

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Learning when effort matters: neural dynamics underlying updating and adaptation to changes in performance efficacy.

### Permalink

<https://escholarship.org/uc/item/2674k6c9>

### Journal

Cerebral Cortex, 33(5)

### Authors

Grahek, Ivan

Frömer, Romy

Prater Fahey, Mahalia

et al.

### Publication Date





2023-02-20

### DOI

10.1093/cercor/bhac215

Peer reviewed

# Learning when effort matters: neural dynamics underlying updating and adaptation to changes in performance efficacy

Ivan Grahek , Romy Frömer , Mahalia Prater Fahey , Amitai Shenhav 

Department of Cognitive, Linguistic, & Psychological Sciences, Carney Institute for Brain Science, Brown University, Box 1821, Providence, RI 02912, United States

\*Corresponding author: Department of Cognitive, Linguistic, & Psychological Sciences, Carney Institute for Brain Science, Brown University, Providence, RI 02912, United States. Email: [ivan\\_grahek@brown.edu](mailto:ivan_grahek@brown.edu)

To determine how much cognitive control to invest in a task, people need to consider whether exerting control matters for obtaining rewards. In particular, they need to account for the efficacy of their performance—the degree to which rewards are determined by performance or by independent factors. Yet it remains unclear how people learn about their performance efficacy in an environment. Here we combined computational modeling with measures of task performance and EEG, to provide a mechanistic account of how people (i) learn and update efficacy expectations in a changing environment and (ii) proactively adjust control allocation based on current efficacy expectations. Across 2 studies, subjects performed an incentivized cognitive control task while their performance efficacy (the likelihood that rewards are performance-contingent or random) varied over time. We show that people update their efficacy beliefs based on prediction errors—leveraging similar neural and computational substrates as those that underpin reward learning—and adjust how much control they allocate according to these beliefs. Using computational modeling, we show that these control adjustments reflect changes in information processing, rather than the speed–accuracy tradeoff. These findings demonstrate the neurocomputational mechanism through which people learn how worthwhile their cognitive control is.

**Key words:** motivation; cognitive control; learning; performance efficacy; EEG.

## Introduction

Cognitive control is critical for achieving most goals, but it is effortful (Botvinick and Cohen 2014; Shenhav et al. 2017). To decide how to invest control into a task (e.g. writing an essay for a competition), a person must weigh these effort costs against the potential benefits of a given type and amount of control (Manohar et al. 2015; Verguts et al. 2015; Kool and Botvinick 2018). One aspect of these benefits is the significance of the expected outcomes, both positive (e.g. a monetary prize, social acclaim) and negative (e.g. lost revenue, social derision) (Atkinson et al. 1966; Leng et al. 2021). An equally important aspect of the expected benefits of control is the extent to which control matters for achieving good outcomes and avoiding bad ones (Shenhav et al. 2021; Frömer, Lin, et al. 2021a). This can in turn be decomposed into the extent to which higher levels of control translate into better performance (e.g. whether writing a good essay will require substantial or only minimal control resources; control efficacy) and the extent to which better performance translates into better outcomes (e.g. whether prizes are determined by the strength of an essay or by arbitrary or even biased metrics unrelated to essay writing performance; performance efficacy). Whereas studies have increasingly characterized the ways in which control allocation is influenced by expected outcomes

(e.g. Parro et al. 2018; Leng et al. 2021) and the expected efficacy of control (e.g. as a function of task difficulty; Krebs et al. 2012; Vassena et al. 2014; Chiu and Egner 2019), much less is known about how people estimate and adjust to the perceived efficacy of their performance in a given environment.

We recently showed that when participants are explicitly instructed about how efficacious their performance will be on an upcoming trial, they exhibit behavioral and neural responses consistent with increased control (Frömer, Lin, et al. 2021a). We had participants perform a standard cognitive control task (Stroop) for potential monetary rewards, and we varied whether obtaining those rewards was contingent on performing well on the task (high performance efficacy) or whether those rewards were determined at random (low efficacy). We showed that people allocate more control when they expect to have high compared to low efficacy, reflected in higher amplitudes of an EEG index of proactive control (the contingent negative variation [CNV]) and in improved behavioral performance. These results demonstrate that participants leverage expectations about the extent to which their performance matters when deciding how much cognitive effort to invest in a task. However, in this work, participants were explicitly cued with the level of performance efficacy to expect on a given

trial, and those predictive cues retained the same meaning across the session. Thus, how it is that people learn these efficacy expectations in environments where contingencies are not instructed, and how they dynamically update their expectations as contingencies change, remains unanswered.

Outside of the domain of cognitive control, a relevant line of work has examined how people learn about the factors that determine future outcomes when selecting between potential courses of action. In particular, work in this area has shown that people are able to learn about and update their expectations of the likelihood that a given action will generate a given outcome (action–outcome contingency; Dickinson and Balleine 1995; Moscarello and Hartley 2017; Ly et al. 2019). People preferentially, and more vigorously, select actions that reliably lead to desired outcomes (i.e. the more contingent those outcomes are on the action in question; Liljeholm et al. 2011; Manohar et al. 2017), and work in both animals (Balleine and O’Doherty 2010) and humans (Norton and Liljeholm 2020; Dorfman et al. 2021; Ligneul et al. 2022; Morris et al. 2022) has helped to characterize the neural systems that support this process of learning and action selection. However, given the focus on discrete actions and their immediate relationship with outcomes, research into these action–outcome contingencies is unable to capture key aspects that are unique to selection of control states. Most notably, cognitive control signals (e.g. attention to one or more features of the environment) are multidimensional; not immediately observable by either the participant or experimenter; and their relationship with potential outcomes is intermediated by the many-to-many relationship between control states and task performance (Ritz et al. 2022). While there has been research into how people learn to adjust cognitive control signals based on changes in their task environment, here again work has focused on how people adapt to changes in outcomes (e.g. Otto and Daw 2019; Bustamante et al. 2021) and changes in the relationship between control and performance (with increasing task difficulty; Bugg et al. 2011; Nigbur et al. 2015; Bejjani et al. 2018; Jiang et al. 2020). The mechanisms by which people learn about the relationship between performance and outcomes (performance efficacy), and how they adjust their control allocation accordingly, remain largely unexplored.

Here, we seek to fill this gap by studying the mechanisms through which expectations of performance efficacy are formed, updated, and used to guide control allocation. To do so, we extend our previous approach (Frömer, Lin, et al. 2021a)—which studied how behavioral and neural correlates of control allocation vary when performance efficacy is explicitly cued—to examine how participants learn and adapt their control under conditions where efficacy was un-cued (having to instead be learned from feedback) and was gradually varied across a wide range of potential efficacy values over the course of the session. We use computational reinforcement

learning models to show that expected efficacy can be learned from feedback through iterative updating based on weighted prediction errors (Sutton and Barto 2018), and use model-based single-trial EEG analyses to show that these efficacy prediction errors modulate a canonical neural marker of reward-based learning and behavioral adjustment (Fischer and Ullsperger 2013). We further provide evidence that efficacy estimates learned in this way are used to guide the allocation of control—in our EEG study and a second behavioral study, participants tended to perform better when efficacy was higher. We also provide evidence that a neural marker of control allocation (Schevernels et al. 2014) tends to increase with increasing model-based efficacy estimates. Using a drift diffusion model (DDM; Ratcliff and McKoon 2008; Wiecki et al. 2013), in Study 2, we further show that the performance improvements related to increased performance efficacy reflect facilitation of task-related information processing (reflected in increased drift rates), rather than changes in the speed–accuracy tradeoff (i.e. thresholds). Taken together, these results show that efficacy estimates can be learned and updated based on feedback, leveraging general cognitive and neural mechanisms of predictive inference.

## Materials and methods

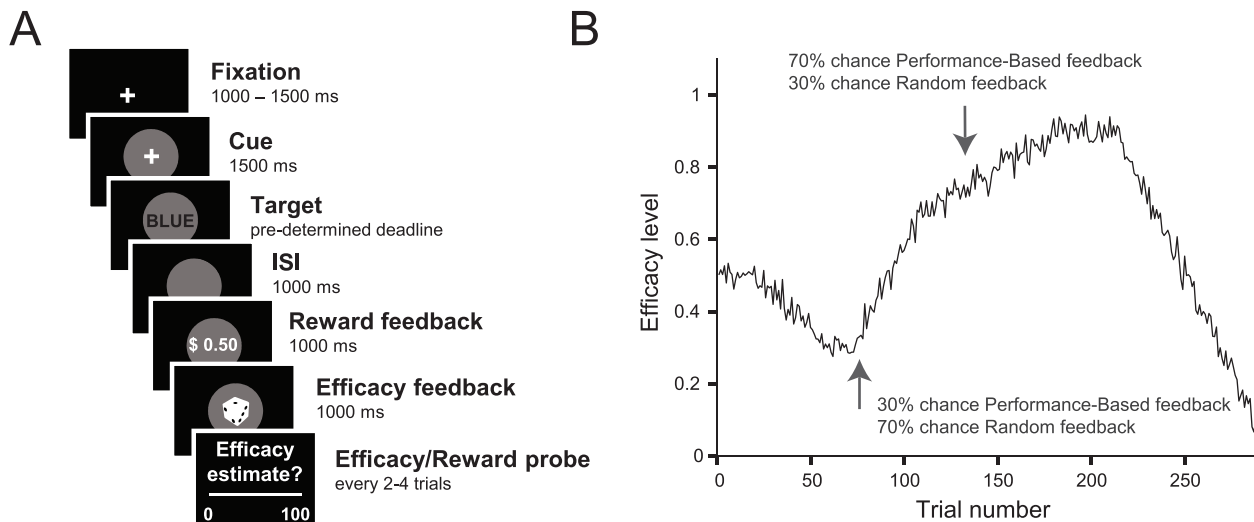
### Study 1

#### Participants

We recruited 41 participants with normal or corrected-to-normal vision from the Brown University subject pool. One participant was excluded due to technical issues. The final data set included 40 participants (24 females, 16 males; median age = 19). Participants gave informed consent and were compensated with course credits or a fixed payoff of \$20. In addition, they received up to \$5 bonus that depended on their task performance (\$3.25 on average). The research protocol was approved by Brown University’s Institutional Review Board.

#### Experimental design

In the main task, taking approximately 45 min, participants performed 288 Stroop trials (Fig. 1A). Each trial started with the presentation of a cue (gray circle) that remained on the screen throughout the trial. After a period of 1,500 ms, a Stroop stimulus was superimposed until a response was made or 1,000 ms elapsed, at which time it was sequentially replaced with 2 types of feedback presented for 1,000 ms each. Each trial onset was preceded by a fixation cross (randomly jittered between 1,000 and 1,500 ms). Participants responded to the ink color of the Stroop stimulus (red, green, blue, or yellow) by pressing one of 4 keyboard keys (D, F, J, and K). Stimuli were either color words same as the ink color (congruent,  $n = 108$ ), or different (incongruent,  $n = 108$ ), or a string of letters “XXXXX” (neutral,  $n = 72$ ). Feedback informed them whether they obtained a reward (reward, “\$50c” or no reward, “\$0c”) and whether the reward they



**Fig. 1.** Manipulating expected efficacy and assessing learning in Study 1. A) Trial Schematic. On each trial, participants saw a cue (gray circle), predicting the onset of a Stroop stimulus (target), and were then sequentially presented with reward and efficacy feedback. On half of the trials, efficacy feedback was presented first, and on the other half, reward feedback was presented first. Every 2–4 trials, participants were subsequently asked to estimate their current efficacy level (“How much do you think your rewards currently depend on your performance?”) or reward rate (“How often do you think you are currently being rewarded?”). B) Efficacy manipulation. We let the probability of performance-based versus random feedback continuously drift over the course of the experiment (inversed for one half of the sample). Arrows mark time points with low and high efficacy, respectively. When efficacy was low, rewards were more likely to be random, whereas when efficacy was high, rewards were more likely to be performance-based.

received depended on their performance (performance-based feedback, a button graphic) or not (random feedback, a dice graphic). In order to earn rewards in the performance-based case, participants had to be both accurate and respond within an individually calibrated response deadline (see details below). The order of the 2 types of feedback was pseudo-randomized with half of the trials showing reward feedback first and the other half efficacy feedback. Every 2–4 trials the feedback was followed by a probe of efficacy (“How much do you think your rewards currently depend on your performance?”) or reward rate (“How often do you think you are currently being rewarded?”) to which participants responded on a visual analog scale ranging from 0 to 100. The number and timing of the probes were randomized per subject resulting in a median of 45 efficacy probes ( $SD = 3.38$ ) and 47 reward probes ( $SD = 2.81$ ).

