

Modulation of neural activity in frontopolar cortex drives reward-based motor learning

Maria Herrojo Ruiz (✉ M.Herrojo-Ruiz@gold.ac.uk)

Goldsmiths University of London

Tom Maudrich

Max Planck Institute for Human Cognitive and Brain Sciences

Benjamin Kalloch

Max Planck Institute for Human Cognitive and Brain Sciences

Daniela Sammler

Max Planck Institute for Human Cognitive and Brain Sciences

Rouven Kenville

Max Planck Institute for Human Cognitive and Brain Sciences

Arno Villringer

Max Planck Institute for Human Cognitive and Brain Sciences

Bernhard Sehm

Max Planck Institute for Human Cognitive and Brain Sciences

Vadim Nikulin

Max Planck Institute for Human Cognitive and Brain Sciences

Research Article

Keywords: frontopolar cortex, tDCS, motor learning, reward

Posted Date: March 18th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-318224/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Modulation of neural activity in frontopolar cortex drives reward-based motor learning

M Herrojo Ruiz^{a,b,*}, T Maudrich^c, B Kalloch^c, D Sammler^c, R Kenville^c, A Villringer^c, B Sehm^c, and V Nikulin^{b,c}

^aPsychology Department, Goldsmiths University of London, UK

^bCenter for Cognition and Decision Making, National Research University Higher School of Economics, Russian Federation

^cDepartment of Neurology, Max Planck Institute for Human Cognitive and Brain Sciences, Germany

March 11, 2021

*To whom correspondence should be addressed. E-mail: M.Herrojo-Ruiz@gold.ac.uk

Abstract

The frontopolar cortex (FPC) contributes to tracking the reward of alternative choices during decision making, as well as their reliability. Whether this FPC function extends to reward gradients associated with continuous movements during motor learning remains unknown. We used anodal transcranial direct current stimulation (tDCS) over the right FPC to investigate its role in reward-based motor learning. Nineteen healthy human participants completed a motor sequence learning task using trialwise reward feedback to discover a hidden goal along a continuous dimension: timing. As additional conditions, we modulated the contralateral motor cortex (left M1) activity, and included a control sham stimulation. Right FPC-tDCS led to faster learning compared to LM1-tDCS and sham through regulation of motor variability. Computational modelling revealed that in all stimulation protocols, an increase in the trialwise expectation of reward was followed by greater exploitation, as shown previously. Yet, this association was weaker in LM1-tDCS suggesting a less efficient learning strategy. The effects of frontopolar stimulation were dissociated from those induced by LM1-tDCS and sham, as motor exploration was more sensitive to inferred changes in the reward tendency (volatility). The findings suggest that rFPC-tDCS increases the sensitivity of motor exploration to updates in reward volatility, accelerating reward-based motor learning.

Keywords: frontopolar cortex | tDCS | motor learning | reward

Introduction

One of the hallmarks of motor skill learning is the reduction in movement variability [1, 2]. As the dancer learns to perform pirouettes, the irregularity in the movement decreases and the turns become smoother. In this context, movement variability is regarded as motor noise, and is dissociated from the intentional use of motor variability, termed motor exploration [3, 4]. Motor learning can also involve motor exploration, particularly when learning from reinforcement, such as feedback about success or failure [5, 2, 4]. In this scenario, initial exploration of movement variables in a continuous space is followed by the exploitation of inferred optimal movements and gradually refined by reducing motor noise [6]. How agents decide about motor exploration and exploitation is the subject of an increasing number of studies on motor learning as a decision-making process [7, 4, 8, 9]. But if motor learning depends on making the right decisions about a movement, one can hypothesize that brain regions involved in regulating the exploration-exploitation tradeoff in cognitive or economic decision-making tasks would also modulate the use of variability during motor decision making.

Here, we postulate that the human frontopolar cortex (FPC) may have a crucial role in driving motor decision making when external reward signals are available to drive exploration-exploitation. In decision-making tasks involving multiple choices, the right FPC has been identified as promoting exploration [10, 11], and tracking the reward value associated with competing options, strategies or goals [12]. Furthermore, FPC accelerates the learning of novel rules [13]. Because motor learning takes place in a continuous movement space [7], we reasoned that the capability of FPC to monitor multiple discrete choices and their reward would make it an ideal candidate to track the reward associated with continuous movement parameters.

To identify the role of the FPC in reward-based motor learning, we used anodal transcranial direct current stimulation (tDCS) over the right FPC while participants completed our recently developed reward-dependent motor sequence learning paradigm [14]. This task required initial exploration along a continuous dimension (timing). Trial-to-trial exploratory behavioural changes were assessed using computational modelling. The *right* FPC was selected due to its greater engagement in regulating the exploration-exploitation balance [15, 11]. We also took an independent measurement of each participant's baseline motor noise [4].

As control tDCS condition we modulated motor cortex activity using contralateral (left) M1 tDCS [16, 17]. Animal and human neurophysiology studies highlight a crucial role of M1 in modulating motor variability and in processing reward [18, 6, 14]. Additional evidence comes from transcranial magnetic stimulation (TMS) and tDCS studies [19, 20, 21, 22, 23], with investigations also linking M1 to decision making [24, 25, 26]. In our study, active tDCS conditions were contrasted to sham stimulation. The effect of the tDCS stimulation protocols on

the individual brain was further assessed using simulations of the electric field strength guided by individual T1-weighted anatomical magnetic resonance images (MRI).

The central hypothesis was that rFPC-tDCS improves learning of the continuous reward landscape associated with movement parameters by balancing the exploration-exploitation tradeoff. Furthermore, we predicted that lM1-tDCS also benefits learning of action-reward associations. However, because the abovementioned studies linking M1 to motor decision making used binary feedback, categorical decisions or simple reward functions, we expected a benefit of rFPC over lM1-tDCS when motor learning occurs in a continuous reward landscape governed by uncertainty. Last, lM1-tDCS was hypothesized to reduce motor noise relative to rFPC stimulation.

Results

Nineteen right-handed participants took part in our study implementing a sham-controlled, double-blinded, cross-over design. They underwent each of three types of a tDCS protocol (rFPC, lM1, sham condition) over three separate weeks in a pseudo-randomized counterbalanced order across participants. Study procedure for each session was identical, with active or sham tDCS applied to the target area for a period of 20 min during task performance (Figure 1). One exception was the last block of the learning phase, which was initiated 5 min after the cessation of tDCS and thus served to assess offline effects on learning. Notably, however, anodal tDCS stimulation effects on motor learning have been shown to last for at least 30 minutes after halting tDCS stimulation [17, 27], but see [28] showing null results on immediate offline motor effects. Accordingly, we assumed that the recent stimulation would strongly influence performance during the last block (~ 5 minutes after tDCS cessation). However, to account for a possible differential effect of tDCS protocols on the offline (3) and online (1) blocks, we performed a specific additional offline minus online contrast in the statistical analysis (see *Materials and Methods*).

Participants received either right FPC or left M1 anodal tDCS, or sham tDCS (targetting either rFPC or lM1, 50% – 50% split across participants). In each tDCS session, they completed a motor task with their right hand on a digital piano (Yamaha Clavinova CLP-150). The task consisted of an initial baseline phase of 20 trials of regular isochronous performance, followed by three blocks of 30 trials of reward-based sequence learning (Figure 1A). The pitch content of the sequences being used for the reward-based learning blocks is displayed in Figure 2 (see details in *Materials and Methods*). The baseline phase allowed us to obtain a measure of motor noise, which represents the residual variability that is expressed when aiming to accurately reproduce

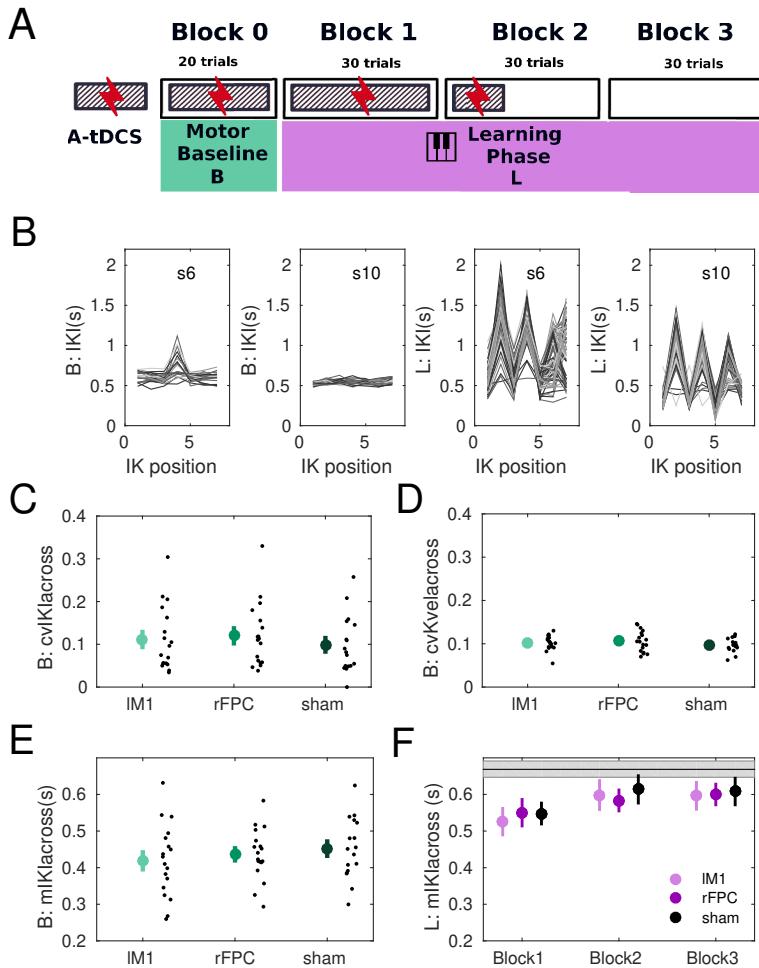


Figure 1: Experimental design and behavioural results. (A) Participants were tested on three separate weeks during which either an active tDCS protocol over the IM1, or the rFPC, or a sham stimulation condition were applied. All tDCS protocols extended for 20 minutes, which included (i) an initial 3 minute resting phase, (ii) a baseline phase of regular isochronous motor performance, and (iii) part of the reward-based learning blocks (1.333). (B) Illustration of timing performance with sham stimulation during the baseline phase (B; panels 1-2) and learning phase (L; panels 3-4) in two participants. Timing was measured using the inter-keystroke-interval (IKI) in seconds, and shown for each inter-keystroke position (1 to 7 for sequences of 8 key presses). Different trajectories denote performance in different trials. (C-E) Effects of stimulation conditions on performance variables during the baseline phase (B). Colored big dots indicate means, with error bars denoting \pm SEM. (C) Across-trials temporal variability, measured with the coefficient of variation of IKI or cvIKlacross in each block; (D) across-trials variability in keystroke velocity or loudness, cvKvelacross; (E) mean performance tempo, mIKI (s). (F) Average performance tempo during learning (L). The gray horizontal line indicates the mean tempo of the hidden target solutions.

the same action [29, 4]. During reward-based learning blocks, participants completed our recently developed reward-based motor sequence learning task [14]. In this task, participants received continuous reward in the form of a trialwise feedback score (range 0-100) to discover a hidden performance goal: a timing pattern. The use of continuous feedback scores to guide learning was based on previous studies investigating motor variability during reward-based learning [5, 4] and motor decision making [30]. Continuous feedback has been shown to be more informative than binary signals (success/failure), contributing to faster learning [31].

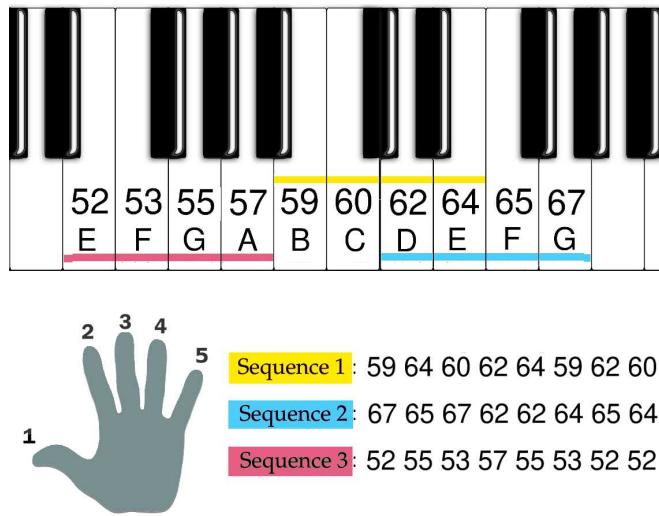


Figure 2: Stimulus materials. The pitch content of the sequences being used for the reward-based learning blocks. Sequences 1,2,3 consisted of different combinations of a set of four pitch values from four neighbouring white keys on the piano. During the preceding baseline phase, participants had to press those same keys but in a successive order (four notes upwards + same four notes downwards) regularly at a self-paced tempo and with their index to little fingers (digits 2-5 shown in the figure).

