

## Cost-Benefit Arbitration Between Multiple Reinforcement-Learning Systems

Wouter Kool, Samuel J. Gershman, & Fiery A. Cushman

Human behavior is sometimes determined by habit, and other times by goal-directed planning. Modern reinforcement learning theories formalize this distinction as a competition between a computationally cheap but inaccurate “model-free” system that gives rise to habits, and a computationally expensive but accurate “model-based” system that implements planning. It is unclear, however, how we choose to allocate control between these systems. Here, we propose that arbitration occurs by comparing each system’s task-specific costs and benefits. Consistent with this proposal, we report two experiments showing that people increase model-based control when it achieves greater accuracy than model-free control, and especially when the rewards of accurate performance are amplified. In contrast, they are insensitive to reward amplification when model-based and model-free control yield equivalent accuracy. This suggests that humans adaptively balance habitual and planned action through online cost-benefit analysis.

Diverse traditions of behavioral research distinguish between two systems for choosing actions: An automatic system that relies on habit, and a controlled system that plans towards goals (Dickinson, 1985; Kahneman, 2003; Sloman, 1996). These systems embody different accuracy-demand tradeoffs: The habit system has low computational demands but is often less accurate, whereas the planning system achieves greater accuracy with greater computational demands. We ask a simple question: How do people decide from moment to moment which system to use?

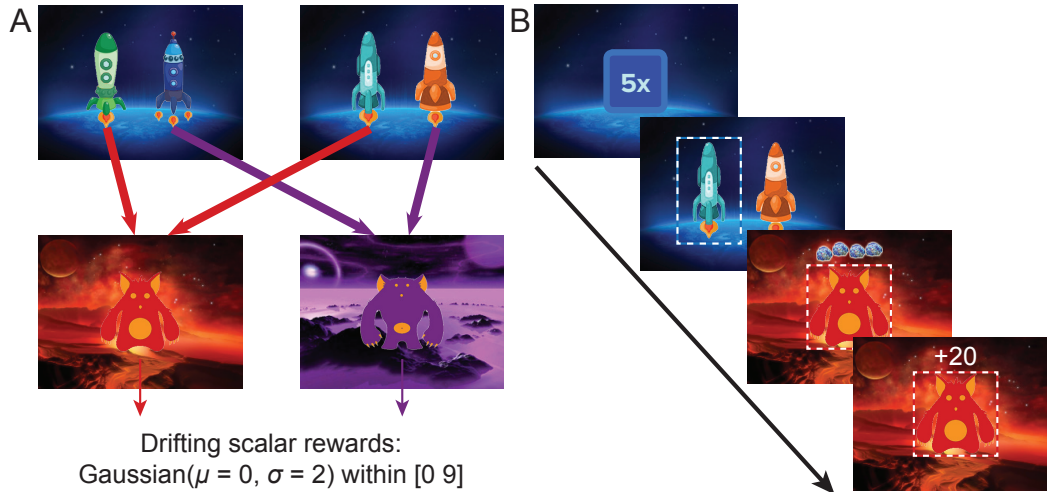
To provide a precise answer, we formalize the notions of habit versus planning in the reinforcement learning (RL) setting. “Model-free” RL chooses actions that previously led to reward (Thorndike, 1911), a computationally efficient but inflexible strategy. “Model-based” RL achieves flexibility by planning in a causal model of the environment, a comparatively accurate but inefficient strategy. This formalization has facilitated research on each systems’ neural basis (Dolan & Dayan, 2013; Gershman, Markman, & Otto, 2014), dependence on executive control (Otto, Gershman, Markman, & Daw, 2013), and contribution to clinical disorders (Gillan, Kosinski, Whelan, Phelps, & Daw, 2016).

Several theories aim to explain how people allocate control between these systems (Gershman, Horvitz, & Tenenbaum, 2015; Griffiths, Lieder, & Goodman, 2015; Keramati, Dezfouli, & Piray, 2011; Pezzulo, Rigoli, & Chersi, 2013), although surprisingly little experimental research targets this question directly (but see Lee, Shimojo, & O’Doherty, 2014). In principle, people could choose to employ the more accurate model-based approach whenever cognitive resources are available, relying on habit only when those resources are already occupied (e.g., under cognitive load; Otto et al., 2013). This “first-come-first-served” approach does not, however, attempt to allocate limited computational resources towards tasks offering the maximum benefit. More sophisticated approaches would be sensitive to two task-specific variables: the amount of reward at stake, and the size of the model-based advantage in

accuracy. Either of these can be incorporated in isolation; people might increase model-based control whenever stakes are high (ignoring whether model-based control is more accurate), or they might increase model-based control when it enjoys the greatest accuracy advantage (ignoring the amount of reward at stake; Daw, Niv, & Dayan, 2005; Rieskamp & Otto, 2006). Finally, people might adaptively integrate both of these pieces of information in order to estimate the comparative reward advantage of model-based control.

Consistent with this final suggestion, we propose that allocation of control is based on the estimated benefits associated with each system in a given task, weighed against the cost of computational demand. Thus, model-based control would be deployed when the combination of high stakes and a robust accuracy advantage are sufficient to offset the cost implied by its reliance on limited cognitive resources. We explore two untested predictions of this proposal: First, that people will increase model-based control when there is a heightened opportunity for reward and, second, that this effect will be eliminated when model-based control cannot reliably outperform model-free control in its accuracy.

Our proposal requires that people assign a cost to model-based control; otherwise, there is nothing to “trade off” against its potential for a reward advantage over model-free control. Two literatures provide strong circumstantial evidence that people represent such a cost. First, model-based control depends on the capacity for cognitive control. Cognitive load, which decreases the capacity for cognitive control, increases the influence of the model-free system (Otto et al., 2013). Additionally, the model-based system depends on prefrontal structures closely associated with cognitive control (Gläscher, Daw, Dayan, & O’Doherty, 2010; Lee et al., 2014; Smittenaar, FitzGerald, Romei, Wright, & Dolan, 2013). Second, people assign an intrinsic cost to exercising cognitive control (Botvinick, 2007; Botvinick & Braver, 2015; Kool, McGuire, Rosen, & Botvinick, 2010; Kool, Shenhav, & Botvinick, 2017). Thus, people



**Fig 1. Design of Experiment 1.** (A) State transition structure. Each first-stage choice deterministically transitions to one of two second-stage states. Each second-stage is associated with a scalar reward (between 0 and 9) that changed across the duration of the experiment according to a random Gaussian walk with  $\sigma = 2$ . (B) Timeline of events. At the start of each trial, a cue indicated whether the points on this trial would be multiplied by 1 (low stakes) or 5 (high stakes). On this high-stakes trial, the alien gives 4 pieces of treasure, and so the total points earned is 20.

avoid tasks that demand cognitive control unless its cost is offset by task-specific rewards (Kool et al., 2010; Westbrook, Kester, & Braver, 2013). This intrinsic cost would presumably serve as a heuristic representation of the opportunity cost associated with deploying limited cognitive resources (Kool & Botvinick, 2014; Kurzban, Duckworth, Kable, & Myers, 2013). Combining these insights, we posit that people assign an intrinsic cost to model-based planning due to its reliance on cognitive control, which is balanced against a task-specific representation of its potential reward advantage.