Efficacy (performance-based or random rewards) on each trial was sampled from a binomial distribution with probabilities ranging between 0.1 and 0.9 that drifted over the course of the experiment and were predetermined (Fig. 1B). In order to ensure that the performance-based and random trials did not differ in reward rate, reward feedback for the random trials was sampled from the moving window of the reward feedback of the previous 10 performance-based trials. At the beginning of the experiment, a window with 8 rewards and 2 no rewards was used to reflect the pre-calibrated reward rate (details below), and this moving window was then updated after every trial. Thus, reward rate was not experimentally manipulated in the experiment and remained constant. We confirmed that reward rate was matched across performance-based and random trials by comparing reward probability between these trial

types ( $b = 0.01$ ; 95% credible intervals [CrI]  $[-0.07, 0.10]$ ;  $P_{b>0} = 0.38$ ).

Prior to the main task, participants performed several practice phases of the Stroop task (~15 min). First, they practiced the mappings between colors and keyboard keys (80 trials). Then, they completed a short practice of the Stroop task with written feedback (“correct” or “incorrect”) on each trial (30 trials). Participants then completed 100 more of such trials during which we individually calibrated the reaction time deadline such that participants yielded approximately 80% reward rate. The reaction time calibration started from a fixed deadline of 750 ms and increased or decreased this threshold in order to ensure that participants earn rewards on 80% of trials (i.e. that they are both accurate and below the deadline). The deadline obtained in this way ( $M = 796$  ms;  $SD = 73$  ms) was used in the main experiment and was not further adjusted. In the final practice, phase participants were introduced to the 2 types of feedback which they would see in the main experiment (30 trials).

The experimental task was implemented in Psychophysics Toolbox (Brainard 1997; Pelli 1997; Kleiner et al. 2007) for Matlab (MathWorks Inc.) and presented on a 23-inch screen with a  $1,920 \times 1,080$  resolution. All of the stimuli were presented centrally while the participants were seated 80 cm away from the screen.

#### Psychophysiological recording and preprocessing

EEG data were recorded at a sampling rate of 500 Hz from 64 Ag/AgCl electrodes mounted in an electrode cap (ECI Inc.), referenced against Cz, using Brain Vision Recorder (Brain Products, München, Germany). Vertical and horizontal ocular activity was recorded from below both eyes (IO1, IO2) and the outer canthi (LO1, LO2), respectively.

Impedances were kept below 10 k $\Omega$ . Offline, data were processed using custom-made Matlab scripts (Frömer et al. 2018) employing EEGLab functions (Delorme and Makeig 2004). Data were re-referenced to average reference and ocular artifacts were corrected using brain electric source analyses (Ille et al. 2002) based on separately recorded prototypical eye movements. The cleaned continuous EEG was then low pass filtered at 40 Hz and segmented into epochs around cue onset (–200 to 1,500 ms), stimulus onset, and both efficacy and reward feedback (–200 to 800 ms). Baselines were corrected to the average of each 200 ms pre-stimulus interval. Segments containing artifacts, values exceeding  $\pm 150 \mu\text{V}$  or gradients larger than  $50 \mu\text{V}$ , were excluded from further analyses.

Single trial event-related potentials (ERPs) were then exported for further analyses in R (R Core Team 2020). The late CNV was quantified between 1,000 and 1,500 ms post neutral cue onset (Schevernels et al. 2014; Frömer et al. 2016; Frömer, Lin, et al. 2021a) as the average activity over 9 fronto-central electrodes (Fz, F1, F2, FCz, FC1, FC2, Cz, C1, and C2). The P3b was quantified between 350 and 500 ms (Fischer and Ullsperger 2013) for both reward and efficacy feedback and calculated as the averaged activity over 9 centroparietal electrodes (Pz, P1, P2, POz, PO1, PO2, CPz, CP1, and CP2).

## Learning models and statistical analyses

### Learning models

Participants provided their subjective estimates of efficacy and reward every 4–8 trials (a total of 45 estimates), and we sought to fit a learning model to these estimates to be able to predict trial-by-trial adjustments in performance and neural markers of learning and cognitive control allocation. In order to obtain trial-by-trial estimates of efficacy and reward rate, we fitted 2 temporal difference learning models (Gläscher et al. 2010; Sutton and Barto 2018) to the continuous subjective estimates of efficacy and reward rate (Rutledge et al. 2014; Eldar et al. 2016; Nagase et al. 2018). The first model (1 learning rate efficacy model) assumed that the estimate of efficacy for the next trial ( $E_{t+1}$ ) depended on the current efficacy estimate ( $E_t$ ) and the prediction error ( $\delta_t$ ) weighted by a constant learning rate ( $\alpha$ ):

$$E_{t+1} = E_t + \alpha * \delta_t$$

Where  $0 \leq \alpha \leq 1$ , and the prediction error is calculated as the difference between the contingency feedback on the current trial ( $e_t$ ) and the efficacy estimate on that trial:  $\delta_t = e_t - E_t$ . The model started from an initial value (free parameter) and updated the model-based efficacy estimate based on the binary efficacy feedback on each trial. For example, assuming a learning rate of 0.5 and the initial value of 0.5, the model would update the initial estimate following efficacy feedback signaling performance-based ( $e_t = 1$ ) to 0.75. If on the next trial

contingency feedback was random ( $e_{t+1} = 0$ ), the model-based efficacy estimate would drop to 0.6. The model was fitted separately to the subjective estimates of efficacy with only the learning rate as a free parameter. The second model (2 learning rates efficacy model) was the same as the first model, but it included 2 learning rates: one learning rate for learning from the performance-based feedback, and another for learning from the random feedback. Finally, as a baseline, we also included the intercept model, which did not update the efficacy estimate throughout the experiment but just assumed that the estimate took one constant value. Importantly, the same models were fitted to obtain the model-based estimates of reward on each trial (1 learning rate reward model and the 2 learning rate reward model). These models were fitted using trial-by-trial reward feedback and the subjective estimates of reward. The models were fit hierarchically to the data using maximum likelihood estimation (using mfit (<https://github.com/sjgershm/mfit>)). To calculate the likelihood of each data point, model-based estimates (0–1 range) were compared to the subjective efficacy estimates (range normalized to 0–1 range for each participant). Likelihood was evaluated on trials which included a subjective estimate, as the likelihood that the difference between the model-based and the empirical estimate comes from a Gaussian distribution centered on 0 with a variance that was fitted as a free parameter for each subject. This variance parameter served as the noise in the estimates. Likelihoods were log transformed, summed, and then maximized using the fmincon function in MATLAB.

We performed a parameter recovery study to show that the most complex model (the 2 learning rates model) can be successfully recovered. We simulated a dataset with the same number of trials and subjective efficacy or reward probes as in the actual experiment. We used the efficacy drifts presented to the actual subjects (half of the simulated subjects saw one drift, and half its inverse), and we used the reward feedback sequences of 2 actual subjects from our experiment. We simulated 200 agents which learned both efficacy and reward with the noise parameter fixed to 0.2, intercept fixed to 0.5, and the positive and negative learning rates sampled from a uniform distribution ranging from 0.001 to 0.5. These parameters were matched based on the range of values obtained from the empirical fits to our data. As shown in [Supplementary Fig. S1](#), we were able to very reliably recover the simulated parameters for both efficacy and reward rate learning.

### Statistical analyses

The efficacy and reward rate estimates obtained through fitting the learning model were then used to analyze the behavioral and EEG data. To this end, we fitted Bayesian multilevel regressions to predict subjective estimates of efficacy and reward rates, reaction times, accuracy, as well as the CNV and P3b amplitudes. Subjective estimates of efficacy and reward rate were regressed onto

efficacy or reward feedback. Reaction times and accuracies were regressed onto trial-by-trial model-based estimates of efficacy and reward rate, as well as trial-by-trial CNV amplitude, while controlling for congruency. The P3b amplitudes were analyzed in 2 ways: with trial-by-trial model-based estimates of efficacy and reward rate and current feedback as predictors, and with model-based prediction errors and learning rates for each feedback type. CNV amplitudes were regressed onto trial-by-trial model-based estimates of efficacy. All of the fitted models controlled for the influence of the reward rate estimates. Parallel analyses were done to predict the P3b in response to reward feedback, while controlling for the efficacy estimates.

The regression models were fitted in R with the *brms* package (Bürkner 2016), which relies on the probabilistic programming language *Stan* (Carpenter et al. 2016) to implement Markov Chain Monte Carlo (MCMC) algorithms and estimate posterior distributions of model parameters. The analyses were done based on the recommendations for Bayesian multilevel modeling using *brms* (Bürkner 2016; Bürkner 2017; Nalborczyk and Bürkner 2019). The fitted models included constant and varying effects (also known as fixed and random) with weakly informative priors (except for the behavioral and CNV analyses, see below for details) for the intercept and the slopes of fixed effects and the likelihoods appropriate for the modeled data (Ex-Gaussian for reaction times, Bernoulli for accuracy, and Gaussian for the subjective estimates and the EEG amplitudes). The fitted models included all of the fixed effects as varying effects. All of the continuous predictors in the model were centered and the categorical predictors were contrast coded. Four MCMC simulations (“chains”; 20,000 iterations; 19,000 warmup) were run to estimate the parameters of each of the fitted models. The convergence of the models was confirmed by examining trace plots, autocorrelation, and variance between chains (Gelman and Rubin 1992). After convergence was confirmed, we analyzed the posterior distributions of the parameters of interest and summarized them by reporting the means of the distribution for the given parameter ( $b$ ) and the 95% CrI of the posterior distributions of that model. We report the proportion of the posterior samples on the relevant side of 0 (e.g.  $P_{b < 0} = 0.9$ ), which represents the probability that the estimate is below or above 0. We also report Bayes factors (BFs) calculated using the Savage–Dickey method (Wagenmakers et al. 2010). We report the BFs in support of the alternative hypothesis against the null ( $BF_{10}$ ), except for the analyses of accurate RT, accuracies, and CNV amplitude in which we have informative priors based on our previous study (Frömer, Lin, et al. 2021a), and in which case we support the evidence in favor of the null ( $BF_{01}$ ).

To compare the positive and negative learning rates, we fitted a model in which the learning rates were predicted by the learning rate type (Kruschke 2013). In this model, we used Gaussian distributions (mean,

standard deviation) as priors (intercept: (0.5,0.5); slopes: (0,0.5)).

We fitted 2 separate models to predict the subjective estimates of efficacy and reward rate based on previous feedbacks. At each timepoint, the estimates were predicted by the current, and previous 4 feedbacks. The feedback on each of the trials (performance-based vs. random or reward vs. no reward) was entered as a constant effect and the models also included the intercept as a varying effect. As the subjective estimates could vary between 0 and 1, we used Gaussian distributions (intercept: (0.5,0.2); slopes: (0,0.2)) as priors.

For predicting the P3b amplitude in response to the onset of the efficacy feedback, we fitted 2 models. First, we fitted a model which included the model-based estimate of efficacy (prior to observing the current feedback), the actual feedback, and the interaction between the expected efficacy and the observed efficacy feedback. Additionally, we controlled for the reward rate estimate. Second, we fitted separate models which included the model-based prediction errors, the influence of the between-subject learning rates, and their interaction with the prediction errors, while controlling for the estimate of the reward rate. In this analysis, the learning rates (one for feedback type for each subject) were mean-centered within subjects, and thus, any effect of the learning rates is driven by the difference between the random and the performance-based learning rate. For these models, we selected wide Gaussian priors (intercept: (5,3); slopes (0,3)). The same logic in building models was applied for the analyses of the reward feedback. In these analyses, we focused on the reward feedback processing and how it interacted with the model-based estimates of reward rates, while controlling for the model-based estimates of efficacy. We analyzed only the trials with correct responses for both the efficacy and the reward feedback analyses.