Behavioral changes across blocks

Analysis of the behavioural data focused on the assessment of motor variability across trials and along two different dimensions: time, which was the *instructed* task-related dimension and measured here using the inter-keystroke-interval (IKI, seconds) index; and keystroke velocity (Kvel, a.u.), which was the non-task-related dimension and is associated with the loudness of the key press. We used the coefficient of variation ($cv = sd / mean$) across trials ($cvacross$) to assess the extent of variability within a block in relation to the mean of the sample. The achieved scores and other general performance variables were also evaluated (see *Materials and Methods*).

The different tDCS protocols did not have dissociable effects on baseline motor variability

for timing (cvIKIacross ; $P = 0.91$, one-way factorial analysis with synchronized rearrangements; Figure 1). Neither was there a significant main effect of factor Stimulation on cvKvelacross at baseline ($P = 0.85$). General performance parameters did not differ in this phase of the experiment as a function of the stimulation protocol either (*Supplementary Results*), suggesting that any differential stimulation effects on the subsequent learning phases are not confounded by baseline effects.

During reward-based learning participants demonstrated an initial tendency to explore timing solutions that were close-by but also far away in the movement space (Figure S1). This indicated that initially participants explored the task-related temporal dimension. They explored timing patterns by (i) changing the ratio of IKI values across neighboring keystrokes (different shape of IKI patterns) but also by (ii) changing IKI values at individual keystrokes, while keeping neighbouring IKI values unchanged (same shape of IKI patterns). This second scenario corresponded with timing solutions that were close-by in the movement space (Figure S2). During the last learning block, participants consistently exploited the inferred rewarded solution (Figure S3). In a minority of the cases, however, exploration of different timing solutions was manifest even at the end of the experimental session.

Statistical analysis during reward-based learning demonstrated that participants improved their scores across blocks in all tDCS conditions (Figure 3A; main effect Block, $P = 0.0002$, full 3×3 non-parametric factorial analysis). These data support that participants successfully used the trialwise feedback in all tDCS conditions to learn about the hidden goal. There was no significant main effect for Stimulation or interaction effect ($P > 0.05$).

A separate one-way factorial analysis on the change in scores from online to offline blocks (3 minus 1) across stimulation conditions demonstrated a significant effect of factor Stimulation ($P = 0.0355$). *Post hoc* pair-wise comparisons for each pair of tDCS stimulation conditions revealed a significantly larger increase following rFPC-tDCS relative to LM1-tDCS (Figure 3B; increase of 19.9 [standard error of the mean or SEM 3.77] for rFPC-tDCS; increase of 12.6 [3.8] for LM1-tDCS; $P \leq P_{FDR} = 0.0260$; moderate effect size, assessed with a non-parametric effect size estimator for dependent samples [32], $\Delta_{dep} = 0.72$, confidence interval or CI = [0.61, 0.83]). There was also a significantly larger increase for rFPC-tDCS relative to sham ($P \leq P_{FDR} = 0.0260$; moderate effect size, $\Delta_{dep} = 0.62$, CI = [0.50, 0.80]; the average score change for sham was 12.3 [2.40]). No differences between LM1 and sham were found ($P > 0.05$).

The general increase in scores across blocks was paralleled by a reduction in the expression of task-related motor variability, measured with cvIKIacross (Figure 3C; significant main effect of Block in the full 3×3 non-parametric factorial analysis; $P = 0.0156$). Notably, we also found a significant main effect of factor Stimulation ($P = 0.0260$), but no interaction effect. *Post*

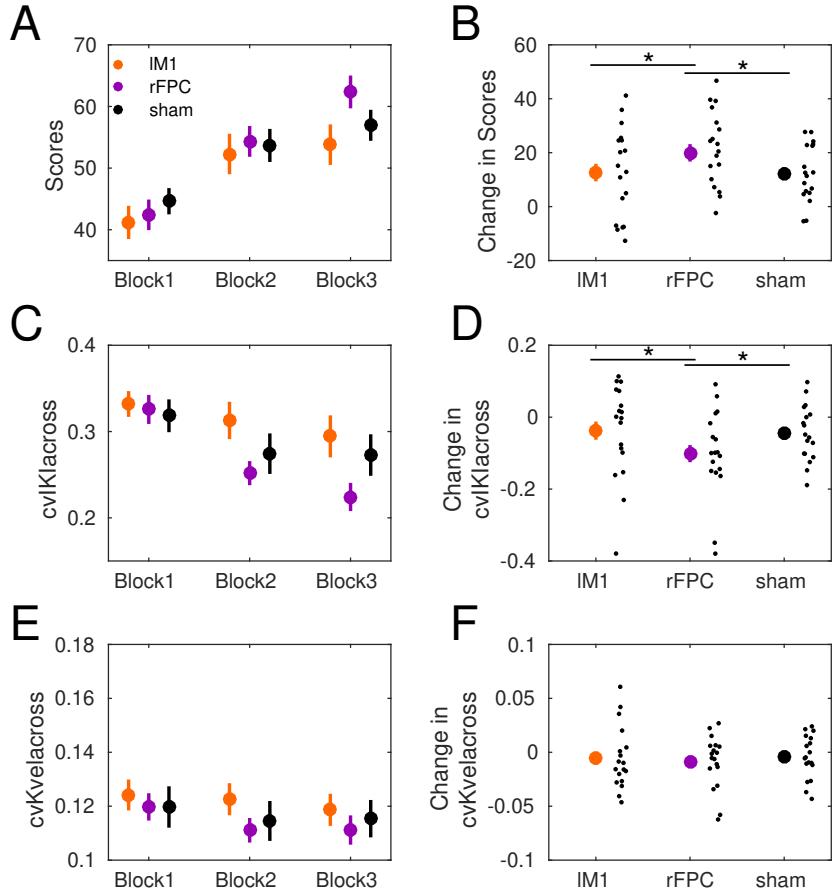


Figure 3: Behavioral results during reward-based learning. **(A)** Participants increased their scores across blocks (significant main effect of Block in full factorial analysis using factors Block (1-3) and Stimulation (IM1, rFPC, sham), supporting they successfully used the trial-by-trial feedback to approach the hidden performance goal. **(B)** The change in scores from online to offline blocks (3 minus 1) was significantly larger in rFPC-tDCS than sham, and in rFPC-tDCS relative to IM1-tDCS (denoted by the asterisk; $P \leq P_{FDR} = 0.0260$; moderate effect sizes: $\Delta_{dep} = 0.68$, CI = [0.50, 0.83] for rFPC-tDCS and IM1-tDCS; $\Delta_{dep} = 0.62$, CI = [0.50, 0.80] for rFPC-tDCS and sham). **(C)** Same as (A) but for the degree of temporal variability, cvIKlacross. The full factorial analysis demonstrated significant main effects of Block and Stimulation. **(D)** The reduction from block 1 to 3 in cvIKlacross was significantly more pronounced in rFPC-tDCS than IM1-tDCS ($P \leq P_{FDR} = 0.0178$, $\Delta_{dep} = 0.68$, CI = [0.50, 0.79]), and also in rFPC-tDCS relative to sham ($P \leq P_{FDR} = 0.0178$, $\Delta_{dep} = 0.72$, CI = [0.61, 0.83]). **(E-F)** Same as (C-D) but for keystroke velocity. No significant main effects or interactions were found when assessing cvKvelacross. Neither were there differential effects of stimulation on the change in cvKvelacross from block 1 to 3. Small black dots represent individual participant data. Colored dots display mean values with error bars denoting \pm SEM.

hoc analyses using 2×3 factorial analyses separately for each pair of stimulation protocols demonstrated a significant main effect of factor Block for pairs {rFPC, sham} and {rFPC, IM1} ($P \leq P_{FDR} = 0.0178$: sham versus rFPC and IM1 versus rFPC). For stimulation pair {IM1, sham} we did not find a significant main effect of factor Block ($P > 0.05$). In addition, there was a significant main effect of Stimulation when considering both active stimulation conditions IM1 and rFPC ($P \leq P_{FDR} = 0.0178$). No other significant main effects for Stimulation or interaction effects were found in the 2×3 *post hoc* analyses.

A separate one-way factorial analysis on the difference between offline and online blocks in cvIKIacross showed a significant effect of Stimulation ($P = 0.0246$, Figure 3D). *Post hoc* analyses on this difference measure revealed that following rFPC-tDCS the drop in temporal variability from online to offline learning blocks was more pronounced than following sham ($P \leq P_{FDR} = 0.0112$, $\Delta_{dep} = 0.72$, CI = [0.61, 0.83]) and also relative to IM1 ($P \leq P_{FDR} = 0.0112$, $\Delta_{dep} = 0.68$, CI = [0.50, 0.79]). When comparing IM1-tDCS to sham, however, the reduction in motor variability did not differ statistically ($P > 0.05$).

Finally, control analyses carried out on non-task-related variables, revealed no significant main or interaction effects on the variability in keystroke velocity (cvKvelacross, Figure 3E, full 3×3 factorial analysis). Thus, variability in the non-task related dimension, Kvel, was not significantly modulated by stimulation or learning block (see also Figure 3F). Additional details on general performance variables are presented in *Supplementary Results*.

Modelling Results

We investigated how individuals adapted their task-related behaviour as a function of the expectation of reward using a hierarchical Bayesian model, the Hierarchical Gaussian Filter for continuous inputs (HGF, [33, 34]). The HGF was adapted to model participants' beliefs about the reward on the current trial k , x_1^k , and about its rate of change, termed environmental volatility, x_2^k . Volatility emerged from the multiplicity of performance-to-score mappings (Figure S2; section *Reward function*), with timing patterns close-by in the movement space leading to different rewards. Volatility was not experimentally manipulated as in previous decision-making studies, where the reward mapping changed every block [35, 34, 36].

In the HGF, the update equations for the mean of the posterior distribution of beliefs on reward and log-volatility depend on the corresponding prediction errors (PE) weighted by precision (pwPE; precision being the inverse variance or uncertainty of the posterior distribution *Materials and Methods*). The perceptual HGF model was complemented with a response model, which defines the mapping from the trajectories of perceptual beliefs onto the observed responses in each participant. We were interested in assessing how belief trajectories or related

computational quantities (e.g. pwPEs, termed $\epsilon_i, i = 1, 2$) influenced subsequent behavioural changes, such as trial-to-trial timing variability or average tempo. Among different alternative response models, random effects Bayesian model selection provided stronger evidence for the response model that explained changes in timing *exploration* (logarithm of *unsigned changes* in cvIKI: $\log(|\Delta \text{cvIKI}|)$) as a linear function of pwPEs updating estimates on reward and log-volatility on the preceeding trial (see *Materials and Methods*, which include a sketch of the modelling approach). Simulations of this model revealed that agents observing a broader range of outcomes or introducing higher trial-to-trial exploration (larger absolute behavioural changes) have greater expectation on volatility, as both contributed towards an increased rate of change in the expectation on reward (Figure S4).

In the winning model, all β coefficients of the multiple linear regression response model for $\log(|\Delta \text{cvIKI}|)$ were significantly different from zero in each stimulation condition ($P \leq P_{FDR} = 0.001$; Figure 4A-B). On average β_1 was negative. This outcome indicated that larger lower-level pwPEs updating reward estimates on the previous trial (ϵ_1 ; increasing the expectation of reward) promoted an attenuation in motor exploration (reduced trial-to-trial unsigned changes in cvIKI). That is, increases in the expectation of reward were followed by exploitative behaviour in the relevant variable (Figure 4C-D), as expected [37, 38]. A non-parametric one-way factorial analysis demonstrated a significant effect of Stimulation on the β_1 coefficients ($P = 0.004$). *Post hoc* analyses further revealed that IM1-tDCS decreased the sensitivity of this association relative to sham and also when compared to rFPC-tDCS (reduced “negative” slope, larger β_1 values in IM1-tDCS than sham, $P \leq P_{FDR} = 0.003, \Delta_{dep} = 0.68, \text{CI} = [0.57, 0.83]$; similar outcome for the IM1-tDCS and rFPC-tDCS comparison: $P \leq P_{FDR} = 0.003, \Delta_{dep} = 0.62, \text{CI} = [0.55, 0.80]$; Figure 4A).

The effect of higher-level pwPEs updating log-volatility (ϵ_2) on trial-to-trial task-related exploration was also dissociated between stimulation conditions (significant effect of Stimulation in one-way factorial analysis, $P = 0.012$; Figure 4B). *Post hoc* analyses additionally demonstrated a dissociation between rFPC-tDCS and sham stimulation in this parameter, as β_2 coefficients were positive and significantly larger for rFPC relative to sham stimulation ($P \leq P_{FDR} = 0.003; \Delta_{dep} = 0.68, \text{CI} = [0.55, 0.83]$). Accordingly, a larger pwPE updating log-volatility on the previous trial—related to an increase in the expectation of volatility—was followed by more pronounced exploration for rFPC-tDCS than sham. Conversely, negative ϵ_2 promoted greater exploitation in rFPC-tDCS. A similar dissociation was found when comparing both active stimulation conditions, rFPC-tDCS and IM1-tDCS, due to significantly larger β_2 coefficients in rFPC-tDCS than in IM1-tDCS (positive coefficients; $P \leq P_{FDR} = 0.003; \Delta_{dep} = 0.71, \text{CI} = [0.58, 0.89]$). These results indicate that for rFPC-tDCS the sensitivity (slope) of this

association was greater than for sham or lM1. Thus, the winning response model identified different behavioural strategies in response to updates in reward and volatility as a function of the tDCS stimulation condition.