Some form of cost-benefit tradeoff occurs in several theoretical models of metacontrol (Gershman et al., 2015; Griffiths et al., 2015; Keramati et al., 2011; Payne, Bettman, & Johnson, 1988; Rieskamp & Otto, 2006). Some of these theories formalize competing control systems within the RL framework, but are not supported by direct experimental evidence (Gershman et al., 2015; Griffiths et al., 2015; Keramati et al., 2011). Others have generated data consistent with a cost-benefit tradeoff (Payne et al., 1988; Rieskamp & Otto, 2006), but do not formalize control systems in terms of dual-system RL models. Building on this background, our study is motivated by three goals: First, to provide experimental evidence for cost-benefit analyses in metacontrol; second, to accomplish this in a setting amenable to formal analysis in the RL framework; third, leveraging this formalization, to assess the adequacy of current theoretical proposals to capture the precise form of cost-benefit analysis. Below, we show how our approach provides new traction in distinguishing among contemporary theories of metacontrol.

## Experiment 1

### Method

**Participants.** One hundred and one participants (range: 21–61 years of age; mean: 32 years of age; 47 female) were recruited on Amazon Mechanical Turk to participate in the experiment. This sample size was chosen so that we would achieve approximately 90% power to detect a true medium effect (Cohen's  $d = 0.3$ ) with a two-tailed  $\alpha = 0.05$ . Participants gave informed consent, and the Harvard Committee on the Use of Human Subjects approved the study.

Participants were excluded from analysis if they timed out on more than 20% of all trials (more than 40), and we excluded all trials on which participants timed out (average 4.4%). After applying these criteria, data from 98 participants were used in subsequent analysis.

**Materials and procedures.** The first experiment was designed to test whether choice behavior shows increased model-based control in the face of increased incentives—that is, when they stand to gain the most from superior accuracy in performance, offsetting the putative subjective cost of executive control. We used a recently developed two-step task (Kool, Cushman, & Gershman, 2016) based, in part, on work by Doll and colleagues (2015) (Figure 1A). In short, this task dissociates model-free from model-based control by exploiting the ability of the model-based system to plan using an explicit model of the task structure, which contrasts with the model-free reliance on the direct experience of action-reward associations.

Each trial started randomly in one of two possible first-stage states, where participants made an initial choice between a pair of spaceships that appeared side by side on a blue earth-like planet background. The spaceships had an equal probability of appearing on the left or right side of the screen. The choice between the left- and right-hand spaceships had to be made using the “F” or “J” button keys within a response deadline of 1,500ms.

This first-stage choice deterministically controlled which second-stage state, a purple or a red planet, would then be encountered. Importantly, each first-stage state afforded the possibility of transitioning to either planet, with one spaceship always leading to the purple planet and the other always to the red planet. Each second-stage state was associated with a scalar reward. Specifically, on each planet, participants found a single alien, and they were told that this alien ‘worked at a space mine’. They were instructed to press the space bar within the time limit in order to receive the reward. Participants were told that sometimes the aliens were in a good part of the mine and they paid off a high number of points or ‘space treasure’, whereas at other times the aliens were mining in a bad spot, and this yielded fewer pieces of space treasure. The payoffs of these mines slowly changed over the course of the experiment according to independent random walks, encouraging learning throughout the entire session. Note that without this slow change in the scalar rewards, a model-free, trial-and-error, strategy would eventually converge to the model-based strategy. The continuous change over time guarantees that the value representations remain different between systems throughout the experimental session.

One of the alien’s reward distributions was initialized randomly within a range of 0 points to 4 points, and the other within a range of 5 to 9 points. They then varied according to a Gaussian random walk ( $\sigma = 2$ ) with reflecting bounds at 0 and 9. A new set of randomly drifting reward distributions was generated for each participant. At the end of the experiment, participants were given 1¢ for every thirty-six points they earned (i.e., 0.25¢ for a maximal win of 9 points on a given trial).

In order to equate the time demands for the model-based and model-free strategies, both the selected spaceship and alien were highlighted for the remainder of the response period. This meant that participants were not able to increase their rate of reward over time by responding more rapidly.

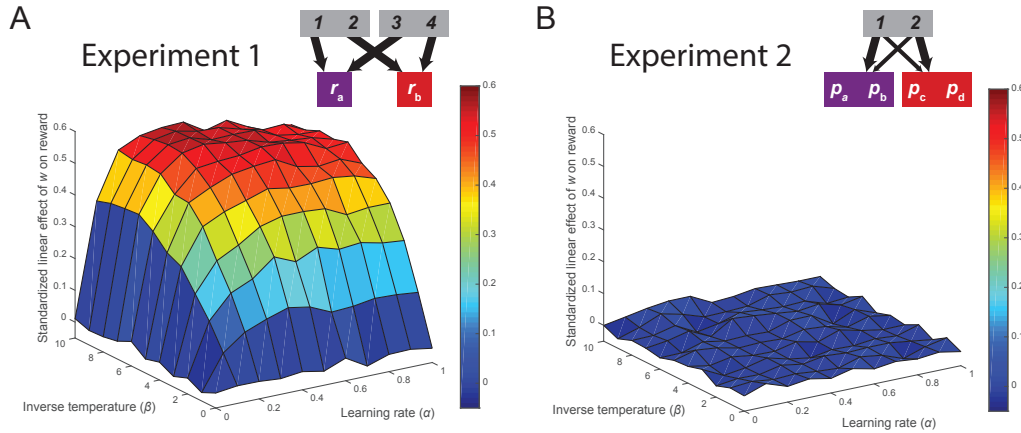
The most important feature of the task is that the choices between spaceships are equivalent between the first-stage states. For each pair, one spaceship always led to the red planet and alien, whereas the other always led to the purple planet and alien. Because of this equivalence, this task can distinguish model-based and model-free strategies, since only the model-based system can transfer experiences learned in one starting state to

the other starting state. For a pure model-based strategy, each second-stage outcome should always affect first-stage preferences on the next trial, regardless of whether this trial started with the same or the other pair of spaceships, because it plans towards the second-stage goals. In contrast, a pure model-free learner does not transfer experiences obtained after one pair of spaceships to the other pair, since it only learns action-reward associations (Doll et al., 2015).

