest that the most likely neuromodulatory mechanism of DBS on enhanced exploration is a partial inhibition of GPI firing and accordingly electrically lesioning the basal ganglia’s output (Wu *et al.*, 2001; Dostrovsky *et al.*, 2000; McCairn and Turner, 2009; Boraud *et al.*, 1996).

Previous studies have confirmed that during acute changes in GPI-DBS (comparable to the ON-OFF conditions tested in our patients) a reduced inhibitory output induced by DBS in GPI leads to increased excitability of the motor cortex (Tisch *et al.*, 2007; Kühn *et al.*, 2003). Motor cortex activity correlates negatively with decision value during explore-exploit decisions (Tomov *et al.*, 2020). This means that within a sequential sampling decision mechanism, increased motor cortex activity correlates with small differences in the value of the two choices and corresponding higher decision uncertainty and choice exploration. This was indeed the effect predicted by the RLDDM in the ON-DBS state. A reduction in the drift-rate scaling parameter (m) has a corresponding effect on decision value as this parameter multiplies the difference in expected value in the two choices (Pedersen *et al.*, 2017). The exploratory influence of pallidal DBS could therefore be explained by thalamocortical disinhibition, increased motor cortex excitability and its activity inversely encoding decision value.

Alternatively, rather than an influence on a downstream cortical decision mechanism, GPI-DBS may have influenced implementation of this process *within* the BG circuit. The simulations of Dunovan *et al.*, 2019 propose that activity in the direct (striato-nigral) and indirect (striato-pallidal) BG pathways differentially encode the rate of evidence accumulation and the amount of information needed to reach decision threshold. We found that GPI-DBS reduced the drift rate scaling parameter during the post-reversal phase of the task. As this parameter determines the rate of evidence accumulation, our results would be consistent with this predicted function of the direct pathway being neuromodulated by pallidal DBS. However, because the GPI sits anatomically at the point of convergence of both direct and indirect pathways, we cannot distinguish this mechanism from a downstream influence on thalamo-cortical excitability. Future studies examining neuromodulation of a basal ganglia nucleus

that lies within the indirect pathway, such as the subthalamic nucleus (STN), could be used to delineate this mechanism further, as this might be predicted to selectively influence the RLDDM boundary separation parameter which determines the decision threshold.

The basal ganglia's role in reinforcement learning (information seeking vs. uncertainty)

During the early stages of learning, firing rates of GPI neurons show transient inhibitory pauses in firing rate which encode exploratory choices. These are replaced with higher firing rates as choices become habitual and stimulus response mapping is obtained (Sheth et al., 2011). These findings are consistent with physiological changes in the direct and indirect pathways during different stages of learning. (Yin et al., 2009) identified increased excitability in the striato-nigral (D1R-expressing) direct pathway early in learning with corresponding inhibition of the GPI. With subsequent consolidation of learning and habituation of the stimulus response relationship, D2R dependent increases in striato-pallidal excitability would dis-inhibit the GPI (Kravitz et al., 2010), increasing its excitability and promoting exploitation. Recent data from Lee and Sabatini, 2021 would however contradict this interpretation as they found exploration increased in optogenetically stimulated indirect pathway striato-pallidal neurons. This discrepancy might be explained by careful distinction between basal ganglia contributions to actively promoting exploratory choices during learning, versus active inhibition of choices from negative outcomes, that may be more relevant to choice extinction (Tecuapetla et al., 2016) and shifting (Kravitz et al., 2012).

An alternative explanation (to the idea of active exploratory information seeking) is that enhanced non-greedy choices ON-DBS represent increased random variability in the update of action values due to computational noise inherent to the learning process (Findling and Wyart, 2021). Examining the effects of GPI-DBS in a non-stationary 'restless' bandit task could delineate the contribution of decision noise as these models that take noise into account of the decision process have been successfully fitted to such tasks (Findling et al., 2019).

The basal ganglia's role in reinforcement learning (learning from previous outcomes)

We found no effect of DBS on the learning rate to either positive or negative outcomes in the task. This is surprising given evidence from Local Field Potential (LFP) recordings from GPI-DBS electrodes which demonstrate that these encode choice outcome (Eisinger et al., 2022). Reward Prediction Error (RPE) encoding in LFP's from GPI electrodes have also been identified and are most marked during exploratory decisions and absent when choices are exploitative (Schroll et al., 2015). Pallidal error-related activity has been identified from LFP recordings. This precedes the cortical error-related negativity (ERN), which is, in turn, closely related to the negative prediction error signal (Herrojo Ruiz et al., 2014; Holroyd et al., 2003). A plausible a priori mechanism for increased exploration would be DBS induced blunting of these learning signals. The absence of any effect of DBS on the RLDDM learning rate parameters results would suggest that outcome signalling including RPE-like and ERN signals in the GPI are relatively unaffected by DBS. Therefore, rather than modify intrinsic encoding of RPE-like signals or through distal connections with other brainstem structures, (Hong and Hikosaka, 2008) DBS enhanced exploration appears to be mediated by a neuromodulation of the decision process, consistent with theoretical accounts proposing the basal ganglia's excitability as an arbitrator of action selection (Gilbertson and Steele, 2021; Humphries et al., 2012; Dunovan et al., 2019). An equally parsimonious explanation for the absence of any impact of DBS on the learning rate, is that by targeting of the (dorsal) motor segment of GPI, there was no effect on the excitability of the ventrally located GPI – lateral habenula pathway (Hong and Hikosaka, 2008) which facilitates negative prediction error signalling in the Ventral Tegmental Area (VTA).

Integrating the results in a network perspective

The postero-lateral portion of the GPI targeted in the group of patients studied here is chosen in routine clinical practice because of its position within the sensorimotor loop of the basal ganglia (Alexander et al., 1986) to relieve motor symptoms of dystonia. It is noteworthy that despite evidence for anatomically segmented regions of the GPI (based on white matter connectivity Draganski et al., 2008) stimulating the sensorimotor pallidum in this study could alter decision-making and

therefore modify 'cognitive' function. At first glance, this result might appear incongruous with its established role in motor control. However, our functional connectivity analysis supports the conclusion that sensorimotor GPI function extends beyond that traditionally attributed to it on the basis of anatomical connections. The finding of a neuromodulatory influence of sensorimotor pallidum on decision-making is consistent with fMRI activation of pre-and post-central gyri, during exploratory choices (*Chakroun et al., 2020*) and the encoding of decision value by motor cortex (*Hare et al., 2011*).