Next, we examined the effects of Stimulation and Block on the perceptual component of the Bayesian model by submitting the relevant model parameters to a full 3×3 factorial analysis (see also *Supplementary Results*). Concerning the mean trajectories of perceptual beliefs, we found a significant main effect of factor Block on the expectation of reward (μ_1 , $P = 0.002$) and volatility (μ_2 , $P = 0.0012$; Figure 4E-F). Across blocks, the expectation of reward increased, whereas the mean log-volatility estimate decreased, as participants gradually shifted to exploitation and approached the hidden performance goal (Figure 3). There was no main Stimulation or interaction effect for μ_1 ($P > 0.05$), which converges with the results from the model-free analysis of the scores. Neither was there an interaction or Stimulation effect for μ_2 ($P > 0.05$). Combining these findings with the results of the β_1 coefficients in the response model (Figure 4A), it follows that, as the scores (and μ_1) increased, rFPC-tDCS and sham promoted more exploitative behaviour than lM1-tDCS. In parallel, as μ_2 decreased (negative ϵ_2), rFPC-tDCS was more sensitive to changes in this estimate than lM1-tDCS and sham; rFPC-tDCS was therefore associated with a more pronounced tendency to exploit the inferred solution (smaller exploration of cvIKI).

Electric field distribution of tDCS

To control for the confound that anodal tDCS likely increases cortical excitability but could also lead to the opposite polarity of the effect [39], we complemented the main analysis with a simulation of the electric field induced by tDCS in each participant using SimNIBS ([40, 41]; see *Materials and Methods: tDCS*). This analysis focused on the focality and magnitude of the neuromodulatory effects induced by the active tDCS protocols. The results revealed that the focus of the induced electric field was within the targeted regions and had a similar magnitude in both structures (Figure 5). The peak values of the vector norm of the electric field (normE) did not differ between active stimulation conditions (99.9% percentile: mean and SEM for lM1-tDCS = 0.132 [0.006] V/m; for rFPC-tDCS = 0.134 [0.010] V/m; permutation test, $P > 0.05$). In addition, the volume corresponding with the 99.9% percentile of the field strength was not significantly different between active tDCS conditions (focality: $1.22 [0.11] \times 10^4$ mm 3 for lM1-tDCS; $1.14 [0.08] \times 10^4$ mm 3 for rFPC-tDCS; $P > 0.05$). Notwithstanding the similarity in peak and focality of the simulated normE values for lM1 and rFPC-tDCS, in both cases the electric field spread to neighboring areas beyond the target coordinate. Under lM1-tDCS, the induced electric field was maximum in the left M1 (area 4 of the human connectome project

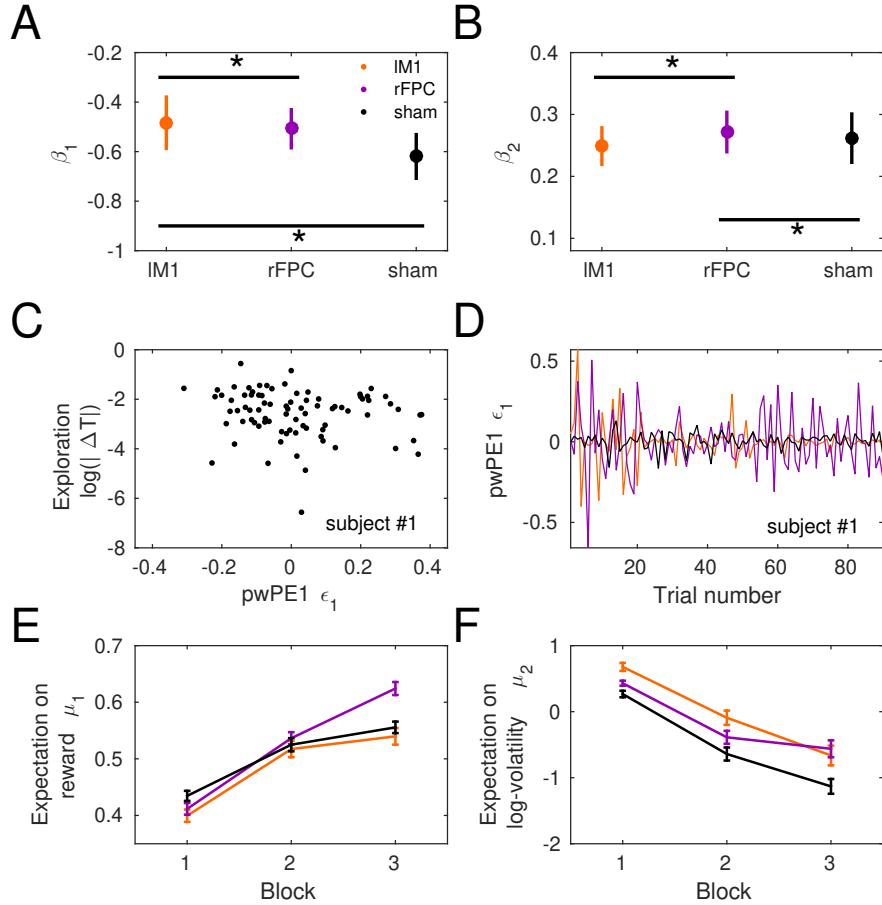


Figure 4: Computational modelling analysis. Data shown as mean and \pm SEM. **(A-B)**. β coefficients of the response model that explain the behavioural changes in trial k as a linear function of the precision-weighted prediction errors (pwPE, termed ϵ) in reward (ϵ_1) and volatility (ϵ_2) on the previous trial, $k - 1$. The performance measure, $\log(|\Delta cvIKI|)$, was the unsigned *change* from trial $k - 1$ to k in the degree of timing variation across keystroke positions. **(A)**. There was a significant main effect of factor Stimulation on β_1 (one-way non-parametric factorial analysis, $P = 0.004$). **(B)**. Coefficient β_2 was also modulated significantly with the factor Stimulation (one-way non-parametric factorial analysis, $P = 0.012$). Differences between pairs of stimulation conditions in β coefficients, as found in *post-hoc* analyses, are denoted by the horizontal black line and the asterisk ($P \leq P_{FDR}$, see main text). **(C)**. Illustration of the association between trialwise pwPE on reward and the subsequent task-related exploration, $\log(|\Delta cvIKI|)$, in one subject during sham. **(D)**. Illustration of the trajectories of pwPE on reward across all stimulation conditions in the same subject. **(E)**. Block-wise average of the expectation on reward, μ_1 . Significant main effect of Block ($P = 0.002$). **(F)** Mean expectation on environmental log-volatility, μ_2 , across blocks. This parameter was significantly modulated by the Block factor ($P = 0.0012$).

multi-modal parcellation, HCP-MMP1; [42]), followed by the premotor cortex (6), prefrontal areas (8Av and 8C) and somatosensory cortex (3). Under rFPC-tDCS, the peak of the electric field corresponded with the rFPC (areas 10p and 10pp), followed by regions in the medial prefrontal cortex (mPFC; 9) and orbitofrontal cortex (OFC; 11). Lastly, the variability in the electric field strength (standard deviation) did not differ between tDCS targets ($P > 0.05$, Figure S6), supporting the comparable effects of both stimulation protocols in our sample.

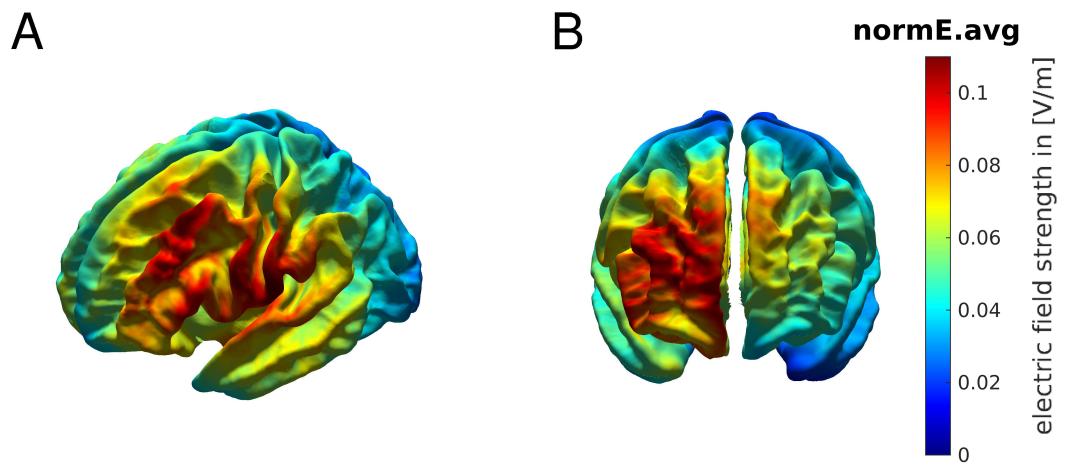


Figure 5: Electric field distribution for anodal IM1-tDCS (A) and rFPC-tDCS (B). Norm of the electric field strength (normE) derived from FEM calculations using SimNIBS, and averaged across participants using the fsaverage-transformed surface.

Discussion

In this study, we identify a potential role of rFPC in reward-based motor learning by using a motor task that requires a shift from exploration to exploitation, a Bayesian computational model of the behaviour [34], and simulations of the electric field induced by tDCS. The results indicated that rFPC-tDCS relative to sham and IM1-tDCS accelerated the increment in scores from online to offline learning blocks primarily through regulation of task-related motor variability. Trial-to-trial analyses using computational modelling further demonstrated that across all stimulation protocols increased expectation of reward led to subsequent exploitation, as expected [4, 38]. Because the sensitivity (slope) of this association was greater for rFPC-tDCS and sham relative to IM1-tDCS, these results suggest that behavioural changes following pwPEs updating reward estimates were enhanced for rFPC-tDCS and sham. Frontopolar stimulation was however dissociated from IM1 and sham stimulation with regards to the effects of trial-

to-trial volatility estimates on exploration. While LM1-tDCS and sham were less sensitive to changes in environmental volatility, rFPC-tDCS promoted greater exploitation of the rewarding timing pattern as the expectation on log-volatility progressively decreased—with increasing expectation on reward. These findings suggest that rFPC contributes to reward-based motor learning by promoting a shift from exploration towards successful exploitation. This shift is mediated by an enhanced sensitivity to environmental volatility, which is associated with changes in the reward structure over time. Our results extend findings in the area of decision-making [13, 12, 11, 10] to that of motor skill learning: Brain regions previously linked to decision-making in the cognitive or perceptual domain, such as the FPC, could be relevant in the motor domain. The findings also complement recent tDCS work associating the dorsolateral PFC to motor decision-making [43].

Neuromodulation of the rFPC via anodal tDCS reduced motor variability during learning blocks. This effect was specific to the learning phase, as rFPC-tDCS did not modulate baseline motor noise, similarly to LM1-tDCS and sham. Recent work has demonstrated that the inhibition of the rFPC via TMS during binary choices decreased directed exploration driven by information seeking [10]. In line with that result, cathodal versus anodal rFPC-tDCS have been shown to decrease or increase exploration, respectively, in a three-armed bandit task [11]. Thus, a surprising outcome was that, compared to LM1-tDCS and sham, anodal rFPC-tDCS did not increase motor variability during learning but instead reduced it towards the third block. This apparent discrepancy can be better understood by comparing our motor task with the decision-making tasks used in FPC studies. Binary or multi-armed bandit tasks as used in decision-making studies generally use a probabilistic reward function that changes over time [10, 11, 35, 34]. In this context, an optimal policy should continuously balance the exploration-exploitation tradeoff [44], while FPC might drive exploration as needed for the task demands [10, 11]. On the other hand, motor tasks that use continuous reward signals to guide learning typically maintain the same the reward structure over time [5, 4, 31]. Here, movement variability is initially high but progressively decreases as participants approach the hidden solution. In this scenario, our study demonstrated that rFPC-tDCS facilitates the decline in motor variability and an increase in scores. This is consistent with the evidence from human and non-human primates, supporting that FPC might have evolved to monitor multiple competing choices and direct the exploratory tendency towards the most rewarding one [12]. This function of FPC thus makes it particularly suitable to track the reward associated with (multiple) continuous movement parameters, such as timing or force in our study.