As an example, imagine an agent starting in the left panel in Figure 1A. There, she chooses the green rocket ship and transitions to the red planet. She receives a very large reward on the red planet. On the next trial, the agent starts in the other panel in Figure 1A, namely the one the right. A model-free algorithm will not have updated the value of either of these rocket ships—only the green ship in the left panel will have received a boost in value. In contrast, a model-based algorithm will use its model of the transition structure of the task to assign value to all ships based on their probability of reaching the red planet. Thus, the model-based algorithm alone predicts that she will exhibit an increased probability of selecting the turquoise rocket ship in the right panel. The analyses reported below exploit this distinction between decision-making strategies.

In order to test our motivational cost-benefit account of metacontrol, we introduced a ‘stakes’ manipulation into this two-step paradigm (Figure 1B). Specifically, at the start of each trial, a cue indicated a multiplication factor for the subsequent points earned. With 50% probability, this cue indicated that all points would be multiplied by 5. For example, earning 4 pieces of space treasure would increase the overall point total by 20 points, earning 9 pieces would increase the total by 45 points. On the other trials, the cue indicated that the points would be multiplied by 1, so that each piece of space treasure would be worth only 1 point. After participants observed the reward outcome on the trial, they were also given an explicit computation of the number of points gained on that trial multiplied by the trial’s stake. The running score was always available in the top-right corner of the screen. We predicted that behavioral performance would show increased contributions of model-based choice on high-stakes trials, since this is the reward-maximizing strategy in this task (see ‘simulations’ section below; Kool et al., 2016).

Each participant completed 25 practice trials followed by 200 rewarded trials. Before these practice trials, participants were instructed extensively about the transition structure, the reward distributions of the aliens, and how the stakes manipulation worked. These sections were designed to make sure participants fully understood every task element, introducing one at a time and assessing understanding at several points during the practice session. In all these practice and instructional sessions, there was no time limit for any of the responses.



**Fig 2. The strength of the control-reward trade-off in the experimental paradigms of Experiments 1 and 2.** (A) The novel version of the two-step task used in Experiment 1 embodies a tradeoff between model-based control and reward. For virtually all combinations of learning rate (the degree to which new information is incorporated) and inverse temperature (the randomness of choice), the task shows a strong positive relationship between the degree of model-based control and reward, as measured by linear regression. (B) In the Daw two-step task used in Experiment 2, increased model-based control does not yield increased reward. This surface map plots the strength of this relationship as a function of the learning rate and inverse temperature. The plot is uniformly flat around zero across the entire range.

### Dual-system RL model

In order to estimate the degree of model-based control on high- and low-stakes trials, we employed a dual-system RL model (see Supporting Information; Daw, Gershman, Seymour, Dayan, & Dolan, 2011). This model consists of a model-free system and a model-based system that both represent values for the actions at the first-stage state. The model-free system learns state-action values for all first- and second-stage states through a simple temporal difference learning algorithm (Sutton & Barto, 1998). In essence, this system simply increases the value of actions that lead to outcomes that are more positive than expected (i.e., that produce a positive prediction error), and decreases the value of actions that lead to outcomes that are less positive than expected (a negative prediction error). The model-based system, on the other hand, computes the values of available actions on the fly, by combining the transition structure of the task with the second-stage model-free values to plan towards goals. In other words, this system plans through its internal model of the experiment to find the expected second-stage outcomes for each first-stage action. Our model included two weighting parameters ( $w_{\text{low}}$  and  $w_{\text{high}}$ ) that determine the relative contribution of the model-based and model-free system on low- and high-stakes trials, respectively. Model-based control is induced by weights closer to 1, whereas model-free control is induced by weights closer to 0.

The model also includes an ‘inverse temperature’ parameter  $\beta$  which controls the exploitation-exploration tradeoff between two choice options given their difference in value. This parameter dictates choice

probability through a logistic function ranging from uniform likelihood across actions (pure exploration, insensitive to the value of actions) towards always picking the action with the highest value, regardless of the difference between options (pure exploitation, never exploring lower value options). In addition to these two choice-related parameters, the model also included a learning rate parameter  $\alpha$  that governs the degree to which action values are updated after a reward outcome, an eligibility trace parameter  $\lambda$  that controls the degree to which outcome information at the second stage transfers to the start stage, and ‘stickiness’ parameters  $\pi$  and  $\rho$  that capture perseveration on either the response or the stimulus choice.

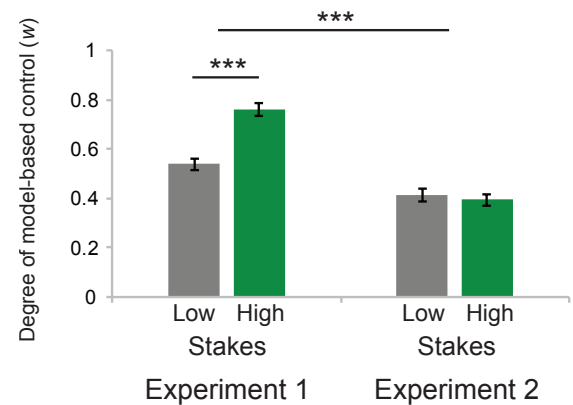
**Simulations.** In order to confirm that the model-based strategy is indeed the reward-maximizing strategy on this task, we used RL simulations to estimate the relationship between the degree of model-based control and reward. Specifically, we used the dual-system RL model to simulate performance on the two-step task, without the stakes manipulation, for agents varying from completely model-free ( $w = 0$ ) to completely model-based ( $w = 1$ ). We recorded the reward rate, the average number of points collected per trial, for each these allocations between RL strategies. We then estimated the strength of the relationship between model-based control and reward through linear regression. This process was repeated 1000 times across a range of inverse temperatures ( $\beta$ ) and learning rates ( $\alpha$ ), resulting in a surface of standardized regression coefficients in a space of these RL parameters (for details, see Kool et al., 2016).