To test the influence of efficacy on the late CNV, we fitted a model which predicted the CNV based on the model-based efficacy estimates, while controlling for the effect of the reward rate estimates. Drawing on the results of our previous study (Frömer, Lin, et al. 2021a), this model included Gaussian priors for the intercept (−0.16, 1.59) and the efficacy (−0.30, 0.73) and reward (0,0.73) slopes.

For predicting reaction times and accuracy, we fitted models which included congruency (Facilitation: difference between neutral and congruent trials; Interference: difference between incongruent and neutral trials) and the model-based efficacy estimates, while controlling for the reward rate estimates. We used Gaussian distributions as informative priors based on our previous study (Frömer, Lin, et al. 2021a), for both the reaction times (intercept: (624, 58.69); facilitation (15.54, 21.93), interference (61.12, 37.49); efficacy (−10.26, 19.51); reward (0, 19.51)) and accuracy (Note that the prior distributions are set in log-odds.) analyses (intercept: (2.11, 0.81);

facilitation ( $-0.45, 0.64$ ), interference ( $-0.53, 0.81$ ); efficacy ( $0.09, 0.32$ ); reward ( $0, 0.32$ )).

To investigate how the late CNV influences the behavior, we fitted 2 models in which we predicted the reaction times and accuracy based on the CNV amplitude. The prior distributions for these models were weakly informative Gaussian distributions for predicting both the reaction times (intercept:  $(650, 200)$ ; slope:  $(0, 50)$ ) and accuracy (intercept:  $(0.7, 0.2)$ ; slope:  $(0, 0.2)$ ).

To visualize the topographies of the relevant ERP effects, we fitted the relevant models to all 64 channels and then plotted the posterior estimates of the effects of interest at each electrode (cf. Frömer, Nassar, et al. 2021b).

## Study 2

### Participants

We recruited 87 participants residing in the United States from Prolific—an online platform for data collection. Participants had normal or corrected-to-normal vision and gave informed consents. They were compensated with a fixed payoff of \$8 per hour (median completion time of 74 minutes) plus a monetary bonus based on points earned during the task (\$1 on average). The research protocol was approved by Brown University's Institutional Review Board.

We a priori excluded participants who did not pass attention checks ( $N=8$ ) or who took substantially longer than the average participant to complete the study ( $N=2$  participants who took over 130 min), suggesting that they did not sustain attention to the experiment over that time. We fit our learning models to data from the remaining 77 subjects and then excluded participants whose performance suggested inattention to the overall task (based on accuracies  $<70\%$  across all trials—including the trials in which performance efficacy was low,  $N=6$ ) or inattention to the task feedbacks and efficacy probes (based on low learning rates ( $N=19$ ), and one subject with no variance in responses to reward probes). To identify participants with exceedingly low learning rates, we submitted all positive and negative efficacy learning rates to unsupervised Gaussian mixture models (as implemented in the Mclust package; Scrucca et al. 2016) to determine the best fitting number and shape of clusters (model comparison via Bayesian information criterion (BIC)). This procedure identified 4 clusters of subjects with different overall learning rates (Supplementary Fig. S2B and C), and we excluded subjects from the first cluster as they all had very low learning rates relative to the other participants (both learning rates  $<0.03$ ). The subjects excluded based on low learning rates were most likely not paying attention to efficacy feedback or were always giving very similar responses to the efficacy probes (Supplementary Fig. S3). The final sample included 51 participants (31 females, 20 males; median age = 29).

### Experimental design

In order to better understand the computational mechanisms that lead to improved behavioral performance in

high efficacy states (Study 1), we wanted to fit a DDM (Ratcliff and McKoon 2008) to our behavioral data. If people allocate more attention when they think they have high performance efficacy, this should be observed as an increase in the drift rate (speed of evidence accumulation). However, Study 1 included a tight respond deadline for earning a reward, making it more challenging to fit the DDM. To avoid this issue, in Study 2, we adjusted the task to a free response paradigm which allowed us to investigate drift rate and threshold adjustments (cf. Leng et al. 2021). We used this design to test the hypothesis that higher efficacy estimates should predict increased drift rates.

Instead of single trials, participants now completed 288 intervals during which they could respond to as many trials (congruent and incongruent, removing the neutral condition) as they wished within a fixed time window (randomly selected between 2,000, 3,000, or 4,000 ms). Apart from this, the structure of the task remained the same: participants saw a fixation cross (1,000, 1,500, or 2,000 ms), then completed as many trials as they wished during a fixed interval, followed by the feedback (1,500 ms) on how many points they earned and whether this was based on their performance or awarded to them based on random chance. Note that participants now received continuous reward feedback (10 points per correct response instead of the binary reward-no reward in Study 1). For example, if participants completed 4 trials correctly and 2 trials incorrectly within a performance-based interval they would receive 40 points and see feedback informing them that the points were based on their performance. To determine the number of points on random intervals, the same yoking procedure as in Study 1 was employed, ensuring that the amount of reward was matched between performance-based and random intervals (reward amounts on random intervals were sampled from the moving average window of the past 10 performance-based intervals). We confirmed that the yoking procedure was successful by comparing the reward amounts on the 2 interval types ( $b = 0.00$ ; 95% CrI  $[-0.00, 0.02]$ ;  $P_{b>0} = 0.16$ ). As in Study 1, participants were probed every 2–4 intervals to estimate either how much they thought their rewards depended on their performance, or how often they were rewarded. We again implemented an efficacy drift (modified, but comparable to the drift in Study 1; Supplementary Fig. S2A), now across the 288 intervals of the task.

We gamified the task in order to make it more appealing for the participants. Instead of the Stroop task, we used a picture-word interference task in which 4 gray-scaled images of fruit (apple, pear, lemon, and peach) were overlaid by those fruit words. Participants responded to the image, while ignoring the word, by pressing one of the 4 corresponding keys. They first practiced this task and then were introduced with a cover story telling them that they are in the garden and need to water the fruit in little patches by pressing the correct keys. They were instructed that they will be moving through a garden and that in some patches watering

will directly translate into how many points they will be earning, while in the others that will not be the case (the efficacy drift). The experiment was implemented in Psiturk (Gureckis et al. 2016) and the participants performed the task on their own computers and were required to have a keyboard.

### Learning models, DDM, and statistical analyses

#### Learning models

We fitted the same set of learning models as in Study 1, performed model comparison, and got the interval-by-interval model-based estimates of performance efficacy and reward. Note that in this version of the task, participants earned points in each interval, unlike the binary rewards (reward vs. no reward) in Study 1. This meant that the reward learning model learned reward magnitudes rather than reward probabilities. However, model fitting and the further analyses were the same as in Study 1.

#### Statistical analyses

We fitted the same Bayesian multilevel models as in Study 1 to predict the influence of previous efficacy feedbacks on the subjective efficacy estimates, as well as the influence of the model-based efficacy estimates on reaction times and accuracies. For the analyses of the subjective efficacy estimates, we used the same priors as in Study 1. For the analyses of the reaction times and accuracies, we used the posterior distributions obtained in Study 1 as the informative priors for the congruency, efficacy, and reward effects. The reaction times and accuracy models also controlled for the effects of the average congruency level in the interval and the interval length.

#### Drift diffusion model

This model decomposes participant's behavior into drift rate (the speed of evidence accumulation) and response threshold (the level of caution), allowing us to investigate which of these 2 components is affected by the efficacy estimates. We fitted the model using Bayesian hierarchical estimation as implemented in the HDDM package (Wiecki et al. 2013). The fitted model included the main effects of efficacy and reward rate estimates onto both drift rate and threshold and included the effect of congruency on the drift rate. The responses were coded as correct or incorrect, and trials with reaction times below 250 ms were excluded. All of the effects were allowed to vary across subjects, and we ran 5 MCMC chains (12,000 iterations; 10,000 warmup). We confirmed convergence by examining trace plots and variance between chains.

## Results

### Study 1

To investigate how efficacy estimates are learned, and how they affect control allocation, in Study 1 we recorded EEG while 40 participants performed a modified version of the Stroop task (Fig. 1A). Across trials, we varied whether reward outcomes (\$0.50 vs. \$0.00) were

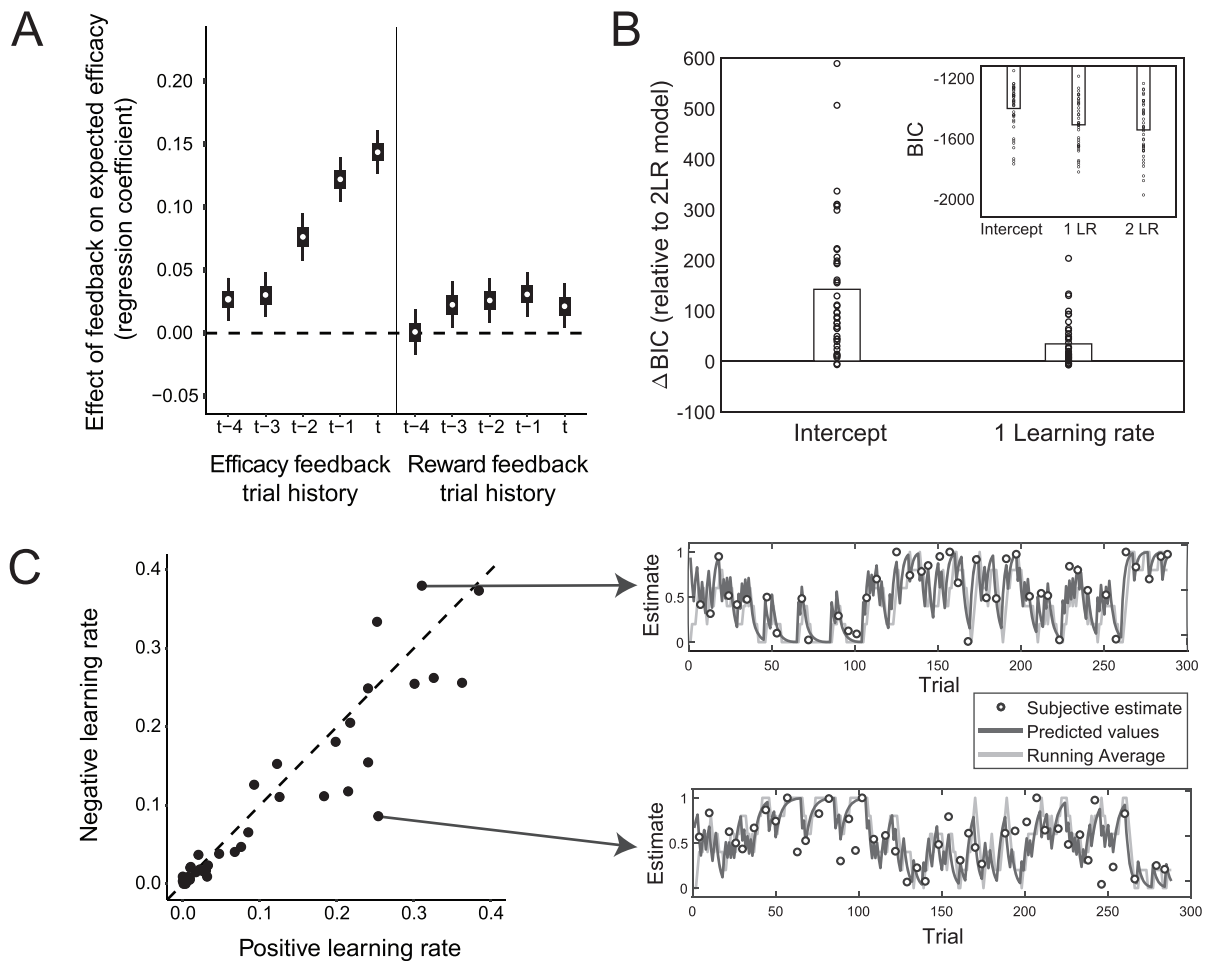
determined by a participant's performance on a given trial (responding accurately and below a pre-determined response time criterion; i.e. performance-based trials) or whether those outcomes were determined independent of performance (based on a weighted coin flip; i.e. random trials). Over the course of the session, we gradually varied the likelihood that a given trial would be performance-based or random such that, at some points in the experiment, most trials were performance-based (high efficacy level), and at other points, most trials had random outcomes (low efficacy level) (Fig. 1B). Importantly, unlike in our previous study (Frömer, Lin, et al. 2021a), participants were not told whether a given trial would be performance-based (maximal efficacy) or random (minimal efficacy), but instead had to estimate their current efficacy level based on recent trial feedback, which indicated both the reward outcome (\$0.50 vs. \$0.00) and how that outcome was determined (performance-based [button symbol] vs. random [dice symbol]). We held reward rate constant across both feedback types by yoking reward rate on random trials to the reward rate on performance-based trials and counter-balanced the time-course of the gradual change in efficacy (see Methods for details). To capture changes in efficacy expectations over the course of the session, we probed participants every 4–8 trials (averaging 44.3 probes per participant) to report their current estimates of efficacy. These efficacy probes were interleaved with probes asking participants to estimate the current reward rate, serving as foils and for control analyses.