Our connectivity analysis is noteworthy for the extent of the connectivity pattern across much of the brain. The largest within session effect of DBS was seen in patients where the electrode stimulation volume shared connectivity with brain regions which would be considered functionally relevant to our task (e.g. pre-frontal regions). Equally, this connectivity profile included regions with little to no recognised functional role in reward-reversal learning. Therefore, caution is required in interpreting the anatomical detail of this network given the limitations of using imaging data with no individualised anatomical specificity. On a more general level, this analysis provides support to the idea that discrete functions in local brain networks (in our case the GPI), are integrated into a broader scale whole-brain network. This leads to co-variation in neural activity between regions which are not conventionally attributed a functional role in a specific motor or sensory act (*Kauvar et al., 2020*) and might explain how DBS exerts influence on remote areas beyond the stimulation site (*Horn et al., 2017*). For instance, such as modulation of distributed networks through subthalamic DBS affects motor learning and risk-taking behaviour in Parkinson's disease patients (*de Almeida Marcelino et al., 2019; Irmen et al., 2019*).

Clinical relevance - apathy

From a clinical perspective, these findings may be relevant to understanding the mechanisms which worsen apathy in Parkinson's disease (PD) following DBS targeting the basal ganglia (*Zoon et al., 2021*). Patients with PD-Apathy tend to preferentially choose high value outcomes over options with lower payouts (*Le Heron et al., 2018*). In the ON-DBS state, heightened exploratory choices observed in our data would lead to poorer encoding of the difference in value of decisions and incentivise only actions which lead to large differences in expected outcomes. This assumes that the mechanism of DBS-induced apathy is analogous to apathy acquired from lesions of the bilateral globus pallidus, where the basal ganglia's output is nullified, and a decrease in the signal to noise ratio that encodes value to incentivise action has been proposed (*Levy and Dubois, 2006*).

Study limitations

Patients with primary forms of dystonia who have not undergone DBS exhibit abnormal reinforcement learning (*Gilbertson et al., 2019; Arkadir et al., 2016; Gilbertson et al., 2020*). Comparison of decision making between our patients and a healthy control group replicated previously described RL abnormal bias towards exploitative choices in these patients (*Gilbertson et al., 2019; Arkadir et al., 2016; Gilbertson et al., 2020*). It seems likely that the observed enhancement of exploration by pallidal neuromodulation of GPI is from a high level of exploitation that is intrinsic to an imbalanced basal ganglia circuit related to the disease process in dystonia patients. We therefore cannot discount the possibility that DBS enhanced exploration was a consequence of neuromodulating a pathologically biased decision making circuit. Given the previous experimental (*Sheth et al., 2011*) and theoretical (*Humphries et al., 2012*) support which motivated testing the hypothesis of this study, it seems likely that the same function is subserved by a physiologically intact basal ganglia circuit. Furthermore, allowing for the uniqueness of this clinical indication for DBS in this patient group, these interpretational limitations are offset by the infrequent opportunities to study the role of this brain region in human explore-exploit decision-making.

The use of normative connectomes instead of patient-specific data enhance the signal to noise ratio and image quality (*Horn et al., 2017*) but prevents individual quantification of connectivity strength in relation to the behavioural effect of DBS. Lastly, computational models account for a heuristic explanation of complex and dynamic neural networks and serve mainly as a theoretical support for the observed experimental results. The transferability of computationally inferred mechanisms remains a limitation of this study.

Conclusions

A recent increase in interest in the explore-exploit dilemma has significantly advanced understanding of the functional neuroanatomy, with computationally demanding cortical algorithms resolving multiple-alternative decision problems (*Gershman, 2018; Schulz and Gershman, 2019*). The experimental results presented here demonstrate that a understanding of a unified account of the brain's approach to adaptive behavioural control benefits from inclusion of subcortical circuits. Future research should aim to delineate the differential contributions of subcortical and cortical circuits to explore-exploit decision-making as well as their interaction, leading to insights into disordered decision-making in psychiatric and neurologic conditions (*Addicott et al., 2017*).

Materials and methods

Participants

Nineteen patients with isolated dystonia (14 f, 59.79 ± 1.93 years old; mean \pm S.E.M) who had undergone Deep Brain Stimulation (DBS) surgery targeting the bilateral Globus Pallidus interna (GPI) were enrolled in the study. One participant was unable to complete more than ~50% of trials in the ON-DBS state and was excluded from the final analysis. Three of the remaining 18 patients were able to complete reversal-learning task in the ON-DBS state but were unable to tolerate testing in the OFF-DBS state. One patient's OFF-DBS testing file was inadvertently overwritten and was therefore not available for analysis. For more clinical details of the patients and specifications of the inserted electrodes please see *Supplementary file 1*.

Eighteen healthy controls (14 f, 57.31 ± 1.13 years old) with no previous diagnosis of a neurological or psychiatric disorder were enrolled.

In accordance with the declaration of Helsinki, participants gave written informed consent to participate in the study, which was approved by the local ethics committee (Charité – Universitätsmedizin Berlin, EA1/179/20).

Experimental task

We used a modified version of a two-choice reward reversal-learning task based on *Pessiglione et al., 2006*. The task was presented on a laptop screen using Psychtoolbox v3.0 (*Brainard and Vision, 1997*) running on MATLAB (R2019, The MathWorks, Natick, MA, USA). The laptop was positioned in front of the patient so that responses could be made on the keyboard. Participants were instructed to try to win as many 'vouchers' as possible. Printed screen text was in German. At the beginning of each trial, a fixation cross presented at the centre of the screen indicated that a new trial had begun. A pair of fractal images were presented and subjects were expected to indicate their choice by a keyboard button press. The order of the fractals was randomly assigned to either the right or left of the fixation cross. In the event that they did not make a response within 3 s the fixation cross was re-presented and a new trial began. A choice was highlighted for 0.5 s by a red circle around the chosen fractal. The choice screen was followed by feedback screen for a further 1 s. Feedback consisted of either – 'you win' or a neutral feedback condition with a screen with the words 'nothing'. Patients performed three sessions of 40 trials with 3–4 min breaks between sessions to improve task compliance. The visual stimuli were associated with a fixed probability of rewarding outcome of 80:20%. Reward contingencies were reversed after 60 trial presentations (i.e. the image associated with a 80% reward probability in the first 60 trials was associated with 20% reward probability in the last 60 trials and vice-versa). Participants were not informed about the existence of a contingency reversal. A short training session of 10 trials with a novel pair of fractals was performed before formal testing began to familiarise the patients with the task.

The initial DBS condition was randomised. Patients repeated the same task OFF-DBS in the opposite DBS state after a 20 min break. The fractal images differed between both task versions to avoid learning effects.