The results of this analysis are shown in Figure 2A. The key feature of this surface map is that there is a strong relationship between model-based control and reward across a large region of the sample parameter space. This analysis confirms that the model-based strategy is associated with increased reward on this two-step task. Therefore, if model-based planning is costly, it would be rational to shift allocation towards the model-based system when potential incentives are high, since the cost-benefit tradeoff is more advantageous under those circumstances.

## Results

**Computational model analyses.** We used maximum *a posteriori* estimation with empirical priors on the parameters to fit the free parameters to our participant's choices in this task (Gershman, 2016). Table 1 reports the estimated parameters. Replicating previous work, we found that the weighting parameters indicated a mix of model-free and model-based action values (mean  $w = 0.65$ ). We also confirmed that the model-based strategy is associated with increased accuracy in this task (Kool et al., 2016), since we found that individual differences in the model-based weighting parameters predicted the reward rate, the average number of points per trial (corrected by the average reward value across each participant's reward distribution),  $p_s < 0.001$  (see Table 2). Most importantly, we found that the degree of model-based control was significantly increased on high-stakes trials (mean  $w = 0.76$ ) compared to low-stakes trials (mean  $w = 0.54$ ) [ $t(97) = 4.67$ ,  $p < 0.001$ , Cohen's  $d = 0.47$ ] (Figure 3)<sup>1</sup>. The Supporting Information reports additional analyses that replicate these findings using a multilevel logistic regression model for both this and the following experiment.

One potential concern in the model-fitting procedure above is that we only varied the weighting parameter between conditions, forcing any behavioral changes induced by the stake manipulation on this parameter. Therefore, we also fit a version of the RL model that varied all parameters between the high- and low-stakes conditions<sup>2</sup>. These analyses replicated the effect that the weighting parameter was larger when the stakes were high (mean = 0.70) compared to when they were low (mean = 0.55) [ $t(97) = 3.15$ ,  $p = 0.002$ ,  $d = 0.32$ ]. In addition, they revealed that the inverse temperature was larger on high-stakes trials compared to low-stakes trials [ $t(97) = 3.95$ ,  $p < 0.001$ ,  $d = 0.40$ ], indicating that participants were also more likely to pick



**Fig 3. Degree of model-based control in low- and high-stakes conditions for Experiments 1 and 2.** We observed an increase in model-based control in the high-stakes condition in Experiment 1, but not in Experiment 2. Error bars indicate within-subject SEM.

\*\*\*  $p < 0.001$

the high-value action at the first stage when the stakes were high, rather than exploring the alternative choice option. This is a rational pattern of behavior on this task, since the cost of exploration is higher when the amount of potential reward increases. The other parameters of the model did not show a significant difference between stake size conditions,  $p_s > 0.10$  (see Supporting Information for more detail).

**Behavioral performance.** Although our model-fitting procedure has the virtue of precision, it has the drawback of being relatively opaque. In order to provide a more intuitive description and statistical test of our data, we analyzed choice probabilities as a function of the previous trial's reward history and the relation between current and previous first-stage state. The basic rationale for this analysis is that for any model-based strategy, the sign of the second-stage prediction error on the previous trial has to influence the likelihood of retaining the same goal on the next trial. For a model-free strategy, this relationship only holds when a reinforced action is presented in consecutive trials.

In this paradigm, the implicit equivalence between the two first-stage states allows for such a quantification of the goal-directed component (Doll et al., 2015). The model-based strategy constructs action values by planning towards the second-stage model-free values, allowing it to generalize knowledge learned from both starting states. Thus, prediction errors at the second

<sup>1</sup> We replicated this result in an experiment ( $n = 93$ ) that was identical to this one, except with negative reward (reward range -4 to +5). We again found the degree of model-based control was higher in the high-stakes (mean  $w = 0.62$ ) compared to low-stakes trials (mean  $w = 0.51$ ) [ $t(93) = 2.81$ ,  $p = 0.006$ ,  $d = 0.29$ ].

<sup>2</sup> These results should be interpreted with some caution because the inverse temperature parameter (exploration vs. exploitation) interacts multiplicatively with the weighting parameter (model-based vs. model-free) to determine choice probabilities. This potentially creates a non-identifiability issue: different combinations of parameter values can result in the same likelihood (Gershman, 2016).



**Table 1.** Best-fitting parameter estimates shown as median plus quartiles across participants and Experiments.

	Predictor	$\beta$	$a$	$\lambda$	$\pi$	$\varrho$	$w_{\text{low}}$	$w_{\text{high}}$
Experiment 1	25 <sup>th</sup> percentile	0.46	0.05	0.00	-0.04	-0.34	0.00	0.63
	Median	0.64	0.82	0.25	0.17	-0.12	0.62	0.95
	75 <sup>th</sup> percentile	2.97	1.00	0.85	0.68	0.05	0.86	1.00
Experiment 2	25 <sup>th</sup> percentile	2.10	0.02	0.39	0.04	-0.02	0.00	0.00
	Median	3.36	0.29	0.69	0.17	0.05	0.43	0.44
	75 <sup>th</sup> percentile	3.95	0.56	1.00	0.37	0.13	0.74	0.62

stage equally affect first-stage preferences, regardless of whether this trial starts with the same starting state as the previous trial. In other words, if the outcome of the previous trial is better than expected (a positive prediction error), the model-based system will be more likely to revisit the previous second-stage state, independent of which first-stage state is presented. The opposite is true when the previous outcome was worse than expected (a negative prediction error). In that case, the model-based system will reduce the likelihood of revisiting that second-stage state. Thus, the effect of the sign of the previous reward prediction error on the probability that the previous second-stage state is revisited represents the model-based component, since it carries over to the next trial even when the first-stage states are different.