#### Participants dynamically update efficacy expectations based on feedback

To determine whether and how participants learned their current efficacy levels, we first analyzed the influence of previous efficacy feedback (whether outcomes on a given previous trial had been performance-based or random) on one's current subjective estimates of efficacy. We found very strong evidence that participants adjusted their subjective efficacy upward or downward depending on whether the most recent trial was performance-based or random ( $b = 0.14$ ; 95% CrI [0.12, 0.16];  $P_{b < 0} = 0$ ;  $BF_{10} > 100$ ) and that this remained true (but to diminishing degrees) up to 5 trials back (all  $P_{b < 0} < 0.01$ ; Fig. 2A; Supplementary Table S1). This effect of feedback on efficacy estimates was present only for the efficacy feedback, while reward feedback did not predict the subjective estimates of efficacy (Fig. 2A). These results suggest that participants were dynamically updating their efficacy estimates based on efficacy feedback.

This pattern of learning was accounted for by a standard reinforcement learning (RL) algorithm, the delta-rule, according to which efficacy estimates are updated in proportion to the prediction error between the expected efficacy and the efficacy feedback on a given trial (i.e. whether a given outcome was determined by performance or not) (Sutton and Barto 2018).





**Fig. 2.** Efficacy learning is captured by a reinforcement learning model with separate learning rates for performance-based and random feedback in Study 1. A) Efficacy estimates track recent efficacy feedback. Regression weights for the influence of the current ( $t$ ) and previous contingent versus random feedback, as well as reward feedback, on the subjective efficacy estimate. Error bars represent 50% and 95% highest density intervals. B) A separate learning rate model captures efficacy learning best. BICs of fitted intercept-only model and one learning rate model relative to two learning rate models are plotted for each participant and favor the two learning rate model. C) Learning rate biases vary between participants. Positive and negative learning rate estimates are plotted for all participants (left). Points below the diagonal indicate higher learning rates for performance-based compared to random feedback, and points above the opposite. Example learning trajectories for two participants (right). Subjective and model-based efficacy estimates, and a running average of the previous 5 efficacy feedbacks, for a participant with a higher learning rate for the contingent compared to random feedback (upper) and a participant with the reverse bias (lower).

Interestingly, consistent with studies of reward learning (Niv et al. 2012; Collins and Frank 2014; Lefebvre et al. 2017; Chambon et al. 2020; Garrett and Daw 2020), we found that the RL model that best accounted for our data allowed efficacy estimates to be updated differently from trials that were more efficacious than expected (Performance-Based feedback) than from trials that were less efficacious than expected (Random feedback), by having separate learning rates scaling prediction errors in the 2 cases. Even when accounting for model complexity, this Two Learning Rate Efficacy model outperformed a One Learning Rate Efficacy model as well as a baseline model that fits a single constant efficacy estimate and no learning (Intercept Model) (Fig. 2B). We were able to successfully recover the parameters of this model from a simulated dataset matched to our data (see Methods and Supplementary Fig. S1). We found that the 2 learning rates for this best-fit model varied across the group (Fig. 2C) but did not

find that one was consistently higher than the other ( $b = 0.02$ ; 95% CrI  $[-0.04, 0.08]$ ;  $P_{b < 0} = 0.260$ ;  $BF_{01} = 13.55$ ; Supplementary Fig. S4), despite most participants (80%) tending to learn more from performance-based than random trials. Finally, model-based efficacy estimates were strongly related to the raw subjective estimates on trials on which participants reported efficacy ( $b = 0.77$ ; 95% CrI  $[0.62, 0.91]$ ;  $P_{b < 0} < 0.001$ ;  $BF_{01} > 100$ ;  $R^2 = 0.50$ ), demonstrating that the model successfully captured the raw estimates. Taken together, these results show that participants dynamically updated their expected efficacy based on trial-by-trial feedback and that they did so differentially based on whether the trial was more or less efficacious than expected. The fitted models further enable us to generate trial-by-trial estimates of expected efficacy and efficacy prediction errors, which we use in model-based analyses of behavior and neural activity below. Note that in all the following analyses, we control for reward estimates

obtained from models fit to reward feedback (for details see below).

### The feedback-related P3b indexes updating of efficacy expectations

To investigate the neural mechanism underlying feedback-based learning of efficacy, we probed the centroparietal P3b ERP component (Fig. 3A), an established neural correlate of prediction-error based learning (Fischer and Ullsperger 2013; Nassar et al. 2019). If the P3b indexes feedback-based updating of efficacy predictions, as it does for reward predictions, we would expect this ERP to demonstrate several key markers of such a learning signal. First, we would expect the P3b to reflect the extent to which efficacy feedback (performance-based vs. random) deviates from the current level of expected efficacy. In other words, the P3b should track the magnitude of the unsigned efficacy prediction error (PE) on a given trial. We tested this by examining how the amplitude of the P3b to a given type of efficacy feedback varied with model-based estimates of the participant's efficacy expectation on that trial, while holding the expected reward rate constant (see Section 2 for the details of the experimental design and the statistical models). If the P3b signaled efficacy PE, then its amplitude should scale inversely with the expected probability of a given type of feedback (i.e. how *unexpected* that feedback is), thus correlating negatively with expected efficacy on trials providing performance-based feedback and correlating positively with expected efficacy on trials providing random feedback. In addition to overall higher P3b to performance-based compared to random feedback ( $b=0.86$ ; 95% CrI [0.42, 1.31];  $P_{b<0}=0$ ;  $BF_{10}=30.86$ ; Fig. 3B), we found exactly this predicted interaction between feedback type and expected efficacy ( $b=-2.40$ ; 95% CrI [-4.07, -0.74];  $P_{b>0}=0$ ;  $BF_{10}=24.40$ ; Fig. 3B; Supplementary Fig. S5A; Supplementary Table S2), with the P3b amplitude decreasing with model-based estimates of expected efficacy on performance-based trials ( $b=-1.49$ ; 95% CrI [-2.72, -0.27];  $P_{b>0}=0.03$ ;  $BF_{10}=1.52$ ) and increasing numerically, but not robustly with expected efficacy on random trials ( $b=0.91$ ; 95% CrI [-0.34, 2.21];  $P_{b<0}=0.12$ ;  $BF_{10}=0.44$ ). Accordingly, when we regressed P3b amplitude on our model-based estimates of trial-to-trial efficacy PE, we found a positive relationship ( $b=1.25$ ; 95% CrI [0.35, 2.15];  $P_{b<0}=0.01$ ;  $BF_{10}=5.84$ ; Fig. 3C-left; Supplementary Table S2).

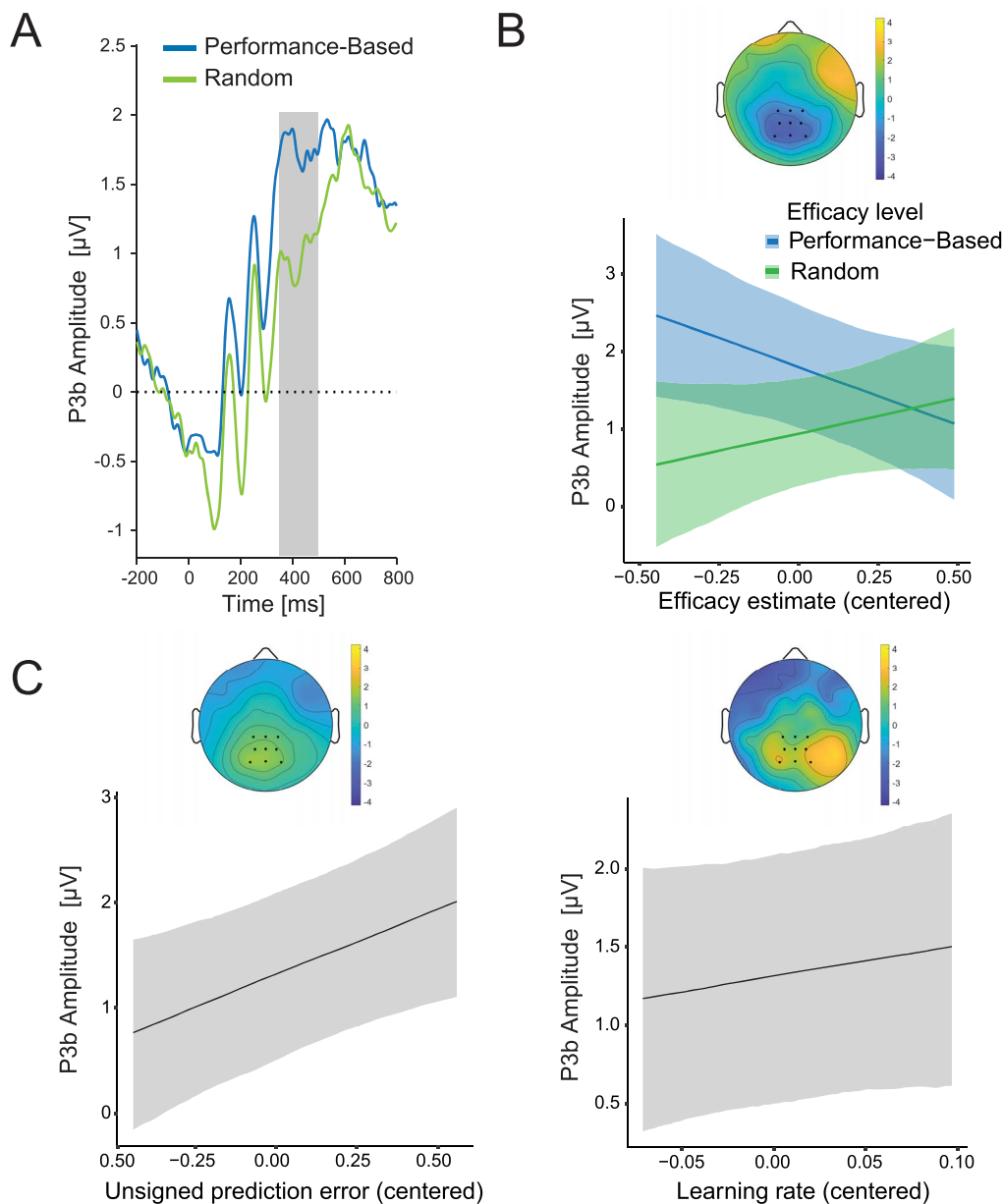
In addition to tracking efficacy PEs, the second key prediction for the P3b if it indexes efficacy learning is that it should track the extent to which PEs are used to update estimates of efficacy (i.e. the learning rate). In the current study, we found that participants differed in their learning rates for the 2 forms of efficacy feedback (performance-based vs. random), providing us with a unique opportunity to perform a within-subject test of whether the P3b tracked differences in learning rate across these 2 conditions. Specifically, we could

test whether a given subject's P3b was greater in the feedback condition for which they demonstrated a higher learning rate. We found only suggestive evidence for such an effect, with the P3b numerically but not robustly higher for the within-subject feedback condition with the higher learning rate ( $b=2.00$ ; 95% CrI [-2.21, 6.04];  $P_{b<0}=0.17$ ;  $BF=1.08$ ; Fig. 3C-right). While, theoretically, prediction error and learning rate might also interact in predicting the P3b amplitude, we did not observe such an interaction here. This finding is in line with previous work on reward processing (Fischer and Ullsperger 2013), which has found additive effects of prediction errors and learning rate on P3b.