Reinforcement learning drift diffusion model

The RLDDM model consists of a hybrid of the sequential sampling model choice rule based on the DDM with the drift rate, ν modified by the difference in the expected values of the two choices. The

DDM calculates the likelihood of the decision time (DT) of choice, x , on trial, t , with the Weiner first-passage time (WFPT) distribution;

$$DT(x) \sim WFPT[a, T_{er}, z, v(t)], \quad (1)$$

The non-decision time and bias (starting point) of the diffusion process are represented by T_{er} and z . The boundary separation parameter a defines the point at which the diffusion process reaches the decision threshold and a response is estimated. The expected values, Q_{upper} and Q_{lower} of the choices represented by the upper and lower response boundaries, are estimated for each choice using the temporal difference learning rule (Rescorla, 1972).

$$Q_t = Q_{t-1} + \alpha * (R - Q_{t-1}), \quad (2)$$

where α is the learning rate and R is the outcome of the choice on trial t . We fitted two variations of the RLDDM model, the first with a single learning rate, the second with dual learning rates, where the Q values were modified separately by a learning rate α^+ that multiplied the positive prediction errors (i.e. outcomes that were better than expected) and by α^- , which multiplied negative prediction errors (when outcomes were worse than expected). The drift rate on each trial is then influenced by the difference in the expected values in the two choices by the drift rate scaling parameter m , where,

$$v(t) = [Q_{upper}(t) - Q_{lower}(t)] * m. \quad (3)$$

Given our hypothesis and prior knowledge that the learning rate, scaling parameter and boundary separation parameters on explore-exploit choices (Pedersen et al., 2017), we estimated values of T_{er} and z with the a priori assumption of no effect of DBS and estimated values of α^+ , α^- , m and a as dependant variables for both on and OFF-DBS states.

Model parameter estimation was performed using the Hierarchical Bayesian framework implemented in Python (3.6.12) which uses the Hierarchical estimation of DDM (HDDM) module (0.8.0; <https://hddm.readthedocs.io/en/latest/>; Wiecki et al., 2013). All models were run with four chains with 7,000 burn-in samples and 15000 posterior samples each. Convergence between the Markov chain Monte Carlo (MCMC) chains was assessed using the Gelman Rubin statistic (Gelman, 2013). For all estimates including those for parameter recovery the Gelman Rubin values were all less than 1.1 indicating successful convergence.

Posterior distributions in the differences in the parameter estimates were estimated by subtracting the posterior estimates for each parameter in the ON and OFF-DBS states. The posterior means and highest density intervals (HDI) were calculated using the 'bayestestR' and 'tidybays' packages in R version 4.1.1 (<http://www.r-project.org>) (Kay, 2019; Makowski et al., 2019).

Mapping of the influence of the difference parameters on the exploratory choice preference in the RLDDM model was performed by simulating 20 synthetic data sets at each parameter combination and estimating the P(Explore) values. Parameter space mapping was performed for the positive learning rate, drift rate scaling and boundary separation parameters (α^+ , m and a). The parameter values for each simulation were centred around the OFF state group mean estimates with each combination of parameters derived from a range across of values relevant to the parameter space estimated from fitting the experimental data.

Lead localisation and connectivity analysis

DBS Electrodes were localised for 14/19 of the patients whom both testing data was available ON and OFF-DBS using the software 'Lead-DBS' V2.5 (Horn et al., 2019) (<https://www.lead-dbs.org/>) as described (de Almeida Marcelino et al., 2019; Neumann et al., 2018) and mapped in MNI template space. Stimulation volumes were modelled based on SimBio-FieldTrip pipeline incorporated in Lead DBS (Bhatia et al., 2018). Stimulation volume modelling was informed by DBS parameters used to control patients' symptoms during behavioural task (Supplementary file 1). Each patient's stimulation volumes were then used as seed regions to estimate individual functional connectivity profile maps using an openly available group connectome derived from resting state functional connectivity MRI images of 1000 neurologically healthy volunteers (Holmes et al., 2015). Then, voxel-wise correlations of individual connectivity to DBS-induced changes in the probability of exploring the lower value choice 'P(Explore)' were calculated and visualised as R-maps (Figure 4B) after thresholding R

(Horn et al., 2017; de Almeida Marcelino et al., 2019; Neumann et al., 2018). The DBS induced change in P(Explore) was defined as the maximum increase in P(Explore) ON-DBS minus P(Explore) OFF-DBS in one of the three sessions. The similarity between the 'ideal' connectivity profile of this R-map (Figure 4B) and the connectivity profile of each individual patient was mathematically assessed by calculating spatial correlation coefficients between the R-map and the individual non-behavioural connectivity maps (Horn et al., 2017; AlFatly et al., 2019). Finally, the predictive potential of the DBS connectivity maps for the increase in exploration was estimated by correlating the individuals increase in P(Explore) with DBS to the spatial correlation coefficient of each patient (Figure 4A). A permutation distribution was constructed from the same analysis performed 1000 times and the correlation coefficient re-estimated with the P(Explore) values for each subject randomly re-ordered.

Generating a heat-map of stimulation volumes

This section describes a method to highlight the subcortical cluster stimulated by most of stimulation volumes in our cohort and its relation to GPi/GPe nuclei (extracted from the DISTAL atlas Ewert et al., 2018). This is important to explore the extent of stimulation and possible differential modulation of the basal ganglia nuclei in the vicinity of the stimulated area. The DISTAL atlas is a probabilistic atlas (which means it offers voxel-wise probability that a voxel is belonging to a specific structure). Therefore, we first thresholded bilateral GPi and GPe images to include only voxels with 50% probability of belonging to the GPi or the GPe and extracted binary masks of them. Later, stimulation volumes were extracted and saved as binary Nifty images as has been described in the main manuscript method section. We then overlapped all binary stimulation volumes (n=14) to extract a total stimulation volume representative of the full cohort. The latter has been explored and depicted as a 3D volume in relation to the GPi/GPe (Supplementary file 1) to visualise the extent of overlap with each of these structures. Next, we summed up the number of stimulation volumes contributing to each voxel in the total stimulation volume. This helped extracting an n-map (or heat-map) which can illustrate the frequency with which each voxel is being stimulated by stimulation volumes in the cohort. As a final step, we have calculated the extent of overlap between the heat-map and the GPi or GPe nuclei as the weighted sum of overlapping voxels. This means that we summed the values of the heat-map intersecting voxels with the binary mask of each of GPi or GPe nucleus. The latter method ensures that the voxel-wise frequency information would not be lost compared to simple binary sum of the overlap of total stimulation volume with GPi/GPe binary masks. All analyses mentioned in this section were performed in a grid of $0.5 \times 0.5 \times 0.5$ mm resolution using bihemispheric information in respect of GPi/GPe or stimulation volumes (since all patients have been bilaterally stimulated).

Behavioural data analysis

Prior to statistical analyses, trials without responses (errors of omission) were excluded ON-DBS mean 1.06 ± 2.26 (range 0–8) trials, OFF-DBS 1.78 ± 2.5 (range 0–7 trials). Accordingly, DT and task performance (number of rewards won in a session) were analyzed using ANOVA with the MATLAB function *anovan*. In case of significant interactions, post-hoc tests were conducted using paired samples t-tests. Normality assumptions were tested using Kolmogorov Smirnov tests (all $p > 0.05$). All results are reported as mean values \pm S.E.M.

Acknowledgements

TPG was funded by a NRS Fellowship from the Chief Scientist Office (Scotland). Additional funding for this research was provided by the University of Dundee/NHS Tayside Movement disorders research Endowment Fund. AL de A Marcelino is a fellow in the BIH Charité Junior Clinician Scientist Program funded by Stiftung Charité.