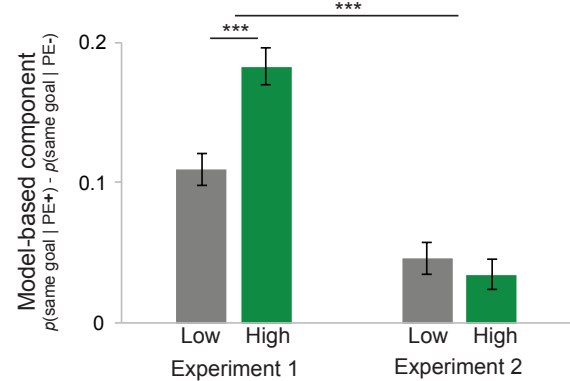
We used this logic to test our cost-benefit hypothesis by analyzing the proportion of trials on which participants revisited the previous second stage state as a function of the sign of the prediction on the previous trial (positive vs. negative), separately for the high- and low-stakes trials (Figure 4). The trial-by-trial estimates for the second-stage prediction errors for these analyses were derived from the computational model described above. Consistent with our previous results, we found that the model-based choice component, i.e., the effect of the previous trial's prediction error's sign on the probability of revisiting the second-stage state, was significantly higher on high-stakes trials compared to low-stakes trials [ $t(97) = 3.70, p < 0.001, d = 0.37$ ].

## Discussion

These findings suggest that metacontrol is governed by a cost-benefit analysis. Participants exerted increased model-based control on high-stakes trials compared to low-stakes trials, indicating increased willingness to engage in effortful planning.

## Experiment 2

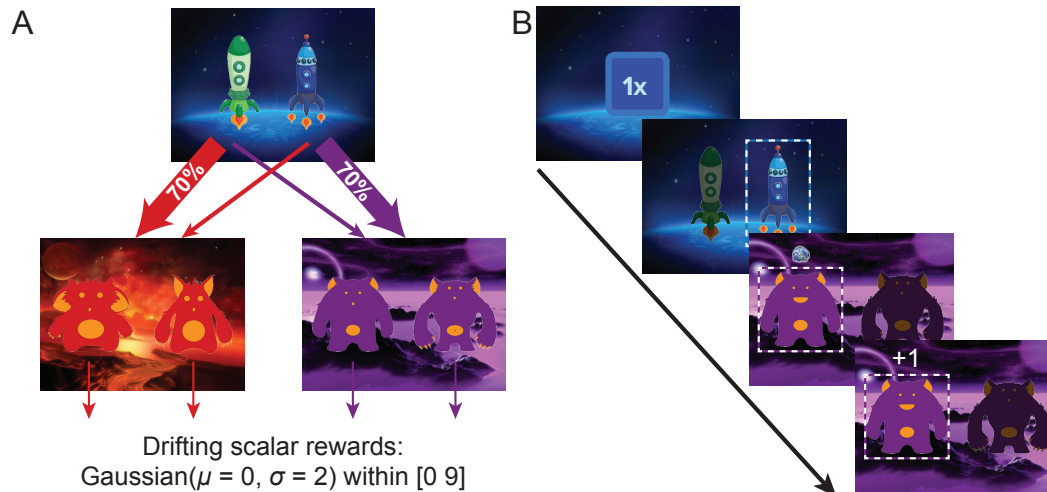
The results from Experiment 1 are consistent with our cost-benefit hypothesis, but also with a simple decision heuristic that reflexively increases model-based control in any context of enhanced reward opportunity, but without assessing the task-specific advantage of either control mechanism as we propose. This alternative heuristic model predicts that we should still observe an

**Fig 4. Model-based choice component for both stake-size conditions and experiments.** We calculated the model-based choice component as the increase in probability of choosing the same goal as on the previous trial after positive outcomes compared to negative outcomes.\*\*\*  $p < 0.001$ **Table 2.** Correlation coefficients between model-based weighting parameter and reward rate, the average number of points per trial (corrected for differences in chance performance) in Experiments 1 and 2.

	Parameter	$r$	$p$
Experiment 1	$w_{\text{low}}$	0.54	$< 0.001$
	$w_{\text{high}}$	0.32	0.001
Experiment 2	$w_{\text{low}}$	-0.01	0.93
	$w_{\text{high}}$	-0.02	0.81

increase in model-based control on high-stakes trials even when model-based control yields no accuracy advantage.

In prior work (Kool et al., 2016), we found that the original two-step task developed by Daw and colleagues (2011) embodies precisely this detachment of model-based control and performance accuracy. As described in more detail below, this task comprises only one first-stage state with two choices that lead to the second-stage states with different probabilities: Each choice leads to one state more frequently than the other. Importantly, the model-free, but not the model-based, system, is affected by these low-probability transitions, since the latter uses the transition structure to discount those associations. Crucially, and in contrast to the two-step task employed in Experiment 1, model-based control in the Daw task does not result in increased reward (see also Akam, Costa, & Dayan, 2015). This difference



**Fig 5. Design of Experiment 2.** (A) State transition structure. Each first-stage choice has a high probability of transitioning to one of two second-stage states and a low probability of transitioning to the other. Each second-stage choice is associated with a probability of obtaining a binary reward. (B) Timeline of events. At the start of each trial, a cue indicated whether 1 point (low stakes) or 5 points (high stakes) could be won. On this low-stakes trial, the alien gives a piece of treasure, and so the total number points earned is 1.

between the tasks allows us to discriminate between the two hypotheses about the nature of cost-benefit arbitration. If the stake-size effect is driven by an incentive-based heuristic, then high stakes should increase model-based control in both tasks. However, if the brain estimates task-specific values for both systems in order to guide arbitration, then high stakes should not increase model-based control in the Daw task. Here, we test these contrasting predictions.

### Method

**Participants.** One hundred participants (range: 18–58 years of age; mean: 34 years of age; 43 female) were recruited on Amazon Mechanical Turk to participate in the experiment. Participants gave informed consent, and the Harvard Committee on the Use of Human Subjects approved the study. No participants timed out on more than 20% of all trials. We excluded all trials on which participants timed out (average 4.1%).

