

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/317379316>

Rational metareasoning and the plasticity of cognitive control

Preprint · February 2018

DOI: 10.13140/RG.2.2.24500.14721

CITATIONS

3

READS

346

4 authors, including:



Falk Lieder

University of California, Berkeley

43 PUBLICATIONS 310 CITATIONS

[SEE PROFILE](#)



Amitai Shenhav

Brown University

34 PUBLICATIONS 1,164 CITATIONS

[SEE PROFILE](#)



Sebastian Musslick

Princeton University

12 PUBLICATIONS 41 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Goal Setting [View project](#)



Developing tools and theories for helping people make better decisions [View project](#)

All content following this page was uploaded by [Falk Lieder](#) on 16 February 2018.

The user has requested enhancement of the downloaded file.

Rational metareasoning and the plasticity of cognitive control

Falk Lieder¹, Amitai Shenhav², Sebastian Musslick³, Thomas L. Griffiths⁴

¹ Helen Wills Neuroscience Institute, University of California, Berkeley

² Department of Cognitive, Linguistic, and Psychological Sciences,
Brown Institute for Brain Science, Brown University

³ Princeton Neuroscience Institute, Princeton University

⁴ Department of Psychology, University of California, Berkeley

Abstract

The human brain has the impressive capacity to adapt how it processes information to high-level goals. While it is known that these cognitive control skills are malleable and can be improved through training, the underlying plasticity mechanisms are not well understood. Here, we develop and evaluate a model of how people learn when to exert cognitive control, which controlled process to use, and how much effort to exert. We derive this model from a general theory according to which the function of cognitive control is to select and configure neural pathways so as to make optimal use of finite time and limited computational resources. The central idea of our Learned Value of Control model is that people use reinforcement learning to predict the value of candidate control signals of different types and intensities based on stimulus features. This model correctly predicts the learning and transfer effects underlying the adaptive control-demanding

1 behavior observed in an experiment on visual attention and four experiments on
2 interference control in Stroop and Flanker paradigms. Moreover, our model explained
3 these findings significantly better than an associative learning model and a Win-Stay
4 Lose-Shift model. Our findings elucidate how learning and experience might shape
5 people's ability and propensity to adaptively control their minds and behavior. We
6 conclude by predicting under which circumstances these learning mechanisms might lead
7 to self-control failure.

8 Author Summary

9 The human brain has the impressive ability to adapt how it processes information to high
10 level goals. While it is known that these cognitive control skills are malleable and can be
11 improved through training, the underlying plasticity mechanisms are not well understood.
12 Here, we derive a computational model of how people learn when to exert cognitive
13 control, which controlled process to use, and how much effort to exert from a formal
14 theory of the function of cognitive control. Across five experiments, we find that our
15 model correctly predicts that people learn to adaptively regulate their attention and
16 decision-making and how these learning effects transfer to novel situations. Our findings
17 elucidate how learning and experience might shape people's ability and propensity to
18 adaptively control their minds and behavior. We conclude by predicting under which
19 circumstances these learning mechanisms might lead to self-control failure.

1 Introduction

2 The human brain has the impressive ability to adapt how it processes information
3 and responds to stimuli in the service of high level goals, such as writing an article [1].

4 The mechanisms underlying this behavioral flexibility range from seemingly simple
5 processes, such as inhibiting the impulse to browse your Facebook feed, to very complex
6 processes such as orchestrating your thoughts to reach a solid conclusion. Our capacity
7 for *cognitive control* enables us to override automatic processes when they are
8 inappropriate for the current situation or misaligned with our current goals. One of the
9 paradigms used to study cognitive control is the Stroop task, where participants are
10 instructed to name the hue of a color word (e.g., respond “green” when seeing the
11 stimulus RED) while inhibiting their automatic tendency to read the word (“red”) [2].

12 Similarly, in the Eriksen flanker task, participants are asked to report the identity of a
13 target stimulus surrounded by multiple distractors while overcoming their automatic
14 tendency to respond instead to the distractors. Individual differences in the capacity for
15 cognitive control are highly predictive of academic achievement, interpersonal success,
16 and many other important life outcomes [3,4].

17 While exerting cognitive control improves people’s performance in these tasks, it
18 is also effortful and appears to be intrinsically costly [5,6]. The Expected Value of
19 Control (EVC) theory maintains that the brain therefore specifies how much control to
20 exert according to a rational cost-benefit analysis, weighing these effort costs against
21 attendant rewards for achieving one’s goals [7]. In broad accord with the predictions of
22 the EVC theory, previous research has found that control specification is context-
23 sensitive [8,9] and modulated by reward across multiple domains [10,11], such as

1 attention, response inhibition, interference control, and task switching. While previous
2 theories account for that fact that people's performance in these task is sensitive to
3 reward [7,12–14], it remains unclear *how* these dependencies arise from people's
4 experience. Recently, it has been proposed that the underlying mechanism is associative
5 learning [15,16]. Indeed, a number of studies have demonstrated that cognitive control
6 specification is plastic: whether people exert cognitive control in a given situation, which
7 controlled processes they employ, and how much control they allocate to them is learned
8 from experience. For instance, it has been demonstrated that participants in visual search
9 tasks gradually learn to allocate their attention to locations whose features predict the
10 appearance of a target [17], and a recent study found that learning continuously adjusts
11 how much cognitive control people exert in a Stroop task with changing difficulty [18].
12 Furthermore, it has been shown that people learn to exert more cognitive control after
13 their performance on a control-demanding task was rewarded [10] and learn to exert
14 more control in response to potentially control-demanding stimuli that are associated with
15 reward than to those that are not [11].

16 These studies provide evidence that people can use information from their
17 environment (e.g., stimulus features) to learn when to exert cognitive control and how to
18 exert control, and it has recently been suggested that this can be thought of in terms of
19 associative learning [15,16]. Other studies suggested that cognitive control can be
20 improved through training [19–21]. However, achieving transfer remains challenging
21 [22–25], the underlying learning mechanisms are poorly understood, and there is
22 currently no theory that could be used to determine which training regimens will be most
23 effective and which real-life situations the training will transfer to. Developing precise

1 computational models of the plasticity of cognitive control may be a promising way to
2 address these problems and to enable more effective training programs for remediating
3 executive dysfunctions and enabling people to pursue their goals more effectively.

4 In this article, we extend the EVC theory to develop a theoretical framework for
5 modeling the function and plasticity of cognitive control specification. This extension
6 incorporates recent theoretical advances inspired by the rational metareasoning
7 framework developed in the artificial intelligence literature [26,27]. We leverage the
8 resulting framework to derive the Learned Value of Control (LVOC) model which can
9 learn to efficiently select control signals based on features of the task environment. The
10 LVOC model can be used to simulate cognitive control (e.g., responding to a goal-
11 relevant target that competes with distractors) and, more importantly, how it is shaped by
12 learning. According to the LVOC model, people learn the value of different cognitive
13 control signals (e.g., how much to attend one stimulus or another). A key strength of this
14 model is that it is very general and can be applied to phenomena ranging from simple
15 learning effects in the Stroop task to the acquisition of complex strategies for reasoning
16 and problem-solving. In order to demonstrate the validity and generality of this model,
17 we show that it can capture the empirical findings of five cognitive control experiments
18 on the plasticity of visual attention [17], the interacting effects of reward and task
19 difficulty on the plasticity of interference control [10,11], and the transfer of such
20 learning to novel stimuli [8,9]. Moreover, the LVOC model outperforms alternate models
21 of such learning processes that rely only on associative learning or a basic win-lose-stay-
22 shift strategy. Our findings shed light on how learning and experience might shape
23 people's ability and propensity to adaptively control their minds and behavior, and the

1 LVOC model predicts under which circumstances these mechanisms might lead to self-
2 control failure.

3 Models

4 Formalizing the function of cognitive control

5 At an abstract level, all cognitive control processes serve the same function: to adapt
6 neural information processing to achieve a goal [28]. At this abstract level, neural
7 information processing can be characterized by the computations being performed, and
8 the extent to which the brain achieves its goals can be quantified by the expected utility
9 of the resulting actions. From this perspective, an important function of cognitive control
10 is to select computations so as to maximize the agent's reward rate (i.e., reward per unit
11 time). This problem is formally equivalent to the *rational metareasoning* [26,29] problem
12 studied in computer science: selecting computations so as to make optimal use of the
13 controlled system's limited computational resources (i.e., to achieve the highest possible
14 sum of rewards with a limited amount of computation).

15 Thus, rational metareasoning suggests that the specification of cognitive control is
16 a metacognitive decision problem. In reinforcement learning [30], decision problems are
17 typically defined by a set of possible actions, the set of possible states, an initial state, the
18 conditional probabilities of transitioning from one state to another depending on the
19 action taken by the agent, and a reward function. Together these five components define a
20 Markov decision process (MDP [30]). In a typical application of this framework the agent
21 is an animal, robot, or computer program, actions are behaviors (e.g., pressing a lever),
22 the state characterizes the external environment \mathcal{E} (e.g., the rat's location in the maze),

1 and the rewards are obtained from the environment (e.g., pressing a lever dispenses
2 cheese). In general, the agent cannot observe the state of the environment directly; for
3 instance, the rat running through a maze does not have direct access to its location but has
4 to infer this from sensory observations. The decision problems posed by an environment
5 that is only partially observable can be modelled as a partially observable MDP (POMDP
6 [31]). For each POMDP there is an equivalent MDP whose state encodes what the agent
7 knows about the environment and is thus fully observable; this is known as the belief-
8 MDP [31].

9 Critically, the belief-MDP formalism can also be applied to the choice of internal
10 computations [27] – such as allocating attention [32] or gating information into working
11 memory [33,34] – rather than only physical actions. In the rational metareasoning
12 framework, the agent is the cognitive control system whose actions are control signals
13 that specify which computations the controlled systems should perform. The internal state
14 of the controlled systems is only partially observable. We can formally define the
15 problem of optimal cognitive control specification as maximizing reward the in the meta-
16 level MDP

$$17 \quad M = (\mathcal{S}, s_0, \mathcal{C}, T, r), \quad (1)$$

18 where \mathcal{S} is the set of possible information states, comprising beliefs about the external
19 environment (e.g., the choices afforded by the current situation) and beliefs about the
20 agent’s internal state (e.g., the decision system’s estimates of the choices’ utilities), s_0
21 denotes the initial information state, \mathcal{C} is the set of possible control signals that may be
22 discrete (e.g., “Simulate action 1.”) or continuous (e.g., “Increase the decision threshold
23 by 0.175.” or “Suppress the activity of the word-reading pathway by 75%.”), T is a

1 transition model, and r is the reward function that cognitive control seeks to maximize.
2 The transition model specifies the conditional probability of transitioning from belief
3 state s to belief state s' if the control signal is c by $T(s, c, s')$. The meta-level reward
4 function r combines the utility of outcome X (of actions resulting from control signal c in
5 belief state s) with the computational cost associated with exerting cognitive control:

$$6 \quad r(s, c) = u(X) - \text{cost}(s, c), \quad (2)$$

7 where X is the outcome of the resulting action, u is utility function of the brain's reward
8 system, and $\text{cost}(s, c)$ is the cost of implementing the controlled process.

9 Within this framework, we can define a cognitive control strategy $\pi: \mathcal{S} \rightarrow \mathcal{C}$ as a
10 mapping from belief states $s \in \mathcal{S}$ to control signals $c \in \mathcal{C}$. The optimal cognitive control
11 strategy π^* is the one that always chooses the computation with the highest expected
12 value of computation (EVOC):

$$13 \quad \pi^*: s \mapsto \operatorname{argmax}_c \text{EVOC}(c, s). \quad (3)$$

14 The EVOC is the expected sum of computational costs and benefits of performing the
15 computation specified by the control signal c and continuing optimally from there on:

$$16 \quad \text{EVOC}(c, s) = Q_{\pi^*}(s, c) = \mathbb{E}[r(s, c) + V_{\pi^*}(S_{t+1}) \mid S_t = s, C_k = c, T], \quad (4)$$

17 where Q_{π^*} is known as the Q-function of the optimal control strategy π^* , and $V_{\pi^*}(S_{t+1})$
18 is the expected sum of meta-level rewards of starting π^* in state S_{t+1} .

19 In summary, cognitive control specification selects the sequence of cognitive
20 control signals that maximizes the expected sum of rewards of the resulting actions minus
21 the cost of the controlled process. The optimal solution to this problem is given by the
22 optimal control policy π^* .

