Behavioral/Cognitive

# Causal Evidence for Learning-Dependent Frontal Lobe Contributions to Cognitive Control

Paul S. Muhle-Karbe,[1,2,3] Jiefeng Jiang,[1,4] and Tobias Egner[1]

[1]Center for Cognitive Neuroscience, Duke University, Durham, North Carolina 27708, [2]Department of Experimental Psychology, University of Oxford, Oxford OX1 3UB, United Kingdom, [3]Department of Experimental Psychology, Ghent University, Ghent, Belgium 9000, and [4]Department of Psychology, Stanford University, Stanford, California 94305

The lateral prefrontal cortex (LPFC) plays a central role in the prioritization of sensory input based on task relevance. Such top-down control of perception is of fundamental importance in goal-directed behavior, but can also be costly when deployed excessively, necessitating a mechanism that regulates control engagement to align it with changing environmental demands. We have recently introduced the "flexible control model" (FCM), which explains this regulation as resulting from a self-adjusting reinforcement-learning mechanism that infers latent statistical structure in dynamic task environments to predict forthcoming states. From this perspective, LPFC-based control is engaged as a function of anticipated cognitive demand, a notion for which we previously obtained correlative neuroimaging evidence. Here, we put this hypothesis to a rigorous, causal test by combining the FCM with a transcranial magnetic stimulation (TMS) intervention that transiently perturbed the LPFC. Human participants (male and female) completed a nonstationary version of the Stroop task with dynamically changing probabilities of conflict between task-relevant and task-irrelevant stimulus features. TMS was given on each trial before stimulus onset either over the LPFC or over a control site. In the control condition, we observed adaptive performance fluctuations consistent with demand predictions that were inferred from recent and remote trial history and effectively captured by our model. Critically, TMS over the LPFC eliminated these fluctuations while leaving basic cognitive and motor functions intact. These results provide causal evidence for a learning-based account of cognitive control and delineate the nature of the signals that regulate top-down biases over stimulus processing.

Key words: cognitive control; computational modeling; prefrontal cortex; reinforcement learning; transcranial magnetic stimulation

## Significance Statement

A core function of the human prefrontal cortex is to control the signal flow in sensory brain regions to prioritize processing of task-relevant information. Abundant work suggests that such control is flexibly recruited to accommodate dynamically changing environmental demands, yet the nature of the signals that serve to engage control remains unknown. Here, we combined computational modeling with noninvasive brain stimulation to show that changes in control engagement are captured by a self-adjusting reinforcement-learning mechanism that tracks changing environmental statistics to predict forthcoming processing demands and that transient perturbation of the prefrontal cortex abolishes these adjustments. These findings delineate the learning signals that underpin adaptive engagement of prefrontal control functions and provide causal evidence for their relevance in behavioral control.

## Introduction

The prefrontal cortex is known to play a vital role in cognitive control (Miller and Cohen, 2001). In particular, the lateral prefrontal

cortex (LPFC) is commonly conceived as a source of top-down signals that amplify the processing of task-relevant information in posterior cortices to support focused and distractor-resistant cognition (Egner and Hirsch, 2005; Zanto et al., 2011). Importantly, however, although such control signals are of fundamental importance in the pursuit of goals, their excessive deployment can

also be costly because it can, for example, hinder the discovery of unattended but valuable information (Bocanegra and Hommel, 2014, Schuck et al., 2015). Adaptive cognition therefore requires a careful regulation of LPFC-based control to align it with fluctuating levels of environmental demand (Amer et al., 2016; Shenhav et al., 2017).

Empirically, such regulation can be witnessed in classic selective attention interference tasks such as the Stroop or Flanker protocols, in which participants typically enhance their attentional focus on task-relevant stimulus features in response to conflict induced by incongruent task-irrelevant features (Gratton et al., 1992; Botvinick et al., 2001). Intriguingly, such adaptation reflects both short-term (phasic) and long-term (tonic) trial history, suggesting that the brain effectively synthesizes conflict experiences over different time scales to establish optimal levels of cognitive focus (Torres-Quesada et al., 2013; Egner, 2014).

We have recently shown that both types of adaptation are captured by a single reinforcement-learning mechanism that minimizes the mismatch ("prediction error") between exerted and required levels of control (Jiang et al., 2014, 2015). The integration of recent and remote conflict experiences is realized via a flexible, self-adjusting learning rate that changes based on the inferred volatility of the task environment (Behrens et al., 2007). In stable environments, the learning rate is low so that predictions are based on a large trial history and the effect of occasional noise is minimal. In contrast, in volatile environments, the learning rate rises to ensure that predictions are based only on the recent trial history while discarding older, outdated evidence. This "flexible control model" (FCM) captures phasic and tonic adaptation in nonstationary tasks in which conflict probabilities change dynamically within runs (Jiang et al., 2014). Moreover, by combining the model with neuroimaging, we were able to reveal that changes in the model's learning rate and predicted control demand were encoded in the anterior insula and dorsal striatum, respectively. Conversely, LPFC activity tracked the extent to which striatal control predictions were used to adjust performance (Jiang et al., 2015).

Collectively, these findings suggest that the regulation of prefrontal top-down control is guided by a learning mechanism that infers changing environmental statistics to predict a task's forthcoming cognitive demand (see also Botvinick et al., 2001, Shenhav et al., 2013). Here, we aimed to put this hypothesis to a stringent, causal test by combining the FCM with a transcranial magnetic stimulation (TMS) intervention that transiently perturbed the LPFC during task performance. This setup created a powerful and novel window with which to study how the LPFC's role in performance changes due to learning.

Previous research has focused primarily on modulations of LPFC activity based on phasic (MacDonald et al., 2000; Kerns et al., 2004; Egner and Hirsch, 2005) and/or tonic changes in control demand (Carter et al., 2000; Braver et al., 2003; De Baene and Brass, 2013). Although informative, these approaches can neither reveal the nature of the learning signals that drive control engagement nor establish a causal link between changing levels of LPFC activation and performance. Our model-based approach to TMS permitted us to overcome these limitations by formulating an explicit learning mechanism (via the model) and directly probing the LPFC's causal role in the task (via TMS). To the extent that model-based predictions of control demand reflect varying levels of LPFC engagement, we expected them to capture fluctuations in the disruptive effects of TMS over time, thereby providing causal evidence for learning-dependent LPFC engagement.

## Materials and Methods

### Participants

Thirty-six healthy adults (9 male, mean age = 27.1 years, range = 18–45) were invited to participate in the study, which consisted of 2 experimental sessions taking place on separate days. Participants were recruited from a local database of former fMRI study participants and all had normal or corrected-to-normal visual acuity. Before participation, they were screened extensively for the presence of TMS contraindications based on guidelines by Rossi et al. (2009). Moreover, each participant's tolerance to the TMS intervention was tested at the start of the first session (see section below for details on TMS parameters). In the course of this procedure, seven participants withdrew from further participation due to unpleasant side effects of the stimulation (e.g., twitches of jaw or eye muscles). Two further participants had to terminate the study prematurely, one due to difficulty with positioning of the TMS coil and the other one due to technical failure in response collection during the first session. Therefore, the final sample consisted of 27 participants (6 male, mean age = 27.2 years, range = 18–39). Approval for all procedures was obtained from the Duke University Health System Institutional Review Board and participants gave written informed consent before each experimental session.

### Apparatus and stimuli

Stimulus presentation and response collection was controlled via Psychtoolbox in MATLAB (The MathWorks) on a laptop placed at a distance of ∼50 cm from the participants. Responses were collected with the "Z" and "N" keys of a QWERTY keyboard that was placed on the participant's lap. Stimuli consisted of a set of 24 grayscale photographs of faces (12 female, and 12 male, taken from the Cohn–Kanade face image database) of neutral expression that were overlaid with gender word labels (i.e., the words "man," "woman," "male," "female") printed in red font in either lowercase or uppercase letters. Compound face–word stimuli were presented centrally on a gray screen (RGB values = 100, 100, 100) subtending at an ∼4.6 × 5.7° visual angle. Between trials, a central fixation cross was shown (2.3 × 2.3°) in either white or red font (see next section for details).

### Experimental task

The experimental task required participants to categorize the compound face–word stimuli based on the gender of the face via button press while ignoring the gender of the overlaid word. Each target was paired with each distracter so that the gender of the face and the word could either correspond (congruent trials) or diverge (incongruent trials). Trials started with the presentation of a white crosshair for a randomly jittered duration taken from a uniform distribution of either 3000, 4000, or 5000 ms. Thereafter, a warning period was inserted, during which a red crosshair was shown for 500 ms, followed by a face–word target stimulus. Targets were shown for 250 ms and then replaced by a white crosshair, and responses were recorded for a duration of 2000 ms after target onset until the next trial started (Fig. 1).