**Materials and Procedures.** The second experiment employed the Daw two-step decision making task (Figure 5A) (Daw et al., 2011; Decker, Otto, Daw, & Hartley, 2016). The response buttons and response deadlines were identical to those used in Experiment 1. Participants made an initial choice between a pair of spaceships. This choice then led probabilistically to one of the two second-stage states. Each spaceship led to one planet more frequently than the other (70% vs. 30%). On each planet, the participant made a second choice between two aliens that work at different space mines (the second-stage states), and offered them a chance to win a piece of space treasure (in contrast with

the scalar reward in Experiment 1). Participants were told that sometimes the aliens were in a good part of the mine, where they were more likely to deliver a piece of space treasure. At other times, the aliens were mining in a bad spot, and they were less likely to deliver space treasure. The reward probabilities of these mines changed slowly over the course of the experiment, encouraging learning throughout the entire session. One pair of aliens was initialized with probabilities of 0.25 and 0.75, and the other pair with probabilities of 0.4 and 0.6, after which they changed according to a Gaussian random walk ( $\sigma = 0.025$ ) with reflecting bounds at 0.25 and 0.75 for the remainder of the experiment. A new set of randomly drifting reward distributions was generated for each participant. At the end of the experiment, participants were given 1¢ for every four points they earned, so that the maximal reward was worth the same in dollar cents for both experiments (Experiment 1: 9 points; Experiment 2: 1 point; both 0.25¢).

In this task, a pure model-free strategy learns about the spaceships' action values in an associative manner, increasing the probability of choosing a spaceship if it previously led to reward, independent of the type of transition that preceded this reward. Choice under a model-based strategy, however, takes into account the type of transition that was preceded by the reward, since the values of the first-stage actions are computed in a prospective fashion based on the transition structure and the learned value of each of the second-stage aliens.

As an example, imagine an agent picking the green spaceship, transitioning through a rare transition to the purple planet, and then receiving a reward. What spaceship should the agent choose on the following trial? A pure model-free agent would be more likely to repeat

the previous trial's choice (green), since this is the choice that led to the positive reward outcome. However, this choice would be less likely to lead to the purple planet where reward was received given the transition structure of the task. In light of this, a model-based agent would be more likely to switch spaceship choices (to blue) on the next trial, thus maximizing the expectation of reward given its model of the environment.

We implemented the same stakes manipulation used in Experiment 1 (Figure 5B). At the start of each trial, a cue indicated how many points could be won on a trial. With 50% probability, this cue indicated that space treasure would only be worth 1 point, otherwise it indicated that space treasure was worth 5 points. After participants observed the reward outcome on the trial, they were also given an explicit computation of the number of points gained on that trial multiplied by the trial's stake. The running score was always available in the top-right corner of the screen.

Participants again completed 25 practice trials followed by 200 rewarded trials. Before these practice trials, participants were instructed extensively on the task, similarly to Experiment 1.

### **Dual-system RL model**

We again employed the dual-system RL model to capture the degree of model-based control on high- and low-stakes trials. This model was identical to the one described in Experiment 1, but with the number of states and actions changed so as to accommodate the new task structure.

*Simulations.* We first used RL simulations to show that the model-based strategy is not reward-advantageous in this task, a basic premise of our experimental design. As before, we ran RL simulation to estimate the strength of the relationship between model-based control and reward on this task, across a wide range of inverse temperatures and learning rates (see Kool et al., 2016 for details), by recording the reward rate, the proportion of trials with a positive outcome, of agents varying from completely model-free ( $w = 0$ ) to completely model-based ( $w = 1$ ). We then used linear regression to calculate the strength of the relationship between reward and control for each combination of inverse temperature and learning rate.

The results are shown Figure 2B. Importantly, the regression coefficients of the resulting surface map are uniformly close to zero. This indicates that nowhere across the sample range of RL parameters there was a positive relationship between model-based control and reward. This confirms that the model-based strategy is not reward-advantageous in the Daw two-step task. Therefore, a cost-benefit account would not predict an increase in control in response to high stakes. However, an incentive-heuristic account, which does not rely on

the explicit representation of costs and rewards, would still predict an increase in model-based control on the high-stakes trials.

The absence of this tradeoff is produced by the interaction of several factors (for a detailed description, see Kool et al., 2016). First, the model-based strategy is weakened in this task, because the first-stage choices carry relatively decreased importance due to the rare transitions, the existence of a second-stage choice, and the low distinguishability between second-stage reward distributions. Furthermore, it is much harder for the model-based system to establish accurate representations of the second-stage probabilistic reward outcomes, compared to the scalar (point value) outcomes used in the novel two-step task. To see this, note that one needs multiple reward observations from the same alien in this Experiment, integrating across outcomes, whereas in Experiment 2 a single observation theoretically provides full information about that alien's value.

### **Results**

*Computational model analyses.* We fit the dual-system RL model to the participants' choices. This model was largely similar to the model used in Experiment 1, except the number of states and actions were changed to reflect the structure of the novel paradigm. Most importantly, the model again includes weighting parameters that determine the contributions of the model-based and model-free system on the high- and low-stakes trials separately.

The estimated parameters for Experiment 2 are reported in Table 1. As before, we found that the weighting parameters suggested that both model-based and model-free strategies were mixed in the population (mean  $w = 0.40$ ). Also, consistent with our previous findings and RL simulations, we found that individual differences in the model-based weighting parameters did not predict the reward rate, the proportion of trials with a positive outcome,  $p_s > 0.80$  (see Table 2). Most importantly, we observed no difference in model-based control in the high-stakes trials (mean  $w = 0.39$ ) compared to low-stakes trials (mean  $w = 0.41$ ), [ $t(99) = -0.38, p = 0.71, d = -0.04$ ] (Figure 3), in contrast with the increase in model-based control in Experiment 1.

As before, we estimated a model that varied all parameters between stake-size conditions. This model replicated the finding that model-based control did not differ on high-stakes trials (mean  $w = 0.36$ ) compared to low-stakes trials (mean  $w = 0.38$ ), [ $t(99) = -0.48, p = 0.63, d = -0.05$ ]. However, we found that the inverse temperature, controlling the exploration-exploitation tradeoff, show a significant increase under high-stakes compared to low-stakes trials [ $t(99) = 4.00, p < 0.001, d = 0.40$ ]. This result suggests that, independent of their use of RL strategy, participants showed more exploiting



behavior under high-stakes trials. As noted before, this is a rational decision in the current task. Furthermore, this result rules out the potential concern that participants were simply not paying attention to the stake-size cue selectively for this task. Rather, they show that the lack of a stake-size effect on model-based control was the result of a cost-benefit analysis that takes into account the costs and benefits of either strategy. The other parameters of the model did not show a significant difference between stake-size conditions,  $p$ s  $> 0.50$  (see Supporting Information for more detail).

**Behavioral performance.** In addition to the computational modeling analysis, we again analyzed choice probabilities by computing a metric of model-based influence. As before, the rationale is that for the model-based strategy, the outcome on the previous trial has to influence the likelihood of retaining the identical goal state on the next trial.