So far, we have assumed that the cognitive control system chooses one control signal at a time, but \mathbf{c} could also be a vector comprising multiple control signals (e.g., one that increases the rate at which evidence is accumulated towards the correct decision by via an attentional mechanism and a second one that adjusts the decision threshold). Furthermore, overriding a habit by a well-reasoned decision also requires executing a coordinated sequence of cognitive operations for planning and reasoning. Instead of specifying each of these operations by a separate control signal, the cognitive control system might sometimes use a single control signal to instruct the decision system to execute an entire planning strategy. The rational metareasoning framework allows us to model cognitive strategies as options [35–38]. An option is a policy combined with an initiation set and a termination condition [38]. Options can be treated as if they were elementary computations and elementary computations can be interpreted as options that terminate after the first step. With this extension, the optimal solution to the cognitive control specification problem becomes

$$\pi^*(s) = \arg \max_{o \in \mathcal{O}} Q^*(s, o), \quad (5)$$

where the set of options \mathcal{O} may include control strategies and elementary control signals.

Critically, this rational metareasoning perspective on cognitive control covers not only simple phenomena, such as inhibiting a pre-potent automatic response in the Stroop task, but also more complex ones, such as sequencing one’s thoughts so as to follow a good decision strategy, and very complex phenomena such as reasoning about how to best solve a complex problem.

1 The LVOC model of the plasticity of cognitive control specification

2 The computations required to determine the expected value of control may themselves be
3 costly and time consuming. Yet, in some situations cognitive control has to be engaged
4 very rapidly, because maladaptive reflexes, impulses, and habitual responses have to be
5 inhibited before the triggered response has been executed. In such situations, there is
6 simply not enough time to compute the expected value of control on the fly. Fortunately,
7 this may not be necessary because an approximation to the EVOC can be learned from
8 experience. We therefore hypothesize that the cognitive control system learns to predict
9 the context-dependent value of alternative control signals. By understanding how this
10 learning occurs, we might be able to explain the experience-dependent changes in how
11 people use their capacity for cognitive control which we will refer to as the plasticity of
12 cognitive control specification. In addition to these systematic, experience-driven
13 changes cognitive control is also intrinsically variable. To model the plasticity and the
14 variability of cognitive control, this section develops a model that combines a novel
15 feature-based learning mechanism with a new control specification mechanism that
16 explores promising control signals probabilistically to accelerate learning which of them
17 is most effective.

18 The previous section characterized the problem of cognitive control specification
19 as a sequential meta-decision problem. This makes reinforcement learning algorithms
20 [39] a natural starting point for exploring how the cognitive control systems learns the
21 EVOC from experience. Approximate Q-learning appears particularly suitable because
22 the optimal control strategy can be expressed in terms of the optimal Q-function
23 (Equations 3-5). From this perspective, the plasticity mechanisms of cognitive control

1 specification serve to learn an approximation to the value $Q_t(s, c)$ of selecting control
2 signal c in state s based on one's experience with selecting control signals $\mathbf{c} =$
3 (c_1, \dots, c_t) in states $\mathbf{s} = (s_1, \dots, s_t)$ and receiving the meta-level rewards $\mathbf{r} = (r_1, \dots, r_t)$.
4 Learning an approximate Q-function Q_t from this information could enable the cognitive
5 control system to efficiently select a control strategy by comparing learned values rather
6 than reasoning about their effects.

7 Learning the optimal meta-level state-value function Q^* can be challenging
8 because the value of each control signal may depend on the outcomes of the control
9 signals selected afterwards. Furthermore, the state space of the meta-level MDP has a
10 very high dimensionality as it comprises all possible states that the controlled system
11 could be in. To overcome these challenges, a neural system like the brain might learn a
12 linear approximation to the meta-level state value function instead of estimating each of
13 its entries separately. Concretely, the cognitive control system might learn to predict the
14 value of selecting a control strategy (e.g., focusing on the presenting speaker instead of
15 attending to an incoming phone call) by a weighted sum of features of the internal state
16 and the current context (e.g. being in a conference room). For instance, the value $Q^*(s, c)$
17 of choosing control signal c in the internal state s can be predicted from the features
18 $f_k(s)$, the implied control signal intensities \mathbf{c} , their interactions with the features, that is
19 $f_k(s) \cdot c_i$, and their costs. Concretely, the EVOC of selecting control signal c in state s is
20 approximated by the *Learned Value of Control* (LVOC),

$$\begin{aligned}
& \text{LVOC}(s, c; \mathbf{w}) = w_0 + \left(\sum_{k=1}^K w_k^{(f)} \cdot f_k(s) \right) + \left(\sum_{l=1}^L w_l^{(c)} \cdot c_l \right) \\
& + \left(\sum_{k=1}^K \sum_{l=1}^L w_{k,l}^{(f \times c)} \cdot f_k(s) \cdot c_l \right) - \text{cost}(c) - w^{(T)} \cdot T, \quad (6)
\end{aligned}$$

where the weight vector \mathbf{w} includes the offset w_0 , the weights $\mathbf{w}_k^{(f)}$ of the states' features, the weights $\mathbf{w}^{(c)}$ of the control signal intensities, the weights $\mathbf{w}_{k,l}^{(f \times c)}$ of their interaction terms, the weight $w^{(T)}$ of the response time T , and $\text{cost}(c)$ is the intrinsic cost of control which scales with the amount of cognitive control applied to the task.

The optimal way to update the weights based on experience in a stationary environment is given by Bayes rule. Our model therefore maintains and continues to update an approximation to the posterior distribution

$$P(\mathbf{w} | e_{1,\dots,t}) \propto P(\mathbf{w} | e_{1,\dots,t-1}) \cdot P(e_t | \mathbf{w}), \quad (7)$$

on the weight vector \mathbf{w} given its experience $e_{1,\dots,t}$ up until the present time t , where each experience $e_i = (s_i, c_i, r_i, T_i, s_{i+1})$ comprises the state, the selected control signal, the reward, the response time, and the next state. In simple settings where a single control signal determines a single reward our model's learning mechanism is equivalent to Bayesian linear regression [40,41]. In more complex settings involving a series of control signals or delayed rewards the learning rule approximates the Bayesian update by substituting the delayed costs and benefits of control by the model's predictions. For more details, see S1 Text.

If the value of control is initially unknown, the optimal way to select control signals is to balance exploiting previous experience to maximize the expected immediate performance with exploring alternative control allocations that might prove even more

1 effective. Our model solves this dilemma by an exploration strategy similar to Thompson
 2 sampling: It draws k samples from the posterior distribution on the weights and averages
 3 them, that is

$$4 \quad \tilde{w}_1, \dots, \tilde{w}_k \sim P(\mathbf{w} | e_{1,\dots,t}), \quad \tilde{w} = \frac{1}{k} \cdot \sum_{i=1}^k \tilde{w}_i. \quad (8)$$

5 According to the LVOC model the brain then selects a control signal by maximizing the
 6 EVOC predicted by the average weight \tilde{w} , that is

$$7 \quad c_t \approx \arg \max_c \text{LVOC}(s_t, c; \tilde{w}). \quad (9)$$

8 Together, Equations 6-9 define the LVOC model of the plasticity of cognitive
 9 control. The LVOC model extends the EVC theory [7] which defines optimal control
 10 signals in terms of the EVOC (Equation 3), by proposing two mechanisms through which
 11 the brain might be able to approximate this normative ideal: learning a feature-based,
 12 probabilistic model of the EVOC (Equations 6-7) and selecting control signals by
 13 sampling from this model (Equations 8-9). This model is very general and can be applied
 14 to model cognitive control of many different processes (e.g., which location to saccade to
 15 vs. how strongly to inhibit the word-reading pathway) and different components of the
 16 same process (e.g., rate of evidence accumulation towards the correct decision vs. the
 17 decision threshold). The LVOC model's core assumptions are that the brain learns to
 18 predict the EVOC of alternative control specifications from features of the situation and
 19 the control signals, and that the brain then probabilistically selects the control
 20 specification with the highest predicted value of control. Both of these components could
 21 be implemented by many different mechanisms. For instance, instead of implementing
 22 the proposed approximation to Bayesian regression, the brain might learn to predict the

1 EVOC through the rewarded-modulated associative plasticity mechanism outlined in the
2 SI. We are therefore not committed to the specific instantiation we used for the purpose
3 of the simulations reported below (Equations 7-9).

4 The LVOC model instantiates the very general theory that the brain learns how to
5 process information via metacognitive reinforcement learning. This includes not only the
6 plasticity of cognitive control but also how people might discover cognitive strategies for
7 reasoning and decision-making and how they learn to regulate their mental activities
8 during problem solving. As a proof of concept, the following sections validate the LVOC
9 model against five experiments on the plasticity of attention and interference control.

10 [Alternative models: Associative learning and Win-Stay Lose-Shift](#)

11 In principle, the control-demanding behavior considered in this paper could result from
12 simpler mechanisms than the ones proposed here. In this section, we consider two simple
13 models that we use as alternatives to compare against the more complex LVOC model.

14 The first model relies on the assumption that the plasticity of cognitive control can be
15 understood in terms of associative learning [15,16]. We therefore evaluate our model
16 against an associative learning model based on the Rescorla-Wagner learning rule [42].

17 This model forms stimulus-control associations based on the resulting reward. The
18 association $A_{s,c}$ between a stimulus s and a control signal c is strengthened when it is
19 accompanied by (intrinsic or extrinsic) reward and weakened otherwise. Concretely, the
20 association strengths involving the chosen response were updated according to the
21 Rescorla-Wagner rule, that is

22
$$A_{s,c} = A_{s,c} + \alpha \cdot \left(R - \sum_s I_s \cdot A_{s,c} \right), \quad (10)$$

1 where α is the learning rate, R is the reward and the indicator variable I_s is 1 when the
2 stimulus s was present and 0 else. Given the learned associations, the control signal is
3 chosen probabilistically according to the exponentiated Luce's choice rule, that is each
4 control signal c is selected with probability

$$5 \quad p(c) = \frac{\exp(A_{s,c})}{\sum_c \exp(A_{s,c})}. \quad (11)$$

6 The second alternative model is based on previous research suggesting that people
7 sequentially adjust their strategy through a simple Win-Stay Lose-Shift mechanism [43].
8 This mechanism starts with a random strategy on the first trial, and on each subsequent
9 trial it either repeats the previous strategy when it was successful or switches to a
10 different strategy when the current strategy failed. Here, we apply this idea to model how
11 the brain learns which control signal to select. Concretely, our WSLS model repeats the
12 previous control signal (e.g., "Attend to green.") when it leads to a positive outcome
13 (Win-Stay) and randomly selects a different control signal (e.g., "Attend to red.")
14 otherwise (Lose-Shift).

15 Simulations of learning and transfer effects in cognitive control paradigms

16 To evaluate the proposed models, we used them to simulate the plasticity of
17 attentional control in a visual search task [17] as well as learning and transfer effects in
18 Stroop and Flanker paradigms [8–11]. Table 1 summarizes the simulated phenomena and
19 how the LVOC model explains each at a conceptual level.

20 Learning to control visual attention

21 Previous research has shown that how people allocate their attention is shaped by
22 learning [8–11,15–17]. For instance, Lin and colleagues [17] had participants perform a

1 visual search task in which they gradually learned to allocate their attention to locations
2 whose color predicted the appearance of the target (Figure 1a). In this task, participants
3 viewed an array of four rotated letters (one T and three L's), each encompassed by a
4 different colored circle. They were instructed to report the orientation of the T. The
5 circles appeared before the letters allowing participants to allocate their attention by
6 saccading to a promising location before the letters appeared. In the training phase, the
7 target always appeared within the green circle, but in the test phase it was equally likely
8 to appear in any of the four circles.

9 Visual search entails sequentially allocating cognitive control to different
10 locations based on their visual features. Since attention is can be understood as an
11 instance of cognitive control, this problem is naturally modeled as a meta-level MDP. We
12 therefore applied our LVOC theory to predict the dynamics and consequences of learning
13 which locations to attend to based on their features (the colored circles) in this paradigm.
14 Since the stimuli were presented along a circle, approximating locations people might
15 naturally attend around a clock, we assumed that the control signal $c \in \{1,2,3, \dots, 12\}$
16 specifies which of 12 locations to attend, the state s_t encodes which of the 12 locations
17 were highlighted by a colored circle (see Figure 1a), the circles' colors, the unknown
18 position of the target, and the list of locations that have already been inspected on the
19 current trial. Since, the set of possible control signals is small, our simulation assumes
20 that the brain always finds the control signal that maximizes the predicted EVOC
21 (Equation 9).