Both experimental sessions began with a brief practice block of 12 trials in which performance feedback was presented centrally on the screen for 500 ms immediately after each response. Subsequently, participants worked through six runs of the task, each of which contained 96 trials. The first two runs served as behavioral training and did not entail TMS. The final four runs were the experimental blocks of interest and were completed with concurrent TMS (see section below for details). During the initial 16 trials of each run, half of the trials were congruent and the other half incongruent. These trials served as a "burn-in" period to establish equivalent and neutral (0.5) expectations of conflict across runs and participants. During the subsequent 80 trials, the probability of incongruent trials was varied in 4 phases of 20 trials with alternating probabilities of 0.2 (low proportion conflict) and 0.8 (high proportion conflict). This volatility served to encourage continuous learning of changing control–demand throughout the experiment. The sequence of low- and high-conflict phases was counterbalanced across runs and participants (Fig. 1). Stimulus presentation was random with the two constraints that face identities never repeated across consecutive trials and that distracter words always alternated from trial to trial between words with lowercase
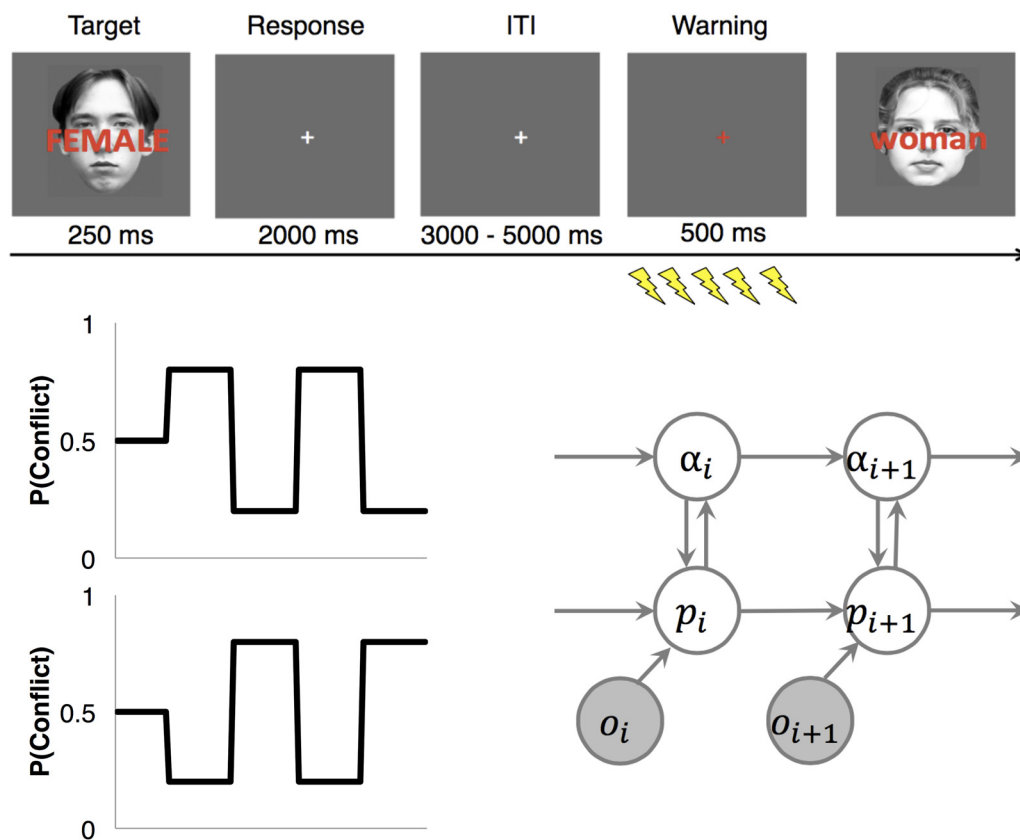
**Figure 1.** Task design and model architecture. Top, Single-trial structure of the Stroop task. The timing of the TMS pulses is displayed via the flash symbols. Bottom left, Two conflict probability distributions that were randomly assigned to each run. Bottom right, Structure of the FCM. The model consists of three variables: $\alpha$, $p$, and $o$. The colors indicate the difference between latent variables (white circles) and observable variables (gray circles). Arrows indicate the direction of information flow. On each trial ($i$), the observed congruency serves as input to the model and is used to update the state of the latent variables for the subsequent trial ($i + 1$; see Materials and Methods for details).

and uppercase letters. These constraints served to circumvent priming effects based on low-level stimulus feature repetitions (Mayr et al., 2003).

*TMS protocol*
The two experimental sessions involved identical procedures except that TMS was applied over different target sites. TMS pulses were delivered with a Magstim Rapid$^2$ stimulator via a Double 70 mm Air Film Coil. The navigation of the coil was guided by a frameless stereotaxic neuronavigation system (Brainsight) that located target sites onto individual anatomical MR images (Fig. 2). Target sites were selected based on anatomical criteria. One site was located in the LPFC and was localized at the posterior end of the left inferior frontal sulcus (average coordinates in MNI space: 35, 12, 24). Previous work has shown that this region is robustly activated by conflict in the Stroop task (Derrfuss et al., 2005, 2009), that damage to this region is associated with enhanced performance costs of conflict (Gläscher et al., 2012; Schroeter et al., 2012), and that its activity scales parametrically with behavioral adjustments due to learned control predictions (Jiang et al., 2015). Accordingly, this region is a likely candidate for the translation of conflict anticipation into cognitive control over stimulus processing. The other target site was located in the secondary somatosensory cortex and was localized by placing the coil on the interhemispheric midline and moving it beyond the central sulcus (average coordinate in MNI space: 0, −33, 57). This region is not implicated in cognitive control and TMS over this site was used as a control condition for nonspecific effects of TMS such as the discharge sound or the somatosensory sensation of the pulses (Clerget et al., 2013; Muhle-Karbe et al., 2014).

During the TMS runs, five pulses were delivered on each trial at a frequency of 10 Hz and an intensity corresponding to 60% of the maximum stimulator output. TMS trains started with the onset of the warning signal and ended 100 ms before the onset of the target stimulus. The rationale for this timing was to affect preparatory top-down control

mechanisms while leaving basic perceptual and motor processes unaffected. All parameters were modeled closely after one of our previous studies, in which a very similar protocol proved effective in disrupting control functions of the LPFC during context-based decision making (Muhle-Karbe et al., 2014). We chose not to calibrate TMS intensities based on participants' resting motor thresholds because the excitability of the primary motor cortex does not provide a reliable index for cortical excitability elsewhere in the brain (Stewart et al., 2001; Antal et al., 2004) and we considered a fixed stimulation intensity to provide a less arbitrary criterion (D'Ardenne et al., 2012).

*Statistical analyses*
*General linear model (GLM) analyses.* In the first set of analyses, we aimed to measure the effects of TMS over the two target regions on global indices of phasic and tonic conflict adaptation that are commonly used in the literature. In addition to reaction time (RT) and error rates, we also computed inverse efficiency scores (IES) for each design cell by dividing the respective RTs by the corresponding percentage of correct responses (Townsend and Ashby, 1983). IES serve to integrate speed and accuracy into a single index and were used to maximize the power of our analyses. The effects of TMS on phasic adaptation were analyzed in 2 (current trial congruency) × 2 (previous trial congruency) × 2 (TMS site) repeated-measures ANOVAs separately for each performance index. The effects of TMS on tonic adaptation were analyzed via 2 (current trial congruency) × 2 (proportion conflict) × 2 (TMS site) repeated-measures ANOVAs. Significant interaction terms were unpacked via planned paired-samples *t* tests. Phasic adaptation should be reflected in a significant interaction between the factors "current trial congruency" and "previous trial congruency," also known as conflict adaptation effect (Gratton et al., 1992; Botvinick et al., 2001). Tonic adaptation should be reflected in a significant interaction between the factors "current trial congruency" and "proportion conflict," also known as proportion congruency effect (Bugg and Crump,
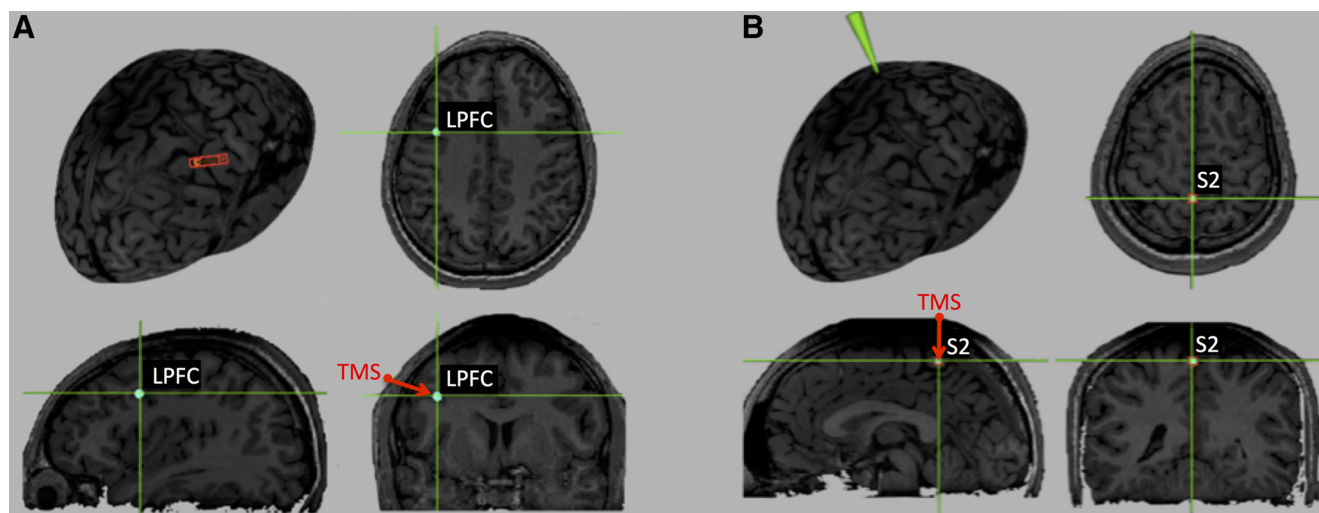
**Figure 2.** Illustration of the TMS navigation in the LPFC session (**A**) and in the active control condition (**B**). The top left image of each panel illustrates the position of the TMS coil center relative to a 3D reconstruction of an example participant's brain. The other images display the location of the target site in the axial, sagittal, and coronal planes. Average MNI coordinates of the target sites were −35, 12, 24 in the LPFC session and 0, −33, 57 in the control (S2) session. Red circles illustrate the central position of the TMS coil where the magnetic field was maximal and red arrows indicate the direction of the induced current flow. Note that the application of TMS also affects the brain tissue between the skull location of the coil and the respective brain target.