Additional information

Competing interests

Andrea A Kühn: has received from honoraria from Boston Scientific, Medtronic and Teva. The other authors declare that no competing interests exist.

Funding

Funder	Grant reference number	Author
Chief Scientist Office		Tom Gilbertson
NHS Tayside Movement disorders research Endowment Fund		Tom Gilbertson
Stiftung Charité		Ana Luisa de A Marcelino

The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

Author contributions

Ana Luisa de A Marcelino, Conceptualization, Formal analysis, Investigation, Writing – original draft, Writing – review and editing; Owen Gray, Formal analysis, Methodology, Writing – review and editing; Bassam Al-Fatly, Formal analysis, Visualization, Methodology, Writing – review and editing; William Gilmour, Data curation, Investigation, Writing – review and editing; J Douglas Steele, Supervision, Investigation, Methodology, Writing – review and editing; Andrea A Kühn, Conceptualization, Resources, Supervision, Writing – review and editing; Tom Gilbertson, Conceptualization, Formal analysis, Supervision, Funding acquisition, Investigation, Visualization, Methodology, Writing – original draft, Writing – review and editing

Author ORCIDs

Ana Luisa de A Marcelino  <http://orcid.org/0000-0002-3291-7222>

Bassam Al-Fatly  <http://orcid.org/0000-0003-0067-6177>

Tom Gilbertson  <http://orcid.org/0000-0002-9866-1565>

Ethics

Human subjects: The which was approved by the local ethics committee (Charité - Universitätsmedizin Berlin, EA1/179/20).

Decision letter and Author response

Decision letter <https://doi.org/10.7554/eLife.79642.sa1>

Author response <https://doi.org/10.7554/eLife.79642.sa2>

Additional files

Supplementary files

- MDAR checklist
- Supplementary file 1. Clinical Demographics. TWSTRS: Toronto Western Spasmodic Torticollis Rating Scale; BDI: Beck's Depression Inventory; F: Female; M: Male; S.E.M: standard error of mean; L: left; R: right. *These patients only performed the task in one stimulation condition.
- Supplementary file 2. Summary of Connectivity (R-map) analysis. ACC, anterior cingulate cortex; BA, Brodmann area; CBM, cerebellum; IFG, inferior frontal gyrus, ins., insula; ITG, inferior temporal gyrus; MCC, midcingulate cortex; MedFG, medial frontal gyrus; MFG, middle frontal gyrus; MTG, middle temporal gyrus; OG, orbital gyrus; OL, occipital lobule; PCC, posterior cingulate cortex; PreCG, precentral gyrus; Prec, precuneus; PostCG, postcentral gyrus; SFG, superior frontal gyrus; SMG, supramarginal gyrus; SNr, substantia nigra; STG, superior temporal gyrus.

Data availability

Raw choice and reaction time data, computational model parameter estimates, simulated data and r-maps from connectivity analysis are available via the Open Science Framework <https://osf.io/fs36g/>.

The following dataset was generated:

Author(s)	Year	Dataset title	Dataset URL	Database and Identifier
Marcelino A, Gray O, Al-Fatly B, Gilmour W, Steele D, Kühn AA, Gilbertson T	2022	Explore-Exploit DBS	https://osf.io/fs36g/	Open Science Framework, fs36g

References

- Addicott MA**, Pearson JM, Sweitzer MM, Barack DL, Platt ML. 2017. A primer on foraging and the explore/exploit trade-off for psychiatry research. *Neuropsychopharmacology* **42**:1931–1939. DOI: <https://doi.org/10.1038/npp.2017.108>, PMID: 28553839
- Alexander GE**, DeLong MR, Strick PL. 1986. Parallel organization of functionally segregated circuits linking basal ganglia and cortex. *Annual Review of Neuroscience* **9**:357–381. DOI: <https://doi.org/10.1146/annurev.ne.09.030186.002041>, PMID: 3085570
- AlFatly B**, Ewert S, Kübler D, Kroneberg D, Horn A, Kühn AA. 2019. Connectivity profile of thalamic deep brain stimulation to effectively treat essential tremor. *Brain* **142**:3086–3098. DOI: <https://doi.org/10.1093/brain/awz236>, PMID: 31377766
- Arkadir D**, Radulescu A, Raymond D, Lubarr N, Bressman SB, Mazzoni P, Niv Y. 2016. DYT1 dystonia increases risk taking in humans. *eLife* **5**:e14155. DOI: <https://doi.org/10.7554/eLife.14155>, PMID: 27249418
- Badre D**, Doll BB, Long NM, Frank MJ. 2012. Rostrolateral prefrontal cortex and individual differences in uncertainty-driven exploration. *Neuron* **73**:595–607. DOI: <https://doi.org/10.1016/j.neuron.2011.12.025>, PMID: 22325209
- Bar-Gad I**, Elias S, Vaadia E, Bergman H. 2004. Complex locking rather than complete cessation of neuronal activity in the globus pallidus of a 1-methyl-4-phenyl-1,2,3,6-tetrahydropyridine-treated primate in response to pallidal microstimulation. *The Journal of Neuroscience* **24**:7410–7419. DOI: <https://doi.org/10.1523/JNEUROSCI.1691-04.2004>, PMID: 15317866
- Bartolo R**, Averbeck BB. 2020. Prefrontal cortex predicts state switches during reversal learning. *Neuron* **106**:1044–1054. DOI: <https://doi.org/10.1016/j.neuron.2020.03.024>, PMID: 32315603
- Bhatia KP**, Bain P, Bajaj N, Elble RJ, Hallett M, Louis ED, Raethjen J, Stamelou M, Testa CM, Deuschl G, Tremor Task Force of the International Parkinson and Movement Disorder Society. 2018. Consensus statement on the classification of tremors from the task force on tremor of the international parkinson and movement disorder society. *Movement Disorders* **33**:75–87. DOI: <https://doi.org/10.1002/mds.27121>, PMID: 29193359
- Boes AD**, Prasad S, Liu H, Liu Q, Pascual-Leone A, Caviness VS Jr, Fox MD. 2015. Network localization of neurological symptoms from focal brain lesions. *Brain* **138**:3061–3075. DOI: <https://doi.org/10.1093/brain/awv228>, PMID: 26264514
- Bogacz R**, Gurney K. 2007. The basal ganglia and cortex implement optimal decision making between alternative actions. *Neural Computation* **19**:442–477. DOI: <https://doi.org/10.1162/neco.2007.19.2.442>, PMID: 17206871
- Boorman ED**, Behrens TEJ, Woolrich MW, Rushworth MFS. 2009. How green is the grass on the other side? frontopolar cortex and the evidence in favor of alternative courses of action. *Neuron* **62**:733–743. DOI: <https://doi.org/10.1016/j.neuron.2009.05.014>, PMID: 19524531
- Boraud T**, Bezard E, Bioulac B, Gross C. 1996. High frequency stimulation of the internal globus pallidus (gpi) simultaneously improves parkinsonian symptoms and reduces the firing frequency of gpi neurons in the MPTP-treated monkey. *Neuroscience Letters* **215**:17–20. DOI: [https://doi.org/10.1016/s0304-3940\(96\)12943-8](https://doi.org/10.1016/s0304-3940(96)12943-8), PMID: 8880743
- Brainard DH**, Vision S. 1997. The psychophysics toolbox. *Spatial Vision* **10**:433–436. DOI: <https://doi.org/10.1163/156856897X00357>, PMID: 9176952
- Chakravarthy VS**, Joseph D, Bapi RS. 2010. What do the basal ganglia do? A modeling perspective. *Biol Cybern* **103**:237–253. DOI: <https://doi.org/10.1007/s00422-010-0401-y>, PMID: 20644953
- Chakroun K**, Mathar D, Wiehler A, Ganzer F, Peters J. 2020. Dopaminergic modulation of the exploration/exploitation trade-off in human decision-making. *eLife* **9**:e51260. DOI: <https://doi.org/10.7554/eLife.51260>, PMID: 32484779
- Cleary DR**, Raslan AM, Rubin JE, Bahgat D, Viswanathan A, Heinricher MM, Burchiel KJ. 2013. Deep brain stimulation entrains local neuronal firing in human globus pallidus internus. *Journal of Neurophysiology* **109**:978–987. DOI: <https://doi.org/10.1152/jn.00420.2012>, PMID: 23197451
- Cohen JD**, McClure SM, Yu AJ. 2007. Should I stay or should I go? how the human brain manages the trade-off between exploitation and exploration. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* **362**:933–942. DOI: <https://doi.org/10.1098/rstb.2007.2098>, PMID: 17395573
- Cools R**, Clark L, Owen AM, Robbins TW. 2002. Defining the neural mechanisms of probabilistic reversal learning using event-related functional magnetic resonance imaging. *The Journal of Neuroscience* **22**:4563–4567. DOI: <https://doi.org/10.1523/JNEUROSCI.22-11-04563.2002>, PMID: 12040063