22 The features $f(s, c)$ encode only observable aspects of the state s that are relevant
23 to the value of the control signal c . Concretely, our simulations assumed that the features

1 encode whether the attended location was highlighted by a colored circle, the color of
2 that circle (one binary indicator variable for each possible color), its position (by one
3 binary indicator variable for each of the four possible locations), and whether or not it has
4 been attended before. To capture people's prior knowledge that attending a location a
5 second time is unlikely to provide new information, we set the prior on the weight of the
6 last feature to -1 ; this captures the well-known inhibition of return mechanism in visual
7 attention [44]. For all other features the mean of the prior on the weights was 0. Based on
8 the results reported by [17], we modeled reaction times as the sum of a non-decision time
9 of 319ms and a decision-time of 98ms per attended location. Our simulation assumed that
10 people incur a fixed cost ($r(s, c) = \text{cost}(c) = -1$ for all $c \in \{1, 2, 3, 4\}$) every time they
11 deploy their attention to a location. For simplicity, we assume that in this simple task
12 people always search until they find the target and that when they attend to a location
13 they always recognize the presence/absence of the target and respond accordingly.
14 Hence, the intuition that people should try to find the target with as few saccades as
15 possible follows directly from the objective of maximizing the sum of meta-level
16 rewards. Applied to this visual search task, the LVOC model offers a mechanism for how
17 people learn where to allocate their attention based on environmental cues in order to find
18 the target as quickly as possible.

19 Our associative learning model assumed that finding the target yields an intrinsic
20 reward of $+1$ and no reward or cost otherwise. The responses C were saccades to one of
21 the 12 locations. The stimuli S comprised indicator variables for each of the four colors,
22 the absence of a circle, and whether the location had been inspected before, and one
23 feature that was always 1. To capture the inhibition of return, the reward associations

1 with the stimulus-feature indicating that a location had been attended previously were
2 initialized by -1 . All other association strengths were initialized as 0 and the learning
3 rate of the Rescorla-Wagner model to the data from [17] using maximum-likelihood
4 estimation.

5 Learning and transfer effects in inhibitory control

6 In Stroop and Flanker paradigms, the cognitive control strategy o is defined by a
7 single control signal $c \in [0,1]$ serves to bias processing away from an automatic
8 mechanism. Following the classic model by Cohen and colleagues [45], we assume that
9 control signals determine the relative contribution of the automatic versus the controlled
10 process to the drift rate d at which evidence is accumulated towards the controlled
11 response [46]:

$$12 \quad d = c \cdot d_{\text{controlled}} + (1 - c) \cdot C \cdot d_{\text{automatic}}, \quad (12)$$

13 where $d_{\text{controlled}}$ and $d_{\text{automatic}}$ are the drift rates of the controlled and the automatic
14 process respectively, and $C = 1$ when the trial is congruent or -1 when the trial is
15 incongruent. The drift rates, in turn, affect the response and response time according to a
16 drift-diffusion model [46]. When the decision variable exceeds the threshold $+\theta$, then the
17 response agrees with the controlled process (equivalent to the correct response for these
18 tasks). When the decision variable falls below $-\theta$, the response is incorrect. To capture
19 sources of error outside of the evidence accumulation process (e.g., motor execution
20 errors), our simulation assumes that people accidentally give the opposite of their
21 intended response on a small fraction of trials ($p_{\text{flip}} < 0.05$).

We model the selection of continuous control signals as a gradient ascent on the EVOC predicted by Thompson sampling (Equations 8-9). Concretely, continuous control signals are selecting by repeatedly applying the update rule

$$c \leftarrow c + \eta \cdot \frac{d \text{LVOC}(s, c; \tilde{w})}{dc} \quad (13)$$

until the change in the Euclidean norm of the control signal intensity vector is less than δ_c . This mechanism starts from the control signal deployed on the previous trial and thereby captures the inertia of control specification and the resulting reconfiguration cost [46], consistent with the task-set inertia hypothesis [47]. Furthermore, it predicts that control intensities are adjusted gradually and continually, thereby allowing control to be exerted while the optimal control signal is still being determined. This feature of our model makes the intuitive prediction that time pressure might reduce the magnitude of control adjustment [cf. 48,49].

We model the cost associated with a continuous control signal c as the sum of the control cost required to exert that amount of control ($\text{cost}(c)$) and the opportunity cost of executing the controlled process ($\omega \cdot t$), that is

$$\text{cost}(s, c) = \omega \cdot t + \text{cost}(c), \quad (14)$$

where ω is the opportunity cost per unit time¹, t is the duration of the controlled process, and $\text{cost}(c)$ is the intrinsic cost of exerting the control signal c . While the first term captures that goal-directed control processes, such as planning, can take significantly longer than automatic processes, such as habits, the second term captures that due to interference between overlapping pathways the cost of a control signal increases with its

¹ In many real-world scenarios and some experiments, the opportunity cost is time-varying. This can be incorporated into our model by adding a learning mechanism that estimates ω from experience [54,86].

1 intensity [13] even when control intensity accelerates the decision process [11,50,51].

2 Following [46], we model the intrinsic cost of control as the implementation cost

3
$$\text{cost}(c) = \exp(a_i \cdot \|c\| + b_i), \quad (15)$$

4 where c is the control signal, a_i specifies how rapidly control cost increases with control

5 intensity, and b_i determines the lowest possible cost.² The monotonic increase of control

6 cost with control signal intensity expressed by this equation models that the more

7 intensely you focus on one process, say color-naming, the less you are able to do other

8 valuable things, such as verbal reasoning. This cognitive opportunity cost of control is a

9 consequence of overlap between neural pathways serving different functions [13,52,53].

10 In all of our simulations, the number of samples drawn from the posterior
11 distribution on the weights was $k = 2$. For simplicity, we modeled control allocation in
12 each trial of the Stroop and Flanker tasks simulated below as an independent, non-
13 sequential, metacognitive control problem. The opportunity cost of time (ω in Equation
14 14) was set to \$8/h [cf. 51]. Model parameters were fitted by maximum likelihood
15 estimation using Bayesian optimization [55]. The drift rates of the controlled and the
16 automatic process (Equation 12) were determined from people's response times on
17 neutral trials. The model's prior precision on the weights was set to assign 95%
18 confidence to the EVOC of a stimulus lying between the equivalent of ± 5 cents per
19 second.

20 In the color-word Stroop task by Krebs et al. [11] the participant's task was to
21 name the font color of a series of color words which were either congruent or incongruent

² This framework can easily be extended to also other sources of control costs, such as reconfiguration costs [46].

with the word itself (Figure 2a). For two of the four colors, giving the correct response yielded a monetary reward whereas responses to other two colors were never rewarded. Our simulation of this experiment assumed that people represent each stimulus by a list of binary features that encode the presence of each possible color and each possible word independently but do not encode their combinations. To capture the contribution of the experiment's financial incentives for correct responses, we assumed that the utility $u(X)$ in Equation 2 is the sum of the financial reward and the intrinsic utility of getting it right [56], that is

$$u(\text{correct}) = r_{\text{external}} + r_{\text{intrinsic}}, \quad (16)$$

$$u(\text{wrong}) = -(r_{\text{external}} + r_{\text{intrinsic}}), \quad (17)$$

where the monetary reward r_{external} was 10 cents on rewarded trials and zero otherwise. The non-decision time was set to 300ms. The implementation cost parameters (a_i and b_i in Equation 15), the probability of accidental response flips (p_{flip}), the intrinsic reward $r_{\text{intrinsic}}$ of responding correctly (Equations 16-17), and the noise parameter σ of the drift diffusion model (Equation 12) were fit to the empirical data shown in Figure 2b-c. To enable a fair comparison, we gave the associative learning model and the Win-Stay Lose-Shift model degrees of freedom similar to those of the LVOC model by adding parameters for the intrinsic reward of being correct, the probability of response error, and the noise of the drift-diffusion process. In addition, the Rescorla-Wagner model was equipped with a learning rate parameter. Each model was fitted using maximum-likelihood estimation using the Bayes adaptive direct search algorithm [57].

In the Flanker task by Braem et al. [10], participants were instructed to name the color of a central square (the *target*) flanked by two other squares (*distractors*) whose

1 color was either the same as the color of the target (congruent trials) or different from it
 2 (incongruent trials) (Figure 3a). On a random 25% of the trials, responding correctly was
 3 rewarded and on the other 75% of the trials it was not. Our simulation assumed that
 4 people predict the EVOC from two features that encode the presence of conflict and
 5 congruency respectively: The conflict feature was +1 when the flankers and the target
 6 differed in color and zero otherwise. Conversely, the value of the congruency feature was
 7 +1 when the flankers had the same color as the target and zero else. To capture that
 8 people exert more cognitive control when they detect conflict [58], the prior mean on
 9 these weights was +1 for the interaction between control signal intensity and
 10 incongruence and -1 for the interaction between control signal intensity and congruence.
 11 Providing our model with these features instantiates our assumption that in the Flanker
 12 task perception is easy but response inhibition can be challenging. In other words, our
 13 model assumes that errors in the Flanker arise from the failure to translate three correct
 14 percepts into one correct response by inhibiting the automatic responses to the other two.
 15 Furthermore, the incongruency feature can also be interpreted as a proxy for the resulting
 16 response conflict that is widely assumed to drive the within-trial adjustment of control
 17 signals in the Flanker task [58].

18 Our simulation assumed that people only learn on trials with feedback. The effect
 19 of control was modelled as inhibiting the interference from the flankers according to

$$d = d_{\text{target}} + (1 - c) \cdot C \cdot d_{\text{flankers}} , \quad (18)$$

21 where $C = 1$ if the distractors are congruent and $C = -1$ when they are incongruent.

22 The drift rates for accumulating information from the target (d_{target}) and the distractors
 23 (d_{flankers}) were assumed to be identical. Their value was fit to the response time for color

1 naming on rewarded neutral trials reported in [11], and the non-decision time was
2 300ms. The perceived reward value of the positive feedback was determined by
3 distributing the prize for high performance (EUR 10) over the 168 rewarded trials of the
4 experiment ($z = 7.5$ US cents per correct response). Braem et al. [10] found that the
5 effect of reward increased with people's reward sensitivity. To capture individual
6 differences in reward sensitivity, we modelled people's subjective utility by

$$7 \quad u(\text{correct}) = z^\alpha + r_{\text{intrinsic}}, \quad (19)$$

$$8 \quad u(\text{wrong}) = -r_{\text{intrinsic}}, \quad (20)$$

9 where for $z \geq 0$ is the payoff and $\alpha \in [0,1]$ is the reward sensitivity. The reward
10 sensitivity was set to 1, and the intrinsic reward of being correct ($r_{\text{intrinsic}}$), the standard
11 deviation of the noise (σ), the threshold of the drift-diffusion model (θ), the
12 implementation cost parameters (a_i, b_i) were fit to the effects of reward on the reaction
13 times on congruent trials (Figure 3c), the average reaction time, and the effect of reward
14 sensitivity on conflict adaptation reported by [10]. The probability of accidentally giving
15 the opposite of the intended response was set to zero.

16 To enable a fair comparison between LVOC model and the two simpler models,
17 we equipped the associative learning model and the Win-Stay Lose-Shift model with the
18 same assumptions and degrees of freedom as the LVOC model. Equivalently to the
19 LVOC model of this task, they included a bias against exerting control was instantiated
20 by an association of -1 between either stimulus feature and control exertion. The effect of
21 control was modeled using the same drift-diffusion model with same set of free
22 parameters, and like the LVOC model they also included a free parameter for the intrinsic
23 reward of being correct and the probability of response error. Furthermore, these models

1 included a free parameter for the cost of exerting control that is equivalent to two
2 parameters of the LVOC model's implementation and reconfiguration cost parameters,
3 because their control signal was either 1 or 0. Furthermore, the Rescorla-Wagner model
4 included an additional parameter for its learning rate, giving it the same number of
5 parameters as the LVOC model.