### The CNV indexes control allocation based on updated expectations of efficacy

Thus far, our findings suggest that participants dynamically updated expectations of their performance efficacy based on feedback, and that the P3b played a role in prediction error-based updating of these expectations. Next, we tested the prediction that learned efficacy estimates determine the expected benefits of control and thus influence how much control is allocated (Shenhav et al. 2013). We have previously shown that people exert more control when expecting their performance to be more rather than less efficacious on the upcoming trial (Frömer, Lin, et al. 2021a). This was reflected in better behavioral performance and higher amplitudes of the CNV (Fig. 4B-left)—a slow negative fronto-central waveform preceding target onset, which is increasingly negative as the amount of control allocated in preparation for the task increases (Grent-'t-Jong and Woldorff 2007; Schevernels et al. 2014). Here, we used the same marker of control, but, unlike in our previous study, efficacy expectations were (i) learned rather than explicitly instructed; (ii) varied over time rather than having a fixed association with a cue; and (iii) varied continuously across the range of efficacy probabilities rather than alternating between maximal and minimal efficacy. We were therefore interested in testing whether these dynamically varying expectations of efficacy, as estimated by our model, would still exert the same influence on behavior and neural activity.

Consistent with our predictions and our previous findings, participants tended to perform better when they expected performance to be more efficacious, responding faster on correct trials with increasing model-based estimates of efficacy (Fig. 4A; Supplementary Fig. S6A; Supplementary Table S3). This finding replicates the performance effects we observed using instructed cues, albeit with only modest evidence ( $b=-16.00$ ; 95% CrI [-34.91, 2.96];  $P_{b>0}=0.05$ ;  $BF_{01}=1.68$ ). Like in our previous studies, faster responding was not explained by a change in speed-accuracy trade-off, as accuracy did not decrease with increasing efficacy ( $b=0.12$ ; 95% CrI [-0.20, 0.44];  $P_{b>0}=0.23$ ;  $BF_{01}=1.99$ ; Supplementary Fig. S6A). These analyses

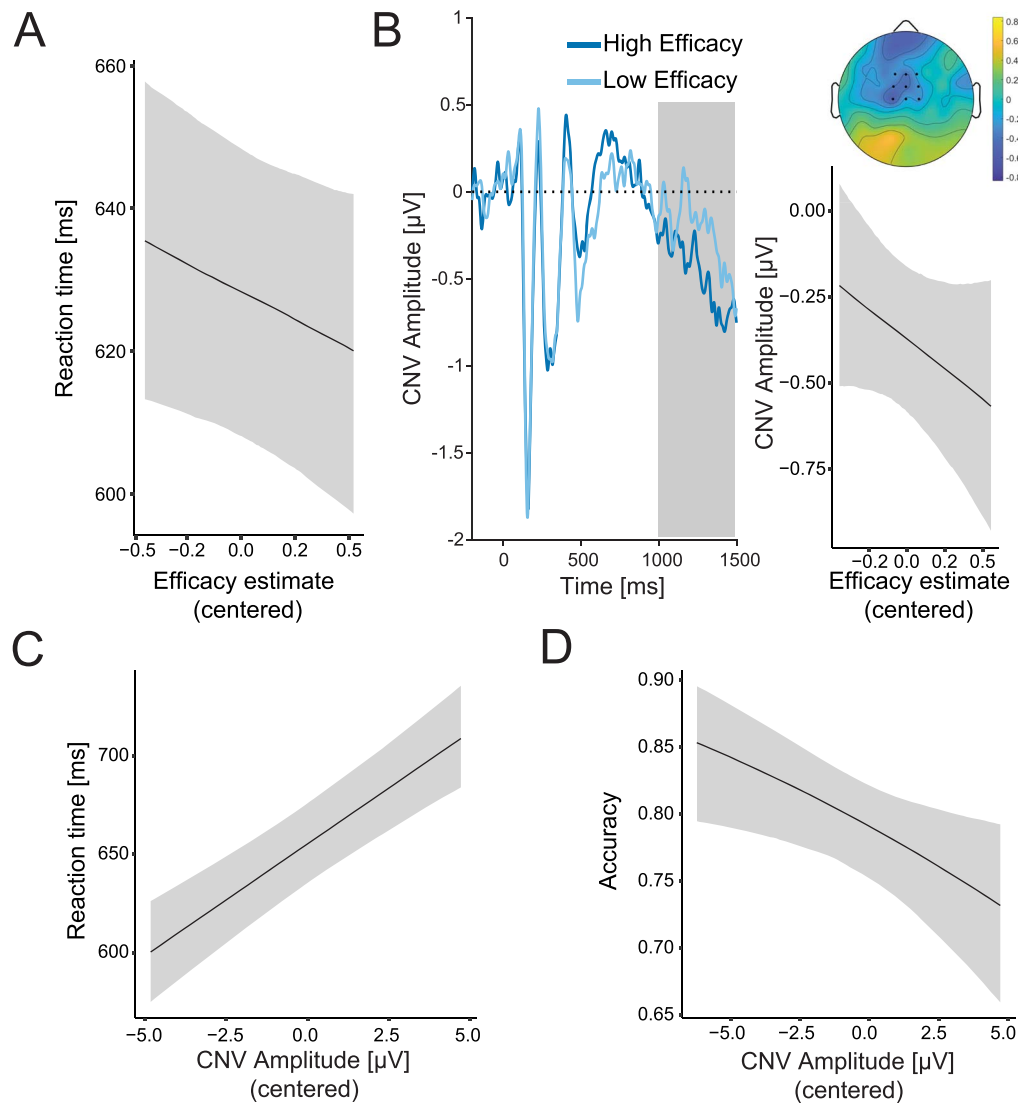


**Fig. 3.** P3b reflects dynamically changing efficacy estimates during processing of efficacy feedback in Study 1. A) ERP average for the P3b locked to the onset of the efficacy feedback separately for performance-based and random feedback. The gray area shows the time window used for quantifying the P3b. B) LMM predicted P3b amplitudes are plotted for performance-based and random feedback as a function of efficacy estimates. The topography shows the interaction of efficacy estimate with efficacy feedback in the P3b time window. C) Predicted (centered) effects of unsigned prediction errors (left) and model-based learning rates (right) on the P3b. Shaded error bars represent 95% confidence intervals. Topographies display fixed-effects estimates.

controlled for the standard behavioral effects related to Stroop congruency (i.e. slower and less accurate responding for incongruent relative to congruent trials; [Supplementary Fig. S7](#)), as well as for the reward rate estimates.

If the CNV provides an index of control allocation based on current incentive expectations, it should both reflect one's latest efficacy estimate and predict performance on the upcoming trial. Our findings support both predictions. Regressing single-trial CNV amplitude onto model-based efficacy estimates, and controlling for expectations of reward (discussed below), we found that CNV amplitudes had a positive relationship (more negative.) with the current efficacy expectations

( $b = -0.35$ ; 95% CrI  $[-0.85, 0.16]$ ;  $P_{b > 0} = 0.09$ ;  $BF_{10} = 2.73$ ; [Fig. 4B-right](#); [Supplementary Fig. S5B](#); [Supplementary Table S4](#)). (Note that the CNV is a negative component, and thus higher CNV amplitudes - i.e., more control allocation - are more negative.) However, this effect was weaker than in the previous experiment with cued efficacy levels, which is to be expected given that in this experiment participants had to learn their efficacy levels. As with the behavioral finding above, this result provides evidence consistent with our previous CNV finding using instructed cues. We further replicate our earlier finding ([Frömer, Lin, et al. 2021a](#)) that larger CNV amplitude in turn predicted better performance on the upcoming trial, with participants responding



**Fig. 4.** Efficacy estimates influence allocation of control and behavior in Study 1. A) Higher efficacy predicts faster accurate responses. B) CNV increases with higher efficacy. Left: Grand-average ERP for high and low efficacy estimates (median split used for plotting). The shaded area marks the time window used for CNV quantification. Time 0 corresponds to the onset of the neutral cue. Right: Predicted CNV amplitudes as a function of efficacy estimates. The topography shows the fixed effect of the efficacy estimate from the fitted linear model. C, D) Larger CNV amplitude predicts better performance. Predicted accurate RT (C) and accuracy (D) as a function of efficacy estimates. Shaded error bars indicate 95% confidence intervals.

faster ( $b = 11.41$ ; 95% CrI [8.09, 14.72];  $P_{b < 0} = 0$ ;  $\text{BF}_{10} > 100$ ; Fig. 4C; Supplementary Table S5) and more accurately ( $b = -0.07$ ; 95% CrI [-0.12, -0.01];  $P_{b > 0} = 0.01$ ;  $\text{BF}_{10} = 3.14$ ; Fig. 4D; Supplementary Table S5) as CNV increased.

#### Parallel learning of efficacy and reward rate

We held the amount of reward at stake constant over the course of the experiment, but the frequency of reward (reward rate) varied over the course of the session based on a participant's recent performance, and participants were sporadically probed for estimates of their reward rate (interleaved with trials that were followed by efficacy probes). Our efficacy manipulation explicitly controlled for this variability by yoking random-outcome feedback to a participant's recent reward rate (see Methods for details). However, this additional source of variability also provided an

opportunity to examine additional mechanisms of learning and adaptation in our task. As in the case of efficacy estimates, reward rate estimates were robustly predicted by reward feedback on previous trials (Supplementary Table S1; Supplementary Fig. S8A), and this reward rate learning process was well captured by a two learning rate reward rate model (Garrett and Daw 2020; Supplementary Fig. S8B and C), with the model-based estimates successfully predicting the reported subjective estimates ( $b = 0.82$ ; 95% CrI [0.60, 1.02];  $P_{b < 0} < 0.001$ ;  $\text{BF}_{01} > 100$ ;  $R^2 = 0.58$ ; Supplementary Fig. S8B and C). Updates to these reward rate estimates were reflected in P3b modulations of (unsigned) reward prediction errors and associated learning rates (Fischer and Ullsperger 2013; Supplementary Fig. S9; Supplementary Table S6). This pattern of results provides additional evidence that efficacy learning involves

similar neural and computational mechanisms as reward-based learning.

## Study 2

### *Learned efficacy modulates control over information processing*

Our findings suggest that people rely on domain-general mechanisms to learn about their performance efficacy in a given environment and use these learned estimates of efficacy to optimize performance on their task. Specifically, in Study 1, we found that higher levels of learned efficacy are associated with faster responses, albeit with modest evidence ( $b = -16.00$ ; 95% CrI  $[-34.91, 2.96]$ ;  $P_{b > 0} = 0.05$ ;  $BF_{01} = 1.68$ ). We also found that this speeding occurred on correct but not incorrect trials, suggesting that these performance adjustments reflected attentional control rather than adjustments to speed–accuracy tradeoffs. However, these findings remain only suggestive in the absence of a formal model, and the presence of a stringent response deadline in this study (individually calibrated for each subject during the practice phase to ensure the reward rate of 80%;  $M = 796$  ms;  $SD = 73$  ms) presented an obstacle to fitting our behavioral data to such a model without additional assumptions (e.g. regarding the form of a collapsing threshold).