- Corp DT**, Joutsa J, Darby RR, Delnooz CCS, van de Warrenburg BPC, Cooke D, Prudente CN, Ren J, Reich MM, Batla A, Bhatia KP, Jinnah HA, Liu H, Fox MD. 2019. Network localization of cervical dystonia based on causal brain lesions. *Brain* **142**:1660–1674. DOI: <https://doi.org/10.1093/brain/awz112>, PMID: [31099831](https://pubmed.ncbi.nlm.nih.gov/31099831/)
- Costa VD**, Tran VL, Turchi J, Averbeck BB. 2014. Dopamine modulates novelty seeking behavior during decision making. *Behavioral Neuroscience* **128**:556–566. DOI: <https://doi.org/10.1037/a0037128>, PMID: [24911320](https://pubmed.ncbi.nlm.nih.gov/24911320/)
- Costa VD**, Tran VL, Turchi J, Averbeck BB. 2015. Reversal learning and dopamine: a Bayesian perspective. *The Journal of Neuroscience* **35**:2407–2416. DOI: <https://doi.org/10.1523/JNEUROSCI.1989-14.2015>, PMID: [25673835](https://pubmed.ncbi.nlm.nih.gov/25673835/)
- Costa VD**, Mitz AR, Averbeck BB. 2019. Subcortical substrates of explore-exploit decisions in primates. *Neuron* **103**:533–545. DOI: <https://doi.org/10.1016/j.neuron.2019.05.017>, PMID: [31196672](https://pubmed.ncbi.nlm.nih.gov/31196672/)
- Daw ND**, O’Doherty JP, Dayan P, Seymour B, Dolan RJ. 2006. Cortical substrates for exploratory decisions in humans. *Nature* **441**:876–879. DOI: <https://doi.org/10.1038/nature04766>, PMID: [16778890](https://pubmed.ncbi.nlm.nih.gov/16778890/)
- Daw ND**, Gershman SJ, Seymour B, Dayan P, Dolan RJ. 2011. Model-Based influences on humans’ choices and striatal prediction errors. *Neuron* **69**:1204–1215. DOI: <https://doi.org/10.1016/j.neuron.2011.02.027>, PMID: [21435563](https://pubmed.ncbi.nlm.nih.gov/21435563/)
- de Almeida Marcelino AL**, Horn A, Krause P, Kühn AA, Neumann W-J. 2019. Subthalamic neuromodulation improves short-term motor learning in Parkinson’s disease. *Brain* **142**:2198–2206. DOI: <https://doi.org/10.1093/brain/awz152>, PMID: [31169872](https://pubmed.ncbi.nlm.nih.gov/31169872/)
- Dostrovsky JO**, Levy R, Wu JP, Hutchison WD, Tasker RR, Lozano AM. 2000. Microstimulation-induced inhibition of neuronal firing in human globus pallidus. *Journal of Neurophysiology* **84**:570–574. DOI: <https://doi.org/10.1152/jn.2000.84.1.570>, PMID: [10899228](https://pubmed.ncbi.nlm.nih.gov/10899228/)
- Draganski B**, Kherif F, Klöppel S, Cook PA, Alexander DC, Parker GJM, Deichmann R, Ashburner J, Frackowiak RSJ. 2008. Evidence for segregated and integrative connectivity patterns in the human basal ganglia. *The Journal of Neuroscience* **28**:7143–7152. DOI: <https://doi.org/10.1523/JNEUROSCI.1486-08.2008>, PMID: [18614684](https://pubmed.ncbi.nlm.nih.gov/18614684/)
- Dunovan K**, Vich C, Clapp M, Verstynen T, Rubin J, Bogacz R. 2019. Reward-driven changes in striatal pathway competition shape evidence evaluation in decision-making. *PLOS Computational Biology* **15**:e1006998. DOI: <https://doi.org/10.1371/journal.pcbi.1006998>, PMID: [31060045](https://pubmed.ncbi.nlm.nih.gov/31060045/)
- Edlow BL**, Mareyam A, Horn A, Polimeni JR, Witzel T, Tisdall MD, Augustinack JC, Stockmann JP, Diamond BR, Stevens A, Tirrell LS, Folkerth RD, Wald LL, Fischl B, van der Kouwe A. 2019. 7 tesla MRI of the ex vivo human brain at 100 micron resolution. *Scientific Data* **6**:244. DOI: <https://doi.org/10.1038/s41597-019-0254-8>, PMID: [31666530](https://pubmed.ncbi.nlm.nih.gov/31666530/)
- Eisinger RS**, Cagle JN, Alcantara JD, Opri E, Cernera S, Le A, Torres Ponce EM, Lanese J, Nelson B, Lopes J, Hundley C, Ravy T, Wu SS, Foote KD, Okun MS, Gunduz A. 2022. Distinct roles of the human subthalamic nucleus and dorsal pallidum in Parkinson’s disease impulsivity. *Biological Psychiatry* **91**:370–379. DOI: <https://doi.org/10.1016/j.biopsych.2021.03.002>, PMID: [33993998](https://pubmed.ncbi.nlm.nih.gov/33993998/)
- Ewert S**, Plettig P, Li N, Chakravarty MM, Collins DL, Herrington TM, Kühn AA, Horn A. 2018. Toward defining deep brain stimulation targets in MNI space: a subcortical atlas based on multimodal MRI, histology and structural connectivity. *NeuroImage* **170**:271–282. DOI: <https://doi.org/10.1016/j.neuroimage.2017.05.015>, PMID: [28536045](https://pubmed.ncbi.nlm.nih.gov/28536045/)
- Findling C**, Skvortsova V, Dromnelle R, Palminteri S, Wyart V. 2019. Computational noise in reward-guided learning drives behavioral variability in volatile environments. *Nature Neuroscience* **22**:2066–2077. DOI: <https://doi.org/10.1038/s41593-019-0518-9>, PMID: [31659343](https://pubmed.ncbi.nlm.nih.gov/31659343/)
- Findling C**, Wyart V. 2021. Computation noise in human learning and decision-making: origin, impact, function. *Current Opinion in Behavioral Sciences* **38**:124–132. DOI: <https://doi.org/10.1016/j.cobeha.2021.02.018>
- Gelman A**. 2013. Two simple examples for understanding posterior p-values whose distributions are far from uniform. *Electronic Journal of Statistics* **7**:2595–2602. DOI: <https://doi.org/10.1214/13-EJS854>
- Gershman SJ**. 2018. Deconstructing the human algorithms for exploration. *Cognition* **173**:34–42. DOI: <https://doi.org/10.1016/j.cognition.2017.12.014>, PMID: [29289795](https://pubmed.ncbi.nlm.nih.gov/29289795/)
- Ghahremani DG**, Monterosso J, Jentsch JD, Bilder RM, Poldrack RA. 2010. Neural components underlying behavioral flexibility in human reversal learning. *Cerebral Cortex* **20**:1843–1852. DOI: <https://doi.org/10.1093/cercor/bhp247>, PMID: [19915091](https://pubmed.ncbi.nlm.nih.gov/19915091/)
- Gilbertson T**, Humphries M, Steele JD. 2019. Maladaptive striatal plasticity and abnormal reward-learning in cervical dystonia. *The European Journal of Neuroscience* **50**:3191–3204. DOI: <https://doi.org/10.1111/ejn.14414>, PMID: [30955204](https://pubmed.ncbi.nlm.nih.gov/30955204/)
- Gilbertson T**, Arkadir D, Steele JD. 2020. Opposing patterns of abnormal D1 and D2 receptor dependent cortico-striatal plasticity explain increased risk taking in patients with DYT1 dystonia. *PLOS ONE* **15**:e0226790. DOI: <https://doi.org/10.1371/journal.pone.0226790>, PMID: [32365120](https://pubmed.ncbi.nlm.nih.gov/32365120/)
- Gilbertson T**, Steele D. 2021. Tonic dopamine, uncertainty and basal ganglia action selection. *Neuroscience* **466**:109–124. DOI: <https://doi.org/10.1016/j.neuroscience.2021.05.010>, PMID: [34015370](https://pubmed.ncbi.nlm.nih.gov/34015370/)
- Gurney K**, Prescott TJ, Redgrave P. 2001. A computational model of action selection in the basal ganglia. I. A new functional anatomy. *Biol Cybern* **84**:401–410. DOI: <https://doi.org/10.1007/PL00007984>, PMID: [11417052](https://pubmed.ncbi.nlm.nih.gov/11417052/)
- Hampshire A**, Chaudhry AM, Owen AM, Roberts AC. 2012. Dissociable roles for lateral orbitofrontal cortex and lateral prefrontal cortex during preference driven reversal learning. *NeuroImage* **59**:4102–4112. DOI: <https://doi.org/10.1016/j.neuroimage.2011.10.072>, PMID: [22075266](https://pubmed.ncbi.nlm.nih.gov/22075266/)