6 Experiment 2 by Bugg et al. [8] asked participants to name the color of Stroop
7 stimuli like those used by Krebs et al. [11]. Critically, some of the color words were
8 printed in color that appeared on congruent trials 80% of the time whereas other color
9 words were printed in a color that appeared on incongruent trials 80% of the time (Figure
10 4a). Each word was written either in cursive or standard font. We modeled the stimuli by
11 four binary features indicating the presence of each of the four possible words (1 if the
12 feature is present and 0 otherwise), and a fifth feature indicating the font type (0 for
13 regular and 1 for cursive). The non-decision time was set to 400ms. Since there were no
14 external rewards for good performance, the utility of correct/incorrect responses was
15 $\pm r_{\text{intrinsic}}$. The implementation cost parameters (a_i and b_i), the probability of accidental
16 response flips (p_{flip}), the intrinsic reward of being correct ($r_{\text{intrinsic}}$), and the standard
17 deviation of the noise (σ) were fit to the empirical data shown in Figure 4b and Figure 4c.
18 Given these parameters, the drift rate for color naming and reading were determined to
19 match the reaction times on unrewarded neutral trials reported in [11].

20 Finally, Bugg et al. [9] presented their participants with pictures of animals
21 overlaid by animal names (Figure 4d). The participants' task was to name the animal
22 shown in the picture. Critically, for some animals, the picture and the word were usually
23 congruent whereas for other the picture and the word were usually incongruent. The

1 training phase was followed by a test phase that used novel pictures of the same animal
2 species. We modelled this picture-word Stroop by representing each stimulus by a vector
3 of binary indicator variables. Concretely, our representation assumed one binary indicator
4 variable for each word (i.e., BIRD, DOG, CAT, FISH) and one indicator variable for
5 each image category (i.e., bird, dog, cat, fish). The non-decision time was set to 400ms.
6 The implementation cost parameters (a_i and b_i), the intrinsic reward of being correct
7 ($r_{\text{intrinsic}}$), the standard deviation of the noise (σ), and the probability of accidental
8 response flips (p_{flip}) were fit to the empirical data shown in Figure 4e-f. Given these
9 parameters, the drift rate for word reading was fit as above and the drift rate for picture
10 naming was fit to a response time of 750ms.

	Phenomenon	Explanation of the LVOC model
Lin et al. (2016), Exp. 1	In the training block, participants learn to find the target increasingly faster when it always appears in a location with a certain color. In the test block, participants are significantly slower on trials that violate this regularity.	People learn to predict the value of attending to different locations from their color.
Krebs et al. (2010), Exp. 1	People come to name the color of incongruent words faster and more accurately for colors for which performance is rewarded.	People learn to predict the value of increasing control intensity from the color of the word.
Braem et al. (2012), Exp. 1	On a congruent Flanker trial, people are faster when the previous trial was rewarded and congruent than when it was unrewarded and congruent, but the opposite holds when the previous trial was incongruent. These effects are amplified in people with high reward sensitivity.	People learn to exert more control on incongruent trials. Thus, rewarded incongruent trials tend to reinforce higher control signals while rewarded congruent trials tend to reinforce low control signals. Thus, people increase control after the former and lower control after the latter.
Bugg et al. (2008), Exp. 2	People become faster and more accurate at naming the color of an incongruently colored word when it is usually incongruent than when it is usually congruent.	People learn that exerting more control is more valuable when the color or word is predictive of incongruence.

Bugg et al. (2011), Exp. 2	People are faster at naming animals in novel, incongruently labelled images when that species was mostly incongruently labelled in the training phase than when it was mostly congruently labelled.	People learn that exerting more control is more valuable when the semantic category of the picture is predictive of incongruence.
----------------------------	---	---

Table 1: The core assumption of the LVOC model explains the learning effects observed in five different cognitive control experiments.

a) Visual search task of Lin et al. (2016)

T: left or right rotated?

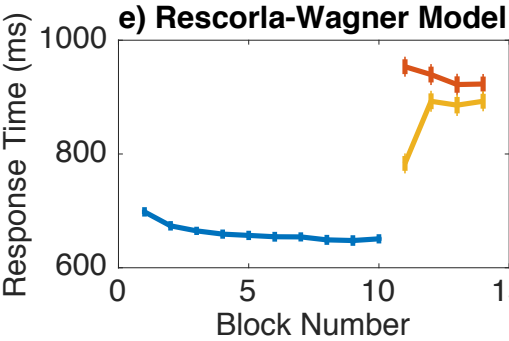
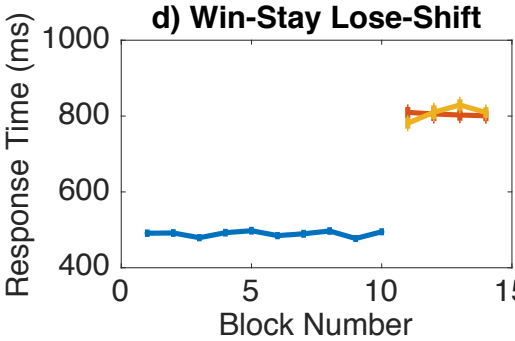
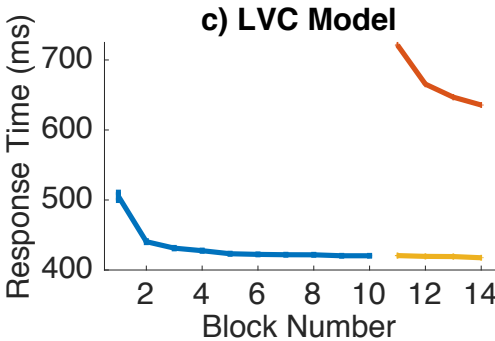
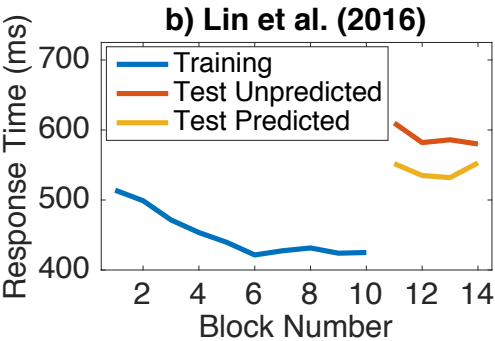
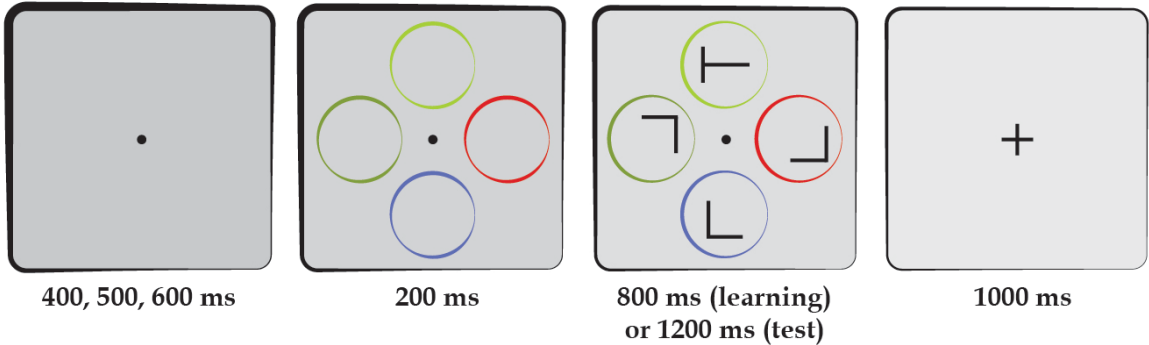


Figure 1: Learning to control the allocation of attention. a) Visual search task used by Lin et al. (2016). b) Predictions of the LVOC model. c) Human data from Experiment 1 of Lin et al. (2016). d) Fit of Win-Stay Lose-Shift model. e) Fit of Rescorla-Wagner model.

1 Results

2 We found that our model correctly predicted the learning effects observed in five
3 different cognitive control experiments by virtue of its fundamental assumption that
4 people reinforcement-learn to predict the value of potential control signals and control
5 signal intensities from situational features (see Table 1). The following sections describe
6 these findings in detail.

7 Plasticity of attentional control in visual search

8 Lin et al. [17] had participants perform a visual search task for which the target of
9 attention could either be predicted (training and predictable test trials) or not
10 (unpredictable test trials) (Figure 1a). For this task, given its core reinforcement learning
11 assumption (Table 1), the LVOC model predicts that 1) people should learn to attend to
12 the circle with the predictive color and thus become faster at finding the target over the
13 course of training, 2) continue to use the learned attentional control strategy in the test
14 block and hence be significantly slower when the target appears in a circle of a different
15 color during the test block, and 3) gradually unlearn their attentional bias during the test
16 block (Figure 1c). As shown Figure 1b, all three predictions were confirmed by Lin and
17 colleagues [17].

18 We compared the performance of LVOC to two plausible alternative models of
19 these control adjustments: a Win-Stay Lose-Shift model and a simple associative learning
20 model based on the Rescorla-Wagner learning rule. We found that the Win-Stay Lose-
21 Shift model failed to capture that people's performance improved gradually during
22 training, and it also failed to capture the difference between people's response times to
23 predicted versus unpredicted target locations in the test block (see Figure 1d). As Figure

1 le shows, the fit of the associative learning model (estimated learning rate: 0.0927)
2 captures that after learning to exploit the predictive regularity in the training block
3 participants were significantly slower in the test block. However, this simple model
4 predicted significantly less learning induced improvement and significantly slower
5 reaction times than was evident from the data by [17]. A quantitative model comparisons
6 using the Bayesian Information Criterion [59,60] provided very strong evidence that the
7 LVOC model explains the data by [17] better than the Rescorla-Wagner model or the
8 Win-Stay Lose-Shift model ($BIC_{LVOC} = 1817.8$, $BIC_{RW} = 9763.2$, $BIC_{WSLS} = 3449.9$).
9 This reflects that our model was able to accurately predict the data from [17] without any
10 free parameters being fitted to those data. In conclusion, findings suggest that the LVOC
11 model correctly predicted essential learning effects observed by [17] and explains these
12 data significantly better than a simple associative learning model and a Win-Stay Lose-
13 Shift model.

14 To more accurately capture both the slow improvement in the training block and
15 the rapid unlearning in the test block simultaneously, the LVOC model could be extended
16 by including a mechanism that discounts what has been learned or increases the learning
17 rate when a change is detected [61,62]. Next, we evaluate the LVOC model against
18 empirical data on the plasticity of inhibitory control.

19 Plasticity of Inhibitory Control

20 We found that our model can capture reward-driven learning effects in Stroop and
21 Flanker tasks, as well as how people learn to adjust their control allocation based on
22 features that predict incongruence and the transfer of these learning effects to novel
23 stimuli. In each case, the LVOC model captured the empirical phenomenon more

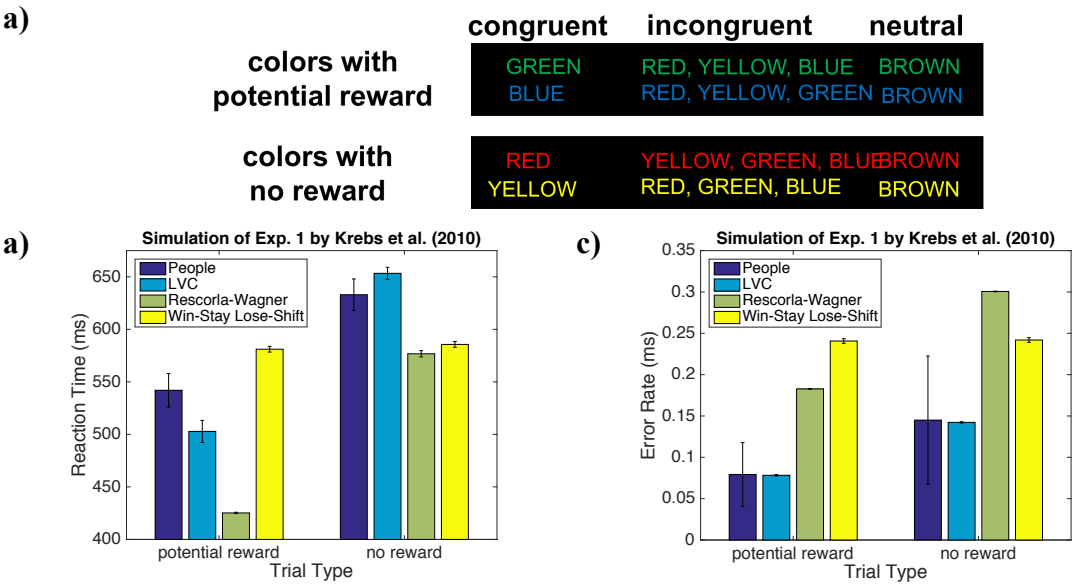
1 accurately than either a simple Win-Stay Lose-Shift model or a simple associative
 2 learning model. The following two sections present these results in turn.