2012). An influence of the TMS intervention on either form of adaptation should be reflected in a significant three-way interaction term. Note that only trials from TMS runs after the burn-in phase were subjected to all analyses because only those trials could be meaningfully analyzed in terms of TMS sites and proportion conflict. Moreover, error trials, trials subsequent to errors, and outlier trials (i.e., trials with RTs deviating >2 SDs from the grand median of the respective session) were removed from all RT analyses. We conducted two additional follow-up analyses to evaluate the robustness of obtained effects. First, to judge the impact of our data-trimming procedure, we also compared adaptation scores between TMS sites with data in which outliers were identified based on the pooled SDs across both sessions. Second, to attenuate the impact of extreme participant values, we also conducted nonparametric comparisons using the Wilcoxon rank-sum test to evaluate whether there were significant differences between TMS sites.

*Response hand analyses.* We performed another set of validation analyses to ascertain the nature of TMS-induced changes in behavioral adaptation. These analyses compared adaptation effects between the left and right response hand to evaluate the possibility that behavioral effects of TMS over the LPFC were due to spread of the induced currents to the adjacent premotor cortex. In that case, TMS-induced performance modulation should be stronger with right-hand responses, due to the target site in the left hemisphere. We evaluated this possibility by including the additional factor response hand in the foregoing ANOVAs and by comparing adaptation scores directly between left-hand and right-hand responses.

*Model-based analyses: rationale.* The traditional behavioral indices of phasic and tonic adaptation, described above, reflect rather static measures of adaptation that conceive of participants as relying either exclusively on a very long-term (block-wise) or a very short-term (previous trial) trial history to adjust top-down control. In dynamic task environments, however, control engagement is likely more flexible, taking into account a variable trial history based on the inferred rate of change (volatility) of the environment (Behrens et al., 2007). The FCM (Jiang et al., 2014, 2015) seeks to account for this flexibility by estimating the optimal way of combining recent and remote trial history on each trial to determine the predicted conflict level (i.e., the probability of encountering an incongruent trial). In the central part of our analyses, we applied this model to evaluate directly our hypothesis that phasic and tonic adaptation are both expressions of this common learning mechanism that determines the relative engagement of LPFC-based top-down control.

The model's architecture has been described in detail previously (Jiang et al., 2014, 2015) and is therefore only briefly summarized here (Fig. 1).

Overall, the FCM consists of three key variables: the adaptive learning rate ($\alpha$), the predicted conflict level ($p$), and the observed trial congruency ($o$). The model takes trial-by-trial congruency as input and infers trialwise states of the first two variables as output. In its implementation, the model tracks a joint probabilistic distribution of the learning rate and the predicted conflict level. At the beginning of each trial before congruency is observed, this joint distribution is smoothed using a $\beta$ distribution to account for the approximated nature of the model predictions (for details, see Jiang et al., 2015). The predicted conflict level is then updated based on a reinforcement learning rule [i.e., $p \leftarrow p + \alpha(o - p)$]. The marginal mean of $\alpha$ and $p$ serve as estimates of learning rate and predicted conflict level on the current trial, respectively. After congruency is observed, the joint distribution is updated via Bayes' rule. The updated joint distribution is then used to simulate the next trial.

To link the model to behavior directly, a continuous variable of "control prediction error" (CPE) is computed, which is defined as the trialwise discrepancy between $p$ and $o$. Similar to other accounts of behavioral adaptation (Botvinick et al., 2001), the FCM assumes that a mismatch between anticipated and actual control demand leads to less efficient task information processing and thus to slower and less accurate performance. Within the model, the amount of mismatch is reflected in the CPE variable. To underscore the virtue of this model architecture with a flexible learning rate, we have demonstrated previously that it can simultaneously account for tonic and phasic adaptation effects (Jiang et al., 2014) and that it outperforms traditional models that rely on fixed learning rates, even when the latter are optimized *post hoc*; that is, after the whole trial sequence has been observed, rather than "on the fly" as in the FCM (Jiang et al., 2014, 2015). Note that, in one of our recent studies, the outlined architecture was amended by a RT variable that was used, in addition to the observed trial congruency, to update the model's latent variables and to reveal individual differences in learning-based control engagement (Jiang et al., 2015). Here, we chose not to include this variable because we focused on average TMS-induced changes in behavioral adaptation between the two sessions with different TMS target sites. Potential individual differences in adaptation were accounted for by our within-subjects design. Moreover, we expected prefrontal TMS to interfere with behavioral adaptation based on predicted control demand (see below). Therefore, including an RT variable would add noise selectively to the model estimation in the LPFC TMS session, thus biasing the comparison of latent model states.

When mapping the core components of the FCM onto the brain via model-based fMRI, we observed previously a segregation between brain regions involved in the learning of the predicted control level (anterior
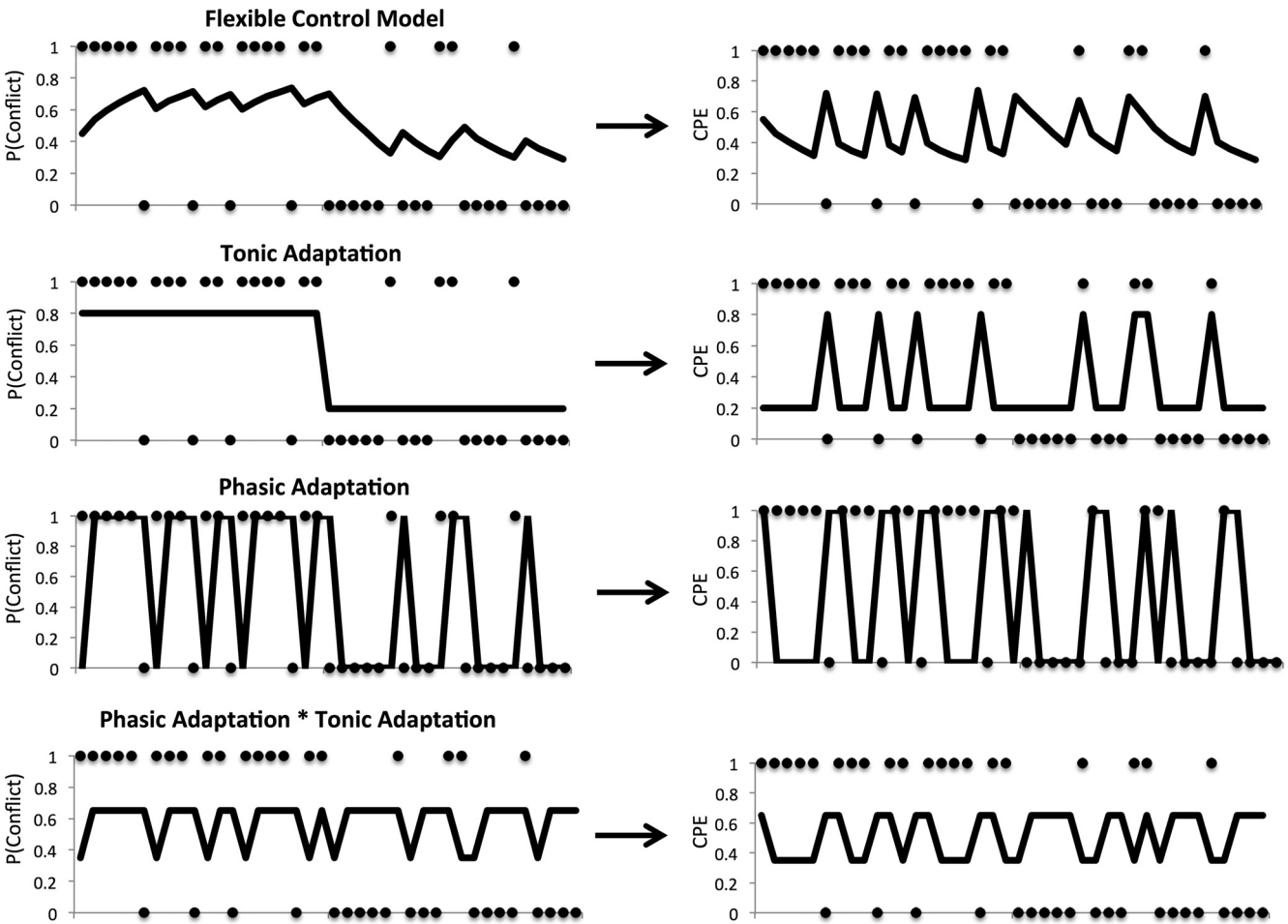
**Figure 3.** Illustration of the regressors used in the generation of models for the model comparison. The left panel illustrates how an example sequence of 40 trials (displayed as dots) encompassing one block phase with a high proportion of conflict and one block phase with a low proportion of conflict translates into estimated levels of conflict likelihood for each regressor. The right panel illustrates the resulting mismatch between the predicted and the observed level of conflict (i.e., the CPE of the respective regressor). The AMs for the model comparison were constructed by combining the different main effects and interactions shown here (see Materials and Methods section for details). Please note that CPE estimates are only included in the analysis for the FCM (top row) and are merely shown for illustrative purposes for the other models (bottom three rows).

insula and caudate nucleus) and those involved in implementing cognitive control based on model predictions (anterior cingulate cortex and LPFC). For the latter, we observed that individuals who exhibited a stronger neural representation of predicted conflict level in the LPFC also displayed a stronger correlation between CPE scores and response speed (Jiang et al., 2015). This suggests that the level of predicted conflict serves as a signal that engages LPFC-based top-down control in anticipation of the forthcoming trial. Critically, this negative correlation between CPE and behavioral performance relies on proactive control being guided by the predicted conflict level. Therefore, we expected that disrupting the LPFC via TMS before stimulus onset would perturb the modulation of predicted conflict level on control, thereby diminishing the correlation between CPE and performance.