Here, the estimation of the model-based influence on choice follows a slightly different logic than in Experiment 1. The model-based strategy uses the second-stage model-free values and the experiment's transition structure to compute the expected values of the first-stage actions. Therefore, if a positive outcome is obtained at the second stage, the model-based system will increase the likelihood of planning towards that goal on the next trial. After a rare transition, this means that the system will decrease the probability of repeating the previous choice after a reward, because this achieves a higher likelihood to get to the previously rewarded second-stage state. After a rare transition and a loss, the model-based strategy is more likely to stick with the original action, since this decreases the likelihood of getting to the unrewarded state.

We tested our cost-benefit analysis by analyzing the proportion of trials on which participants chose the first-stage action that would most likely lead to the previous second-stage state as a function of the outcome on the previous trial, separately for the high- and low-stakes trials (Figure 4). Consistent with the computational modeling results, we found that this model-based component was not affected by the stakes manipulation [ $t(99) = -0.72, p = 0.470, d = -0.07$ ].

### Comparison between experiments

We directly compared behavior between the two experiments in order to test whether the different task structures yielded reliable differences. First, we found that the increase in model-based control induced by increased incentives was significantly larger in Experiment 1 compared to Experiment 2, both for the weighting parameter estimated from the computational model (Figure 3) [ $t(196) = 3.56, p < 0.001, d = 0.51$ ], as well as for the more direct behavioral estimation of the model-based component (Figure 4) [ $t(196) = 3.29, p =$

$0.001, d = 0.47$ ]. Second, confirming earlier results (Kool et al., 2016), we found a reliable shift in model-based control between experiments, with the average weighting parameter significantly higher in Experiment 1 (mean  $w = 0.65$ ) compared to Experiment 2 (mean  $w = 0.40$ ) [ $t(196) = 6.21, p < 0.001, d = 0.88$ ]. Third, consistent with our previous findings (Kool et al., 2016), multiple regression analyses confirmed that the relationships between the weighting parameter and average reward rate were significantly stronger in the novel paradigm compared to the original paradigm for the low-stakes trials [ $t(194) = 3.97, p < 0.001$ ], as well as on the high-stakes trials [ $t(194) = 2.41, p < 0.05$ ].

### Discussion

In contrast with Experiment 1, participants did not increase model-based control when this strategy was not associated with superior performance. This finding is consistent with a cost-benefit analysis, but not with an incentive-heuristic in which high stakes always trigger increased model-based control.

If model-based control did not yield increased reward, then why did we still observe a mixture of model-based and model-free strategies in Experiment 2? As we have previously noted (Kool et al., 2016), one possibility is that behavior on this task reflects a prior belief that model-based control is associated with increased reward in the real world (where this is presumably valid). Furthermore, the extensive training on the transition structure may have induced an assumption that it should be employed during task performance.

Of course, the two-step tasks used in Experiment 1 and 2 differed in several respects. We have shown that each of these is necessary to yield a robust accuracy advantage for model-based control (Kool et al., 2016). Although there is a strong *a priori* basis for predicting that the relationship between stake size and metacognitive control depends on our intended manipulation of the model-based accuracy advantage for each task, it may instead be moderated incidentally by some other, correlated difference between the tasks. This is a potential topic for further study.

### General Discussion

We find that arbitration between model-based and model-free control is sensitive to the task-specific costs and benefits associated with each system. Participants relied more on model-based control on trials with larger incentives (Experiment 1), but only when this yielded more accurate performance (Experiment 2). This implies that participants estimate the expected reward for each system, and then weigh this against the increased costs of model-based control.

Several contemporary models of metacontrol are broadly consistent with the idea of calculating the costs and benefits of each system, but differ in their details. For instance, some models posit that control is eventually allocated to whichever system yields greater accuracy (Daw et al., 2005), or reward (Rieskamp & Otto, 2006). This cannot explain our stake-size effect on metacontrol, because the relative accuracy of two strategies will not change under a monotonic scaling of reward. Rather, our data favor models that balance accuracy against the cost of cognitive control (Gershman et al., 2015; Griffiths et al., 2015), such as increased decision time (Keramati et al., 2011) or limited cognitive resources (Kurzban et al., 2013).

How might this cost be assigned? In principle, it could be computed from a model of opportunity costs. For real world tasks, however, this is likely to be prohibitively demanding. One way around this problem is to assign model-based control an intrinsic subjective cost. Consistent with the observation that model-based control relies on cognitive control (Otto et al., 2013), several prior proposals that cognitive control involves an intrinsic cost (Botvinick & Braver, 2015). For example, Shenhav and colleagues (2013) propose that the brain selects actions based on the ‘expected value of control’, the expected reward associated with exerting cognitive control discounted by the cost of this exertion. Similarly, Motivational Intensity Theory (Brehm, Wright, Solomon, Silka, & Greenberg, 1983) proposes that effort is only invested when success on the task is both possible and worthwhile, and the effort justified. Our results suggest a similar process governs the deployment of model-based control, as formalized within the RL framework.

Our results also demonstrate that the costs of model-based control are weighed against its accuracy benefits in a task-specific manner: High stakes failed to increase model-based control for a task where it provided no advantage. Several current models of metacontrol cannot accommodate this finding, because they do not explicitly compare the rewards obtained by both systems. Rather, these depend on heuristic approximations of model-based advantage: For instance, by assuming perfect model-based accuracy (Keramati et al., 2011), calculating errors in the state transitions predicted by the model-based system (Lee et al., 2014), or assuming a model-based advantage when uncertainty in model-free value estimates is high (Keramati et al., 2011; Lee et al., 2014; Pezzulo et al., 2013). Our data instead favor a mechanism that explicitly compares the rewards obtained by model-based versus model-free control.

Combining these insights, our results favor a model of metacontrol would first learn task-dependent reward history of different control mechanisms, and then integrate these rewards with their unique cognitive control costs to guide controller allocation. We provide

the first direct support for this class of model by explicitly comparing two tasks that (1) both distinguish model-based and model-free control, (2) vary in the benefits of model-based control (Akam et al., 2015; Kool et al., 2016), and (3) both include trial-by-trial variation in the magnitude of rewards at stake. These findings invite a computational formalization, as well as neuroimaging work to establish the locus of metacontrol in the brain.