3

4 **Table 2: Model parameters used in the simulations of empirical findings.**

	a_i	b_i	θ	σ	$r_{\text{intrinsic}}$	p_{flip}
Krebs, et al. (2010)	1.60	-0.01	3	0.05	1.60¢	3.5%
Braem et al. (2012)	4.17	-2	2.75	5	4.17¢	0.8%
Bugg et al. (2008)	1.95	-2.1	2.65	3.01	3.89¢	0.4%
Bugg et al. (2011)	5	-2	2.75	3	18.00¢	0.8%

5
6



7 Figure 2: LVOC model captures that in the paradigm by Krebs et al. (a) people learn to exert more
 8 cognitive control on stimuli whose features predict that performance will be rewarded which manifests in
 9 faster responses (b) and fewer errors (c).

10 Reward-driven plasticity in interference control

11 People learn to allocate more control on rewarded trials

12 To determine whether our integrated theory captures the reward-modulated
 13 plasticity of cognitive control specification, we used the LVOC model to simulate two
 14 sets of experiments that examined the influences of reward on cognitive control. Krebs et

al. [11] found that participants in a color-word Stroop task learned to respond faster and more accurately to incongruently colored color words when their color predicted that performance would be rewarded than when it predicted that performance would not be rewarded. We found that our model can capture these effects with reasonable parameter values (see Table 2). Figure 2 shows that our model captures Krebs et al.'s finding that people learn to exert more control on trials with rewarded colors than on trials with unrewarded colors even though they were interspersed within the same block. Concretely, our model captured that people become faster (541 ± 7 ms vs. 691 ± 8 ms; $t(959998) = -14.6, p < 10^{-15}$) and more accurate ($4.9 \pm 0.02\%$ errors vs. $11.8 \pm 0.03\%$ errors; $t(959998) = -164.7, p < 10^{-15}$) when the color of the word is associated with reward. Critically, the qualitative effects observed in this experiment follow logically from the core assumption of the LVOC model (see Table 1).

We compared the LVOC's model performance to that of an associative learning model with equivalent parameters (see Methods); the maximum likelihood estimates of these parameters were $\alpha = 0.0447$ for the learning rate, $r_{\text{intrinsic}} = 0.1811$ for the intrinsic reward, $\sigma_{\epsilon} = 0.1525$ for the noise of the drift-diffusion process, and $p_{\text{error}} = 0.1799$. While the Rescorla-Wagner model was able to qualitatively capture the effect of potential reward on reaction time and error rate, its quantitative fit was far worse than the fit of the LVOC model (see Figure 2b and Figure 2c); thus, a quantitative model comparison controlling for the number of parameters provided very strong evidence for the LVOC model over the Rescorla-Wagner model ($\text{BIC}_{\text{LVOC}} = 45.3$ vs. $\text{BIC}_{\text{RW}} = 1333.9$). We also fitted the Win-Stay Lose-Shift model and its parameter estimates were $r_{\text{intrinsic}} = 0$ for the intrinsic reward, $\sigma_{\epsilon} = 0$ for the noise of the drift-diffusion process,

1 and $p_{\text{error}} = 0.07$). We found that the WSLS model was unable to capture the effect of
2 reward on response times and error rates (see Figure 2b and Figure 2c) because its control
3 signals are uninformed by the stimulus presented on the current trial. Consequently, a
4 formal model comparison provided strong evidence for the LVOC model over the Win-
5 Stay Lose-Shift model ($\text{BIC}_{\text{LVOC}} = 45.3$ vs. $\text{BIC}_{\text{WSLS}} = 2454.8$).

6 Reward accelerates trial-by-trial learning of how to allocate control

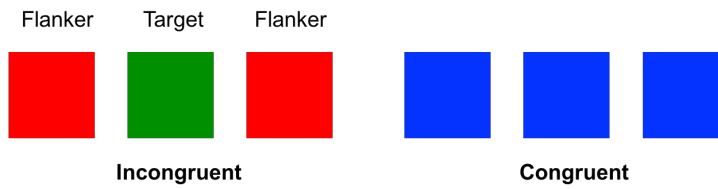
7 Braem et al. [10] found that participants in their Flanker task allocated more
8 cognitive control after rewarded incongruent trials than after rewarded congruent trials or
9 unrewarded trials. As Figure 3b shows, the LVOC model can capture this reward-induced
10 conflict-adaptation effect with a plausible set of parameters (see Table 2). Our model
11 correctly predicted that people's responses on congruent trials are faster when they are
12 preceded by rewarded congruent trials than when they are preceded by rewarded
13 incongruent trials. The predicted difference (7 ms) was smaller than the empirically
14 observed difference (27 ms) but it was statistically significant ($t(99) = 37.99$, $p <$
15 10^{-15}). According to our model, people learn to exert more control on incongruent trials
16 than on congruent trials. Furthermore, being rewarded for exerting a low level of control
17 reduces the control intensity on the subsequent trial, whereas being rewarded for exerting
18 a high level of control increases the control intensity on the subsequent trial. Thus, our
19 model predicts that control intensity should increase after rewarded incongruent trials but
20 decrease after rewarded congruent trials. On congruent trials, more control leads to
21 slower responses because it inhibits the facilitating signal from the flankers (Equation
22 18). This suggests that our model's metacognitive reinforcement learning mechanism
23 correctly predicts the findings of Braem et al. [10] (see Figure 3 and Table 1).

1 The LVOC model's learning increases with the magnitude of the reward.
2 Consequently, the LVOC model predicts that the effect shown in Figure 3b should
3 increase with people's reward sensitivity. Concretely, as we increased the reward
4 sensitivity parameter α from 0 to 1, the predicted reward-driven effect of conflict
5 monotonically increased from 0.6ms to 8.0ms ($t(102) = 4.77, p < 0.0001$). Consistent
6 with this prediction, Braem and colleagues [10] found that the magnitude of reward-
7 driven conflict adaptation effect increased with people's reward sensitivity, suggesting
8 that the reward experienced for exerting cognitive control was the driving force of their
9 adjustments. The significant positive correlation between people's reward sensitivity and
10 the magnitude of their conflict adaptation effect reported by [10] confirms our model's
11 prediction. Our model captures all of these effects because it learns to predict the
12 expected rewards and costs of exerting control from features of the situation and
13 probabilistically chooses the control signal that achieves the best cost-benefit tradeoff.

14 We compared the LVOC's fit to these behaviors with the associative learning and
15 Win-Stay Lose-Shift models. Even though the associative learning model had the same
16 number of parameters as the LVOC model, its fit was substantially worse than the fit of
17 the LVOC model (mean squared errors: $MSE_{RW} = 3.33$ vs. $MSE_{LVOC} = 1.35$), and its
18 best fit (Figure 3d) failed to capture the qualitative effect shown in Figure 3c. Finally, we
19 evaluated a Win-Stay Lose-Shift model. This model was equipped with the same set of
20 parameters as our Rescorla-Wagner model except for the learning rate parameter. We
21 found that the Win-Stay Lose-Shift model was unable to capture the data by Braem et al.
22 Figure 3e) because it stays with the controlled process forever once it has been rewarded
23 for using it ($MSE_{WSLS} = 19.4$). Taken together with the previous results, this suggests

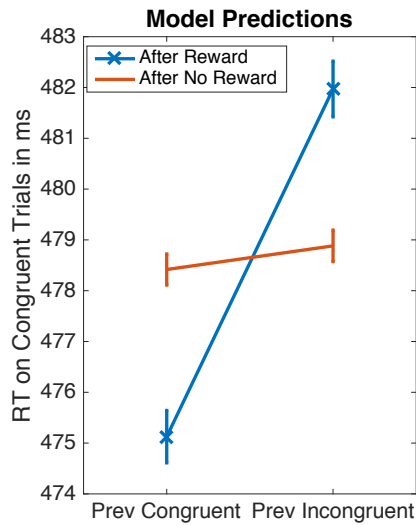
- 1 that the simple mechanisms assumed by the associative learning model and the Win-Stay
- 2 Lose-Shift model are insufficient to explain the complexity of cognitive control plasticity,
- 3 but the LVOC model can capture it.

1 a)

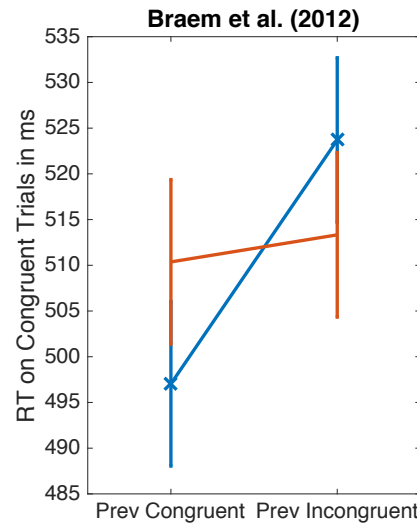


2

3 b)

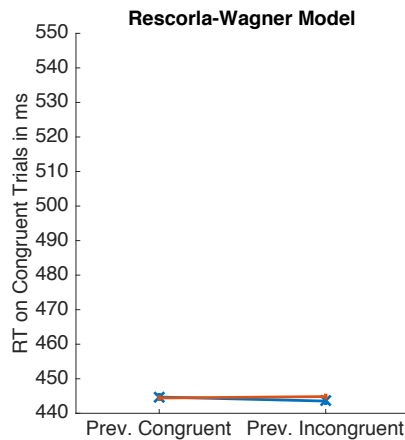


c)



4

5 d)



e)

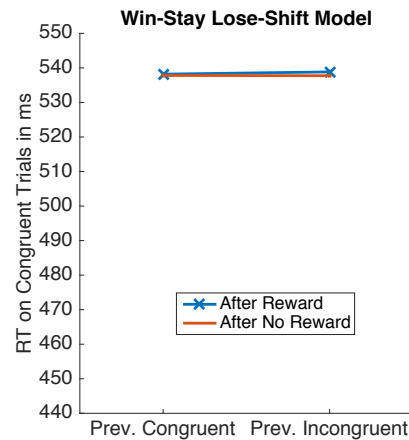
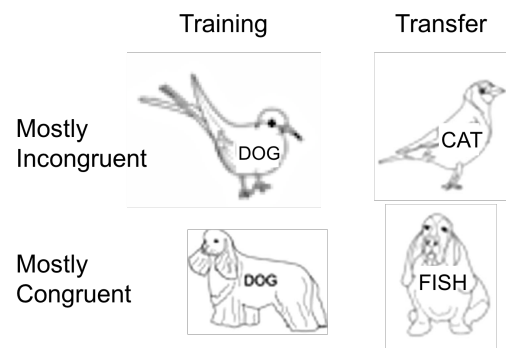


Figure 3: Metacognitive reinforcement learning captures the effect of reward on learning from experienced conflict observed by Braem et al. (2012). a) Illustration of the Flanker task by Braem et al. (2012). b) Fit of LVOC model. c) Human data by Braem et al. (2012). d) Fit of Rescorla-Wagner model. e) Fit of Win-Stay Lose-Shift model.