*Model comparison.* Before evaluating the outcomes of the TMS intervention on model-based performance indices, we sought to replicate our previous findings that the FCM captures performance variability that cannot be accounted for by phasic and tonic adaptation effects alone. Moreover, we also investigated whether the performance variability captured by the FCM could be accounted for by the (weighted) sum of these two effects and/or their interaction. To this end, we conducted a formal model comparison using the behavioral data from the session with TMS over the control site to test the FCM against six alternative models (AMs; Fig. 3, Table 1): a phasic adaptation model (AM1), a tonic adaptation model (AM2), a hybrid model containing both phasic and tonic adaptation effects (AM3), a hybrid model containing both phasic and tonic adaptation effects as well as their interaction (AM4), a hybrid model

**Table 1. Summary of regression models that were submitted to the model comparison**

|  | FCM | AM1 | AM2 | AM3 | AM4 | AM5 | AM6 |
|---|---|---|---|---|---|---|---|
| Constant | X | X | X | X | X | X | X |
| Trial congruency | X | X | X | X | X | X | X |
| Control prediction error | X |  |  |  |  | X | X |
| Tonic adaptation |  |  | X | X | X | X | X |
| Phasic adaptation |  | X |  | X | X | X | X |
| Tonic adaptation * phasic adaptation |  |  |  |  | X |  | X |
| Total number of parameters | 4 | 4 | 4 | 6 | 8 | 8 | 10 |
| Protected exceedance probability (%) | 61.6 | 26.1 | 12.2 | 0.009 | 0.006 | 0.009 | 0.006 |

Columns represent models, the upper six lines represents regressors, and X's indicate that a given regressor was part of the respective model architecture. The bottom two lines indicate the number of free parameters and the protected exceedance probability of each model (see Materials and Methods for details).

containing both phasic and tonic adaptation effects as well as CPE (AM5), and a hybrid model containing both phasic and tonic adaptation effects, their interaction, and CPE (AM6). All seven models were GLMs with each row of a GLM encoding the effect(s) for one trial. Each effect (e.g., phasic adaptation, tonic adaptation, their interaction, and CPE) was represented by two regressors, one for congruent trials and one for incongruent trials, respectively. Two additional regressors were included for each model, a constant regressor and a regressor encoding trialwise congruency to account for the classic congruency effect (see Table 1 for an overview). Please note that an additional regressor for the main effect

of the respective modulator does not change the obtained results because these main effects are already captured by linear combinations of the other model regressors. Moreover, when modeling both trial types separately, the CPE regressor is equivalent to a regressor representing the actual control prediction (with reversed coefficient sign on incongruent trials). Therefore, the FCM performs identically to a model that uses control prediction rather than CPE to explain performance.

For the FCM, the modulator was CPE, as discussed above. For the phasic and tonic adaptation models, the modulator was the previous trial congruency and the average block probability of incongruent trials, respectively, according to the definition of the two adaptation effects. For the hybrid models, performance is affected by a combination of phasic and tonic adaptation effects (Fig. 3). Note that these adaptation effects only represent two special cases of conflict-level predictors. Even their combination provides limited flexibility in how previous trials can influence behavior because an older trial can only have one of three possible weights in predicting conflict level (i.e., one weight for the previous trial, one weight for other trials within the block, and one weight for all other trials). In contrast, the FCM adaptively adjusts the contributions of all previous trials based on the environmental volatility, thus ensuring greater flexibility than provided by the AMs.

For each model and participant, model performance was assessed using a leave-one-run-out cross-validation (Chiu et al., 2017) to control for potential overfitting by candidate models with a larger number of free parameters. That is, the model was fit to the trialwise RT in three of the four runs (training sample) and then used to predict the trialwise RT in the remaining run (test sample). This procedure was repeated four times, with each run serving as the test sample once. Model performance was quantified using the log likelihood, calculated in the following manner:

$$Log - likelihood = nln\left(\frac{\sum_i^n PE_i^2}{n}\right)$$

where $n$ is the number of trials and $PE_i$ is the prediction error of RT at trial $i$. Assuming that a flexible learning rate captures variance in control engagement that is missed by the individual adaptation effects, their weighted sum, and their interactions and their weighted sum, we expected the FCM to outperform all AMs.

Finally, we tested whether the model-based CPE variable could account for unique variance in performance, even when phasic and tonic adaptation are both modeled explicitly. To assess this, we compared model performance between AM4, which includes tonic and phasic adaptation effects and their interaction, and AM6, which includes model-based CPE, in addition to the regressors of AM4. Both models were fit to all trialwise RTs (without cross-validation) for each subject and the model prediction error was used to approximate log likelihood (i.e., the number of trials times the logged mean squared prediction error). Group sums of log likelihood were compared between the two models via a likelihood ratio test.

*Relation between phasic and tonic adaptation.* We conducted another set of analyses to investigate the relationship between phasic and tonic adaptation in more detail. Here, we aimed to clarify to what extent phasic adaptation in our task could be explained merely via the periodic changes in the proportions of congruent and incongruent trials in the different block phases. This could be the case given that the different trial congruency sequences, which are used to compute phasic adaptation scores, are not equally distributed across block phases with low and high proportion conflict (e.g., sequences with two consecutive congruent trials are more common in low proportion conflict phases). Accordingly, a tonic adaptation effect alone that improves performance on congruent trials in low proportion conflict phases and on incongruent trials in high proportion conflict phases could produce a spurious phasic adaptation effect despite no short-term changes in control engagement taking place. To evaluate this possibility, we conducted two types of control analyses. Initially, we performed a follow-up model comparison between the phasic adaptation model (AM1) and the tonic adaptation model (AM2) from the foregoing section. This analysis should reveal directly the extent to which each form of adaptation contributed to performance in our task. In a second step, we reanalyzed data from the session with TMS over the control site in a 2

(proportion conflict) × 2 (previous trial congruency) × 2 (current trial congruency) repeated-measures ANOVA. This ANOVA models both adaptation effects simultaneously and was performed separately for RT, error rates, and IES. Potential influences of tonic changes in conflict probability on phasic adaptation should be reflected in a significant three-way interaction term.

*Model-based analyses of TMS effects.* In the final and central part of our analyses, we applied the FCM to evaluate the effects of the TMS intervention on learning-based control engagement during task performance. As noted above, we expected that TMS over the LPFC would diminish the modulation of performance based on anticipated control demand, which is reflected in the correlation between CPE and performance. To examine a modulation of RT based on CPE, we constructed two GLMs, one for congruent trials and one for incongruent trials. Each GLM encoded a regressor of trialwise CPE levels as well as a constant regressor. GLMs were fit to trialwise RTs and the coefficient of the CPE regressor was considered the modulation of CPE on the GLM's corresponding congruency type. Given the noncontinuous nature of performance accuracy, we could not examine a CPE-modulation of error rates and IES via trialwise regression analyses. Instead, we grouped experimental trials into quartiles based on their estimated level of CPE and calculated error rates and IES for each quartile. This was done separately for both TMS sites (LPFC vs control site) and trial types (congruent vs incongruent). We then computed slopes for each condition, expressing the extent to which performance linearly scaled across CPE quartiles. These slopes served as an index for performance modulation based on control predictions and were analyzed in 2 (TMS site) × 2 (trial type) repeated-measures ANOVAs. Assuming that CPE serves to engage LPFC-based top-down control, we expected slopes to be significantly greater than zero in the active control condition, reflecting a modulation of performance based on the anticipated level of conflict. In contrast, in the LPFC session, we expected this correlation to be weakened, reflecting interference with the learning-dependent engagement of cognitive control.

## Results

### Phasic adaptation

The RT analysis revealed a significant main effect of current trial congruency ($F_{(1,26)} = 79.547$, $p < 0.001$, $\eta^2 = 0.754$), reflecting faster responses on congruent trials (502 ms) than on incongruent trials (526 ms). The three-way interaction involving the factors of previous trial congruency, current trial congruency, and TMS site was marginally significant ($F_{(1,26)} = 3.873$, $p = 0.059$, $\eta^2 = 0.130$), reflecting a trend toward greater phasic adaptation in the control session than in the LPFC session. All other main effects and interactions were nonsignificant.

The analysis of error rates revealed a significant main effect of current trial congruency ($F_{(1,26)} = 29.202$, $p < 0.001$, $\eta^2 = 0.529$), reflecting fewer errors on congruent trials (1.8%) than on incongruent trials (4.7%). The main effect of previous trial congruency was significant as well ($F_{(1,26)} = 6.284$, $p = 0.019$, $\eta^2 = 0.195$), reflecting greater error rates after congruent trials (3.7%) than after incongruent trials (2.9%). These two factors also interacted ($F_{(1,26)} = 6.752$, $p = 0.015$, $\eta^2 = 0.206$), reflecting phasic conflict adaptation; that is, reduced congruency effects after incongruent trials (1.9%) than after congruent trials (3.8%). The three-way interaction involving the TMS site factor was nonsignificant ($F_{(1,26)} = 1.231$, $p = 0.277$, $\eta^2 = 0.045$).