To understand why FPC stimulation led to the largest increase in scores from online to offline blocks, the computational results should be considered. Across blocks, there was an

attenuation of the expectation on log-volatility in all stimulation conditions, which suggests that the reward tendency estimate became more stable over time. This result also indicates that over time, trial-to-trial changes in scores (and expectation on reward) were associated with increasingly more negative update steps ($\text{pwPE } \epsilon_2$) on volatility. Because under rFPC-tDCS the positive slope of the association between exploration and ϵ_2 was greater than for IM1-tDCS and sham—indicating higher sensitivity in this association (larger β_2)—smaller ϵ_2 values over time would lead to more exploitative behaviour in rFPC-tDCS. This finding thus dissociates the effects of IM1-tDCS and sham on reward-dependent motor learning from those of rFPC-tDCS. It hints at an important condition for successful reward-based motor learning: increased sensitivity of task-related exploration to changes in volatility.

How was the mapping between movement and reward acquired in the different stimulation protocols? In the continuous movement space over which our task was defined, timing was the task-related dimension that needed to be mapped to reward. Thus, one possibility is that rFPC-tDCS might have facilitated the acquisition of this complex mapping, increasing the achieved scores and the expectation on reward, and reducing volatility estimates. This remains speculative at this moment as a limitation of this study is that we did not track neural dynamics during task performance. Therefore, we could not assess how rFPC-tDCS modulated the emergence of a neural representation of the mapping between movement parameters and reward. On the behavioural level, however, the results converge in showing that rFPC-tDCS led in parallel to increased task-related exploitation and increased scores. Accordingly, the timing pattern that was exploited under frontopolar stimulation was indeed closer to the hidden target than the solutions inferred under sham and IM1-tDCS. The results are consistent with previous findings indicating that FPC infers the absolute reliability of several alternative goals [45, 46]. Moreover, FPC manages competing goals by keeping track of alternative choices [13, 12]. Notably, our results expand previous findings by revealing that during motor learning, which is generally solved in a continuous movement space [7], learning the mapping between different movement configurations and their associated reward partially relies on FPC, although an additional involvement of the dorsolateral PFC is also likely ([43]).

Stimulation over the contralateral M1 did not modulate motor variability across blocks (cvIKIacross) when compared to sham stimulation. This lack of significant effects should be interpreted with caution as our statistical approach does not allow us to make inferences on null results. However, in the modelling analysis, motor cortex stimulation was associated with a reduced sensitivity to changes in the expectation of reward. Although we had hypothesized that rFPC would accelerate reward-based learning more when compared to IM1-tDCS, which we confirmed, we also predicted that IM1-tDCS would increase exploration and improve

learning from reward signals relative to sham. This prediction was based on the existing TMS and tDCS studies demonstrating the involvement of M1 in reward-based motor learning [47, 26] and decision-making [24, 25]. Furthermore, M1 contributes to reward-based motor learning via neurophysiological plasticity changes, such as long-term potentiation [26]. This would be consistent with the maximum electric field intensity over lM1 in our study indicating excitatory effects, which relate to a decrease in local GABAergic activity and enhanced long-term potentiation-like activity [48, 16, 49]. It is noteworthy that most previous studies linking M1 to reward-based motor learning (see above), reward-guided motor processing [19, 20] and valued-based decision making [21, 22] used TMS protocols. Accordingly, the focality of TMS may be necessary to overcome the large inter-individual variability affecting tDCS studies [39]. Arguing similarly, [23] used smaller electrodes anterior and posterior to M1 to demonstrate a benefit of reward signals and M1-tDCS on motor retention. Thus, resolving the issue of the dissociable role of M1 and FPC in regulating motor variability during reward-based motor learning will require follow-up TMS studies and, additionally, a comparison between initial learning and motor retention [50]. Lastly, because the largest effects of rFPC relative to lM1-tDCS emerged when contrasting the achieved scores in offline versus online blocks, future work should clarify whether the benefits of rFPC over lM1-tDCS during reward-based motor learning are specific to offline stimulation, as M1-tDCS effects on motor learning may be limited to online stimulation [17, 28].

Of note, a limitation of any tDCS study is that the diffuse spatial effect of tDCS does not allow us to determine whether modulation of the target area alone is responsible for the observed behavioural effect [39]. Even with the higher spatial resolution of TMS, it has been argued that any effect of FPC stimulation on behaviour could be mediated by other brain regions coupled with FPC during behavioural exploration [10], such as the inferior parietal cortex or ventral premotor cortex [45] or by regions engaged in other aspects of goal-directed behaviour, such as the ventromedial or dorsolateral PFC and OFC [15, 12, 43].

To mitigate that limitation, we modelled the electric field in the individual anatomy and assessed the strength and focality of the induced electric field for each tDCS target. We found an enhanced focal activation with maxima in the targeted areas for both active stimulation protocols. Recent work demonstrated an association between enhanced electric field strength in SimNIBS due to anodal tDCS and excitatory effects [48], however their results are limited to M1. Accordingly, neurophysiological implications for rFPC stimulation cannot be drawn out at this point. Future studies combining electroencephalography and functional MRI should assess the network of interactions between FPC, other regions in the PFC, and cortical motor regions to determine the precise mechanism underlying the rFPC-tDCS effect on reward-based motor

learning reported here.

Materials and Methods

Participants

Nineteen right-handed participants (10 females, mean = 27.7 yrs, std = 3.3 yrs, range 21-33) with no history of neurological disease or hearing impairment and with no musical training outside of the requirements of the general music curriculum in school were recruited. Laterality quotient was assessed by the Oldfield handedness inventory ([51]; mean = 90, standard error of the mean or SEM = 3.2; values available in 17/19 participants). The sample size is similar to that found in tDCS studies focusing on motor learning [52, 23] and was based on our previous estimation of the minimum sample size required to detect effects of different experimental manipulations in this paradigm (e.g. reward or affective manipulations, [14]) with a statistical power of 0.95. The study protocols were approved by the local ethics committee of the University of Leipzig (277-14-25082014) and agrees with the provisions of the Helsinki Declaration [53]. Participants gave written informed consent before the beginning of the first experimental session. To incentivize participants during completion of the reward-based learning phase of the motor task, they were informed about a €50 voucher for online purchases that would be awarded to the participant scoring the highest average score across the three sessions.

Experimental Design and Procedure

A sham-controlled, double-blinded, cross-over design was implemented. The study was comprised of three sessions with a 7-days interval between them (same time of day) to reduce potential carry-over effects. Selection of target coordinates for each tDCS protocol was guided by individual T1-MRI. minutes after tDCS cessation). However, to account for a possible differential effect of tDCS protocols on the offline (3) and online (1) blocks, we performed a specific additional offline minus online contrast in the statistical analysis, which is described below.

In the learning phase, we implemented a mapping rule between movement and reward governed by uncertainty, as in our task different timing patterns could receive the same reward, whereas similar timing patterns would obtain different rewards (Figure S2). This choice was based on previous research suggesting that the acquisition of complex motor skills in daily life often involves an uncertain or variable mapping between actions and outcomes [54]. Moreover, higher reward uncertainty can be beneficial for motor retention [54]. Accordingly, in our study, higher overall exploration would lead to the perception of higher uncertainty in the environment

(or greater estimated change in the reward tendency, termed environmental volatility, [34]). In this scenario, the optimal strategy to maximize the mean total reward involves two phases: (i) an initial increase in exploration to learn the mapping between reward and movement parameters, followed by (ii) a swift switch to the exploitation of the performance inferred as most rewarding. The quantitative analysis of exploration is shown below.

Prior to receiving tDCS and completing the motor task, participants had to familiarize themselves with the series of tones they had to produce during each phase of the task. The stimulus material for the baseline phase consisted of a series of eight consecutive white piano keys with four fingers (four notes upwards + same four notes downwards; one finger per key with fixed finger-to-key mapping). For the learning phase, the stimulus material was comprised of a sequence of eight notes, which was a combination of the four neighboring white keys they had to press during the baseline phase (Figure 2). Three different types of sequences were used for each stimulation session, with a pseudo-randomized counterbalanced order across participants. Each sequence was defined over a similar range of semitones but the range had a different spatial location on the keyboard (i.e. towards higher or lower pitch values, Figure 2). Participants were explicitly taught the order of the notes for the baseline and reward-based learning phase by one of the experimenters, who played the notes for the participants using an isochronous timing. Participants had to repeat the sequence of notes after the demonstration by the experimenter using a self-paced tempo. Because the stimulus materials for both phases were short and the order of the notes easy to play, all participants demonstrated an error-free performance after just a few repetitions (baseline materials: 2 repetitions on average, range 1-4; learning materials: 4 repetitions on average, range 2-6).

Once under tDCS, the baseline phase required participants to press the corresponding series of white keys regularly at a self-paced tempo. This phase allowed us to assess fine motor control during regular performance as a proxy for baseline motor noise.

Next, during the reward-based learning blocks, participants had to play the corresponding sequence of notes (Figure 2). Crucially, however, at this point they were instructed that the timing of the performance target was not isochronous and thus their goal was to use trial-based feedback scores to discover and approach that hidden target. Participants were not aware that different timing solutions could receive the same reward.

Reward function

The performance measure that was rewarded during learning blocks was the Euclidean norm of the vector corresponding to the pattern of temporal differences between adjacent inter-keystroke-intervals (IKI, in s) for a trial-specific performance (as in [14]). To approach the

hidden target performance, participants had to deviate from an isochronous performance and find out the right combination of successive IKIs.

Here we denote the vector norm by $\|\Delta\mathbf{z}\|$, with $\Delta\mathbf{z}$ being the vector of differences, $\Delta\mathbf{z} = (z_2 - z_1, z_3 - z_2, \dots, z_n - z_{n-1})$, and z_i representing the IKI at each keystroke ($i = 1, 2, \dots, n$). Notably, IKI values represent the difference between the onset of consecutive keystrokes, and therefore $\Delta\mathbf{z}$ indicates a vector of differences of differences (put simply: differences of intervals). The target value of the performance measure for each sequence was a vector norm of 1.9596 (e.g. one of the maximally rewarded performances leading to this vector norm of IKI-differences would consist of IKI values: [0.2, 1, 0.2, 1, 0.2, 1, 0.2] s; that is a combination of short and long intervals). The score was computed in each trial using a measure of proximity between the target vector norm $\|\Delta\mathbf{z}^t\|$ and the norm of the performed pattern of IKI differences $\|\Delta\mathbf{z}^p\|$, using the following expression:

$$\text{score} = 100 \exp(-|\|\Delta\mathbf{z}^t\| - \|\Delta\mathbf{z}^p\||). \quad (1)$$

As mentioned above, different combinations of IKIs could lead to the same IKI differences and thus same Euclidean norm. Thus, timing patterns that were far away in the movement space could receive the same reward, whereas timing patterns close-by in the movement space would obtain different rewards (Figure S2). Accordingly, the mapping rule between movement and reward was governed by uncertainty, and higher overall exploration could be associated with the perception that the environment and reward structure was more unstable. This was explicitly assessed in the mathematical model of the behaviour described below.

tDCS

During the experiment, tDCS was applied via saline-soaked sponge electrodes to the individual target coordinate using a battery-driven DC-stimulator (NeurConn, Ilmenau, Germany). tDCS can transiently modulate cortical excitability via application of direct currents, as shown in combined TMS-tDCS studies [16, 55], and further supported by recent simulation studies [48]. Anodal tDCS has been shown to increase cortical excitability, however additional evidence indicates that it could lead to the opposite polarity of the effect, reducing cortical excitability [39]. For instance, extending the stimulation duration beyond 26 min has been shown to result in inhibitory rather than excitatory effects after anodal tDCS [56]. To control for this confound, we complemented the main analysis with a simulation of the electric field induced by tDCS in each participant using SimNIBS (see below). Anodal or sham tDCS was applied to the right FPC or left (contralateral) M1 regions. The target coordinate for rFPC-tDCS was selected from previous tDCS and fMRI work investigating the role of rFPC on exploration (Montreal

Neurological Institute or MNI peak: $x = 27, y = 57, z = 6$; [15, 11]). For lM1-tDCS, we used a target coordinate in the hand area of the left primary motor cortex (MNI peak: $x = -37, y = -21, z = 58$), based on [57]. The target coordinates were transformed to the individual native space using a T1-weighted high-resolution magnetic resonance image from each participant. Specifically, the MNI coordinates were converted into participants's native MNI space using the reverse native-to-MNI transformation from Statistical Parametric Mapping (SMP, version SPM12). The point on the scalp corresponding with each of our targeted brain areas was marked to place the active electrode ($5 \times 5 \text{ cm}^2$). The reference (cathode, $10 \times 10 \text{ cm}^2$) electrode for rFPC-tDCS was placed at the vertex [11], whereas it was located over the frontal orbit for lM1-tDCS [16]. Flexible elastic straps were used to fixate the electrodes on the head. A three-dimensional (3D) neuronavigation device (Brainsight Version 2; Rogue Research, Montreal, Canada) was used to guide positioning of active electrodes.