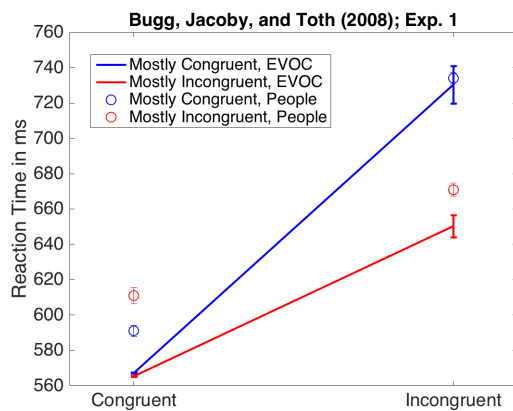
a) Bugg et al. (2008)

Mostly Congruent vs. Mostly Incongruent	
80% BLUE	80% GREEN
20% BLUE	20% GREEN

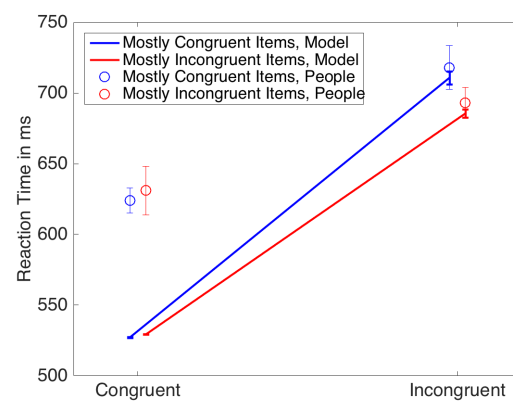
d) Bugg et al. (2011)



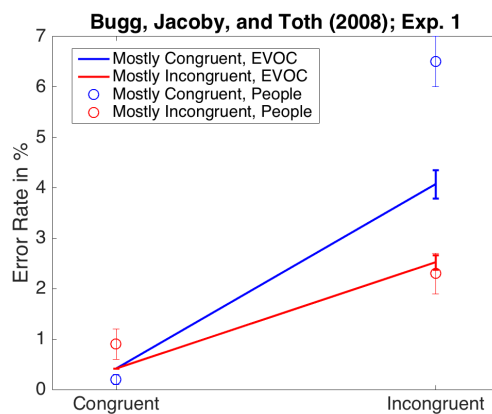
b)



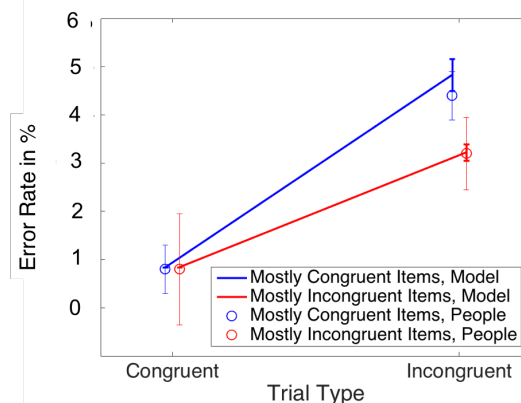
e)



c)



f)



1 Figure 4: The LVOC model captures that people learn to adjust their control intensity based on features that
2 predict incongruence. a) Color-Word Stroop paradigm by Bugg et al. (2008). b-c) LVOC model captures
3 that people learn to exploit features that predict incongruency to respond faster and more accurately on
4 incongruent trial. d) Picture-Word Stroop paradigm by Bugg, Jacoby, and Chanani (2011). e-f) The LVOC
5 model captures that the learning effects enabling people to respond more quickly and more accurately to
6 incongruently labelled images of animals whose labels are usually incongruent transfer to novel images.

1 Transfer of learning effects in interference control

2 The expected value of computation depends not only on the rewards for correct
3 performance but also on the difficulty of the task. In easy situations, such as the
4 congruent trials of the Stroop task, the automatic response can be as accurate, faster, and
5 less costly than the controlled response. In cases like this, the expected value of exerting
6 control is less than the EVOC of exerting no control. By contrast, in more challenging
7 situations, such as incongruent Stroop trials, the controlled process is more accurate and
8 therefore has a positive EVOC as long as accurate performance is sufficiently important.
9 Therefore, on incongruent trials the expected value of control is larger than the EVOC of
10 exerting no control. Our model thus learns to exert control on incongruent trials but not
11 on congruent trials. Our model achieves this by learning to predict the EVOC from
12 features of the stimuli. This predicts that people should learn to exert more control when
13 they encounter a stimulus feature (such as a color or word) that is predictive of
14 incongruence than when they encounter a feature that is predictive of congruence (see
15 Table 1).

16 Consistent with our model's predictions, Bugg and colleagues [8] found that
17 people learn to exert more control in response to stimulus features that predict
18 incongruence than stimulus features that predict congruence. Their participants
19 performed a color-word Stroop task with four colors and their names printed either in
20 cursive or regular font. Our model captured the effects of congruency-predictive features
21 on control allocation with a plausible set of parameters (see Table 2). As shown in Figure
22 4a-b, the LVOC model predicted that responses should be faster (655 ± 9 ms vs. $722 \pm$
23 11 ms; $t(49) = 5.39, p < 0.0001$) and more accurate ($2.85 \pm 0.2\%$ errors vs. $4.3 \pm$

1 0.3% errors; $t(49) = 5.01, p < 0.0001$) on incongruent trials if the word was predictive
2 of incongruence than when it was not. To their surprise, Bugg and colleagues observed
3 that adding an additional feature (font) that conveyed the same information about
4 congruence as the color, did not enhance learning. This is exactly what our model
5 predicted because the presence of a second predictive feature reduces the evidence for the
6 predictive power of the first one and vice versa – this is directly analogous to a
7 phenomenon from the Pavlovian literature known as blocking, whereby an animal fails to
8 learn an association between a stimulus and an outcome that is already perfectly
9 predicted by a second stimulus [63].

10 Since our model learns about the predictive relationship between features and the
11 EVOC, it predicts that all learning effects should transfer to novel stimuli that share the
12 features that were predictive of the expected value of control in the training trials (see
13 Table 1). A separate study by Bugg and colleagues [9] confirmed this prediction. They
14 trained participants in a picture-word Stroop task to associate particular images of certain
15 categories (e.g., cats and dogs) with incongruence and associated particular images of
16 other categories (e.g., fish and birds) with congruence. As expected, participants learned
17 to exert more control when viewing the stimuli associated with incongruence. More
18 importantly, these participants also exerted more control when tested on *novel* instances
19 of the category associated with incongruence (e.g., cats) than on novel instances of the
20 category associated with congruence (e.g., fish). This finding provides strong evidence
21 for the feature-based learning mechanism that is at the core of our model of the plasticity
22 of cognitive control and is entirely accounted for by our model. As shown in Figure 4e
23 and Figure 4f, our model correctly predicted the positive and the negative transfer effects

1 reported by [9] with reasonable parameters (see Table 2): The model's responses were
2 faster (685 ± 2 ms vs. 709 ± 3 ms; $t(99) = -8.13, p < 0.0001$) and more accurate
3 ($3.2 \pm 0.1\%$ errors vs. $4.8 \pm 0.3\%$ errors; $t(99) = -5.06, p < 0.0001$) on incongruent
4 trials if the word was predictive of incongruence than when it was not (positive transfer).
5 Conversely, on congruent trials, the predicted responses were slightly slower when the
6 features wrongly predicted incongruence (530 ± 0.1 ms vs. 527 ± 0.2 ms, $t(99) =$
7 $9.28, p < 0.0001$; negative transfer).

8 Discussion

9 Building on previous work modeling the specification of cognitive control in
10 terms of meta-decision making [12–14,29,33,64,65] and reinforcement learning
11 [33,34,66–68], we have illustrated that at least some of the functions of cognitive control
12 can be characterized using the formal framework of rational metareasoning [26] and
13 meta-level Markov decision processes [27]. Concretely, modeling the function of
14 cognitive control as a meta-level MDP allowed us to derive the first formal
15 computational model of how people learn to specify continuous control signals and how
16 these learning effects transfer to novel situations. This model provides a unifying
17 explanation for how people learn where to attend, the interacting effects of reward and
18 incongruence on interference control, and their transfer to novel stimuli.

19 Our simulations of learning in Stroop and Flanker paradigms illustrate that the
20 LVOC model can account for people's ability to learn when and how intensely to engage
21 controlled processing and inhibit automatic processing. We further found that the LVOC
22 model correctly predicted the learning curve in the visual attention experiment by Lin et

1 al. [17] without any free parameters. Critically, all of our model's qualitative predictions
2 follow directly from our theory's core assumption that people reinforcement-learn to
3 predict the value of alternative control signals and control signal intensities from stimulus
4 features (see Table 1). None of our model's auxiliary assumptions about the cost of
5 control, the reward of being correct, the drift-diffusion model, the details of learning and
6 control signal selection, and the corresponding parameters summarized in Table 2 are
7 necessary to derive these qualitative predictions; instead they only serve to increase the
8 quantitative accuracy of those predictions.

9 While the LVOC model is more complex than basic associative learning and the
10 Win-Stay Lose-Shift mechanism, neither of these simpler models was able to capture
11 human learning in the simulated visual search, Stroop, and Flanker paradigms. This
12 suggests that the complexity of the LVOC model may be currently warranted to capture
13 how people learn when to exert how much cognitive control. Furthermore, the LVOC
14 model's sophistication may be necessary to explain more complex phenomena such as
15 how people learn to orchestrate their thoughts to solve complex problems and acquire
16 sophisticated cognitive strategies. Recent work has indeed shown that the learning
17 mechanism instantiated by the LVOC model can also capture aspects of how people learn
18 how to plan [69] and to flexibly and adaptively choose between alternative cognitive
19 strategies [54]. Testing whether people learn to select sequences of control signals in the
20 way predicted by our model is an interesting direction for future research.

21 [The LVOC model integrates control specification and strategy selection learning](#)

22 The model developed in this article builds on two previous theories: the EVC theory,
23 which offered a normative account of control specification [7], and the rational

1 metareasoning theory of strategy selection [54], which suggested that people acquire the
2 capacity to select heuristics adaptively by learning a predictive model of the execution
3 time and accuracy of those heuristics. The LVOC model synergistically integrates these
4 two theories: it augments the EVC theory with the metacognitive learning and prediction
5 mechanisms identified by [54], and it augments rational metareasoning models of
6 strategy selection with the capacity to specify continuous control signals that gradually
7 adjust parameters of the controlled process (see S2 Text).

8 Empirical predictions

9 All else being equal, the proposed learning rules (see Equation 7, S1 Text Equations 1-7,
10 and S3 Text Equations 13-14) predict that people's propensity to exert cognitive control
11 should increase when the controlled process was less costly (e.g., faster) or generated
12 more reward than expected [19]. The experience that less controlled (more automatic)
13 processing was more costly or less rewarding than expected should also increase our
14 propensity to exert cognitive control [70–72]. Conversely, if a controlled process
15 performed worse than expected or if an automatic process performed better than
16 expected, people's propensity to exert cognitive control should decrease [73].

17 At a more detailed level, our theory predicts that the influence of environmental
18 features on control allocation generalizes across contexts, to the extent that their features
19 are similar. Thus, adding or removing features to the internal predictive model of the
20 EVOC should have a profound effect on the degree to which observed performance of the
21 controlled process in Context A changes people's propensity to select it in Context B, and
22 vice versa. This mechanism can account for empirical evidence that suggests a role for
23 feature-binding in mechanisms of task switching [74–77]. These studies suggest that

1 participants associate the task that they perform on a stimulus with the features of that
2 stimulus. Once they are asked to engage in a new task on that stimulus, the old
3 (associated) task interferes, leading to switch costs.

4 Furthermore, our theory predicts that increasing the rewards and punishments for
5 the outcomes of the controlled or automatic processes should increase the speed at which
6 people's control allocation adapts to new task requirements, because the resulting weight
7 updates will be larger; this becomes especially apparent when the updates are rewritten in
8 terms of prediction errors (see S3 Text, Equations 1-2). Finally, when the assumptions of
9 the internal model are met and its features distinguish between the situations in which
10 each controlled process performs best, then control signal selection should become
11 increasingly more adaptive over time [78,79]. But in situations where the internal
12 model's assumptions are violated, for instance because the value of control is not additive
13 and linear in the features, then the control system's plasticity mechanisms may become
14 maladaptive.

15 This prediction has been confirmed in a recent experiment with a novel color-
16 word Stroop paradigm comprising two association phases and a test phase [80]. In the
17 first association phase, participants learned that color naming was rewarded for certain
18 colors whereas word reading was rewarded for the other colors. In the second association
19 phase, participants learned that color-naming was rewarded for certain words whereas
20 word-reading was rewarded for other words. Critically, in the test phase, naming the
21 color was rewarded if either the word or the color had been associated with color naming
22 (SINGLE trials); but when both the color and the word were associated with color
23 naming then participants had to instead read the word (BOTH trials). This non-linear

1 relationship between stimulus-features and control demands caused mal-transfer from
2 SINGLE trials to BOTH trials that significantly interfered with participants' performance
3 (resulting in participants incorrectly engaging in color-naming, the more control-
4 demanding task which in that context was *also* less rewarding). The LVOC model may
5 thus be able to explain the puzzling phenomenon that people sometimes overexert
6 cognitive control even when it hurts their performance, such as when upon seeing a non-
7 urgent email on a topic you have learned to be careful about you cannot help but compose
8 the perfect response in your mind while trying to finish the talk you have to give in 5
9 minutes.