Finally, we analyzed participants' IES by dividing RT of each design cell by the percentage of correct responses (see Materials and Methods). This analysis revealed a significant main effect of current trial congruency ($F_{(1,26)} = 80.893$, $p < 0.001$, $\eta^2 = 0.757$), reflecting enhanced performance on congruent trials (511 ms) compared with incongruent trials (552 ms). The main effect of previous trial congruency was significant as well ($F_{(1,26)} = 4.269$, $p < 0.049$, $\eta^2 = 0.141$), reflecting enhanced performance after incongruent trials (529 ms) relative to congruent trials (534

ms). These two factors also interacted ($F_{(1,26)} = 6.583$, $p < 0.016$, $\eta^2 = 0.202$), reflecting phasic adaptation due to smaller congruency effects after incongruent trials (33 ms) than after congruent trials (49 ms). Importantly, as shown in the bottom panel in Figure 4A, this effect was further qualified by a significant 3-way interaction ($F_{(1,26)} = 4.395$, $p < 0.046$, $\eta^2 = 0.145$) due to significant adaptation effects in the control session ($t_{(26)} = 3.493$, $p = 0.002$, $d = 0.672$), but not in the LPFC session ($t_{(26)} = 0.534$, $p = 0.598$, $d = 0.103$). Notably, the difference in phasic adaptation scores between target sites was reduced to trend level when RT outliers were identified based on the pooled SDs across experimental sessions ($t_{(26)} = 1.965$ $p = 0.060$). However, the nonparametric Wilcoxon rank-sum test confirmed the significance of the effect ($z = 2.114$, $p = 0.034$), further suggesting that the reduction in phasic adaptation with LPFC stimulation reflected a shift across the whole sample rather than just individual extreme values. Therefore, in sum, we obtained evidence that perturbing the LPFC before stimulus onset diminishes phasic adjustments of cognitive control.

**Tonic adaptation**
Beyond the main effect of current trial congruency reported in the previous section, the RT analysis revealed a significant three-way interaction involving current trial congruency, proportion conflict, and TMS site ($F_{(1,26)} = 4.477$, $p < 0.044$, $\eta^2 = 0.147$). As shown on the top panel in Figure 4B, tonic adaptation was marginally significant in the control session ($t_{(26)} = 2.031$, $p = 0.053$, $d = 0.391$) and nonsignificant in the LPFC session ($t_{(26)} = -1.267$, $p = 0.217$, $d = 0.244$). However, the effect was only at trend level with the alternative trimming procedure ($t_{(26)} = 1.969$, $p = 0.059$), and with the nonparametric Wilcoxon test ($z = 1.946$, $p = 0.052$).

Next to a main effect of current trial congruency, described above, the error analysis revealed a main effect of proportion conflict ($F_{(1,26)} = 7.332$, $p < 0.012$, $\eta^2 = 0.220$), reflecting fewer errors in high conflict phases (2.9%) compared with low conflict phases (4.0%). These two factors also interacted ($F_{(1,26)} = 5.232$, $p < 0.031$, $\eta^2 = 0.168$), reflecting tonic adaptation, i.e., reduced congruency effects in high conflict phases (2.3%), compared with low conflict phases (4.2%). The three-way interaction with the factor TMS site was nonsignificant ($F_{(1,26)} = 1.504$, $p < 0.231$, $\eta^2 = 0.055$).

Finally, beyond the aforementioned main effect of current trial type, the IES analysis revealed a significant two-way interaction between current trial congruency and proportion conflict ($F_{(1,26)} = 4.974$, $p < 0.035$, $\eta^2 = 0.161$), reflecting tonic adaptation with smaller congruency effects in high conflict phases (33 ms) than in low conflict phases (48 ms). Importantly, this effect was further qualified by a significant three-way interaction with the TMS site factor ($F_{(1,26)} = 6.423$, $p < 0.018$, $\eta^2 = 0.198$). As shown in the bottom panel in Figure 4B, this interaction was driven by a significant tonic adaptation effect in the control session ($t_{(26)} = 2.952$, $p = 0.007$, $d = 0.568$) that was abolished in the LPFC session ($t_{(26)} = 0.105$, $p = 0.917$, $d = 0.020$). This effect was replicated with the alternative data trimming procedure ($t = 2.376$, $p = 0.025$), and with the nonparametric Wilcoxon test ($z = 2.643$ $p = 0.008$). Therefore, in sum, we also obtained evidence that perturbing the LPFC before stimulus onset diminishes tonic adjustments of cognitive control.

**Response hand effects**
To determine whether the effects of prefrontal stimulation, observed above, could be ascribed to activation spread to premotor regions rather than genuine effects on cognitive control processes, we first reran the previous ANOVAs by including the additional factor 'response hand' (left vs right). We did not observe any significant effect involving the response hand factor, suggesting that both hands were contributing equally to the observed patterns of results. In a second step, we evaluated adaptation scores (expressed as IES) separately for both response hands, target sites, and time scales. In the control session, phasic adaptation was marginally significant in the left hand ($t_{(26)} = 1.914$, $p = 0.067$, $d = 0.368$) and significant in the right hand ($t_{(26)} = 3.314$, $p = 0.003$, $d = 0.638$), whereas in the LPFC session phasic adaptation was absent for both left-hand ($t_{(26)} = 0.193$, $p = 0.849$, $d = 0.037$) and right-hand responses ($t_{(26)} = 0.660$, $p = 0.515$, $d = 0.127$). Similarly, in the control session, tonic adaptation scores were marginally significant in the left hand ($t_{(26)} = 1.947$, $p = 0.062$, $d = 0.375$) and significant in the right hand ($t_{(26)} = 2.364$, $p = 0.026$, $d = 0.455$), whereas in the LPFC session adaptation was absent both with left-hand ($t_{(26)} = 0.179$, $p = 0.859$, $d = 0.034$) and with right-hand responses ($t_{(26)} = 1.042$, $p = 0.307$, $d = 0.200$). Together, these results provide no evidence for an alternative explanation of our results in terms of activation spread to the premotor cortex.

**Model comparison**
After confirming that the TMS intervention diminished the specific indices of phasic and tonic adaptation to conflict, we next used model-based analyses to study in more detail how TMS over the LPFC affected the learning-based engagement of cognitive control. Initially, we tested whether the FCM accounts for variance in behavioral performance that is not captured by the AMs (see section on model comparison in the Materials and Methods section for details). The FCM (log likelihood = −40737) indeed outperformed all other models under scrutiny (log likelihoods: AM1 = −40728; AM2 = −40670; AM3 = −40630; AM4 = −40505; AM5 = −40567; AM6 = −40437; all $p$-values <0.0001). Following the recommendations by Rigoux et al. (2014), we also calculated the Bayesian omnibus risk to index the statistical risk associated with our model selection, and protected exceedance probabilities to index the relative model likelihood. This confirmed that the FCM performed best among all tested models (protected exceedance probabilities: FCM = 61.6%; AM1 = 26.1%; AM2 = 12.2%, AM3 = 0.009%, AM4 = 0.006%, AM5 = 0.009%, AM6 = 0.006%; Bayesian omnibus risk = 0.000129), though it should be emphasized that the observed exceedance probability provides only modest evidence in favor of the FCM. The superior performance of the FCM, relative to AMs 1–4, suggests that it captures unique behavioral variance that is not accounted for by the individual effects of phasic and tonic adaptation, or by their weighted sum. Moreover, the inferior performance of the more complex hybrid models (i.e., AMs 4 and 6) likely reflects overfitting and consequently poor cross-generalization across experimental runs. Finally, the comparison of the group sums of log likelihood between AM4 and AM6 revealed that the latter performed significantly better ($\chi^2_1 = 86.79$, $p < 0.001$), corroborating that CPE explains unique variance in RT, even when proportion congruency and previous trial congruency are both explicitly modeled. In sum, we obtained evidence that the FCM provided the most effective explanation of our data.

**Relation between phasic and tonic adaptation**
We next conducted two sets of control analyses to examine the extent to which phasic adaption in our task could be explained as a mere side effect of tonic changes in the proportion of congruent
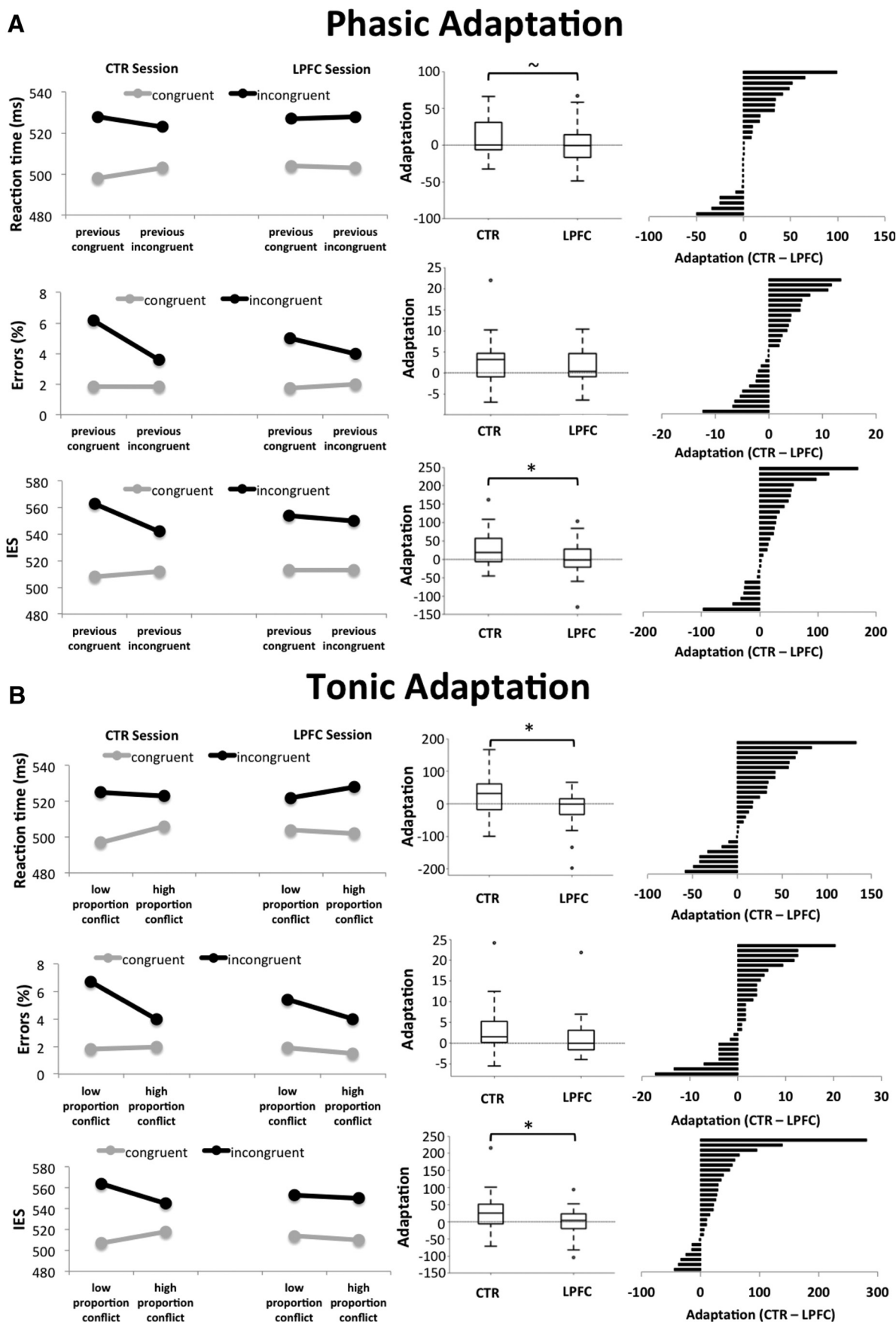
**Figure 4.** Illustration of behavioral performance across TMS sessions. **A**, Performance as a function of TMS target site [control site (CTR) vs LPFC], previous trial congruency (congruent vs incongruent), and current trial congruency (congruent vs incongruent) separately for RT, error rates, and IES. Boxplots in the middle display the distribution (*Figure legend continues.*)