To provide further support for our proposal that learned efficacy influences control over information processing, we ran an additional behavioral study (Study 2). Participants in this study ( $N = 51$ ) performed a web-based version of the task in Study 1, with the biggest modification being that the Stroop trials (now using picture-word rather than color-word interference) were performed over the course of short self-paced time intervals rather than trial-by-trial as in Study 1. Specifically, participants were given limited time windows (2–4 s) to complete as many Stroop trials as they wanted to and were rewarded at the end of each interval (cf. [Leng et al. 2021](#)). When rewards were performance-based, participants received a number of points exactly proportional to the number of correct responses they gave during that window. When rewards were random, the number of points was unrelated to performance on that interval but (as in Study 1) was yoked to their performance in previous performance-contingent intervals. Following our approach in Study 1, we varied the likelihood of a given interval being performance-based or random over the course of the session ([Supplementary Fig. S2A](#)), and sporadically probed participants for their subjective estimates of expected efficacy and reward rate. While in most respects very similar to the paradigm in Study 1, one noteworthy feature of this self-paced design is that it resulted in a less stringent deadline for responding, thus producing reaction time patterns more typical of free-response paradigms for which the traditional (fixed-threshold) DDM was designed. Note that because this was an online sample we also employed additional cutoff criteria to exclude inattentive participants, as detailed in

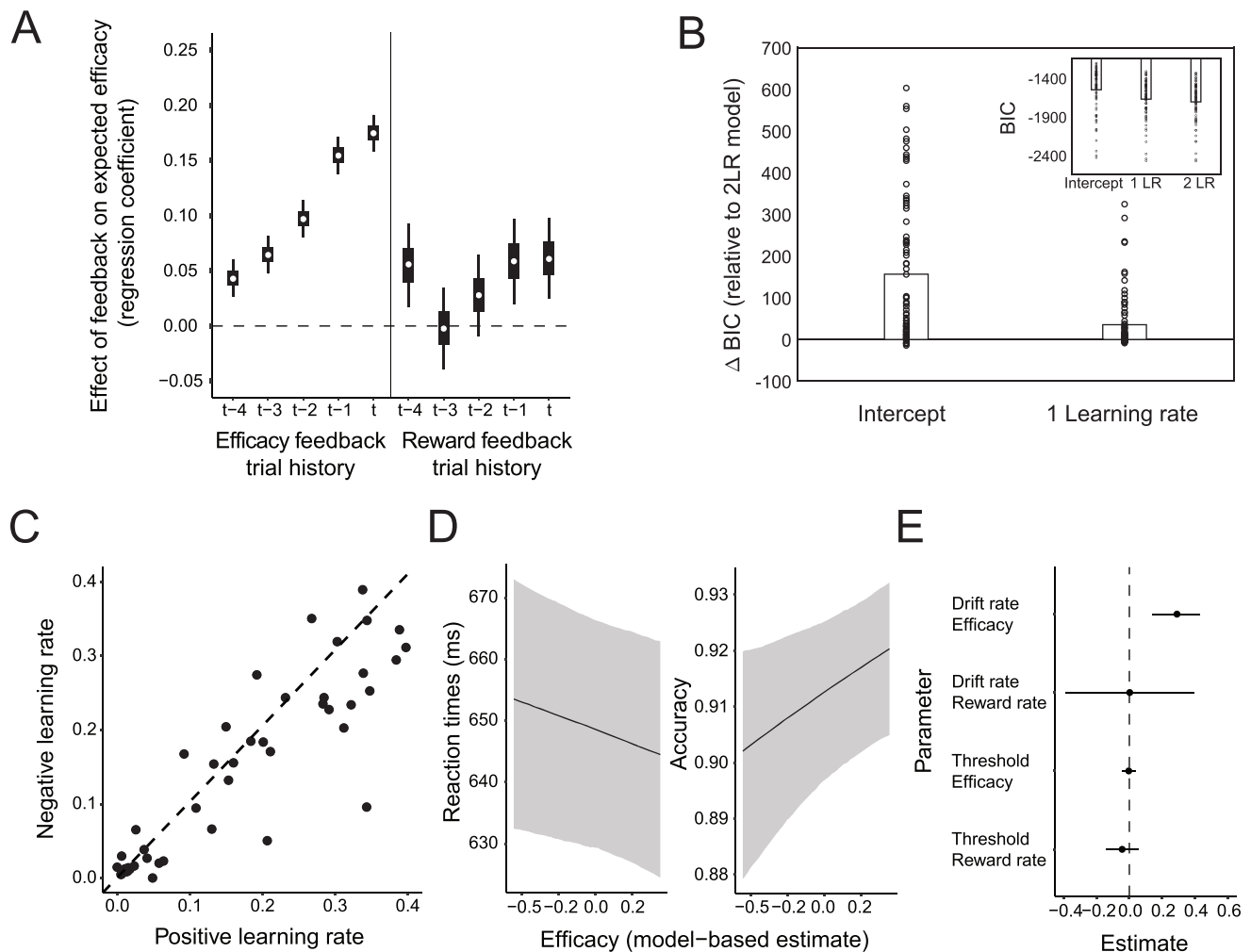
the Methods section and [Supplementary Figs. S2B and C](#) and [S3](#).

Replicating the learning patterns observed in Study 1, subjective estimates of efficacy in Study 2 reflected a running average over recent efficacy feedback ( $b = 0.18$ ; 95% CrI  $[0.16, 0.20]$ ;  $P_{b < 0} = 0$ ;  $BF_{10} > 100$ ). This effect was again weighted towards the most recent feedback but still present up to 5 intervals back (all  $P_{b < 0} < 0.01$ ; [Fig. 5A](#), [Supplementary Fig. S10](#), and [Supplementary Table S7](#)). As in Study 1, this learning pattern was best captured by an RL algorithm with 2 learning rates ([Fig. 5B](#)) and positive and negative efficacy learning rates did not significantly differ from one another on average ( $b = 0.03$ ; 95% CrI  $[-0.03, 0.08]$ ;  $P_{b < 0} = 0.260$ ;  $BF_{01} = 11.74$ ; [Fig. 5C](#)). As in Study 1, the model-based efficacy estimates successfully predicted the raw subjective estimates on intervals on which participants reported their efficacy beliefs ( $b = 0.86$ ; 95% CrI  $[0.81, 0.90]$ ;  $P_{b < 0} < 0.001$ ;  $BF_{01} > 100$ ;  $R^2 = 0.60$ ), and the same was true for the model-based reward estimates predicting the subjective reward estimates ( $b = 1.00$ ; 95% CrI  $[0.95, 1.04]$ ;  $P_{b < 0} < 0.001$ ;  $BF_{01} > 100$ ;  $R^2 = 0.66$ ).

Critically, we once again found that higher model-based estimates of efficacy predicted better performance on the upcoming interval. Participants responded faster ( $b = -10.25$ ; 95% CrI  $[-19.86, -0.20]$ ;  $P_{b > 0} = 0.02$ ;  $BF_{01} = 2.33$ ) and more accurately ( $b = 0.23$ ; 95% CrI  $[0.03, 0.43]$ ;  $P_{b > 0} = 0.01$ ;  $BF_{01} > 100$ ) with increasing model-based efficacy estimates ([Fig. 5D](#), [Supplementary Fig. S6B](#), and [Supplementary Table S8](#)). To formally test whether these behavioral patterns reflected adjustments in information processing (i.e. the rate of evidence accumulation once the stimuli appeared) or instead reflected adjustments in speed–accuracy tradeoffs (i.e. the threshold for responding), we fit these data with the hierarchical drift diffusion model (HDDM; [Wiecki et al. 2013](#)). We tested whether model-based efficacy estimates predicted trial-by-trial changes in drift rate, threshold, or both, while controlling for influences of expected reward rate on those same DDM parameters. We found that higher levels of expected efficacy were associated with higher drift rates ( $b = 0.07$ ; 95% CrI  $[0.14, 0.43]$ ;  $P_{b < 0} = 0.00$ ) but were uncorrelated with threshold levels ( $b = -0.00$ ; 95% CrI  $[-0.04, 0.04]$ ;  $P_{b < 0} = 0.74$ ) ([Fig. 5E](#)). Expected reward rate was not correlated with either drift rate or threshold ([Supplementary Table S9](#)). These results suggest that participants responded to changes in performance efficacy by adjusting their attention to the task, rather than simply adjusting their response threshold (i.e. becoming more or less impulsive).

## Discussion

To evaluate the expected benefits of investing cognitive control into a task, people need to consider whether this control is necessary and/or sufficient for achieving their desired outcome (i.e. whether these efforts are worthwhile). A critical determinant of the worthwhileness of control is performance efficacy, the extent to which



**Fig. 5.** Efficacy is learned in the same way in a modified task and influences a behavioral marker of control allocation. A) Regression weights for the influence of previous feedback type (efficacy and reward) on the subjective efficacy estimate. Error bars represent 50% and 95% highest density intervals. B) Two learning rate model captures efficacy learning best. Differences in BICs between the two learning rate model, and the intercept-only and one learning rate models. C) Efficacy learning rates. Positive and negative efficacy learning rates for each participant. D) Higher model-based efficacy estimates predict better behavioral performance. Higher efficacy estimates reduce reaction times (left) and improve accuracy (right). E) Higher model-based efficacy estimates predict increased allocation of attention. Parameter estimates from the DDM. Higher efficacy estimates increase drift rates, but not response caution (thresholds).

performance on a control-demanding task matters for outcome attainment versus those outcomes being determined by performance-unrelated factors. However, the mechanisms through which people estimate the efficacy of their performance based on previous experience are largely unknown. Here, we identified the computational and neural mechanism through which efficacy is learned and used to guide the allocation of cognitive control. Across 2 experiments, we found that participants dynamically updated expectations of the efficacy of their task performance (i.e. the likelihood that this performance will determine reward attainment) and used those expectations to adjust how much control they allocated. The feedback-based updating of efficacy was well captured by a prediction error-based learning algorithm. Model-based estimates of efficacy and efficacy prediction errors were encoded by canonical neural signatures of effort allocation and belief updating, respectively. Importantly, these findings cannot be

explained by variability in reward, as reward rate was held constant across efficacy levels, and the subjective reward rate was controlled for statistically. Further, our model-based analyses revealed that people allocated more control when they learned that they had more efficacy, extending our previous findings on instructed efficacy (Frömer, Lin, et al. 2021a). Taken together, our results uncover the mechanism through which efficacy estimates are learned and used to inform mental effort investment within a given task environment.

Previous research has characterized the learning algorithms people use to learn the reward value of different states and actions in their environment (Gläscher et al. 2010; Sutton and Barto 2018). Recent theoretical (Jiang et al. 2014; Lieder et al. 2018; Verbeke and Verguts 2019) and empirical (Bejjani et al. 2018; Otto and Daw 2019; Jiang et al. 2020; Bustamante et al. 2021) work has extended this research to show how similar algorithms guide learning and adaptation of cognitive

control under varying rewards and task demands within a given task environment. Our findings extend this work further in several ways. First, we show that people leverage weighted prediction errors when learning about the efficacy of task performance, independently of potential reward and task difficulty. Second, we show that they update their efficacy expectations differently depending on whether performance was more efficacious or less efficacious than they expected, demonstrating a striking parallel with dual learning rate models that have been found to prevail in research on reward learning (Collins and Frank 2013; Lefebvre et al. 2017; Garrett and Daw 2020), including in our own reward rate data (Supplementary Fig. S9). Third, we show that participants dynamically adjust their control allocation based on learned efficacy, just as they do for learned rewards and task demands (Bugg et al. 2011; Jiang et al. 2014; Lieder et al. 2018).

Our neural findings build further on past research on learning and control adaptation. The P3b component has been shown to track prediction error-based learning from action-relevant outcome values (Fischer and Ullsperger 2013; Nassar et al. 2019; Lohse et al. 2020). Here we show that this neural marker tracks learning about efficacy in the same way as it tracks learning about rewards. We found increased P3b amplitudes when people experienced feedback about outcome contingencies that was less expected given their current estimate of efficacy (e.g. expecting low efficacy, but getting performance-contingent feedback), relative to when these contingencies were more expected (e.g. expecting low efficacy and getting random feedback). Our additional finding that P3b amplitude was overall larger for efficacy compared to no efficacy feedback demonstrates that our participants were not just tracking the frequency of the 2 types of feedback, as would be predicted by an oddball account. Instead this finding suggests that they were actively learning from the feedback.