10 According to the LVOC model, control allocation is a process of continuing
11 gradual adjustment (Equation 13). This means that the control intensity for a new
12 situation starts out with the control intensity from the previous situation and is then
13 gradually adjusted towards its optimal value (Equation 13)—just like in anchoring-and-
14 adjustment [48,81]. This might provide a mechanism for commonly observed phenomena
15 associated with task set inertia and switch costs [47]. Since control adjustment takes time,
16 this mechanism predicts that increased time pressure could potentially lead to decreased
17 control adjustment, thereby biasing people's control allocation to its value on the
18 previous trial and thus decreasing their cognitive flexibility. Finally, thinking about the
19 neural implementation of the LVOC model leads to additional neural predictions as
20 detailed in the S3 Text.

21 [Avenues for future research](#)

22 We view rational metareasoning as a general theoretical framework for modeling the
23 allocation and plasticity of cognitive control. As such, it could be used to develop

1 unifying models of different manifestations of cognitive control, such as attention,
2 response inhibition, and cognitive flexibility. Furthermore, rational metareasoning can
3 also be used to connect existing models of cognitive control [6,7,32–
4 34,46,64,65,78,79,82,83]. Interpreting previously proposed mechanisms of control
5 allocation as approximations to rational metareasoning and considering how else rational
6 metareasoning could be approximated might facilitate the systematic evaluation of
7 alternative representations and computational mechanisms and inspire new models.
8 While our computational explorations have focused on which control signal the cognitive
9 control system should select, future work might also shed light on how the cognitive
10 control system monitors the state of the controlled system by viewing the problem solved
11 by the cognitive control system as a partially observable MDP. Concretely, the function
12 of cognitive monitoring could be formulated as a meta-level MDP whose computational
13 actions include sensing operations that update the cognitive control system’s beliefs
14 about the state of the monitored system.

15 Future work should further evaluate the proposed computational mechanism and
16 its neural implementation by performing quantitative model comparisons against simpler
17 models across a wider range of cognitive control phenomena. We will investigate the
18 performance of the proposed metacognitive learning mechanism and evaluate it against
19 alternative mechanisms (e.g., temporal difference learning mechanisms with eligibility
20 traces [30]).

21 Another interesting direction will be to use the learning models to investigate the
22 plasticity of people’s cognitive control skills. We are optimistic that this line of work will
23 lead to better quantitative models of control plasticity that can be used to develop

1 interventions to improve people's executive functions via a combination of cognitive
2 training and augmenting environments where people's automatic responses are
3 maladaptive with cues that prime them to employ an appropriate control signal. In
4 addition, future work may also explore model-based metacognitive reinforcement
5 learning [84] as a model of the plasticity of cognitive control specification. Model-based
6 hierarchical reinforcement learning approaches [37], such as option models [38], could be
7 used to integrate the learning mechanisms for the value of individual control signals with
8 the strategy selection model to provide an account of how the brain discovers control
9 strategies. This might explain how people learn to adaptively coordinate their thoughts
10 and actions to pursue increasingly more challenging goals over increasingly longer
11 periods of time.

12 Finally, the rational metareasoning framework can also be used to model how
13 people reason about the costs and benefits of exerting mental effort and to delineate self-
14 control failure from rational resource-preservation through a normative account of effort
15 avoidance [13,85].

16 Conclusion

17 Our simulation results suggested that the LVOC model provides a promising step
18 towards a mathematical theory of cognitive plasticity that can serve as a scientific
19 foundation for designing cognitive training programs for improving people's executive
20 functions. This illustrates the utility of formalizing the function of cognitive control in
21 terms of rational metareasoning. Rational metareasoning provides a unifying framework
22 for modeling executive functions, and thus opens up exciting avenues for future research.
23 We are optimistic that the connection between executive functions and metareasoning

will channel a flow of useful models and productive ideas from artificial intelligence and machine learning into the neuroscience and psychology of cognitive control.

References

1. Diamond A. Executive Functions. *Annu Rev Psychol.* 2013;64: 135–168. doi:doi: 10.1146/annurev-psych-113011-143750
2. Stroop JR. Studies of interference in serial verbal reactions. *J Exp Psychol. Psychological Review Company;* 1935;18: 643.
3. Tangney JP, Baumeister RF, Boone AL. High self-control predicts good adjustment, less pathology, better grades, and interpersonal success. *J Pers. Wiley Online Library;* 2004;72: 271–324.
4. Moffitt TE, Arseneault L, Belsky D, Dickson N, Hancox RJ, Harrington H, et al. A gradient of childhood self-control predicts health, wealth, and public safety. *Proc Natl Acad Sci. National Acad Sciences;* 2011;108: 2693–2698.
5. Kool W, Botvinick MM. The intrinsic cost of cognitive control. *Behav Brain Sci.* 2013;36: 697–8. doi:10.1017/S0140525X1300109X
6. Kool W, Botvinick MM. A labor/leisure tradeoff in cognitive control. *J Exp Psychol Gen.* 2014;143: 131–141. doi:10.1037/a0031048
7. Shenhav A, Botvinick MM, Cohen J. The Expected Value of Control: An Integrative Theory of Anterior Cingulate Cortex Function. *Neuron. Cell Press,;* 2013;79: 217–240. doi:doi: 10.1016/j.neuron.2013.07.007
8. Bugg JM, Jacoby LL, Toth JP. Multiple levels of control in the Stroop task. *Mem \& Cogn. Springer;* 2008;36: 1484–1494.
9. 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. American Psychological Association;* 2011;37: 844.
10. Braem S, Verguts T, Roggeman C, Notebaert W. Reward modulates adaptations to conflict. *Cognition. Elsevier;* 2012;125: 324–332.
11. Krebs RM, Boehler CN, Woldorff MG. The influence of reward associations on conflict processing in the Stroop task. *Cognition. Elsevier;* 2010;117: 341–347.
12. Kool W, Shenhav A, Botvinick MM. Cognitive Control as Cost-Benefit Decision Making. *Wiley Handb Cogn Control. John Wiley \& Sons, Ltd;* 2017; 167–189.
13. Shenhav A, Musslick S, Lieder F, Kool W, Griffiths TL, Cohen JD, et al. Toward a rational and mechanistic account of mental effort. *Annu Rev Neurosci.* 2017;
14. Boureau Y-L, Sokol-Hessner P, Daw ND. Deciding How To Decide: Self-Control and Meta-Decision Making. *Trends Cogn Sci.* 2015; doi:10.1016/j.tics.2015.08.013

- 1 15. Abrahamse E, Braem S, Notebaert W, Verguts T. Grounding cognitive control in
2 associative learning. *Psychol Bull. American Psychological Association*;
3 2016;142: 693.
- 4 16. Egner T. Creatures of habit (and control): a multi-level learning perspective on the
5 modulation of congruency effects. *Front Psychol. Frontiers Media SA*; 2014;5.
- 6 17. Lin Z, Lu Z-L, He S. Decomposing experience-driven attention: Opposite
7 attentional effects of previously predictive cues. *Attention, Perception, \&
8 Psychophys. Springer*; 2016;78: 2185–2198.
- 9 18. Muhle-Karbe PS, Jiang J, Egner T. Causal Evidence for Learning-Dependent
10 Frontal Lobe Contributions to Cognitive Control. *J Neurosci. Soc Neuroscience*;
11 2018;38: 962–973.
- 12 19. Anguera JA, Boccanfuso J, Rintoul JL, Al-Hashimi O, Faraji F, Janowich J, et al.
13 Video game training enhances cognitive control in older adults. *Nature. Nature
14 Publishing Group*; 2013;501: 97–101. doi:10.1038/nature12486
- 15 20. Karbach J, Kray J. How useful is executive control training? Age differences in
16 near and far transfer of task-switching training. *Dev Sci. Wiley Online Library*;
17 2009;12: 978–990.
- 18 21. Heller SB, Shah AK, Guryan J, Ludwig J, Mullainathan S, Pollack HA. Thinking,
19 fast and slow? Some field experiments to reduce crime and dropout in Chicago.
20 2015.
- 21 22. Melby-Lervåg M, Hulme C. Is working memory training effective? A meta-
22 analytic review. *Dev Psychol. US: American Psychological Association*; 2013;49:
23 270.
- 24 23. Shipstead Z, Redick TS, Engle RW. Is working memory training effective?
25 *Psychol Bull. American Psychological Association*; 2012;138: 628.
- 26 24. Friston K. Functional integration and inference in the brain. *Prog Neurobiol.*
27 2002;68: 113–143. doi:10.1016/S0301-0082(02)00076-X
- 28 25. Owen AM, Hampshire A, Grahn JA, Stenton R, Dajani S, Burns AS, et al. Putting
29 brain training to the test. *Nature. Nature Research*; 2010;465: 775–778.
- 30 26. Russell S, Wefald E. Principles of metareasoning. *Artif Intell.* 1991;49: 361–395.
31 doi:doi: 10.1016/0004-3702(91)90015-c
- 32 27. Hay N, Russell S, Tolpin D, Shimony S. Selecting Computations: Theory and
33 Applications. In: de Freitas N, Murphy K, editors. *Uncertainty in Artificial
34 Intelligence: Proceedings of the Twenty-Eighth Conference.* P.O. Box 866
35 Corvallis, Oregon 97339 USA OR - *Uncertainty in Artificial Intelligence: AUA*
36 *Press*; 2012.
- 37 28. Miller EK, Cohen JD. An integrative theory of prefrontal cortex function. *Annu
38 Rev Neurosci.* 2001;24: 167–202. doi:10.1146/annurev.neuro.24.1.167
- 39 29. Gershman SJ, Horvitz EJ, Tenenbaum JB. Computational rationality: A

- 1 converging paradigm for intelligence in brains, minds, and machines. *Science* (80-
2). American Association for the Advancement of Science; 2015;349: 273–278.
- 3 30. Sutton RS, Barto AG. Reinforcement learning: An introduction. Cambridge, MA,
4 USA: MIT press; 1998.
- 5 31. Kaelbling LP, Littman ML, Cassandra AR. Planning and acting in partially
6 observable stochastic domains. *Artif Intell.* 1998;101: 99–134.
7 doi:[https://doi.org/10.1016/S0004-3702\(98\)00023-X](https://doi.org/10.1016/S0004-3702(98)00023-X)
- 8 32. Gottlieb J, Balan P. Attention as a decision in information space. *Trends Cogn Sci.*
9 2010;14: 240–248. doi:10.1016/j.tics.2010.03.001
- 10 33. Dayan P. How to set the switches on this thing. *Curr Opin Neurobiol.* 2012;22:
11 1068–1074. doi:10.1016/j.conb.2012.05.011
- 12 34. Todd MT, Niv Y, Cohen JD. Learning to Use Working Memory in Partially
13 Observable Environments through Dopaminergic Reinforcement. In: Koller D,
14 Schuurmans D, Bengio Y, Bottou L, editors. *Advances in Neural Information*
15 *Processing Systems 21*. Curran Associates, Inc.; 2009. pp. 1689–1696.
- 16 35. Holroyd CB, Yeung N. Motivation of extended behaviors by anterior cingulate
17 cortex. *Trends Cogn Sci.* Elsevier; 2012;16: 122–128.
- 18 36. Holroyd CB, McClure SM. Hierarchical control over effortful behavior by rodent
19 medial frontal cortex: A computational model. *Psychol Rev.* American
20 Psychological Association; 2015;122: 54.
- 21 37. Botvinick MM, Weinstein A. Model-based hierarchical reinforcement learning and
22 human action control. *Philos Trans R Soc B Biol Sci.* The Royal Society;
23 2014;369: 20130480.
- 24 38. Sutton RS, Precup D, Singh S. Between MDPs and semi-MDPs: A framework for
25 temporal abstraction in reinforcement learning. *Artif Intell.* Essex, UK: Elsevier
26 Science Publishers Ltd.; 1999;112: 181–211. doi:10.1016/s0004-3702(99)00052-1
- 27 39. Watkins CF, Dayan P. Q-learning. *Mach Learn.* Kluwer Academic Publishers;
28 1992;8: 279–292. doi:10.1007/bf00992698
- 29 40. Lindley D V, Smith AFM. Bayes Estimates for the Linear Model. *J R Stat Soc Ser*
30 *B.* Blackwell Publishing for the Royal Statistical Society; 1972;34: 1–41. doi:
31 10.2307/2985048
- 32 41. Kunz S. The Bayesian Linear model with Unknown Variance. 2009.
- 33 42. Rescorla RA, Wagner AR, others. A theory of Pavlovian conditioning: Variations
34 in the effectiveness of reinforcement and nonreinforcement. *Class Cond II Curr*
35 *Res theory.* New-York; 1972;2: 64–99.
- 36 43. Restle F. The selection of strategies in cue learning. *Psychol Rev.* American
37 Psychological Association; 1962;69: 329.
- 38 44. Klein RM. Inhibition of return. *Trends Cogn Sci.* Elsevier; 2000;4: 138–147.