and incongruent trials (see Materials and Methods section for details). First, we performed a follow-up comparison between the phasic adaptation model (AM1) and the tonic adaptation model (AM2) from the model comparison reported above. This revealed that the phasic adaptation model (log likelihood = −40728) tended to outperform the tonic adaptation model (log likelihood = −40670, protected exceedance probablilty = 65%, Bayesian omnibus risk = 0.65), but please note that the observed Basian omnibus risk suggests a similar fit of the two models to our data. Second, we reanalyzed performance in the session with TMS over the control site in a 2 (proportion conflict) × 2 (previous trial congruency) × 2 (current trial congruency) repeated-measures ANOVA. This analysis yielded nonsignificant three-way interactions for RT ($F_{(1,26)} = 0.383, p = 0.541$), error rates ($F_{(1,26)} = 0.873, p = 0.359$), and IES ($F_{(1,26)} = 1.906, p = 0.179$). Hence, there was no evidence in our data for a modulation of phasic adaption based on changing conflict probabilities across block phases (see Torres-Quesada et al., 2013 for similar results). Together, these results corroborate that phasic adaptation in our task was not a mere side effect of tonic changes in in the relative proportion of trial congruency sequences.

**Model-based analyses of TMS effects**
Having established the model's fit to our data, we next examined the effects of the TMS intervention on model-based indices of behavioral adaptation (see Materials and Methods section for details). The RT analysis revealed a marginally significant main effect of TMS site ($F_{(1,26)} = 3.371, p = 0.078, \eta^2 = 0.115$), reflecting a trend toward greater slopes in the control session than in the LPFC session. The main effect of trial congruency, and the interaction term were both nonsignificant. The analysis of error rates revealed main effect of TMS site ($F_{(1,26)} = 9.056, p = 0.006, \eta^2 = 0.258$), reflecting greater slopes in the control session than in the LPFC session. The main effect of trial congruency was significant as well ($F_{(1,26)} = 12.744, p = 0.001, \eta^2 = 0.329$), driven by greater slopes on incongruent trials, relative to congruent trials. These two factors also interacted ($F_{(1,26)} = 4.572, p = 0.042, \eta^2 = 0.150$). As shown in Figure 5, *post hoc* comparisons revealed that, on incongruent trials, slopes were significantly greater in the control session than in the LPFC session ($t_{(26)} = 4.258, p < 0.001, d = 0.819$), whereas slopes did not differ between TMS sites on congruent trials ($t_{(26)} = 0.913, p = 0.369, d = 0.176$). Finally, the analysis of IES revealed nonsignificant main effects of TMS site ($F_{(1,26)} = 2.047, p = 0.164, \eta^2 = 0.073$) and trial congruency ($F_{(1,26)} = 1.551, p = 0.224, \eta^2 = 0.056$). Critically, however, the two factors interacted ($F_{(1,26)} = 4.808, p < 0.037, \eta^2 = 0.156$). As shown in Figure 5, *post hoc* comparisons revealed that on incongruent trials, slopes were significantly larger in the control session than in the LPFC session ($t_{(26)} = 2.668, p = 0.013, d = 0.513$). By

←

(*Figure legend continued.*) of adaptation scores (for each performance index) that were computed by subtracting congruency effects after incongruent trials from congruency effects after congruent trials. Solid lines of the boxplots indicate the median of the distribution, the box outlines indicate the 25[th] and 75[th] percentile, and the whiskers indicate 1.5 times the interquartile range. Extreme values are shown separately as unfilled circles. Bar graphs on the right display the TMS-induced change in phasic adaptation separately for each participant. ***B***, Performance as a function of TMS target site (control site vs LPFC), proportion conflict (low vs high), and current trial congruency (congruent vs incongruent), separately for RT, error rates, and IES. Boxplots in the middle display adaptation scores that were computed by subtracting congruency effects in high proportion conflict from congruency effects with low proportion conflict. Bar graphs on the right display the TMS-induced change in tonic adaptation separately for each participant. ∼$p < 0.10$, *$p < 0.05$.

contrast, on congruent trials, slopes did not differ between sessions ($t_{(26)} = 0.121, p = 0.905, d = 0.023$). Altogether, the pattern of results corroborates that behavioral adaptation in the active control condition reflected the graded engagement of cognitive control based on changing expectations about the probability of conflict. These expectations were inferred from both recent and remote trial history and effectively captured by the FCM. Most importantly, this modulation was abolished after transient perturbation of the left LPFC, consistent with a failure to engage control based on learned anticipations of control demand.

## Discussion
We transiently perturbed the left LPFC during the performance of a nonstationary Stroop task that entailed dynamic shifts in the probability of conflict over time. In the active control condition, participants exhibited adaptive fluctuations in their attentional focus on task-relevant stimulus features, consistent with changing conflict expectations that were inferred from recent and remote experiences. Perturbation of the LPFC abolished these adjustments while leaving basic cognitive and motor functions intact. Below, we discuss the implications of our findings along with potential directions for further inquiry.

**Toward a learning-based neuroscience of cognitive control**
An extensive body of neuroscience literature has documented a central role of the LPFC in the goal-dependent prioritization of sensory input (Desimone and Duncan, 1995; Miller and Cohen, 2001). LPFC neurons exhibit flexible tuning profiles that encode only those aspects of the environment that are used to perform the task at hand (Freedman et al., 2001; Stokes et al., 2013). The LPFC also adapts its functional connectivity with the posterior brain, where it synchronizes with dedicated occipital and temporal lobe regions implicated in processing task-relevant inputs (Zanto et al., 2010; Baldauf and Desimone, 2014). These mechanisms are closely tied with behavioral control because perturbation of the LPFC diminishes feature selectivity in sensory areas during stimulus encoding (Zanto et al., 2011; Lee and D'Esposito, 2012) and impedes performance in categorization tasks that require the flexible use of different stimulus features (Zanto et al., 2011; Muhle-Karbe et al., 2014).

Our study provides an important extension of this rich literature by linking the relative engagement of LPFC-based top-down control over stimulus processing to a learning mechanism that infers latent statistical structure in dynamic environments to predict forthcoming cognitive demand. By flexibly weighting recent and remote experiences, this mechanism permits to accurately anticipate future task states and to regulate the engagement of control accordingly (Jiang et al., 2014, 2015). In the active control condition, we observed a clear signature of this regulation. Even though participants were uninstructed about the changes in conflict likelihood, and knowledge about these changes was not strictly necessary to perform the task correctly, performance clearly scaled with the validity of model-based conflict predictions.