We stimulated with a weak direct current of 1 mA in all conditions for 20 minutes resulting in a current density of 0.04 mA/cm^2 under the target electrode and 0.01 mA/cm^2 under the reference electrode. In general, the modulatory effect of tDCS on brain excitability is more pronounced after several minutes and can subsequently outlast the stimulation duration for up to 1.5 h [16, 55]. To account for the potential delay in the effect of tDCS, we instructed participants to wait for 3 minutes before we initiated the motor task. At the start of the active tDCS stimulation the current was ramped up for 30 s to minimize the tingling sensation on the scalp, which generally fades over seconds [17], and was also ramped down for 30 s at the end. During sham tDCS, the current was ramped-up for 30 s, held constant at 1 mA for 30 s and ramped-down for 30 s. This procedure aimed to induce a similar initial tingling sensation in active and sham protocols, yet without modulation of cortical excitability for sham tDCS [58].

Before and after each tDCS session, participants rated on a 1-10 scale their fatigue, attention and discomfort levels, as well as the sensation with tDCS (post-tDCS, scale 1-5). No significant differences between active and sham tDCS sessions were found in the reported levels of fatigue, discomfort or attention levels ($P > 0.05$, post minus pre changes; paired permutation test). The sensation was not different between rFPC and sham tDCS, either ($P > 0.05$). However, lM1-tDCS induced a higher sensation than sham ($P = 0.0034$, non-parametric effect size $\Delta_{dep} = 0.87$, CI = [0.57, 0.88]). Details on statistical methods are provided below).

Acquisition and analysis of behavioural data

Performance information was saved as MIDI (Musical Instrument Digital Interface) data, which provided the time onsets of keystrokes relative to the previous event (inter-keystroke interval,

IKI, s), MIDI note number that corresponds with the pitch, and MIDI velocity (related to loudness). Behavioral data are available in the Open Science Framework Data Repository: <https://osf.io/zuab8/>

This cvacross measure was computed separately for each keystroke position (cv across trials within the block) and then averaged across keystroke positions. During the baseline phase participants had to accurately reproduce the same action (regular timing and keystroke velocity). In this context, any residual variability can be regarded to reflect motor noise [29, 4], which here was measured using cvIKIacross and cvKvelacross in this initial phase. During learning blocks, the level of task-related motor variability, cvIKIacross, was considered to primarily reflect intentional exploration of this parameter but also some degree of unintentional motor noise —similarly to other studies [5, 4, 14]. In contrast with timing performance, we assumed that participants would not intentionally modify keystroke velocity during learning. Thus, changes in cvKvelacross across blocks —if present—would be an indication of changes in unintentional motor noise with learning. The achieved scores and other general performance variables, such as the block-wise mean tempo (mIKIacross), mean Kvel (mKvelacross) and rate of wrong notes (error rate) were also evaluated.

During the baseline phase we assessed statistically the effects of stimulation conditions on the relevant behavioural variables (see *Results*), excluding the scores. During the learning blocks, statistical analysis focused on the investigation of the effect of stimulation and learning block on all behavioural dependent variables. In addition, we were specifically interested in assessing the influence of tDCS protocols on the change in task-related motor variability (cvIKIacross) and scores from online to offline (after the cessation of tDCS) learning blocks. Thus, additional dependent variables were the difference between blocks 3 and 1 in cvIKIacross and, separately, the scores. Details on statistical testing are provided in section *Statistical Analysis*. When providing mean values on behavioural variables, we also indicate the standard error of the mean or SEM.

Bayesian model of behaviour

In the HGF model for continuous inputs we implemented ([33, 34]), beliefs on x_1 and x_2 were Gaussian distributions and thus fully determined by the sufficient statistics μ_i ($i = 1, 2$, mean of the posterior distribution for x_i , corresponding to participants' expectation) and σ_i (variance of the distribution, representing uncertainty of the estimate). The belief trajectories about the external states x_1 and x_2 (mean, variance) were further used to estimate the most likely response corresponding with those beliefs. A schematic illustrating the model structure and outputs is

shown in Figure 6.

In the HGF, the update equations for the posterior mean at level i and for trial k , μ_i , depend on the prediction errors (PE) weighted by uncertainty σ_i (or its inverse, precision $\pi_i = 1/\sigma_i$) according to the following expression:

$$\Delta\mu_i^k = \mu_i^k - \mu_i^{k-1} \propto \frac{\pi_{i-1}^{k-1}}{\pi_i^k} \delta_{i-1}^k. \quad (2)$$

The change in expectations, $\Delta\mu_i^k$, is proportional to the prediction error of the level below, δ_{i-1}^k , representing the difference between the expectation μ_{i-1}^k and the prediction μ_{i-1}^{k-1} of the level below x_{i-1}^k . The prediction error is weighted by the ratio between the prediction of the precision on the level below, π_{i-1}^{k-1} , and the precision on the current level, π_i^k . The product of the precision weights ratio and the prediction error constitute the precision-weighted prediction error (pwPE), termed ϵ_i . Thus, ϵ_i regulates the update of expectations on trial k : $\Delta\mu_i^k = \epsilon_i$. The pwPE expressions for level 1 and 2 share the general form of Eq. 2, and detailed definitions can be found in [33, 34, 14]. Equation 2 illustrates that higher uncertainty in the current level (larger σ_i^k , smaller π_i^k in the denominator) leads to faster update of expectations; moreover, a smaller prediction of uncertainty on the level below (larger π_{i-1}^{k-1}) also increases the update of expectations. The intuition from this expression is that the more uncertain we are about the level we're trying to estimate (current level), the more we should update that level using new information (prediction errors). On the other hand, the less certain we are about the level below (less precise information), the less that new information should contribute to our update of beliefs.

Details on the free parameters of the HGF model to be estimated in each individual are presented in the Supplementary Materials online. In addition, Table S1 shows our choice of prior values on the HGF parameters that were used to generate belief trajectories.

The response model defines the mapping from the trajectories of perceptual beliefs onto the observed responses in each participant. We were interested in assessing how belief trajectories or related computational quantities influenced subsequent behavioural changes, such as trial-to-trial variability or exploration. Accordingly, we constructed different response models associated with different scenarios in which participants would link a specific performance measure to reward, such as the mean duration of key presses or the degree of timing variation across keystroke positions. Response variables that were bounded to 0-1 in their native space, such as $|\Delta cvIKI^k|$, were transformed into an unbounded variable using the logarithmic transformation. See details on the Supplementary Materials online.

For each performance measure, the corresponding response model explained that variable as a function of (a) the average of the belief estimates μ_1 , or μ_2 ; (b) the precision-weighted PE

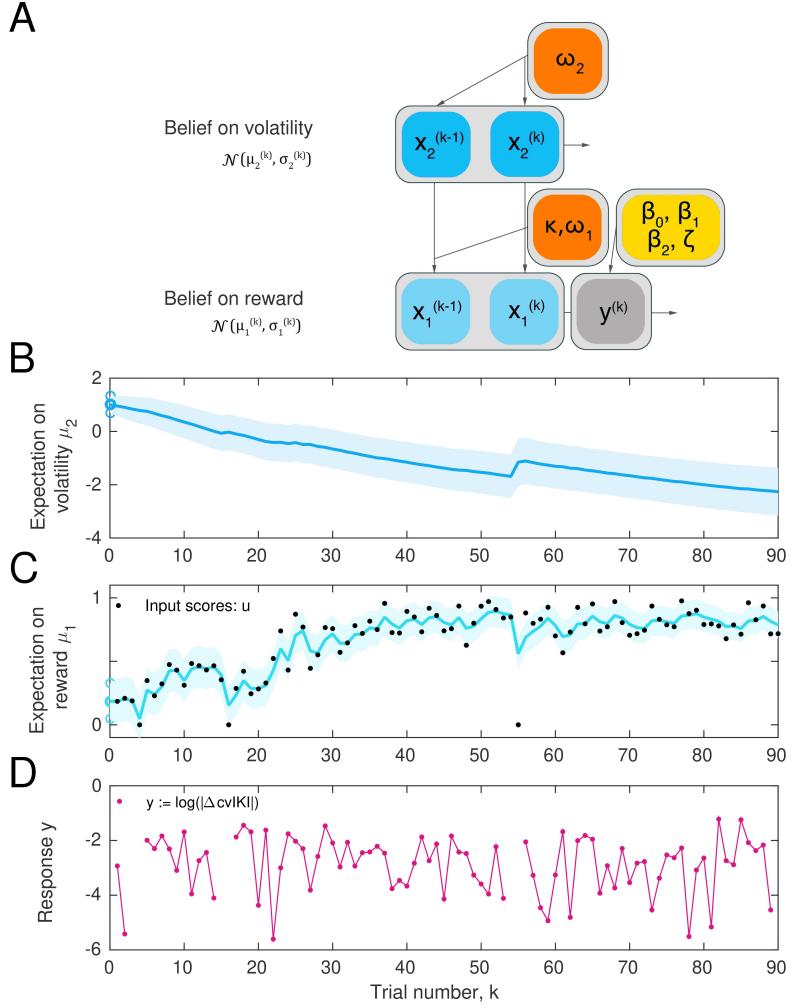


Figure 6: Computational model. Two-level Hierarchical Gaussian Filter for continuous inputs. **(A)** Schematic of the two-level HGF, which models how an agent infers a hidden state in the environment (a random variable), x_1 , as well as its rate of change over time (x_2 , environmental volatility). Beliefs (μ_1, μ_2) about those hierarchically-related hidden states (x_1, x_2) at trial k are updated with the input scores for that trial via prediction errors (PEs). The states x_1 and x_2 are continuous variables evolving as coupled Gaussian random walks, where the step size (variance) of the random walk depends on a set of parameters (shown in orange boxes). The lowest level is coupled to the level above through the variance. The response model generates the most probable response, y , according to the current beliefs about the lower state, x_1 , and is modulated by the response model parameters (yellow box). **(B)** Example of trial-by-trial beliefs about volatility, μ_2 (variance σ_2). **(C)** Belief on the first level, which represents an individual's expectation of reward, μ_1 (variance σ_1). Black dots represent the trialwise input scores (u). **(D)** Performance output (logarithm of the unsigned change in trialwise cvIKI as a proxy for exploration). Shaded areas denote the variance or estimation uncertainty on that level.

(pwPE) about reward, ϵ_1 , or volatility, ϵ_2 ; (c) HGF quantities from the lower level: μ_1, ϵ_1 ; (d) HGF quantities from the higher level: μ_2, ϵ_2 . This led to a total of 16 different models. The rationale for choosing pwPE at level 1 and 2 as predictors in some of the alternative response models was the relevance of PE weighted by uncertainty in current frameworks of Bayesian inference[59, 60]. Moreover, pwPEs determine the step size of the update in the expectation of beliefs (Eq. 2). That is, larger pwPEs about reward increase the expectation of reward, while larger pwPE about volatility increase the corresponding volatility estimate. See [14] for a similar use of the response models defined for this paradigm.

Each model was fitted with the 90 trialwise performance values of the corresponding response variable and with the input scores for each tDCS session. The log model-evidence (LME) was used to optimize the model fit [61]. Random Effects Bayesian Model Selection (BMS; code freely available from the MACS toolbox, [63]) was performed across all 16 models using the LME values. BMS provided stronger evidence for the response model that explained $\log(|\Delta \text{cvIKI}^k|)$ as a linear function of pwPE about reward, ϵ_1 , and volatility, ϵ_2 , on the preceeding trial $k - 1$:

$$\log(|\Delta \text{cvIKI}^k|) = \beta_0 + \beta_1 \epsilon_1^{k-1} + \beta_2 \epsilon_2^{k-1} + \zeta, \quad (3)$$

where ζ is a Gaussian noise variable. The response variable $\log(|\Delta \text{cvIKI}^k|)$ reflected unsigned changes (exploration) from trial $k - 1$ to trial k in cvIKI. In this winning model, the exceedance probability was 0.9888 and the model frequency was 72%.

The HGF with the winning response model provided a good fit to the behavioural data, as the examination of the residuals shows (Figure S5). There were no systematic differences in the model fits across tDCS conditions. The response model noise parameter ζ was not significantly modulated by the stimulation condition ($P > 0.05$; average value $\zeta = 1.3$ [0.08]).

The effect of stimulation on the learning process, as described by the computational model, was assessed by analysing the following dependent variables: (i) The β coefficients of the winning response model (β_1, β_2) that regulate how pwPE on reward and volatility modulate task-related behavioural adaptations, as well as the noise parameter ζ . This analysis thus allowed us to investigate how trial-to-trial update steps in the expectation of reward and volatility (related to changes in observed scores) modulated task-related motor exploration in the next trial. Perceptual HGF model variables, (ii) including the trialwise expectation of reward (μ_1) and volatility (μ_2); and (iii) parameters ω_1 and ω_2 , which modulate the step size of the update equations at each level and characterize the individual learning tendency (See *Supplementary Material* and *Supplementary Results*).

The main variables for the statistical analysis were the mean trajectories of perceptual beliefs (μ_1 , μ_2) and the parameters ω_1 and ω_2 . We found no significant effect of Stimulation on the perceptual model parameters, neither for ω_1 ($P > 0.05$, one-way factorial analysis) nor for ω_2 ($P > 0.05$). On average, ω_1 was -3.5 (0.45), and ω_2 was -4.1 (0.35).