- Hare TA**, Schultz W, Camerer CF, O'Doherty JP, Rangel A. 2011. Transformation of stimulus value signals into motor commands during simple choice. *PNAS* **108**:18120–18125. DOI: <https://doi.org/10.1073/pnas.1109322108>, PMID: 22006321
- Hayden BY**, Pearson JM, Platt ML. 2011. Neuronal basis of sequential foraging decisions in a patchy environment. *Nature Neuroscience* **14**:933–939. DOI: <https://doi.org/10.1038/nn.2856>, PMID: 21642973
- Herrnstein RJ**. 1970. On the law of effect. *Journal of the Experimental Analysis of Behavior* **13**:243–266. DOI: <https://doi.org/10.1901/jeab.1970.13-243>, PMID: 16811440
- Herrojo Ruiz M**, Huebl J, Schönecker T, Kupsch A, Yarrow K, Krauss JK, Schneider GH, Kühn AA. 2014. Involvement of human internal globus pallidus in the early modulation of cortical error-related activity. *Cerebral Cortex* **24**:1502–1517. DOI: <https://doi.org/10.1093/cercor/bht002>, PMID: 23349222
- Holmes AJ**, Hollinshead MO, O'Keefe TM, Petrov VI, Fariello GR, Wald LL, Fischl B, Rosen BR, Mair RW, Roffman JL, Smoller JW, Buckner RL. 2015. Brain genomics superstruct project initial data release with structural, functional, and behavioral measures. *Scientific Data* **2**:150031. DOI: <https://doi.org/10.1038/sdata.2015.31>, PMID: 26175908
- Holroyd CB**, Nieuwenhuis S, Yeung N, Cohen JD. 2003. Errors in reward prediction are reflected in the event-related brain potential. *Neuroreport* **14**:2481–2484. DOI: <https://doi.org/10.1097/00001756-200312190-00037>, PMID: 14663214
- Hong S**, Hikosaka O. 2008. The globus pallidus sends reward-related signals to the lateral habenula. *Neuron* **60**:720–729. DOI: <https://doi.org/10.1016/j.neuron.2008.09.035>, PMID: 19038227
- Horn A**, Reich M, Vorwerk J, Li N, Wenzel G, Fang Q, Schmitz-Hübsch T, Nickl R, Kupsch A, Volkmann J, Kühn AA, Fox MD. 2017. Connectivity predicts deep brain stimulation outcome in Parkinson disease. *Annals of Neurology* **82**:67–78. DOI: <https://doi.org/10.1002/ana.24974>, PMID: 28586141
- Horn A**, Li N, Dembek TA, Kappel A, Boulay C, Ewert S, Tietze A, Husch A, Perera T, Neumann W-J, Reiser M, Si H, Oostenveld R, Rorden C, Yeh F-C, Fang Q, Herrington TM, Vorwerk J, Kühn AA. 2019. Lead-DBS V2: towards a comprehensive pipeline for deep brain stimulation imaging. *NeuroImage* **184**:293–316. DOI: <https://doi.org/10.1016/j.neuroimage.2018.08.068>, PMID: 30179717
- Humphries MD**, Khamassi M, Gurney K. 2012. Dopaminergic control of the exploration-exploitation trade-off via the basal ganglia. *Frontiers in Neuroscience* **6**:9. DOI: <https://doi.org/10.3389/fnins.2012.00009>, PMID: 22347155
- Irmen F**, Horn A, Meder D, Neumann W-J, Plettig P, Schneider G-H, Siebner HR, Kühn AA. 2019. Sensorimotor subthalamic stimulation restores risk-reward trade-off in Parkinson's disease. *Movement Disorders* **34**:366–376. DOI: <https://doi.org/10.1002/mds.27576>, PMID: 30485537
- Izquierdo A**, Brigman JL, Radke AK, Rudebeck PH, Holmes A. 2017. The neural basis of reversal learning: an updated perspective. *Neuroscience* **345**:12–26. DOI: <https://doi.org/10.1016/j.neuroscience.2016.03.021>, PMID: 26979052
- Kauvar IV**, Machado TA, Yuen E, Kochalka J, Choi M, Allen WE, Wetzstein G, Deisseroth K. 2020. Cortical observation by synchronous multifocal optical sampling reveals widespread population encoding of actions. *Neuron* **107**:351–367. DOI: <https://doi.org/10.1016/j.neuron.2020.04.023>, PMID: 32433908
- Kay M**. 2019. Tidybayes: tidy data and geoms for bayesian models. **1.1.1**. R Package.
- Kravitz AV**, Freeze BS, Parker PRL, Kay K, Thwin MT, Deisseroth K, Kreitzer AC. 2010. Regulation of parkinsonian motor behaviours by optogenetic control of basal ganglia circuitry. *Nature* **466**:622–626. DOI: <https://doi.org/10.1038/nature09159>, PMID: 20613723
- Kravitz A.V.**, Tye LD, Kreitzer AC. 2012. Distinct roles for direct and indirect pathway striatal neurons in reinforcement. *Nature Neuroscience* **15**:816–818. DOI: <https://doi.org/10.1038/nn.3100>, PMID: 22544310
- Kühn AA**, Meyer BU, Trottenberg T, Brandt SA, Schneider GH, Kupsch A. 2003. Modulation of motor cortex excitability by pallidal stimulation in patients with severe dystonia. *Neurology* **60**:768–774. DOI: <https://doi.org/10.1212/01.wnl.0000044396.64752.4c>, PMID: 12629231
- LafreniereRoula M**, Kim E, Hutchison WD, Lozano AM, Hodaie M, Dostrovsky JO. 2010. High-frequency microstimulation in human globus pallidus and substantia nigra. *Experimental Brain Research* **205**:251–261. DOI: <https://doi.org/10.1007/s00221-010-2362-8>, PMID: 20640411
- Lee J**, Sabatini BL. 2021. Striatal indirect pathway mediates exploration via collicular competition. *Nature* **599**:645–649. DOI: <https://doi.org/10.1038/s41586-021-04055-4>, PMID: 34732888
- Le Heron C**, Plant O, Manohar S, Ang Y-S, Jackson M, Lennox G, Hu MT, Husain M. 2018. Distinct effects of apathy and dopamine on effort-based decision-making in Parkinson's disease. *Brain* **141**:1455–1469. DOI: <https://doi.org/10.1093/brain/awy110>, PMID: 29672668
- Levy R**, Dubois B. 2006. Apathy and the functional anatomy of the prefrontal cortex-basal ganglia circuits. *Cerebral Cortex* **16**:916–928. DOI: <https://doi.org/10.1093/cercor/bhj043>, PMID: 16207933
- Makowski D**, Ben-Shachar M, Lüdtke D. 2019. BayestestR: describing effects and their uncertainty, existence and significance within the bayesian framework. *Journal of Open Source Software* **4**:1541. DOI: <https://doi.org/10.21105/joss.01541>
- McCairn KW**, Turner RS. 2009. Deep brain stimulation of the globus pallidus internus in the parkinsonian primate: local entrainment and suppression of low-frequency oscillations. *Journal of Neurophysiology* **101**:1941–1960. DOI: <https://doi.org/10.1152/jn.91092.2008>, PMID: 19164104
- Mehlhorn K**, Newell BR, Todd PM, Lee MD, Morgan K, Braithwaite VA, Hausmann D, Fiedler K, Gonzalez C. 2015. Unpacking the exploration–exploitation tradeoff: a synthesis of human and animal literatures. *Decision* **2**:191–215. DOI: <https://doi.org/10.1037/dec0000033>