Extending previous findings on cueing efficacy and/or reward (Schevernels et al. 2014; Frömer, Lin, et al. 2021a), our CNV and behavioral results further show that participants used these learned efficacy estimates to calibrate their effort and their performance. Notably, unlike in previous work, our study shows effort-related changes in CNV amplitude entirely divorced from perceptual cues, providing evidence that this activity truly reflects adjustments in control, rather than reactive processing of features associated with the cue. Taken together, our findings suggest that similar neural mechanisms underlie learning and control adaptation in response to variability in one's efficacy in a task environment, as they do in response to variability in expected rewards (Otto and Daw 2019; Leng et al. 2021). By fitting behavioral data from this task to a DDM (Study 2), we were able to further demonstrate that participants were adapting to changes in expected efficacy by enhancing the processing of stimuli (i.e. increasing their rate of evidence accumulation)—potentially via attentional

control mechanisms—rather than by adjusting their threshold for responding. This particular pattern of control adjustments was predicted for the current task because performance-contingent rewards depended on responses being both fast and accurate (as in Frömer, Lin, et al. 2021a), but future work should test the prediction that different control adjustments should emerge under different performance contingencies (cf. Leng et al. 2021; Ritz et al. 2022).

Our findings build on previous research on how people learn about controllability of their environment. Studies have examined the neural and computational mechanisms by which humans and other animals learn about the contingencies between an action and its associated outcome and demonstrated that these learned action–outcome contingencies influence which actions are selected and how vigorously they are enacted (Dickinson and Balleine 1995; Liljeholm et al. 2011; Manohar et al. 2017; Moscarello and Hartley 2017; Ly et al. 2019). Our work extends this research into the domain of cognitive control, where the contingencies between actions (i.e. control adjustments) and outcomes (e.g. reward) depend both on whether control adjustments predicts better performance and whether better performance predicts better outcomes (Shenhav et al. 2021). Learning control–outcome contingencies therefore requires learning about how control states map onto performance (control efficacy) as well as how performance maps onto outcomes (performance efficacy). By describing the mechanisms by which people solve the latter part of this learning problem and demonstrating that these are comparable to those engaged during action–outcome learning, our study lays critical groundwork for better understanding the links between selection of actions and control states.

Our efficacy-updating results are a reminder that many aspects of feedback are reflected in prediction error signals (Langdon et al. 2018; Frömer, Nassar, et al. 2021b). In the present study, we intentionally separated feedback about reward and efficacy to isolate the cognitive and neural learning mechanisms associated with each. In doing so, we have taken an important first step towards understanding the updating mechanisms underlying each. Further work is needed to understand how these are inferred in more naturalistic settings, in which different forms of feedback are often multiplexed, containing information about the values of actions that were taken as well as about the features and structure of the environment (cf. Dorfman et al. 2019, 2021).

Another distinct element of more complex naturalistic environments is that the same feedback can be used to evaluate multiple targets, including internal ones, such as the selected response and its predicted outcome, and external ones, such as the source of feedback/environment (Carlebach and Yeung 2020). Consistent with such multi-level prediction error signals, recent work has shown that the seemingly consistent relationship between the P3b and behavioral adaptation (Yeung and Sanfey 2004; Chase et al. 2011; Fischer and Ullsperger

2013) in fact depends on the context in which a prediction error occurred (Nassar et al. 2019). Reinforcement learning, and predictive inference more generally, have been proposed to not only support the selection of individual actions but also extended sequences of actions and control signals (Holroyd and Yeung 2012; Lieder et al. 2018). In addition to the evaluation of overt actions and the related feedback, neural signatures of feedback-based learning could also reflect the evaluations of covert actions, such as the intensity of cognitive control signals. Given the many potential causes a given outcome can have, and the flexibility that people have in how they use the feedback, it is easy to see how feedback could be misattributed and lead to inaccurate beliefs about performance efficacy. Such beliefs about environmental statistics can drive changes in feedback processing and behavioral adaptation, above and beyond the true statistics (Schiffer et al. 2017), and are thus of particular importance for understanding some of the cognitive symptoms of mental disorders.

Understanding how efficacy estimates develop based on previous experiences is crucial for understanding why people differ in how much cognitive control they allocate in different situations (Shenhav et al. 2021). People differ in their beliefs about how much control they have over their environments (Leotti et al. 2010; Moscarello and Hartley 2017), and in their ability to estimate their efficacy (Cohen et al. 2020). Further, many mental disorders, including depression and schizophrenia, have been linked with one's estimates of their ability to control potential outcomes in their environment, including through allocation of control (Huys and Dayan 2009; Maier and Seligman 2016). We recently proposed that such changes can drive impairments of motivation and control in those populations (Grahek et al. 2019). As we show in this study, when people have learned to expect low efficacy, they will allocate less cognitive control, which can manifest as apparent control deficits. The experimental and modeling approach taken in our study helps uncover a more fine-grained view of how components of motivation are learned and used to support the allocation of cognitive resources. In this way, our study takes a first step towards a better computational and neural account of efficacy learning, which can aid the understanding of individual differences in the willingness to exert mental effort, as well as the development of interventions aimed at teaching individuals when these efforts truly matter.

## Acknowledgments

We would like to thank Natalie Knowles and Hattie Xu for help with data collection for Study 1, Peyton Strong for help with data collection for Study 2, Carolyn Dean Wolf and Liz Cory for help with programming the Study 1 task, Harrison Ritz for help with constructing the efficacy drifts, and Xiamin Leng for advice on fitting the drift-diffusion models.

## Supplementary material

Supplementary material is available at *Cerebral Cortex* online.

## Funding

This work was supported by the Special Research Fund (BOF) of Ghent University (grant #01D02415 to IG), the Research Foundation Flanders (FWO) travel grant (grant #V432718N to IG), a Center of Biomedical Research Excellence grant P20GM103645 from the National Institute of General Medical Sciences (AS), the Alfred P. Sloan Foundation Research Fellowship in Neuroscience (AS), grant R01MH124849 from the National Institute of Mental Health (AS), and an NSF Graduate Research Fellowship (MPF). The funding sources were not involved in the study design; collection, analysis, and interpretation of data; writing of the report; and decision to submit the article for publication.

*Conflict of interest statement:* The authors have no conflicts of interest to declare.

## Data availability

Analysis scripts are available on this link: [https://github.com/igrahek/LFXC\\_EEG\\_2022.git](https://github.com/igrahek/LFXC_EEG_2022.git). Please contact the authors for raw data.

## References

- Atkinson JW, Feather NT et al. *A theory of achievement motivation*. New York: Wiley; 1966
- Balleine BW, O'Doherty JP. Human and rodent homologies in action control: corticostriatal determinants of goal-directed and habitual action. *Neuropsychopharmacology*. 2010;35(1):48–69.
- Bejjani C, Zhang Z, Egnor T. Control by association: transfer of implicitly primed attentional states across linked stimuli. *Psychon Bull Rev*. 2018;25(2):617–626. <https://doi.org/10.3758/s13423-018-1445-6>.
- Botvinick MM, Cohen JD. The computational and neural basis of cognitive control: charted territory and new frontiers. *Cogn Sci*. 2014;38(6):1249–1285. <https://doi.org/10.1111/cogs.12126>.
- Brainard DH. The Psychophysics Toolbox. *Spat Vis*. 1997;10(4):433–436.
- Bugg JM, Jacoby LL, Chanani S. Why it is too early to lose control in accounts of item-specific proportion congruency effects. *J Exp Psychol Hum Percept Perform*. 2011;37(3):844–859. <https://doi.org/10.1037/a0019957>.
- Bürkner P-C. brms: an R package for Bayesian multilevel models using Stan. *J Stat Softw*. 2016;80(1):1–28.
- Bürkner P-C. Advanced Bayesian multilevel modeling with the R Package brms. 2017: arXiv:170511123.
- Bustamante L, Lieder F, Musslick S, Shenhav A, Cohen J. Learning to overexert cognitive control in a Stroop task. *Cogn Affect Behav Neurosci*. 2021;21(3):453–471. <https://doi.org/10.3758/s13415-020-00845-x>.
- Carlbach N, Yeung N. Flexible use of confidence to guide advice requests. *PsyArXiv*. 2020.