SimNIBS

The electric field distribution induced by each tDCS condition was simulated in each participant with the freely available SimNIBS 2.1 software [40, 41]. SimNIBS integrates different tools, such as FreeSurfer, FMRIB's FSL, MeshFix, and Gmsh [64]. Using the headreco head modelling pipeline of SimNIBS, the electrically most relevant tissue structures (skin, skull, cerebrospinal fluid, gray matter, white matter, eyes, and air) were first segmented from the individual T1-weighted anatomical MRI. The segmentation image of the skin tissue was subsequently smoothed to remove any residual artifact. This was carried out independently from the SimNIBS pipeline with the freely available software MIPAV by applying a spatial Gaussian filter (2 mm in each xyz direction). The creation of the head model was then completed with headreco by generating a tetrahedral mesh as volume conductor model. Next, in the SimNIBS GUI, simulated electrodes were placed manually on the head mesh at their precise position and with the corresponding orientation. Stimulation intensities were selected for anodal and cathodal electrodes and the simulation based on the finite element method (FEM) was initiated. The vector norm of the electric field (normE) was extracted and chosen as dependent variable for subsequent group-level statistical analysis. These steps were repeated separately in each active tDCS condition and in each participant.

The individual normE distribution was transformed to the fsaverage space to create a group average of the mean normE values and their standard deviation (Figure S6). This was carried out with the MATLAB scripts provided in the SimNIBS package. To assess statistical differences between stimulation conditions in the peak values and focality of the induced electric field, we extracted the 99.9 percentile value of the normE distribution, as well as the volume in which the normE values reached the 99.9 strength percentile.

Statistical analysis

Statistical analysis was performed with the use of non-parametric permutation tests with 5000 permutations. During the baseline phase, non-parametric one-way factorial analyses with factor Stimulation (LM1, rFPC, sham) were carried out using synchronized rearrangements [65], which are based on permutations. During learning, full 3×3 factorial analyses with factor Block (1-3

levels during learning or 1 level during baseline) and Stimulation were implemented. Effects were considered significant if ($P \leq 0.05$). Following significant main tests or interactions, *post hoc* analyses during learning focused on 2×3 factorial analyses using pairs of tDCS protocols as factor Stimulation: {lM1, sham}, {rFPC, sham}, {lM1, rFPC}. Follow-up *post hoc* pair-wise comparisons between stimulation conditions or blocks were evaluated using pair-wise permutation tests for matched samples.

In all cases, we addressed the issue of multiple comparisons arising from the implementation of several *post hoc* analyses by controlling the false discovery rate (FDR) at level $q = 0.05$ with an adaptive two-stage linear step-up procedure [66]. Significant effects after FDR-control are reported as $P \leq P_{FDR}$, and providing the explicit adapted value of P_{FDR} .

In addition, separately from the main factorial analyses, we performed analyses of offline (block 3) minus online (block 1) differences in motor variability (cvIKIacross) and scores in the learning phase using one-way factorial analyses with factor Stimulation (lM1, rFPC, sham). Here, *post hoc* analyses of pair-wise contrasts between tDCS conditions also controlled the FDR at level $q = 0.05$.

Throughout the manuscript, non-parametric effect sizes and corresponding confidence intervals are provided along with pair-wise permutation tests. As measure of non-parametric effect size we used the probability of superiority for dependent samples Δ_{dep} , ranging 0-1 [32]. Confidence intervals (CI) for Δ_{dep} were estimated with bootstrap methods [67].

1 Acknowledgments

This research was supported by the International Engagement Fund of Goldsmiths University of London (MHR). MHR and VN were partially supported by the HSE Basic Research Program and the Russian Academic Excellence Project '5-100'. We thank Dennis Maudrich for helping with the experiments.

References

- [1] Wolpert DM, Diedrichsen J, Flanagan JR (2011) Principles of sensorimotor learning. *Nature Reviews Neuroscience* 12(12):739.
- [2] Dhawale AK, Smith MA, Ölveczky BP (2017) The role of variability in motor learning. *Annual review of neuroscience* 40:479–498.

- [3] Therrien AS, Wolpert DM, Bastian AJ (2016) Effective reinforcement learning following cerebellar damage requires a balance between exploration and motor noise. *Brain* 139(1):101–114.
- [4] Chen X, Mohr K, Galea JM (2017) Predicting explorative motor learning using decision-making and motor noise. *PLoS computational biology* 13(4):e1005503.
- [5] Wu HG, Miyamoto YR, Castro LNG, Ölveczky BP, Smith MA (2014) Temporal structure of motor variability is dynamically regulated and predicts motor learning ability. *Nature neuroscience* 17(2):312.
- [6] Santos FJ, Oliveira RF, Jin X, Costa RM (2015) Corticostriatal dynamics encode the refinement of specific behavioral variability during skill learning. *Elife* 4:e09423.
- [7] Wolpert DM, Landy MS (2012) Motor control is decision-making. *Current opinion in neurobiology* 22(6):996–1003.
- [8] Parvin DE, McDougle SD, Taylor JA, Ivry RB (2018) Credit assignment in a motor decision making task is influenced by agency and not sensory prediction errors. *Journal of Neuroscience* 38(19):4521–4530.
- [9] Ota K, Tanae M, Ishii K, Takiyama K (2020) Optimizing motor decision-making through competition with opponents. *Scientific Reports* 10(1):1–14.
- [10] Zajkowski WK, Kossut M, Wilson RC (2017) A causal role for right frontopolar cortex in directed, but not random, exploration. *Elife* 6:e27430.
- [11] Beharelle AR, Polanía R, Hare TA, Ruff CC (2015) Transcranial stimulation over frontopolar cortex elucidates the choice attributes and neural mechanisms used to resolve exploration–exploitation trade-offs. *Journal of Neuroscience* 35(43):14544–14556.
- [12] Mansouri FA, Koechlin E, Rosa MG, Buckley MJ (2017) Managing competing goals. a key role for the frontopolar cortex. *Nature Reviews Neuroscience* 18(11):645.
- [13] Boschin EA, Piekema C, Buckley MJ (2015) Essential functions of primate frontopolar cortex in cognition. *Proceedings of the National Academy of Sciences* 112(9):E1020–E1027.
- [14] Sporn S, Hein T, Ruiz MH (2020) Alterations in the amplitude and burst rate of beta oscillations impair reward-dependent motor learning in anxiety. *Elife* 9:e50654.
- [15] Daw ND, O'doherty JP, Dayan P, Seymour B, Dolan RJ (2006) Cortical substrates for exploratory decisions in humans. *Nature* 441(7095):876.

- [16] Nitsche MA, Paulus W (2000) Excitability changes induced in the human motor cortex by weak transcranial direct current stimulation. *The Journal of physiology* 527(3):633–639.
- [17] Nitsche MA, et al. (2003) Facilitation of implicit motor learning by weak transcranial direct current stimulation of the primary motor cortex in the human. *Journal of cognitive neuroscience* 15(4):619–626.
- [18] Mandelblat-Cerf Y, Paz R, Vaadia E (2009) Trial-to-trial variability of single cells in motor cortices is dynamically modified during visuomotor adaptation. *Journal of Neuroscience* 29(48):15053–15062.
- [19] Klein PA, Olivier E, Duque J (2012) Influence of reward on corticospinal excitability during movement preparation. *Journal of Neuroscience* 32(50):18124–18136.
- [20] Galarraga JK, Celnik P, Chib VS (2019) Motor cortex excitability reflects the subjective value of reward and mediates its effects on incentive-motivated performance. *Journal of Neuroscience* 39(7):1236–1248.
- [21] Thabit MN, et al. (2011) Momentary reward induce changes in excitability of primary motor cortex. *Clinical Neurophysiology* 122(9):1764–1770.
- [22] Klein-Flügge MC, Bestmann S (2012) Time-dependent changes in human corticospinal excitability reveal value-based competition for action during decision processing. *Journal of neuroscience* 32(24):8373–8382.
- [23] Spampinato DA, Satar Z, Rothwell JC (2019) Combining reward and m1 transcranial direct current stimulation enhances the retention of newly learnt sensorimotor mappings. *Brain stimulation*.
- [24] Zénon A, et al. (2015) Increased reliance on value-based decision processes following motor cortex disruption. *Brain stimulation* 8(5):957–964.
- [25] Derosiere G, Vassiliadis P, Demaret S, Zénon A, Duque J (2017) Learning stage-dependent effect of m1 disruption on value-based motor decisions. *Neuroimage* 162:173–185.
- [26] Uehara S, Mawase F, Celnik P (2018) Learning similar actions by reinforcement or sensory-prediction errors rely on distinct physiological mechanisms. *Cerebral cortex* 28(10):3478–3490.
- [27] Galea JM, Jayaram G, Ajagbe L, Celnik P (2009) Modulation of cerebellar excitability by polarity-specific noninvasive direct current stimulation. *Journal of Neuroscience* 29(28):9115–9122.

- [28] Kuo MF, et al. (2008) Limited impact of homeostatic plasticity on motor learning in humans. *Neuropsychologia* 46(8):2122–2128.
- [29] van Beers RJ (2009) Motor learning is optimally tuned to the properties of motor noise. *Neuron* 63(3):406–417.
- [30] Derosiere G, Zénon A, Alamia A, Duque J (2017) Primary motor cortex contributes to the implementation of implicit value-based rules during motor decisions. *Neuroimage* 146:1115–1127.
- [31] Nikooyan AA, Ahmed AA (2015) Reward feedback accelerates motor learning. *Journal of neurophysiology* 113(2):633–646.
- [32] Grissom RJ, Kim JJ (2012) *Effect sizes for research: Univariate and multivariate applications*. (Routledge).
- [33] Mathys C, Daunizeau J, Friston KJ, Stephan KE (2011) A bayesian foundation for individual learning under uncertainty. *Frontiers in human neuroscience* 5:39.
- [34] Mathys CD, et al. (2014) Uncertainty in perception and the hierarchical gaussian filter. *Frontiers in human neuroscience* 8:825.
- [35] Iglesias S, et al. (2013) Hierarchical prediction errors in midbrain and basal forebrain during sensory learning. *Neuron* 80(2):519–530.
- [36] De Berker AO, et al. (2016) Computations of uncertainty mediate acute stress responses in humans. *Nature communications* 7:10996.
- [37] Pekny SE, Izawa J, Shadmehr R (2015) Reward-dependent modulation of movement variability. *Journal of Neuroscience* 35(9):4015–4024.
- [38] Van Mastrigt NM, Smeets JB, Van Der Kooij K (2020) Quantifying exploration in reward-based motor learning. *Plos one* 15(4):e0226789.
- [39] Bestmann S, de Berker AO, Bonaiuto J (2015) Understanding the behavioural consequences of noninvasive brain stimulation. *Trends in cognitive sciences* 19(1):13–20.
- [40] Thielscher A, Antunes A, Saturnino GB (2015) Field modeling for transcranial magnetic stimulation: a useful tool to understand the physiological effects of tms? in *2015 37th annual international conference of the IEEE engineering in medicine and biology society (EMBC)*. (IEEE), pp. 222–225.

- [41] Windhoff M, Opitz A, Thielscher A (2013) Electric field calculations in brain stimulation based on finite elements: an optimized processing pipeline for the generation and usage of accurate individual head models. *Human brain mapping* 34(4):923–935.
- [42] Glasser MF, et al. (2016) A multi-modal parcellation of human cerebral cortex. *Nature* 536(7615):171–178.
- [43] Ota K, Shinya M, Kudo K (2019) Transcranial direct current stimulation over dorsolateral prefrontal cortex modulates risk-attitude in motor decision-making. *Frontiers in human neuroscience* 13:297.
- [44] Sutton RS, Barto AG, , et al. (1998) *Introduction to reinforcement learning*. (MIT press Cambridge) Vol. 135.
- [45] Boorman ED, Behrens TE, Woolrich MW, Rushworth MF (2009) How green is the grass on the other side? frontopolar cortex and the evidence in favor of alternative courses of action. *Neuron* 62(5):733–743.
- [46] Domenech P, Koechlin E (2015) Executive control and decision-making in the prefrontal cortex. *Current opinion in behavioral sciences* 1:101–106.
- [47] Mawase F, Uehara S, Bastian AJ, Celnik P (2017) Motor learning enhances use-dependent plasticity. *Journal of Neuroscience* 37(10):2673–2685.
- [48] Antonenko D, et al. (2019) Towards precise brain stimulation: Is electric field simulation related to neuromodulation? *Brain stimulation* 12(5):1159–1168.
- [49] Stagg CJ, Nitsche MA (2011) Physiological basis of transcranial direct current stimulation. *The Neuroscientist* 17(1):37–53.
- [50] Galea JM, Mallia E, Rothwell J, Diedrichsen J (2015) The dissociable effects of punishment and reward on motor learning. *Nature neuroscience* 18(4):597.
- [51] Oldfield RC, , et al. (1971) The assessment and analysis of handedness: the edinburgh inventory. *Neuropsychologia* 9(1):97–113.
- [52] Seidel O, Ragert P (2019) Effects of transcranial direct current stimulation of primary motor cortex on reaction time and tapping performance: A comparison between athletes and non-athletes. *Frontiers in human neuroscience* 13.
- [53] Association WM, , et al. (2013) World medical association declaration of helsinki: ethical principles for medical research involving human subjects. *Jama* 310(20):2191–2194.