Critically, TMS over the LPFC completely abolished model-based and GLM-based indices of behavioral adaptation. These disruptive effects were independent of the effector that was used for task implementation and specific to changes in control engagement (i.e., TMS did not induce generalized performance deficits). Both observations are consistent with the notion that prefrontal goal representations are abstract in nature and guide the information flow in brain regions that serve task-specific sensory and motor functions (Miller and Cohen, 2001). Previous
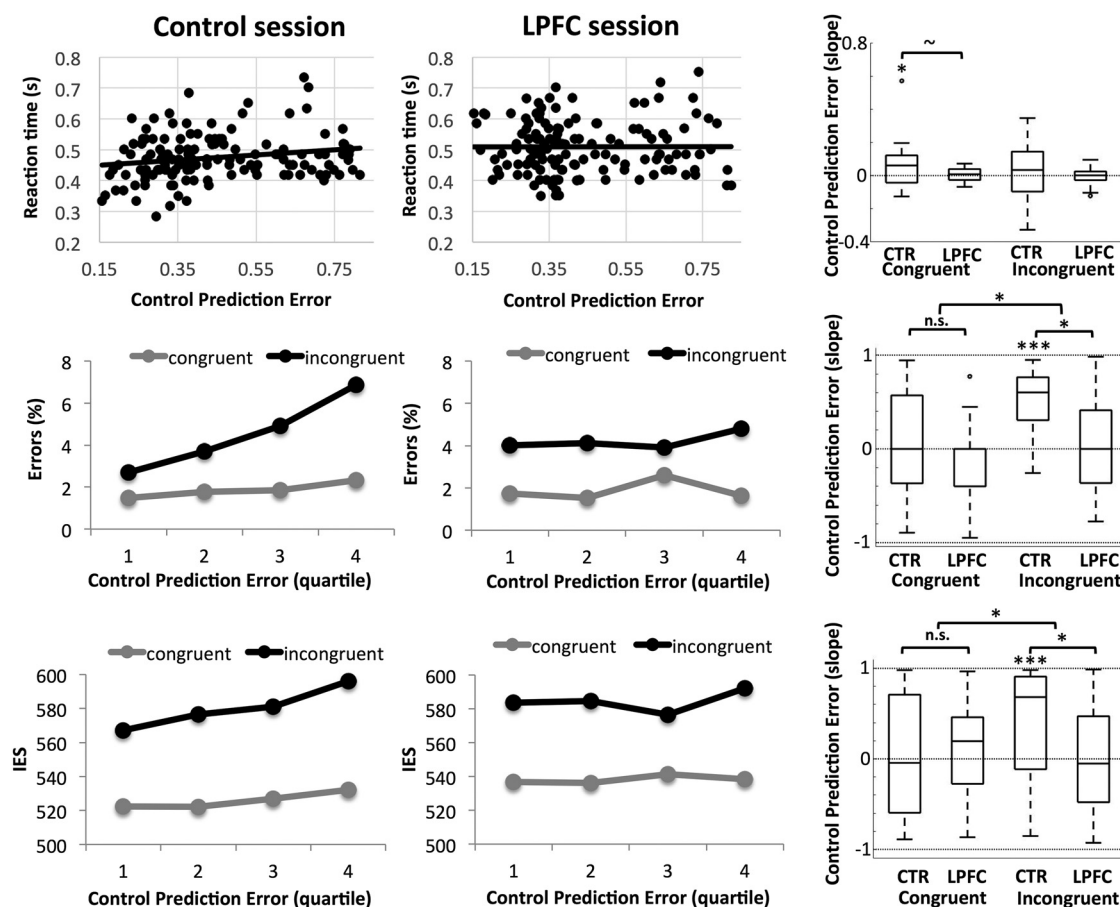
**Figure 5.** Model-based results illustrating performance as a function of CPE, TMS site (control site vs LPFC), and trial type (congruent vs incongruent). Results are displayed separately for each performance index. RT results are shown in the top panel, error rates in the middle panel, and IES in the bottom panel. RT scatter plots display the correlation between CPE and RT (performance on congruent trials of a representative participant) to convey that CPE estimates were obtained on a trial-by-trial basis. In contrast, error rates and IES are plotted as a function of CPE quartile because the computation of these scores required averaging over trials (see Materials and Methods section for details). Boxplots in the right panel display CPE slopes (reflecting the extent to which CPE quartile scores follow a linear trend) as a function of TMS target site (control site vs LPFC) and trial congruency (congruent vs incongruent). Solid lines indicate the median of the distributions, the box outline indicates the 25th and 75th percentile, and the whiskers indicate 1.5× the interquartile range. Extreme values are shown separately as unfilled circles. *$p < 0.05$, ***$p < 0.005$, n.s. not significant.

work has shown repeatedly that conflict adaptation is associated with elevated levels of LPFC activity on postconflict trials (Mac-Donald et al., 2000; Kerns et al., 2004; Egner and Hirsch, 2005) and with enhanced coupling of this region with areas implicated in processing task-relevant stimulus features (Egner and Hirsch, 2005, Morishima et al., 2009). Recent lesion and brain stimulation studies furthermore suggest that this engagement of the LPFC plays a causal role in short-term adaptation to conflict (Boschin et al., 2016, 2017; Gbadeyan et al., 2016). Our results extend the scope of these studies in two important ways. First, they demonstrate a more general role of the LPFC in behavioral adaptation across multiple temporal scales. Second, and more importantly, they provide a detailed and mechanistic framework about the learning signals that underpin multilevel adaptation and instigate the engagement of LPFC-based top-down control.

Such a learning-based perspective on cognitive control has considerable theoretical appeal because it situates the LPFC within a broader neural network that infers regularities in the external world to predict its future state (Vossel et al., 2014; Waskom et al., 2017). Prospective coding is thought to contribute to a variety of mental phenomena (Clark, 2013) and provides an effective means for the brain to establish contextually appropriate levels of task focus that support goal achievement, but minimize the costs

of control engagement (Shenhav et al., 2013, 2017; Amer et al., 2016). Interestingly, within this realm, different types of costs have been distinguished, most prominently intrinsic costs of the cognitive apparatus (e.g., a limited capacity to maintain task-relevant representations or potential metabolic costs of controlled information processing) and opportunity costs of control engagement (e.g., a risk of missing valuable information in the environment due to a selective focus on task-relevant information). An exciting avenue for future research will be to identify the brain mechanisms that support the calculation of these cost parameters and their integration for adjustments in control engagement (see also Shenhav et al., 2013). Clearly, this will require the design of more complex tasks and models, but we are convinced that such a normative, learning-based perspective will ultimately open the door toward a richer and ecologically more valid neuroscience of cognitive control.

### Limitations and future directions

It is worth noting that TMS-induced changes in behavioral adaptation were only of modest magnitude, yielding results at statistical threshold level. This is likely due, at least in part, to our use of long intertrial intervals, which have been shown to diminish behavioral adaptation (Egner et al., 2010), but were imperative in

our protocol to ensure the safe application of TMS (Rossi et al., 2009). This procedure resulted in relatively small adaptation effects in the control session, which also limited the magnitude of potential TMS outcomes. Future studies could allow for a faster trial pacing by delivering TMS only on a subset of all trials (Taylor et al., 2007a, 2007b), though we suspect that this approach could interfere with the learning of task statistics. Specifically, in such a setup, the presence versus absence of TMS would likely become a highly salient task feature that may override the registration of other aspects such as changing conflict likelihood. Beyond considerations about task design, future studies should aim to recruit even larger samples than ours to maximize the robustness of parameter estimation.

Another challenge for future research will be to dissociate the LPFC's role in the implementation of learning-guided top-down control more clearly from a potential role in the underlying learning processes. As noted above, a host of work has implicated the LPFC in control implementation, whereas the monitoring of control demand is typically associated with the medial frontal cortex and subcortical regions such as the thalamus and the striatum (Seifert et al., 2011; Ullsperger et al., 2014; Mansouri et al., 2017). Similarly, our model-based fMRI study associated the LPFC only with the translation of control prediction into behavioral adaptation, but not with the calculation of the underlying learning signals (Jiang et al., 2015). Nonetheless, our paradigm does not permit us to rule out that the TMS intervention also affected the acquisition of predictive knowledge rather than just its usage for behavioral control alone. Interleaving TMS and no-TMS trials could also provide a valuable solution to this challenge, for example, by comparing (short-term) adaptation on trials with and without TMS within the same task blocks. As noted above, however, introducing uncertainty about the application of TMS might cause side effects that diminish the effects of interest and could prove difficult to account for.

On a larger scale, we hope that our study will encourage the field to combine computational modeling and TMS more frequently as a new window to study the dynamic (dis-)engagement of brain regions during task performance. Similar to the productive coalition between modeling and neuroimaging (Forstmann et al., 2011), we believe that a field of model-based TMS could bear mutual benefits for neuroscientists and modelers alike. Whereas modelers would benefit from a causal method that enables to probe the veracity of model predictions with maximal rigor, neuroscientists would gain a toolbox to study dynamic trial-by-trial changes in structure–function relationships instead of treating those relationships as fixed and time stable. Beyond the formulation of more dynamic hypotheses, the prospect of comparing task conditions that are equated in terms of online processing demands, but differ in terms of model-based belief states, could also aid the interpretability of TMS data by dissociating cognitive TMS effects more clearly from nonspecific TMS outcomes (Robertson et al., 2003). In any case, such studies should aim for large sample sizes to account for the considerable variability that is typically observed with TMS effects.

## Conclusion
Our study shows that adaptive fluctuations in the attentional focus on goal-relevant information are captured by a learning mechanism that flexibly integrates recent and remote experiences of conflict between relevant and irrelevant inputs. This mechanism permits to predict a task's forthcoming demand and to engage cognitive control accordingly. TMS-induced perturbation of the LPFC abolished these fluctuations, suggesting a causal

role of this region in translating anticipated cognitive demand into optimal levels of focus. These findings provide causal evidence for the FCM as a mechanistic account for the regulation of cognitive control and emphasize that this flexibility is, at least in part, realized via the dynamic (dis-)engagement of LPFC-based top-down signals.