- Neumann W-J**, Schroll H, de Almeida Marcelino AL, Horn A, Ewert S, Irmen F, Krause P, Schneider G-H, Hamker F, Kühn AA. 2018. Functional segregation of basal ganglia pathways in Parkinson's disease. *Brain* **141**:2655–2669. DOI: <https://doi.org/10.1093/brain/awy206>, PMID: 30084974
- Obeso JA**, Jahanshahi M, Alvarez L, Macias R, Pedroso I, Wilkinson L, Pavon N, Day B, Pinto S, Rodríguez-Oroz MC, Tejeiro J, Artieda J, Talelli P, Swayne O, Rodríguez R, Bhatia K, Rodríguez-Diaz M, Lopez G, Guridi J, Rothwell JC. 2009. What can man do without basal ganglia motor output? the effect of combined unilateral subthalamotomy and pallidotomy in a patient with Parkinson's disease. *Experimental Neurology* **220**:283–292. DOI: <https://doi.org/10.1016/j.expneurol.2009.08.030>, PMID: 19744484
- Pedersen ML**, Frank MJ, Biele G. 2017. The drift diffusion model as the choice rule in reinforcement learning. *Psychonomic Bulletin & Review* **24**:1234–1251. DOI: <https://doi.org/10.3758/s13423-016-1199-y>, PMID: 27966103
- Pessiglione M**, Seymour B, Flandin G, Dolan RJ, Frith CD. 2006. Dopamine-Dependent prediction errors underpin reward-seeking behaviour in humans. *Nature* **442**:1042–1045. DOI: <https://doi.org/10.1038/nature05051>, PMID: 16929307
- Piron C**, Kase D, Topalidou M, Goillandeau M, Orignac H, N'Guyen T-H, Rougier N, Boraud T. 2016. The globus pallidus pars interna in goal-oriented and routine behaviors: resolving a long-standing paradox. *Movement Disorders* **31**:1146–1154. DOI: <https://doi.org/10.1002/mds.26542>, PMID: 26900137
- Ratcliff R**. 1978. A theory of memory retrieval. *Psychological Review* **85**:59–108. DOI: <https://doi.org/10.1037/0033-295X.85.2.59>
- Redgrave P**, Rodriguez M, Smith Y, Rodriguez-Oroz MC, Lehericy S, Bergman H, Agid Y, DeLong MR, Obeso JA. 2010. Goal-Directed and habitual control in the basal ganglia: implications for Parkinson's disease. *Nature Reviews. Neuroscience* **11**:760–772. DOI: <https://doi.org/10.1038/nrn2915>, PMID: 20944662
- Remijne PL**, Nielen MMA, Uylings HBM, Veltman DJ. 2005. Neural correlates of a reversal learning task with an affectively neutral baseline: an event-related fMRI study. *NeuroImage* **26**:609–618. DOI: <https://doi.org/10.1016/j.neuroimage.2005.02.009>, PMID: 15907318
- Rescorla RA**. 1972. A theory of pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement. *Current Research and Theory* **1**:64–99.
- Romano R**, Bertolino A, Gigante A, Martino D, Livrea P, Defazio G. 2014. Impaired cognitive functions in adult-onset primary cranial cervical dystonia. *Parkinsonism & Related Disorders* **20**:162–165. DOI: <https://doi.org/10.1016/j.parkreldis.2013.10.008>, PMID: 24161376
- Rushworth MFS**, Behrens TEJ. 2008. Choice, uncertainty and value in prefrontal and cingulate cortex. *Nature Neuroscience* **11**:389–397. DOI: <https://doi.org/10.1038/nn2066>, PMID: 18368045
- Schroll H**, Horn A, Gröschel C, Brücke C, Lütjens G, Schneider G-H, Krauss JK, Kühn AA, Hamker FH. 2015. Differential contributions of the globus pallidus and ventral thalamus to stimulus-response learning in humans. *NeuroImage* **122**:233–245. DOI: <https://doi.org/10.1016/j.neuroimage.2015.07.061>, PMID: 26220740
- Schulz E**, Gershman SJ. 2019. The algorithmic architecture of exploration in the human brain. *Current Opinion in Neurobiology* **55**:7–14. DOI: <https://doi.org/10.1016/j.conb.2018.11.003>, PMID: 30529148
- Shenhav A**, Cohen JD, Botvinick MM. 2016. Dorsal anterior cingulate cortex and the value of control. *Nature Neuroscience* **19**:1286–1291. DOI: <https://doi.org/10.1038/nn.4384>, PMID: 27669989
- Sheth SA**, Abuelem T, Gale JT, Eskandar EN. 2011. Basal ganglia neurons dynamically facilitate exploration during associative learning. *The Journal of Neuroscience* **31**:4878–4885. DOI: <https://doi.org/10.1523/JNEUROSCI.3658-10.2011>, PMID: 21451026
- Spiegelhalter DJ**, Best NG, Carlin BP, van der Linde A. 2002. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society* **64**:583–639. DOI: <https://doi.org/10.1111/1467-9868.00353>
- Suryanarayana SM**, Hellgren Kotaleski J, Grillner S, Gurney KN. 2019. Roles for globus pallidus externa revealed in a computational model of action selection in the basal ganglia. *Neural Networks* **109**:113–136. DOI: <https://doi.org/10.1016/j.neunet.2018.10.003>, PMID: 30414556
- Sutton RS**, Barto AG. 2018. Reinforcement Learning: An Introduction. MIT press.
- Tecuapetla F**, Jin X, Lima SQ, Costa RM. 2016. Complementary contributions of striatal projection pathways to action initiation and execution. *Cell* **166**:703–715. DOI: <https://doi.org/10.1016/j.cell.2016.06.032>, PMID: 27453468
- Tisch S**, Rothwell JC, Bhatia KP, Quinn N, Zrinzo L, Jahanshahi M, Ashkan K, Hariz M, Limousin P. 2007. Pallidal stimulation modifies after-effects of paired associative stimulation on motor cortex excitability in primary generalised dystonia. *Experimental Neurology* **206**:80–85. DOI: <https://doi.org/10.1016/j.expneurol.2007.03.027>, PMID: 17498697
- Tomov MS**, Truong VQ, Hundia RA, Gershman SJ. 2020. Dissociable neural correlates of uncertainty underlie different exploration strategies. *Nature Communications* **11**:2371. DOI: <https://doi.org/10.1038/s41467-020-15766-z>, PMID: 32398675
- Volkman J**, Mueller J, Deuschl G, Kühn AA, Krauss JK, Poewe W, Timmermann L, Falk D, Kupsch A, Kivi A, Schneider G-H, Schnitzler A, Südmeyer M, Voges J, Wolters A, Wittstock M, Müller J-U, Hering S, Eisner W, Vesper J, et al. 2014. Pallidal neurostimulation in patients with medication-refractory cervical dystonia: a randomised, sham-controlled trial. *The Lancet. Neurology* **13**:875–884. DOI: [https://doi.org/10.1016/S1474-4422\(14\)70143-7](https://doi.org/10.1016/S1474-4422(14)70143-7), PMID: 25127231
- White JK**, Bromberg-Martin ES, Heilbronner SR, Zhang K, Pai J, Haber SN, Monosov IE. 2019. A neural network for information seeking. *Nature Communications* **10**:5168. DOI: <https://doi.org/10.1038/s41467-019-13135-z>, PMID: 31727893