- 1 45. Cohen JD, Dunbar K, McClelland JL. On the control of automatic processes: a
2 parallel distributed processing account of the Stroop effect. *Psychol Rev.*
3 *American Psychological Association*; 1990;97: 332.
- 4 46. Musslick S, Shenhav A, Botvinick MM, Cohen JD. A computational model of
5 control allocation based on the Expected Value of Control. *The 2nd*
6 *Multidisciplinary Conference on Reinforcement Learning and Decision Making.*
7 2015.
- 8 47. Allport DA, Styles EA, Hsieh S. Shifting intentional set: Exploring the dynamic
9 control of tasks. In: Umiltà C, Moscovitch M, editors. *Attention and performance*
10 *XV.* 1994. pp. 421–452.
- 11 48. Lieder F, Griffiths TL, Huys QJM, Goodman ND. Empirical Evidence for
12 Resource-Rational Anchoring and Adjustment. *Psychon Bull \& Rev.* Springer;
13 2017;
- 14 49. Lieder F, Griffiths TL, Huys QJM, Goodman ND. The anchoring bias reflects
15 rational use of cognitive resources. *Psychon Bull \& Rev.* Springer; 2017;
- 16 50. Padmala S, Pessoa L. Reward Reduces Conflict by Enhancing Attentional Control
17 and Biasing Visual Cortical Processing. *J Cogn Neurosci.* MIT Press; 2011;23:
18 3419–3432. doi:10.1162/jocn_a_00011
- 19 51. Manohar SG, Chong TT-J, Apps MAJ, Batla A, Stamelou M, Jarman PR, et al.
20 Reward pays the cost of noise reduction in motor and cognitive control. *Curr Biol.*
21 *Elsevier*; 2015;25: 1707–1716.
- 22 52. Feng SF, Schwemmer M, Gershman SJ, Cohen JD. Multitasking versus
23 multiplexing: Toward a normative account of limitations in the simultaneous
24 execution of control-demanding behaviors. *Cogn Affect \& Behav Neurosci.*
25 Springer; 2014;14: 129–146.
- 26 53. Musslick S, Dey B, Ozcimder K, Patwary MMA, Willke TL, Cohen JD.
27 Controlled vs. Automatic Processing: A Graph-Theoretic Approach to the Analysis
28 of Serial vs. Parallel Processing in Neural Network Architectures. *Proceedings of*
29 *the 38th Annual Conference of the Cognitive Science Society.* 2016. pp. 1547–
30 1552.
- 31 54. Lieder F, Griffiths TL. Strategy selection as rational metareasoning. *Psychol Rev.*
32 *American Psychological Association*; 2017;124: 762.
- 33 55. Kawaguchi K, Kaelbling LP, Lozano-Pérez T. Bayesian Optimization with
34 Exponential Convergence. In: Cortes C, Lawrence ND, Lee DD, Sugiyama M,
35 Garnett R, editors. *Advances in Neural Information Processing Systems 28.* Curran
36 Associates, Inc.; 2015. pp. 2809–2817.
- 37 56. Balci F, Simen P, Niyogi R, Saxe A, Hughes JA, Holmes P, et al. Acquisition of
38 decision making criteria: reward rate ultimately beats accuracy. *Attention,*
39 *Perception, \& Psychophys.* Springer; 2011;73: 640–657.
- 40 57. Acerbi L, Ma WJ. Practical Bayesian Optimization for Model Fitting with

- 1 Bayesian Adaptive Direct Search. In: Guyon I, Luxburg U V, Bengio S, Wallach
2 H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information*
3 *Processing Systems 30*. Curran Associates, Inc.; 2017. pp. 1834–1844. Available:
4 [http://papers.nips.cc/paper/6780-practical-bayesian-optimization-for-model-fitting-](http://papers.nips.cc/paper/6780-practical-bayesian-optimization-for-model-fitting-with-bayesian-adaptive-direct-search.pdf)
5 [with-bayesian-adaptive-direct-search.pdf](http://papers.nips.cc/paper/6780-practical-bayesian-optimization-for-model-fitting-with-bayesian-adaptive-direct-search.pdf)
- 6 58. Botvinick MM, Braver TS, Barch DM, Carter CS, Cohen JD. Conflict monitoring
7 and cognitive control. *Psychol Rev*. Department of Psychology, Carnegie Mellon
8 University, Pennsylvania, USA. mmb@cnbc.cmu.edu; 2001;108: 624–652.
- 9 59. Kass R, Raftery A. Bayes Factors. *J Am Stat Assoc*. American Statistical
10 Association; 1995;90: 773–795. doi:doi: 10.2307/2291091
- 11 60. Schwarz G, others. Estimating the dimension of a model. *Ann Stat*. Institute of
12 Mathematical Statistics; 1978;6: 461–464.
- 13 61. Nassar MR, Rumsey KM, Wilson RC, Parikh K, Heasley B, Gold JI. Rational
14 regulation of learning dynamics by pupil-linked arousal systems. *Nat Neurosci*.
15 Nature Research; 2012;15: 1040–1046.
- 16 62. McGuire JT, Nassar MR, Gold JI, Kable JW. Functionally dissociable influences
17 on learning rate in a dynamic environment. *Neuron*. Elsevier; 2014;84: 870–881.
- 18 63. Kamin LJ. Predictability, surprise, attention, and conditioning. *Punishm aversive*
19 *Behav*. 1969; 279–296.
- 20 64. Keramati M, Dezfouli A, Piray P. Speed/Accuracy Trade-Off between the
21 Habitual and the Goal-Directed Processes. *PLoS Comput Biol*. Public Library of
22 Science; 2011;7: e1002055. doi:doi: 10.1371/journal.pcbi.1002055
- 23 65. Daw N, Niv Y, Dayan P. Uncertainty-based competition between prefrontal and
24 dorsolateral striatal systems for behavioral control. *Nat Neurosci*. Nature
25 Publishing Group; 2005;8: 1704–1711. doi:doi: 10.1038/nn1560
- 26 66. Kool W, Gershman SJ, Cushman FA. Cost-benefit arbitration between multiple
27 reinforcement-learning systems. *Psychol Sci*.
- 28 67. O'Reilly RC, Frank MJ. Making working memory work: a computational model of
29 learning in the prefrontal cortex and basal ganglia. *Neural Comput*. MIT Press;
30 2006;18: 283–328.
- 31 68. Verguts T, Vassena E, Silvetti M. Adaptive effort investment in cognitive and
32 physical tasks: A neurocomputational model. *Front Behav Neurosci*. Frontiers
33 Media SA; 2015;9.
- 34 69. Krueger PM, Lieder F, Griffiths TL. Enhancing Metacognitive Reinforcement
35 learning using reward structures and feedback. *Proceedings of the 39th Annual*
36 *Conference of the Cognitive Science Society*. 2017.
- 37 70. Danielmeier C, Ullsperger M. Post-error adjustments. *Front Psychol*. 2011;2.
38 doi:10.3389/fpsyg.2011.00233
- 39 71. Laming DRJ. *Information theory of choice-reaction times*. Oxford, England:

Academic Press; 1968.

72. Gratton G, Coles MG, Donchin E. Optimizing the use of information: strategic control of activation of responses. *J Exp Psychol Gen.* 1992;121: 480–506.
73. Brown TL, Carr TH. Automaticity in skill acquisition: Mechanisms for reducing interference in concurrent performance. *J Exp Psychol Hum Percept Perform.* American Psychological Association; 1989;15: 686.
74. Waszak F, Hommel B, Allport A. Task-switching and long-term priming: Role of episodic stimulus--task bindings in task-shift costs. *Cogn Psychol.* Elsevier; 2003;46: 361–413.
75. Waszak F, Hommel B, Allport A. Semantic generalization of stimulus-task bindings. *Psychon Bull \& Rev. Springer;* 2004;11: 1027–1033.
76. Waszak F, Hommel B. The costs and benefits of cross-task priming. *Mem \& Cogn.* Springer; 2007;35: 1175–1186.
77. Mayr U, Bryck RL. Outsourcing control to the environment: effects of stimulus/response locations on task selection. *Psychol Res.* Springer; 2007;71: 107–116.
78. Lieder F, Griffiths TL. When to use which heuristic: A rational solution to the strategy selection problem. In: Noelle DC, Dale R, Warlaumont AS, Yoshimi J, Matlock T, Jennings CD, et al., editors. *Proceedings of the 37th Annual Conference of the cognitive science society.* Austin, TX: Cognitive Science Society; 2015.
79. Lieder F, Plunkett D, Hamrick JB, Russell SJ, Hay NJ, Griffiths TL. Algorithm selection by rational metareasoning as a model of human strategy selection. In: Ghahramani Z, Welling M, Weinberger KQ, Cortes C, Lawrence ND, editors. *Advances in Neural Information Processing Systems 27.* Curran Associates, Inc.; 2014.
80. Bustamante L, Lieder F, Musslick S, Shenhav A, Cohen JD. Learning to (mis)allocate control: maltransfer can lead to self-control failure. In: Brunskill E, N. Daw, editors. *The 3rd Multidisciplinary Conference on Reinforcement Learning and Decision Making.* Ann Arbor, Michigan; 2017.
81. Lieder F, Griffiths TL, Huys QJ, Goodman ND. Reinterpreting anchoring-and-adjustment as rational use of cognitive resources.
82. Braver T. The variable nature of cognitive control: a dual mechanisms framework. *Trends Cogn Sci.* 2012;16: 106–113. doi:doi: 10.1016/j.tics.2011.12.010
83. Suchow JW. Measuring, monitoring, and maintaining memories in a partially observable mind. Harvard University. 2014.
84. Sutton RS. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. *Proceedings of the seventh international conference on machine learning.* 1990. pp. 216–224.

85. Inzlicht M, Schmeichel BJ, Macrae CN. Why self-control seems (but may not be) limited. *Trends Cogn Sci*. Elsevier Science,; 2014;18: 127–133. doi:10.1016/j.tics.2013.12.009
86. Niv Y, Daw N, Joel D, Dayan P. Tonic dopamine: opportunity costs and the control of response vigor. *Psychopharmacology (Berl)*. Interdisciplinary Center for Neural Computation, The Hebrew University of Jerusalem, Jerusalem, 91904, Israel.: Springer-Verlag; 2007;191: 507–520. doi:10.1007/s00213-006-0502-4

Supporting Information Legends

S1 Text: Mathematical details of the LVC model's learning mechanism.

S2 Text: Rational metareasoning unifies the EVC theory with the rational metareasoning theory of strategy selection.

S3 Text: Speculations about how the learning mechanism postulated by the LVC model might be implemented in the brain.

Author Contributions:

Conceptualization: FL, TLG, AS

Data Curation: FL, SM, AS

Formal Analysis: FL, SM

Funding Acquisition: TLG

Investigation: FL

Methodology: FL, AS, SM

Project Administration: FL, TLG

Resources: FL, SM, AS

Software: FL, SM, AS

Supervision: TLG, AS, FL

Validation: FL

Visualization: FL, SM

1 Writing – Original Draft Preparation: FL, AS, SM

2 Writing – Review & Editing: FL, AS, TLG, SM

3 **Competing interests statement:**

4 The authors declare no competing interests.

5 **Corresponding author**

6 Correspondence to: Falk Lieder (falk.lieder@berkeley.edu)

7 **Acknowledgements:**

8 The authors thank Matthew Botvinick, Colin Hoy, and S.J. Katarina Slama for comments

9 on an earlier version of the manuscript and Jonathan D. Cohen for useful discussions.