## References
Amer T, Campbell KL, Hasher L (2016) Cognitive control as a double-edged sword. Trends Cogn Sci 20:905–915. CrossRef Medline
Antal A, Nitsche MA, Kincses TZ, Lampe C, Paulus W (2004) No correlation between moving phosphene and motor thresholds: a transcranial magnetic stimulation study. Neuroreport 15:297–302. CrossRef Medline
Baldauf D, Desimone R (2014) Neural mechanisms of object-based attention. Science 344:424–427. CrossRef Medline
Behrens TEJ, Woolrich MW, Walton ME, Rushworth (2007) Learning the value of information in an uncertain world. Nat Neurosci 10:1214–1221.
Bocanegra BR, Hommel B (2014) When cognitive control is not adaptive. Psychol Sci 25:1249–1255. CrossRef Medline
Boschin EA, Brkic MM, Simons JS, Buckley MJ (2016) Distinct roles for the anterior cingulate and dorsolateral prefrontal cortices during conflict between abstract rules. Cereb Cortex 27:34–45. CrossRef Medline
Boschin EA, Mars RB, Buckley MJ (2017) Transcranial magnetic stimulation to dorsolateral prefrontal cortex affects conflict-induced behavioural adaptation in a Wisconsin Card Sorting Test analogue. Neuropsychologia 94:36–43. CrossRef Medline
Botvinick MM, Braver TS, Barch DM, Carter CS, Cohen JD (2001) Conflict monitoring and cognitive control. Psychol Rev 108:624–652. CrossRef Medline
Braver TS, Reynolds JR, Donaldson DI (2003) Neural mechanisms of transient and sustained cognitive control during task switching. Neuron 39:713–726. CrossRef Medline
Bugg JM, Crump MJ (2012) In support of a distinction between voluntary and stimulus-driven control: a review of the literature on proportion congruent effects. Front Psychol 3:367. CrossRef Medline
Carter CS, Macdonald AM, Botvinick M, Ross LL, Stenger VA, Noll D, Cohen JD (2000) Parsing executive processes: strategic vs evaluative funcitons of the anterior cingulate cortex. Proc Natl Acad Sci U S A 97:1944–1948. CrossRef Medline
Chiu YC, Jiang J, Egner T (2017) The caudate nucleus mediates learning of stimulus-control state associations. J Neurosci 37:1028–1038. CrossRef Medline
Clark A (2013) Whatever next? Predictive brains, situated agents, and the future of cognitive science. Behav Brain Sci 36:181–204. CrossRef Medline
Clerget E, Andres M, Olivier E (2013) Deficit in complex sequence processing after a virtual lesion of left BA45. PLoS One 8:e63722. CrossRef Medline
D'Ardenne K, Eshel N, Luka J, Lenartowicz A, Nystrom LE, Cohen JD (2012) Role of prefrontal cortex and the midbrain dopamine system in working memory updating. Proc Natl Acad Sci U S A 109:19900–19909. CrossRef Medline
De Baene W, Brass M (2013) Switch probability context (in)sensitivity within the cognitive control network. Neuroimage 77:207–214. CrossRef Medline
Derrfuss J, Brass M, Neumann J, von Cramon DY (2005) Involvement of the inferior frontal junction in cognitive control: meta-analyses of switching and Stroop studies. Hum Brain Mapp 25:22–34. CrossRef Medline
Derrfuss J, Brass M, von Cramon DY, Lohmann G, Amunts K (2009) Neural activations at the junction of the inferior frontal sulcus and the inferior precentral sulcus: interindividual variability, reliability, and association with sulcal morphology. Hum Brain Mapp 30:299–311. CrossRef Medline
Desimone R, Duncan J (1995) Neural mechanisms of selective visual attention. Annu Rev Neurosci 18:193–222. CrossRef Medline
Egner T (2014) Creatures of habit (and control): a multi-level learning perspective on the modulation of congruency effects. Front Psychol 5.
Egner T, Hirsch J (2005) Cognitive control mechanisms resolve conflict through cortical amplification of task-relevant information. Nat Neurosci 8:1784–1790. CrossRef Medline
Egner T, Ely S, Grinband J (2010) Going, going, gone: characterizing the time-course of congruency sequence effects. Front Psychol 1:154. CrossRef Medline
Forstmann BU, Wagenmakers EJ, Eichele T, Brown S, Serences JT (2011) Reciprocal relations between cognitive neuroscience and formal cognitive models: opposites attract? Trends Cogn Sci 15:272–279. CrossRef Medline

Freedman DJ, Riesenhuber M, Poggio T, Miller EK (2001) Categorical representation of visual stimuli in the primate prefrontal cortex. Science 291:312–316. CrossRef Medline

Gbadeyan O, McMahon K, Steinhauser M, Meinzer M (2016) Stimulation of dorsolateral prefrontal cortex enhances adaptive cognitive control: a high-definition transcranial direct current stimulation study. J Neurosci 36:12530–12536. CrossRef Medline

Gläscher J, Adolphs R, Damasio H, Bechara A, Rudrauf D, Calamia M, Paul LK, Tranel D (2012) Lesion mapping of cognitive control and value-based decision making in the prefrontal cortex. Proc Natl Acad Sci U S A 109:14681–14686. CrossRef Medline

Gratton G, Coles MG, Donchin E (1992) Optimizing the use of information: strategic control of activation of responses. J Exp Psychol Gen 121: 480–506. Medline

Jiang J, Heller K, Egner T (2014) Bayesian modeling of flexible cognitive control. Neurosci Biobehav Rev 46:30–43. CrossRef Medline

Jiang J, Beck J, Heller K, Egner T (2015) An insula-frontostriatal network mediates flexible cognitive control by adaptively predicting changing control demands. Nat Commun 6:8165. CrossRef Medline

Kerns JG, Cohen JD, MacDonald AW 3rd, Cho RY, Stenger VA, Carter CS (2004) Anterior cingulate conflict monitoring and adjustments in control. Science 303:1023–1026. CrossRef Medline

Lee TG, D'Esposito M (2012) The dynamic nature of top-down signals originating from prefrontal cortex: a combined fMRI-TMS study. J Neurosci 32:15458–15466. CrossRef Medline

MacDonald AW 3rd, Cohen JD, Stenger VA, Carter CS (2000) Dissociating the role of the dorsolateral prefrontal and anterior cingulate cortex in cognitive control. Science 288:1835–1838. CrossRef Medline

Mansouri FA, Egner T, Buckley MJ (2017) Monitoring demands for executive control: shared functions between human and nonhuman primates. Trends Neurosci 40:15–27. CrossRef Medline

Mayr U, Awh E, Laurey P (2003) Conflict adaptation effects in the absence of executive control. Nat Neurosci 6:450–452. CrossRef Medline

Miller EK, Cohen JD (2001) An integrative theory of prefrontal cortex function. Annu Rev Neurosci 24:167–202. CrossRef Medline

Morishima Y, Akaishi R, Yamada Y, Okuda J, Toma K, Sakai K (2009) Task-specific signal transmission from prefrontal cortex in visual selective attention. Nat Neurosci 12:85–91. CrossRef Medline

Muhle-Karbe PS, Andres M, Brass M (2014) Transcranial magnetic stimulation dissociates prefrontal and parietal contributions to task preparation. J Neurosci 34:12481–12489.

Rigoux L, Stephan KE, Friston KJ, Daunizeau J (2014) Bayesian model selection for group studies–revisited. Neuroimage 84:971–985. CrossRef Medline

Robertson EM, Théoret H, Pascual-Leone A (2003) Studies in cognition: the problems solved and created by transcranial magnetic stimulation. J Cogn Neurosci 15:948–960. CrossRef Medline

Rossi S, Hallett M, Rossini PM, Pascual-Leone A; Safety of TMS Consensus Group (2009) Safety, ethical considerations, and application guidelines

for the use of transcranial magnetic stimulation in clinical practice and research. Clin Neurophysiol 120:2008–2039. CrossRef Medline

Schroeter ML, Vogt B, Frisch S, Becker G, Barthel H, Mueller K, Villringer A, Sabri O (2012) Executive deficits are related to the inferior frontal junction in early dementia. Brain 135:201–215. CrossRef Medline

Schuck NW, Gaschler R, Wenke D, Heinzle J, Frensch PA, Haynes JD, Reverberi C (2015) Medial prefrontal cortex predicts internally driven strategy shifts. Neuron 86:331–340. CrossRef Medline

Seifert S, von Cramon DY, Imperati D, Tittgemeyer M, Ullsperger M (2011) Thalamocingulate interactions in performance monitoring. J Neurosci 31:3375–3383. CrossRef Medline

Shenhav A, Botvinick MM, Cohen JD (2013) The expected value of control: an integrative theory of anterior cingulate cortex function. Neuron 79: 217–240. CrossRef Medline

Shenhav A, Musslick S, Lieder F, Kool W, Griffiths TL, Cohen JD, Botvinick MM (2017) Toward a rational and mechanistic account of mental effort. Annu Rev Neurosci 40:99–124. CrossRef Medline

Stewart LM, Walsh V, Rothwell JC (2001) Motor and phosphene thresholds: a TMS correlation study. Neuropsychologia 39:415–419. CrossRef Medline

Stokes MG, Kusunoki M, Sigala N, Nili H, Gaffan D, Duncan J (2013) Dynamic coding for cognitive control in prefrontal cortex. Neuron 78:364–375. CrossRef Medline

Taylor PC, Nobre AC, Rushworth MF (2007a) FEF TMS affects visual cortical activity. Cereb Cortex 17:391–399. Medline

Taylor PC, Nobre AC, Rushworth MF (2007b) Subsecond changes in top down control exerted by human medial frontal cortex during conflict and action selection: a combined transcranial magnetic stimulation electroencephalography study. J Neurosci 27:11343–11353. CrossRef Medline

Torres-Quesada M, Funes MJ, Lupiáñez J (2013) Dissociating proportion congruent and conflict adaptation effects in a Simon-Stroop procedure. Acta Psychol (Amst) 142:203–210. CrossRef Medline

Townsend JT, Ashby FG (1983) Stochastic modeling of elementary psychological processes. Cambridge: Cambridge University.

Ullsperger M, Danielmeier C, Jocham G (2014) Neurophysiology of performance monitoring and adaptive behavior. Physiol Rev 94:35–79. CrossRef Medline

Vossel S, Mathys C, Daunizeau J, Bauer M, Driver J, Friston KJ, Stephan KE (2014) Spatial attention, precision, and Bayesian inference: a study of saccadic response speed. Cereb Cortex 24:1436–1450. CrossRef Medline

Waskom ML, Frank MC, Wagner AD (2017) Adaptive engagement of cognitive control in context-dependent decision making. Cereb Cortex 27: 1270–1284. CrossRef Medline

Zanto TP, Rubens MT, Bollinger J, Gazzaley A (2010) Top-down modulation of visual feature processing: the role of the inferior frontal junction. Neuroimage 53:736–745. CrossRef Medline

Zanto TP, Rubens MT, Thangavel A, Gazzaley A (2011) Causal role of the prefrontal cortex in top-down modulation of visual processing and working memory. Nat Neurosci 14:656–661. CrossRef Medline