- Carpenter B, Gelman A, Hoffman M, Lee D, Goodrich B, Betancourt M, Brubaker MA, Guo J, Li P, Riddell A. Stan: a probabilistic programming language. *J Stat Softw.* 2016;2(20):1–37.
- Chambon V, Théro H, Vidal M, Vandendriessche H, Haggard P, Palminteri S. Information about action outcomes differentially affects learning from self-determined versus imposed choices. *Nat Hum Behav.* 2020;4(10):1067–1079. <https://doi.org/10.1038/s41562-020-0919-5>.
- Chase HW, Swainson R, Durham L, Benham L, Cools R. Feedback-related negativity codes prediction error but not behavioral adjustment during probabilistic reversal learning. *J Cogn Neurosci.* 2011;23(4):936–946. <https://doi.org/10.1162/jocn.2010.21456>.
- Chiu YC, Egnér T. Cortical and subcortical contributions to context-control learning. *Neurosci Biobehav Rev.* 2019;99:33–41.
- Cohen A, Nussenbaum K, Dorfman HM, Gershman SJ, Hartley CA. The rational use of causal inference to guide reinforcement learning strengthens with age. *Npj Sci Learn.* 2020;5(1):16.
- Collins AGE, Frank MJ. Cognitive control over learning: creating, clustering, and generalizing task-set structure. *Psychol Rev.* 2013;120(1):190–229. <https://doi.org/10.1037/a0030852>.
- Collins AGE, Frank MJ. Opponent actor learning (OpAL): modeling interactive effects of striatal dopamine on reinforcement learning and choice incentive. *Psychol Rev.* 2014;121(3):337–366. <https://doi.org/10.1037/a0037015>.
- Delorme A, Makeig S. EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J Neurosci Methods.* 2004;134(1):9–21. <https://doi.org/10.1016/j.jneumeth.2003.10.009>.
- Dickinson A, Balleine B. Motivational control of instrumental action. *Curr Dir Psychol Sci.* 1995;4(5):162–167. <https://doi.org/10.1111/1467-8721.ep11512272>.
- Dorfman HM, Bhui R, Hughes BL, Gershman SJ. Causal inference about good and bad outcomes. *Psychol Sci.* 2019;30(4):516–525.
- Dorfman HM, Tomov MS, Cheung B, Clarke D, Gershman SJ, Hughes BL. Causal inference gates corticostriatal learning. *J Neurosci.* 2021;41(32):6892–6904.
- Eldar E, Hauser TU, Dayan P, Dolan RJ. Striatal structure and function predict individual biases in learning to avoid pain. *Proc Natl Acad Sci U S A.* 2016;113(17):4812–4817. <https://doi.org/10.1073/pnas.1519829113>.
- Fischer AG, Ullsperger M. Real and fictive outcomes are processed differently but converge on a common adaptive mechanism. *Neuron.* 2013;79(6):1243–1255. <https://doi.org/10.1016/j.neuron.2013.07.006>.
- Frömer R, Stürmer B, Sommer W. (Don't) mind the effort: effects of contextual interference on ERP indicators of motor preparation. *Psychophysiology.* 2016;53(10):1577–1586. <https://doi.org/10.1111/psyp.12703>.
- Frömer R, Maier M, Rahman RA. Group-level EEG-processing pipeline for flexible single trial-based analyses including linear mixed models. *Front Neurosci.* 2018;12(48):48. <https://doi.org/10.3389/fnins.2018.00048>.
- Frömer R, Lin H, Dean Wolf CK, Inzlicht M, Shenhav A. Expectations of reward and efficacy guide cognitive control allocation. *Nat Commun.* 2021a;12(1):1030. <https://doi.org/10.1038/s41467-021-21315-z>.
- Frömer R, Nassar MR, Bruckner R, Stürmer B, Sommer W, Yeung N. Response-based outcome predictions and confidence regulate feedback processing and learning. *elife.* 2021b;10:e62825.
- Garrett N, Daw ND. Biased belief updating and suboptimal choice in foraging decisions. *Nat Commun.* 2020;11(1):3417. <https://doi.org/10.1038/s41467-020-16964-5>.
- Gelman A, Rubin DB. Inference from iterative simulation using multiple sequences. *Stat Sci.* 1992;7(4):457–472. <https://doi.org/10.1214/ss/1177011136>.
- Gläscher J, Daw N, Dayan P, O'Doherty JP. States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron.* 2010;66(4):585–595. <https://doi.org/10.1016/j.neuron.2010.04.016>.
- Grahek I, Shenhav A, Musslick S, Krebs RM, Koster EHW. Motivation and cognitive control in depression. *Neurosci Biobehav Rev.* 2019;102(March):371–381. <https://doi.org/10.1016/j.neubiorev.2019.04.011>.
- Grent-t-Jong T, Woldorff MG. Timing and sequence of brain activity in top-down control of visual-spatial attention. *PLoS Biol.* 2007;5(1):0114–0126. <https://doi.org/10.1371/journal.pbio.0050012>.
- Gureckis TM, Martin J, McDonnell J, Rich AS, Markant D, Coenen A, Chan P. psiTurk: An open-source framework for conducting replicable behavioral experiments online. *Behav Res Meth.* 2016;48(3):829–842.
- Holroyd CB, Yeung N. Motivation of extended behaviors by anterior cingulate cortex. *Trends Cogn Sci.* 2012;16(2):122–128. <https://doi.org/10.1016/j.tics.2011.12.008>.
- Huys QJM, Dayan P. A Bayesian formulation of behavioral control. *Cognition.* 2009;113(3):314–328. <https://doi.org/10.1016/j.cognition.2009.01.008>.
- Ille N, Berg P, Scherg M. Artifact correction of the ongoing eeg using spatial filters based on artifact and brain signal topographies. *J Clin Neurophysiol.* 2002;19(2):113–124. <https://doi.org/10.1097/00004691-200203000-00002>.
- Jiang J, Heller K, Egnér T. Bayesian modeling of flexible cognitive control. *Neurosci Biobehav Rev.* 2014;46(P1):30–43. <https://doi.org/10.1016/j.neubiorev.2014.06.001>.
- Jiang J, Bramão I, Khazenzon A, Wang SF, Johansson M, Wagner AD. Temporal dynamics of memory-guided cognitive control and generalization of control via overlapping associative memories. *J Neurosci.* 2020;40(10):2343–2356. <https://doi.org/10.1523/JNEUROSCI.1869-19.2020>.
- Kleiner M, Brainard D, Pelli D. What's new in Psychtoolbox-3. In: *Perception 36 ECVF abstract supplement*; 2007
- Kool W, Botvinick M. Mental labour. *Nat Hum Behav.* 2018;2(12):899–908. <https://doi.org/10.1038/s41562-018-0401-9>.
- Krebs RM, Boehler CN, Roberts KC, Song AW, Woldorff MG. The involvement of the dopaminergic midbrain and cortico-striatal-thalamic circuits in the integration of reward prospect and attentional task demands. *Cereb Cortex.* 2012;22(3):607–615.
- Kruschke JK. Bayesian estimation supersedes the T test. *J Exp Psychol Gen.* 2013;142(2):573–588. <https://doi.org/10.1037/a0029177>.
- Langdon AJ, Sharpe MJ, Schoenbaum G, Niv Y. Model-based predictions for dopamine. *Curr Opin Neurobiol.* 2018;49:1–7. <https://doi.org/10.1016/j.conb.2017.10.006>.
- Lefebvre G, Lebreton M, Meyniel F, Bourgeois-Gironde S, Palminteri S. Behavioural and neural characterization of optimistic reinforcement learning. *Nat Hum Behav.* 2017;1(4):1–9. <https://doi.org/10.1038/s41562-017-0067>.
- Leng X, Yee D, Ritz H, Shenhav A. Dissociable influences of reward and punishment on adaptive cognitive control. *PLoS Comput Biol.* 2021;17(12):e1009737.
- Leotti LA, Iyengar SS, Ochsner KN. Born to choose: the origins and value of the need for control. *Trends Cogn Sci.* 2010;14(10):457–463. <https://doi.org/10.1016/j.tics.2010.08.001>.
- Lieder F, Shenhav A, Musslick S, Griffiths TL. Rational metareasoning and the plasticity of cognitive control. *PLoS*

- Comput Biol. 2018;14(4):e1006043. <https://doi.org/10.1371/journal.pcbi.1006043>.
- Ligneul R, Mainen ZF, Ly V, Cools R. Stress-sensitive inference of task controllability. *Nat Hum Behav.* 2022;1–11. <https://doi.org/10.1038/s41562-022-01306-w>.
- Liljeholm M, Tricomi E, O'Doherty JP, Balleine BW. Neural correlates of instrumental contingency learning: differential effects of action-reward conjunction and disjunction. *J Neurosci.* 2011;31(7):2474–2480. <https://doi.org/10.1523/JNEUROSCI.3354-10.2011>.
- Lohse KR, Miller MW, Daou M, Valerius W, Jones M. Dissociating the contributions of reward-prediction errors to trial-level adaptation and long-term learning. *Biol Psychol.* 2020;149(January 2019):107775. <https://doi.org/10.1016/j.biopsycho.2019.107775>.
- Ly V, Wang KS, Bhanji J, Delgado MR. A reward-based framework of perceived control. *Front Neurosci.* 2019;13(February):1–11. <https://doi.org/10.3389/fnins.2019.00065>.
- Maier SF, Seligman MEP. Learned helplessness at fifty: insights from neuroscience. *Psychol Rev.* 2016;123(4):349–367. <https://doi.org/10.1037/rev0000033>.
- Manohar SG, Chong TTJ, Apps MAJ, Batla A, Stamelou M, Jarman PR, Bhatia KP, Husain M. Reward pays the cost of noise reduction in motor and cognitive control. *Curr Biol.* 2015;25(13):1707–1716. <https://doi.org/10.1016/j.cub.2015.05.038>.
- Manohar SG, Finzi RD, Drew D, Husain M. Distinct motivational effects of contingent and noncontingent rewards. *Psychol Sci.* 2017;28(7):1016–1026. <https://doi.org/10.1177/0956797617693326>.
- Morris RW, Dezfouli A, Griffiths KR, Le Pelley ME, Balleine BW. The neural bases of action-outcome learning in humans. *J Neurosci.* 2022;42(17):3636–3647.
- Moscarello JM, Hartley CA. Agency and the calibration of motivated behavior. *Trends Cogn Sci.* 2017;21(10):725–735. <https://doi.org/10.1016/j.tics.2017.06.008>.
- Nagase AM, Onoda K, Foo JC, Haji T, Akaishi R, Yamaguchi S, Sakai K, Morita K. Neural mechanisms for adaptive learned avoidance of mental effort. *J Neurosci.* 2018;38(10):2631–2651. <https://doi.org/10.1523/JNEUROSCI.1995-17.2018>.
- Nalborczyk L, Bürkner P-C. An introduction to Bayesian multilevel models using brms: a case study of gender effects on vowel variability in standard Indonesian. *J Speech Lang Hear Res.* 2019;62(5):1225–1242.
- Nassar MR, Bruckner R, Frank MJ. Statistical context dictates the relationship between feedback-related EEG signals and learning. *elife.* 2019;8:1–26. <https://doi.org/10.7554/eLife.46975>.
- Nigbur R, Schneider J, Sommer W, Dimigen O, Stürmer B. Ad-hoc and context-dependent adjustments of selective attention in conflict control: an ERP study with visual probes. *NeuroImage.* 2015;107:76–84.
- Niv Y, Edlund JA, Dayan P, O'Doherty JP. Neural prediction errors reveal a risk-sensitive reinforcement-learning process in the human brain. *J Neurosci.* 2012;32(2):551–562. <https://doi.org/10.1523/JNEUROSCI.5498-10.2012>.
- Norton KG, Liljeholm M. The rostralateral prefrontal cortex mediates a preference for high-agency environments. *J Neurosci.* 2020;40(22):4401–4409.
- Otto AR, Daw ND. The opportunity cost of time modulates cognitive effort. *Neuropsychologia.* 2019;123:92–105. <https://doi.org/10.1016/j.neuropsychologia.2018.05.006>.
- Parro C, Dixon ML, Christoff K. The neural basis of motivational influences on cognitive control. *Hum Brain Mapp.* 2018;39(12):5097–5111.
- Pelli DD. The VideoToolbox software for visual psychophysics: transforming numbers into movies. *Spat Vis.* 1997;10(4):437–442.
- Ratcliff R, McKoon G. The diffusion decision model: theory and data for two-choice decision tasks. *Neur comp.* 2008;20(4):873–922.
- R Core Team. R: a language and environment for statistical computing, Vienna, Austria. 2020. URL <https://www.R-project.org/>.
- Ritz H, Leng X, Shenhav A. Cognitive control as a multivariate optimization problem. *J Cogn Neurosci.* 2022;34(4):569–591.
- Rutledge RB, Skandali N, Dayan P, Dolan RJ. A computational and neural model of momentary subjective well-being. *Proc Natl Acad Sci.* 2014;111(33):12252–12257. <https://doi.org/10.1073/pnas.1407535111>.
- Schevernels H, Krebs RM, Santens P, Woldorff MG, Boehler CN. Task preparation processes related to reward prediction precede those related to task-difficulty expectation. *NeuroImage.* 2014;84:639–647. <https://doi.org/10.1016/j.neuroimage.2013.09.039>.
- Schiffer AM, Siletti K, Waszak F, Yeung N. Adaptive behaviour and feedback processing integrate experience and instruction in reinforcement learning. *NeuroImage.* 2017;146:626–641. <https://doi.org/10.1016/j.neuroimage.2016.08.057>.
- Shenhav A, Botvinick M, Cohen J. The expected value of control: an integrative theory of anterior cingulate cortex function. *Neuron.* 2013;79(2):217–240. <https://doi.org/10.1016/j.neuron.2013.07.007>.
- Shenhav A, Musslick S, Lieder F, Kool W, Griffiths TL, Cohen JD, Botvinick MM. Toward a rational and mechanistic account of mental effort. *Annu Rev Neurosci.* 2017;40(1):99–124. <https://doi.org/10.1146/annurev-neuro-072116-031526>.
- Shenhav A, Prater Fahey M, Grahek I. Decomposing the motivation to exert mental effort. *Curr Dir Psychol Sci.* 2021;30(4):307–314.
- Sutton RS, Barto AG. *Reinforcement learning: an introduction*. Cambridge, MA: MIT Press; 2018
- Vassena E, Silvetti M, Boehler CN, Achten E, Fias W, Verguts T. Overlapping neural systems represent cognitive effort and reward anticipation. *PLoS One.* 2014;9(3):e91008.
- Verbeke P, Verguts T. Learning to synchronize: how biological agents can couple neural task modules for dealing with the stability-plasticity dilemma. *PLoS Comput Biol.* 2019;15(8):e1006604. <https://doi.org/10.1371/journal.pcbi.1006604>.
- Verguts T, Vassena E, Silvetti M. Adaptive effort investment in cognitive and physical tasks: a neurocomputational model. *Front. Behav Neurosci.* 2015;9(March):57. <https://doi.org/10.3389/fnbeh.2015.00057>.
- Wagenmakers EJ, Lodewyckx T, Kuriyal H, Grasman R. Bayesian hypothesis testing for psychologists: a tutorial on the Savage-Dickey method. *Cogn Psychol.* 2010;60(3):158–189. <https://doi.org/10.1016/j.cogpsych.2009.12.001>.
- Wiecki TV, Sofer I, Frank MJ. HDDM: Hierarchical Bayesian estimation of the drift-diffusion model in Python. *Front in neuroinf.* 2013;14.
- Yeung N, Sanfey AG. Independent coding of reward magnitude and valence in the human brain. *J Neurosci.* 2004;24(28):6258–6264. <https://doi.org/10.1523/JNEUROSCI.4537-03.2004>.