- [54] Dayan E, Averbeck BB, Richmond BJ, Cohen LG (2014) Stochastic reinforcement benefits skill acquisition. *Learning & memory* 21(3):140–142.
- [55] Nitsche MA, Paulus W (2001) Sustained excitability elevations induced by transcranial dc motor cortex stimulation in humans. *Neurology* 57(10):1899–1901.
- [56] Monte-Silva K, et al. (2013) Induction of late ltp-like plasticity in the human motor cortex by repeated non-invasive brain stimulation. *Brain stimulation* 6(3):424–432.
- [57] Mayka MA, Corcos DM, Leurgans SE, Vaillancourt DE (2006) Three-dimensional locations and boundaries of motor and premotor cortices as defined by functional brain imaging: a meta-analysis. *Neuroimage* 31(4):1453–1474.
- [58] Nitsche MA, et al. (2008) Transcranial direct current stimulation: state of the art 2008. *Brain stimulation* 1(3):206–223.
- [59] Friston K, Kiebel S (2009) Predictive coding under the free-energy principle. *Philosophical Transactions of the Royal Society B: Biological Sciences* 364(1521):1211–1221.
- [60] Sedley W, et al. (2016) Neural signatures of perceptual inference. *Elife* 5:e11476.
- [61] Diaconescu AO, et al. (2017) A computational hierarchy in human cortex. *arXiv preprint arXiv:1709.02323*.
- [62] Soch J, Allefeld C (2018) Macs—a new spm toolbox for model assessment, comparison and selection. *Journal of neuroscience methods* 306:19–31.
- [63] Geuzaine C, Remacle JF (2009) Gmsh: A 3-d finite element mesh generator with built-in pre-and post-processing facilities. *International journal for numerical methods in engineering* 79(11):1309–1331.
- [64] Basso D, Chiarandini M, Salmaso L (2007) Synchronized permutation tests in replicated $i \times j$ designs. *Journal of Statistical Planning and Inference* 137(8):2564–2578.
- [65] Benjamini Y, Krieger AM, Yekutieli D (2006) Adaptive linear step-up procedures that control the false discovery rate. *Biometrika* 93(3):491–507.
- [66] Ruscio J, Mullen T (2012) Confidence intervals for the probability of superiority effect size measure and the area under a receiver operating characteristic curve. *Multivariate Behavioral Research* 47(2):201–223.

Figures

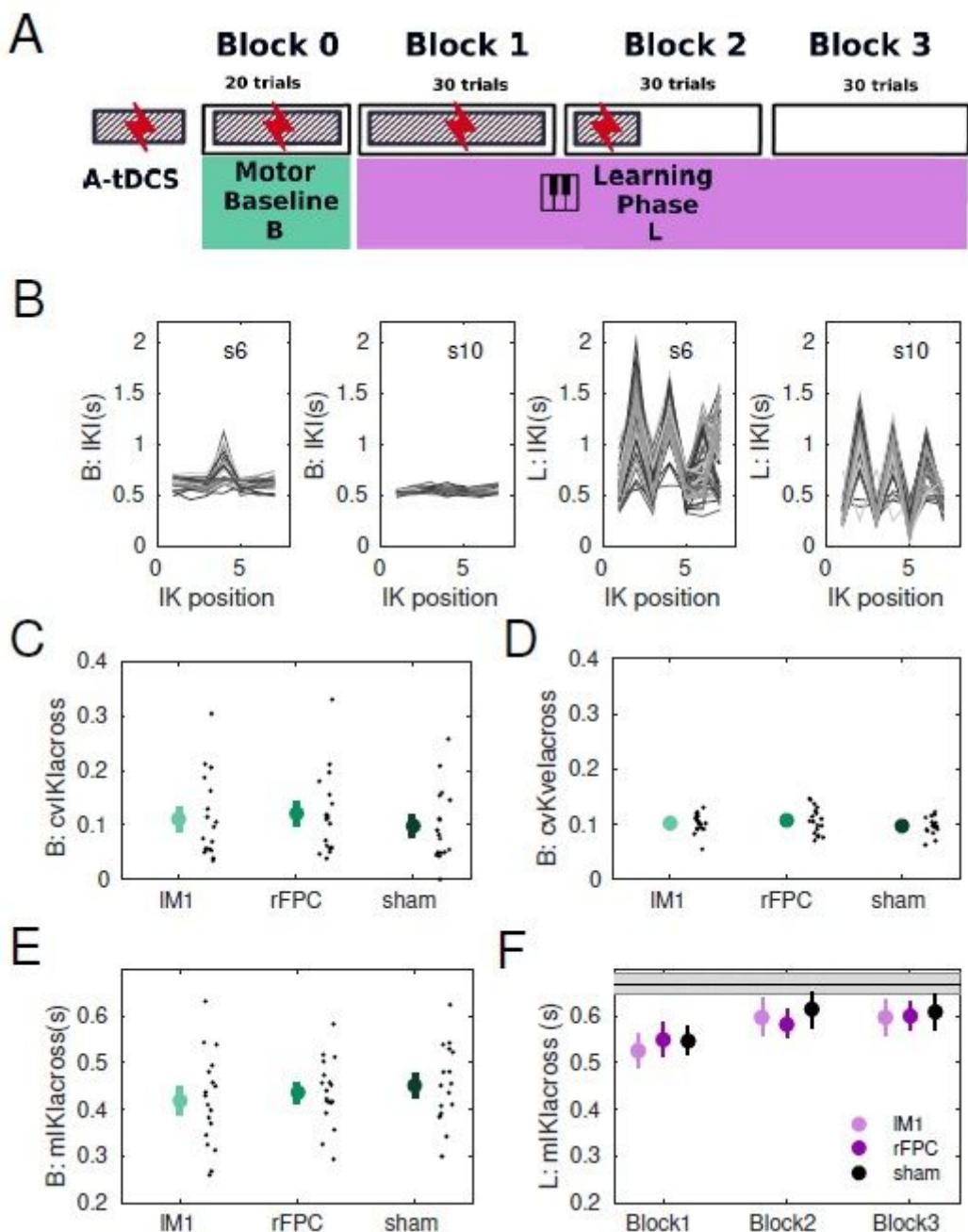
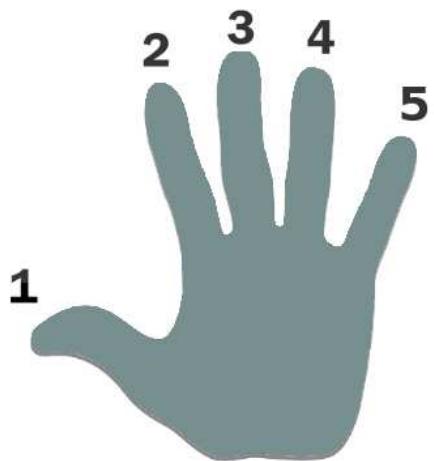
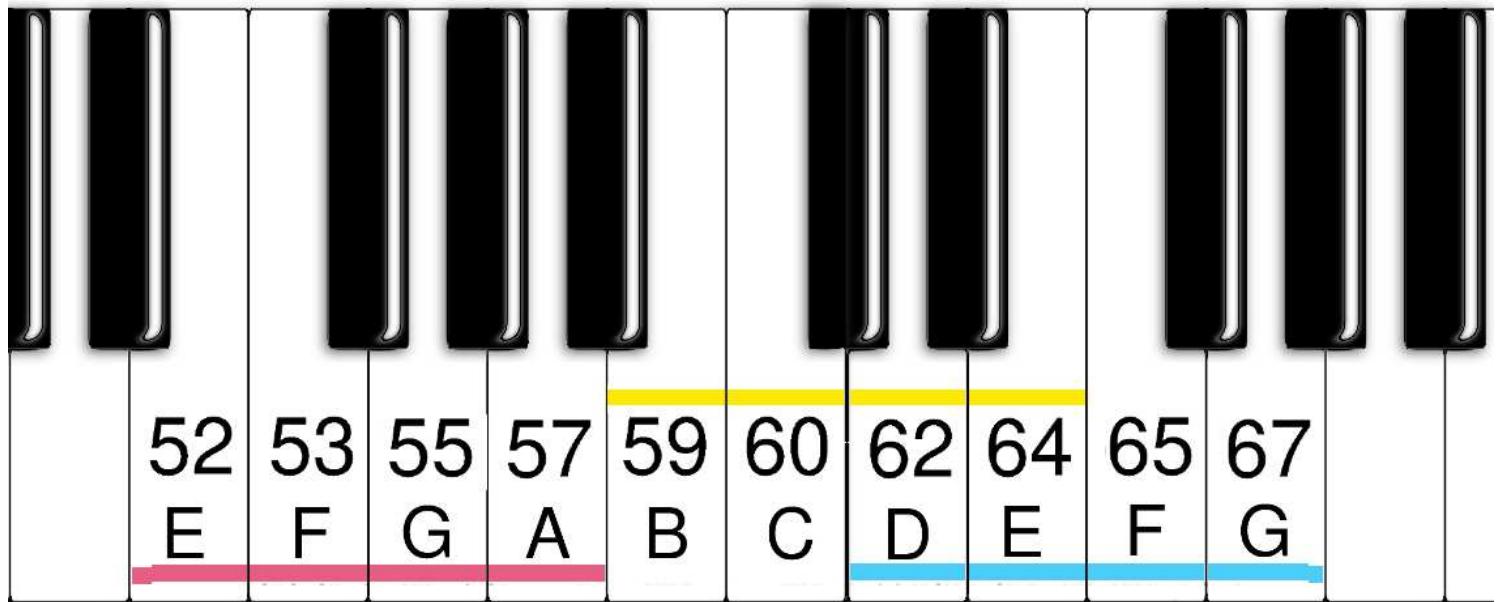


Figure 1

Experimental design and behavioural results. (A) Participants were tested on three separate weeks during which either an active tDCS protocol over the IM1, or the rFPC, or a sham stimulation condition were applied. All tDCS protocols extended for 20 minutes, which included (i) an initial 3 minute resting phase, (ii) a baseline phase of regular isochronous motor performance, and (iii) part of the reward-based learning blocks (1.333). (B) Illustration of timing performance with sham stimulation during the baseline phase (B; panels 1-2) and learning phase (L; panels 3-4) in two participants. Timing was measured using

the inter-keystroke-interval (IKI) in seconds, and shown for each inter-keystroke position (1 to 7 for sequences of 8 key presses). Different trajectories denote performance in different trials. (C-E) Effects of stimulation conditions on performance variables during the baseline phase (B). Colored big dots indicate means, with error bars denoting SEM. (C) Across-trials temporal variability, measured with the coefficient of variation of IKI or cvIKIacross in each block; (D) across-trials variability in keystroke velocity or loudness, cvKvelacross; (E) mean performance tempo, mIKI (s). (F) Average performance tempo during learning (L). The gray horizontal line indicates the mean tempo of the hidden target solutions.



- Sequence 1: 59 64 60 62 64 59 62 60
- Sequence 2: 67 65 67 62 62 64 65 64
- Sequence 3: 52 55 53 57 55 53 52 52

Figure 2

Stimulus materials. (see Manuscript file for full figure legend)