- Wiecki TV**, Sofer I, Frank MJ. 2013. HDDM: hierarchical Bayesian estimation of the drift-diffusion model in python. *Frontiers in Neuroinformatics* **7**:14. DOI: <https://doi.org/10.3389/fninf.2013.00014>, PMID: 23935581
- Wilson RC**, Geana A, White JM, Ludvig EA, Cohen JD. 2014. Humans use directed and random exploration to solve the explore-exploit dilemma. *Journal of Experimental Psychology. General* **143**:2074–2081. DOI: <https://doi.org/10.1037/a0038199>, PMID: 25347535
- Wu YR**, Levy R, Ashby P, Tasker RR, Dostrovsky JO. 2001. Does stimulation of the GPI control dyskinesia by activating inhibitory axons? *Movement Disorders* **16**:208–216. DOI: <https://doi.org/10.1002/mds.1046>, PMID: 11295772
- Yeo BTT**, Krienen FM, Sepulcre J, Sabuncu MR, Lashkari D, Hollinshead M, Roffman JL, Smoller JW, Zöllei L, Polimeni JR, Fischl B, Liu H, Buckner RL. 2011. The organization of the human cerebral cortex estimated by intrinsic functional connectivity. *Journal of Neurophysiology* **106**:1125–1165. DOI: <https://doi.org/10.1152/jn.00338.2011>, PMID: 21653723
- Yin HH**, Mulcare SP, Hilário MRF, Clouse E, Holloway T, Davis MI, Hansson AC, Lovinger DM, Costa RM. 2009. Dynamic reorganization of striatal circuits during the acquisition and consolidation of a skill. *Nature Neuroscience* **12**:333–341. DOI: <https://doi.org/10.1038/nn.2261>, PMID: 19198605
- Zoon TJC**, van Rooijen G, Balm GMFC, Bergfeld IO, Daams JG, Krack P, Denys DAJP, de Bie RMA. 2021. Apathy induced by subthalamic nucleus deep brain stimulation in Parkinson's disease: a meta-analysis. *Movement Disorders* **36**:317–326. DOI: <https://doi.org/10.1002/mds.28390>, PMID: 33331023