## Contributions

All authors contributed to the study design. Testing and data collection were performed by W. Kool. W. Kool performed the data analysis. All authors wrote the paper, and approved the final version for submission.

## Acknowledgements

We thank Catherine Hartley for sharing stimuli, and the Moral Psychology Research Laboratory and Computational Cognitive Neuroscience Laboratory for advice and assistance. This research was supported by grant N00014-14-1-0800 from the Office of Naval Research and by the Center for Brains, Minds and Machines, funded by NSF STC award CCF-1231216.

## Open Practices

All data and materials have been made publicly available at <https://www.github.com/wkool/arbitration>.

## References

- Akam, T., Costa, R., & Dayan, P. (2015). Simple plans or sophisticated habits? State, transition and learning interactions in the two-step task. *PLoS Computational Biology*, 11, e1004648.
- Botvinick, M. M. (2007). Conflict monitoring and decision making: Reconciling two perspectives on anterior cingulate function. *Cognitive, Affective, & Behavioral Neuroscience*, 7, 356-366.
- Botvinick, M. M., & Braver, T. (2015). Motivation and cognitive control: From behavior to neural mechanism. *Annual Review of Psychology*, 66, 83-113.
- Brehm, J. W., Wright, R. A., Solomon, S., Silka, L., & Greenberg, J. (1983). Perceived difficulty, energization, and the magnitude of goal valence. *Journal of Experimental Social Psychology*, 19, 21-48.
- Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-based influences on humans' choices and striatal prediction errors. *Neuron*, 69, 1204-1215.
- Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature Neuroscience*, 8, 1704-1711.

- Decker, J. H., Otto, A. R., Daw, N. D., & Hartley, C. A. (2016). From creatures of habit to goal-directed learners: Tracking the developmental emergence of model-based reinforcement learning. *Psychological Science*, 27, 848-858.
- Dickinson, A. (1985). Actions and habits: The development of behavioural autonomy. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 308, 67-78.
- Dolan, R. J., & Dayan, P. (2013). Goals and habits in the brain. *Neuron*, 80, 312-325.
- Doll, B. B., Duncan, K. D., Simon, D. A., Shohamy, D., & Daw, N. D. (2015). Model-based choices involve prospective neural activity. *Nature Neuroscience*, 18, 767-772.
- Gershman, S. J. (2016). Empirical priors for reinforcement learning models. *Journal of Mathematical Psychology*, 71, 1-6.
- Gershman, S. J., Horvitz, E. J., & Tenenbaum, J. B. (2015). Computational rationality: A converging paradigm for intelligence in brains, minds, and machines. *Science*, 349, 273-278.
- Gershman, S. J., Markman, A. B., & Otto, A. R. (2014). Retrospective revaluation in sequential decision making: A tale of two systems. *Journal of Experimental Psychology: General*, 143, 182-194.
- Gillan, C. M., Kosinski, M., Whelan, R., Phelps, E. A., & Daw, N. D. (2016). Characterizing a psychiatric symptom dimension related to deficits in goal-directed control. *eLife*, 5, e11305.
- Gläscher, J., Daw, N., Dayan, P., & O'Doherty, J. (2010). States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron*, 66, 585-595.
- Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in Cognitive Science*, 7, 217-229.
- Kahneman, D. (2003). A perspective on judgment and choice: Mapping bounded rationality. *American Psychologist*, 58, 697-720.
- Keramati, M., Dezfouli, A., & Piray, P. (2011). Speed/accuracy trade-off between the habitual and the goal-directed processes. *PLoS Computational Biology*, 7, e1002055.
- Kool, W., & Botvinick, M. (2014). A labor/leisure tradeoff in cognitive control. *Journal of Experimental Psychology: General*, 143, 131-141.
- Kool, W., Cushman, F. A., & Gershman, S. J. (2016). When does model-based control pay off? *PLoS Computational Biology*, 12, e1005090.
- Kool, W., McGuire, J. T., Rosen, Z. B., & Botvinick, M. M. (2010). Decision making and the avoidance of cognitive demand. *Journal of Experimental Psychology: General*, 139, 665-682.
- Kool, W., Shenhav, A., & Botvinick, M. (2017). Cognitive control as cost-benefit decision making. In T. Egner (Ed.), *Wiley Handbook of Cognitive Control* (pp. 167-189). Chichester, West Sussex, UK: John Wiley & Sons.
- Kurzban, R., Duckworth, A. L., Kable, J. W., & Myers, J. (2013). An opportunity cost model of subjective effort and task performance. *Behavioral and Brain Sciences*, 36, 661-726.
- Lee, S. W., Shimojo, S., & O'Doherty, J. P. (2014). Neural computations underlying arbitration between model-based and model-free Learning. *Neuron*, 81, 687-699.
- Otto, A. R., Gershman, S. J., Markman, A. B., & Daw, N. D. (2013). The curse of planning: Dissecting multiple reinforcement-learning systems by taxing the central executive. *Psychological Science*, 24, 751-761.
- Payne, J. W., Bettman, J. R., & Johnson, E. J. (1988). Adaptive strategy selection in decision making. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 14, 534-552.
- Pezzulo, G., Rigoli, F., & Chersi, F. (2013). The Mixed Instrumental Controller: Using Value of Information to combine habitual choice and mental simulation. *Frontiers in Psychology*, 4, 92.
- Rieskamp, J., & Otto, P. E. (2006). SSL: A theory of how people learn to select strategies. *Journal of Experimental Psychology: General*, 135, 207-236.
- Shenhav, A., Botvinick, Matthew M., & Cohen, Jonathan D. (2013). The expected value of control: An integrative theory of anterior cingulate cortex function. *Neuron*, 79, 217-240.
- Sloman, S. A. (1996). The empirical case for two systems of reasoning. *Psychological Bulletin*, 119, 3-22.
- Smittenaar, P., FitzGerald, T. H. B., Romei, V., Wright, N. D., & Dolan, R. J. (2013). Disruption of dorsolateral prefrontal cortex decreases model-based in favor of model-free control in humans. *Neuron*, 80, 914-919.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press.
- Thorndike, E. L. (1911). *Animal Intelligence: Experimental Studies*. New York: The Macmillan Company.
- Westbrook, A., Kester, D., & Braver, T. S. (2013). What is the subjective cost of cognitive effort? Load, trait, and aging effects revealed by economic preference. *PLoS ONE*, 22, e68210.