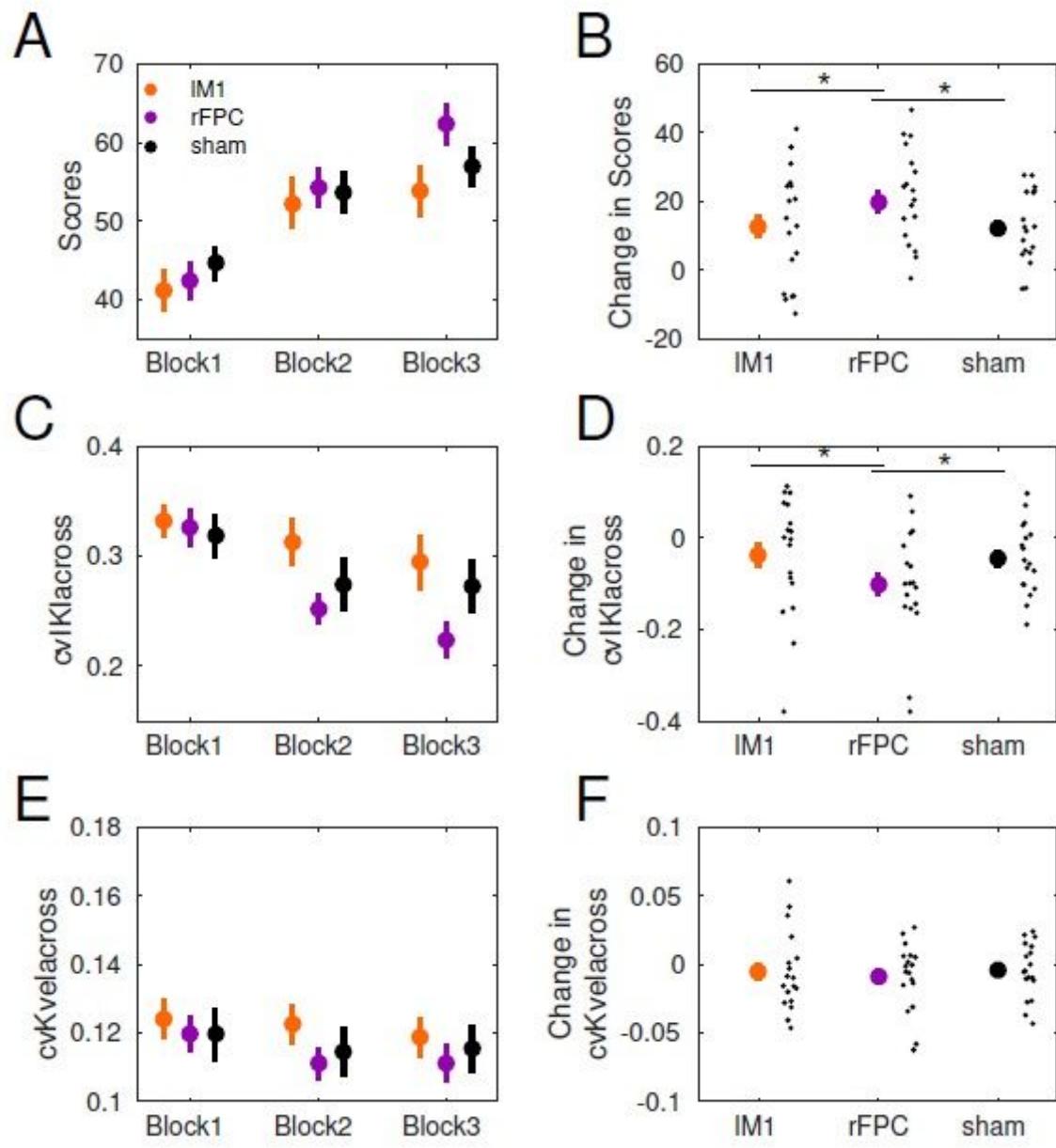


Figure 3

Behavioral results during reward-based learning. (see Manuscript file for full figure legend)

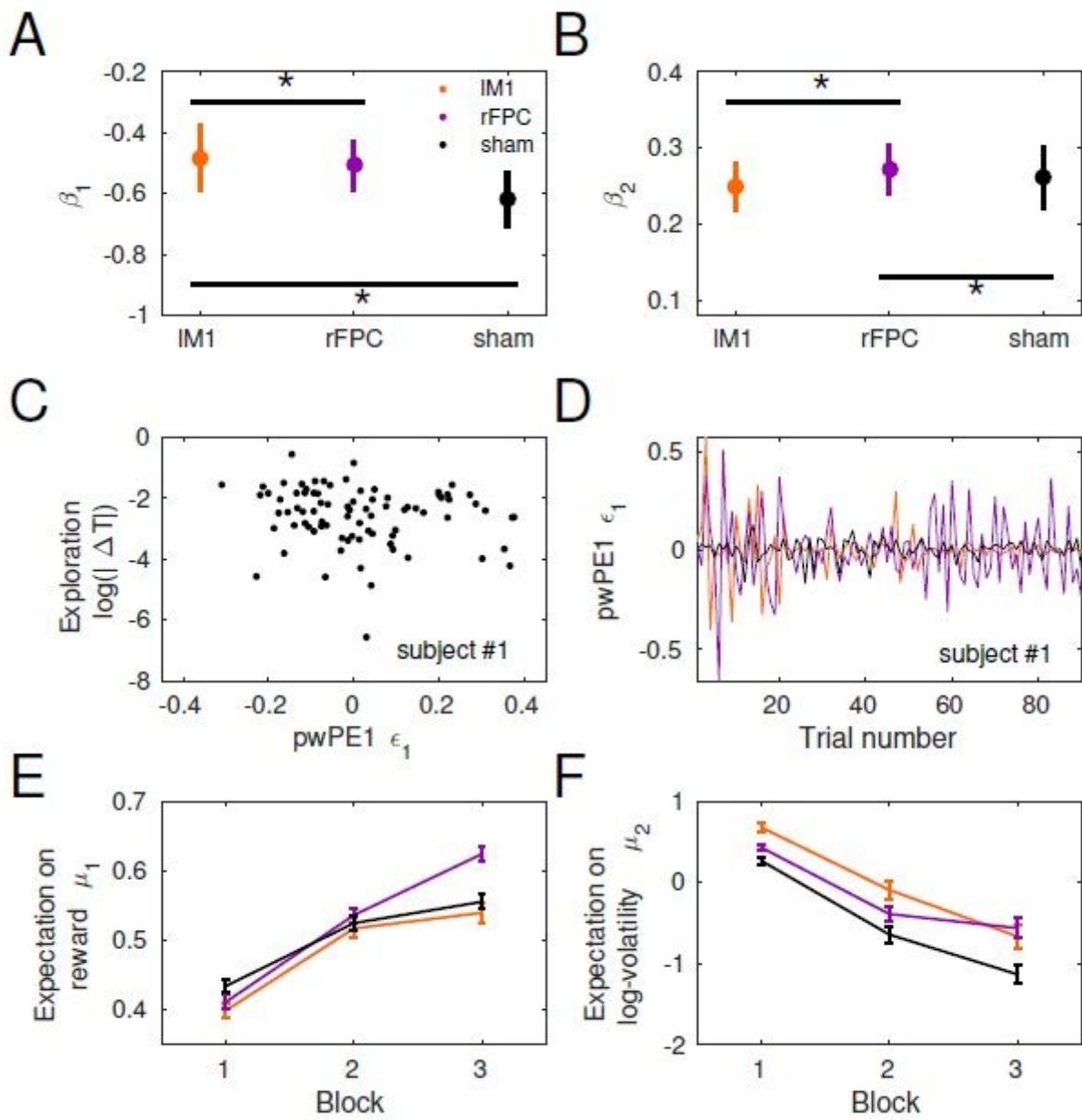
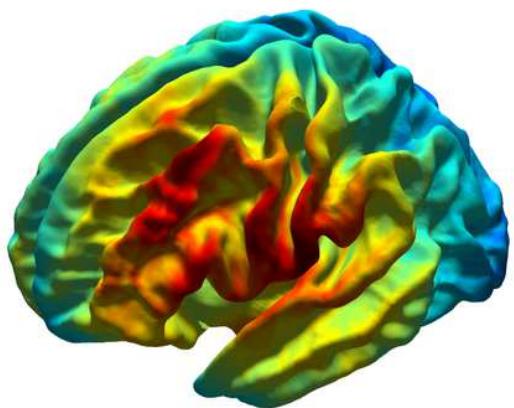
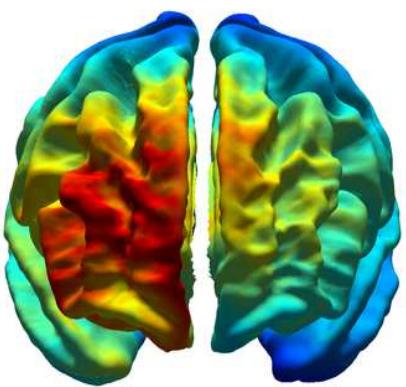
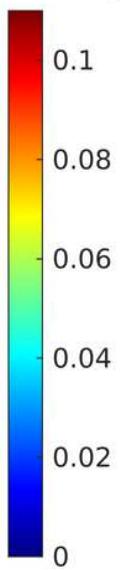


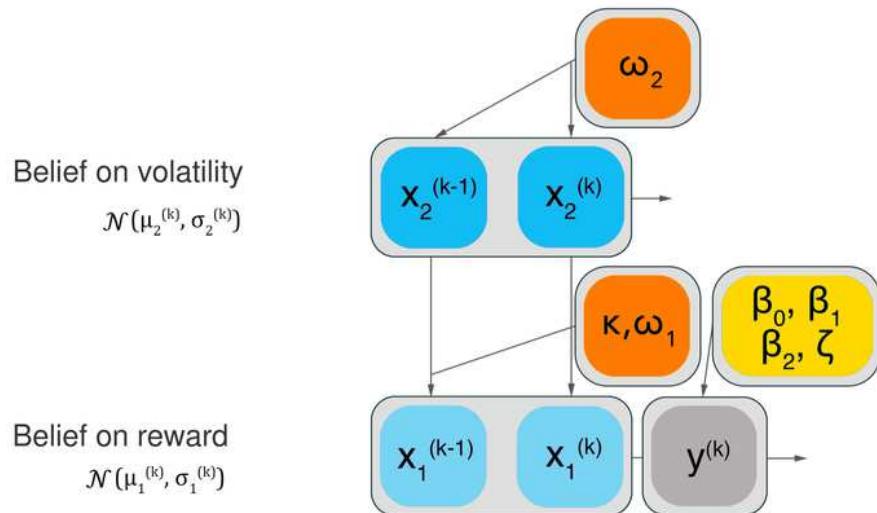
Figure 4

Computational modelling analysis. (see Manuscript file for full figure legend)

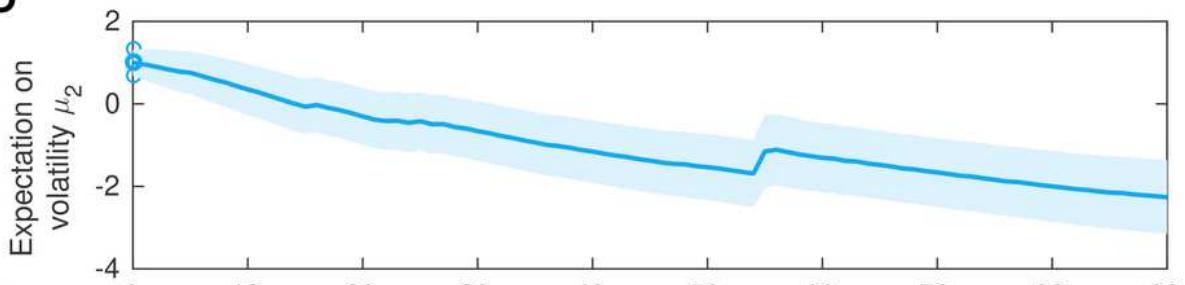
A**B****normE.avg****Figure 5**

Electric field distribution for anodal IM1-tDCS (see Manuscript file for full figure legend)

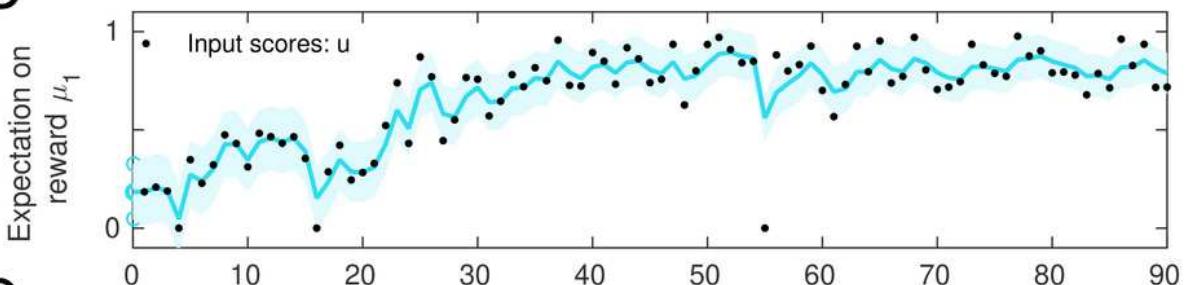
A



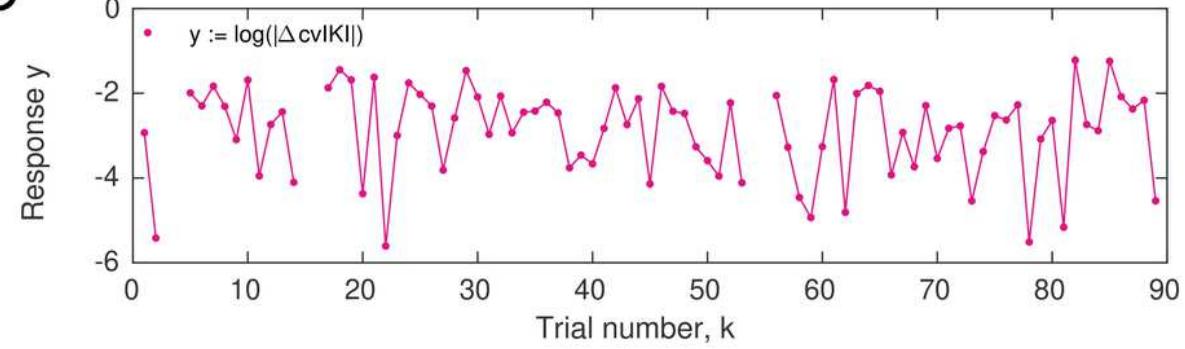
B



C



D

**Figure 6**

Computational model. (see Manuscript file for full figure legend)

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- SupplementaryMaterial.pdf