

Running title: REWARD LEARNING IN THE CEREBELLUM

The role of the cerebellum in learning to predict reward: evidence from cerebellar ataxia

Jonathan Nicholas^{1,2}, Christian Amlang^{3,4}, Chi-Ying R. Lin⁵, Leila Montaser-Kouhsari⁶, Natasha Desai^{3,4}, Ming-Kai Pan^{7,8,9}, Sheng-Han Kuo^{3,4}, Daphna Shohamy^{1,2,10}

Abstract

Recent findings in animals have challenged the traditional view of the cerebellum solely as the site of motor control, suggesting that the cerebellum may also be important for learning to predict reward from trial-and-error feedback. Yet, evidence for the role of the cerebellum in reward learning in humans is lacking. Moreover, open questions remain about which specific aspects of reward learning the cerebellum may contribute to. Here we address this gap through an investigation of multiple forms of reward learning in individuals with cerebellum dysfunction, represented by cerebellar ataxia cases. Nineteen participants with cerebellar ataxia and 57 age- and sex-matched healthy controls completed two separate tasks that required learning about reward contingencies from trial-and-error. To probe the selectivity of reward learning processes, the tasks differed in their underlying structure: while one task measured incremental reward learning ability alone, the other allowed participants to use an alternative learning strategy based on episodic memory alongside incremental reward learning. We found that individuals with cerebellar ataxia were profoundly impaired at reward learning from trial-and-error feedback on both tasks, but retained the ability to learn to predict reward based on episodic memory. These findings provide evidence from humans for a specific and necessary role for the cerebellum in incremental learning of reward associations based on reinforcement. More broadly, the findings suggest that alongside its role in motor learning, the cerebellum likely operates in concert with the basal ganglia to support reinforcement learning from reward.

Running title: REWARD LEARNING IN THE CEREBELLUM

Author affiliations:

1 Department of Psychology, Columbia University, New York, NY, USA

2 Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY, USA

3 Department of Neurology, Columbia University Medical Center, New York, NY, USA

4 Initiative for Columbia Ataxia and Tremor, Columbia University Medical Center, New York, NY, USA

5 Department of Neurology, Baylor College of Medicine, Houston, TX, USA

6 Department of Neurology, Stanford University School of Medicine, Palo Alto, CA, USA

7 Department of Medical Research, National Taiwan University Hospital, 100 Taipei, Taiwan

8 Department and Graduate Institute of Pharmacology, National Taiwan University College of Medicine, 100 Taipei, Taiwan

9 Cerebellar research center, National Taiwan University Hospital, Yun-Lin Branch, Yun-Lin, Taiwan

10 Kavli Institute for Brain Science, Columbia University, New York, NY, USA

Correspondence to:

Daphna Shohamy

Columbia University Zuckerman Institute Quad 3D, 3227 Broadway, New York, NY 10027

ds2619@columbia.edu

Sheng-Han Kuo

650 W. 168th St, Rm 305, New York, NY, 10032

sk3295@cumc.columbia.edu

Running title: REWARD LEARNING IN THE CEREBELLUM

Running title: REWARD LEARNING IN THE CEREBELLUM

Keywords: reward; reinforcement learning; cerebellum; ataxia

Abbreviations: CA – Cerebellar ataxia, HC – Healthy control, SCA – Spinocerebellar ataxias, MSA – Multiple system atrophy, MSA-C – Multiple system atrophy, cerebellar type, FA – Friedreich's ataxia, IMCA – Immune-mediated cerebellar ataxia, ILOCA – Idiopathic late onset cerebellar ataxia, SARA – Scale for the assessment and rating of ataxia, CCAS – Cerebellar cognitive affective/Schmahmann syndrome scale, MoCA – Montreal cognitive assessment, BDI – Beck's depression inventory, QUIP – Questionnaire for impulsive-compulsive disorders in Parkinson's disease, MCI – Mild Cognitive Impairment

Introduction

It is well established that the cerebellum is required for refining movement through supervised motor learning.¹⁻⁴ The cerebellum receives error signals from climbing fiber input which then alters Purkinje cell plasticity to adapt motor behavior in service of minimizing future error.⁵⁻⁷ However, recent findings have challenged the notion that the cerebellum is solely responsible for supervised learning of motor behavior and instead suggest that the cerebellum may also be involved in the processing of reward more generally.⁸⁻¹⁹ In particular, climbing fiber inputs to the cerebellum encode expected reward,^{13,15,17,19} and cerebellar Purkinje cells have been found to report reward-based prediction errors.^{11,12,18} These signals are essential ingredients for reinforcement learning, or learning that allows an organism to determine from trial-and-error feedback which actions should be taken in order to maximize future expected reward. The presence of reward-related processing in the cerebellum suggests that it may play a role in reinforcement learning alongside its capacity for supervised motor learning.²⁰ This proposal challenges not only our current understanding of cerebellar function, but also our understanding of how the brain learns from reward more broadly.^{5,21}

Although research on the cerebellum's function in reward learning is growing, the vast majority of work has been done in animal models,¹⁰⁻¹⁹ and evidence in humans remains limited. Human

Running title: REWARD LEARNING IN THE CEREBELLUM

neuroimaging studies have revealed correlational evidence that the cerebellum is involved in tasks unrelated to movement,²² however, despite some reports of BOLD activity in the cerebellum in response to reward across several early imaging studies,^{23–25} more direct investigations of the role of cerebellum in reward-related behaviors in humans are lacking. The aim of the present study was to fill this gap by testing whether individuals with damage to the cerebellum, as occurs in cerebellar ataxia (CA), are impaired in their ability to acquire stimulus-reward associations.

Our study builds upon a rich literature focused on learning about reward from trial-and-error feedback. This process has been studied extensively using models of incremental learning, which rely on error-driven rules that summarize experiences with a running average.^{26–28} During reward learning of this type, an agent uses the outcome of a recent decision to associate some stimulus with an action. Following successful learning, actions that are more likely to be rewarded are more likely to be repeated. This simple mechanism has been evoked to explain conditioning behavior and is well-captured by reward prediction error signals in midbrain dopamine neurons that project to the striatum.^{27,29} This error signal is also precisely what has been implicated in recent animal models of cerebellar contributions to reward learning,⁹ suggesting an additional, albeit unclear, role for the cerebellum in this process. Whether these cerebellar contributions are actually needed for successful incremental reward learning in humans is at present unknown.

To answer this question, we asked individuals with CA to complete a series of tasks that required them to learn associations between stimuli from trial-and-error feedback in order to maximize expected reward. CA is defined as a lack of coordination caused by disorders that affect cerebellar function.³⁰ A large variety of conditions can cause CA, ranging from immune-mediated disease to genetic and neurodegenerative disorders. Given the presence of cerebellar dysfunction in CA cases, studying individuals with CA is a common method used to investigate the necessary physiological functions of the cerebellum in humans.

Nineteen individuals with CA and 57 age- and sex-matched healthy controls (HC) completed two tasks (**Figure 1**). The first, referred to throughout as the *incremental learning* task, allowed us to measure each participants' ability to learn about reward incrementally. This task was motivated by recent work using a similar simplified paradigm to investigate cerebellar-based incremental learning in non-human primates.^{10,11} The second task, referred to throughout as the *multiple*

learning strategies task, allowed us to measure whether any impairments were specific to incremental learning alone. In the multiple learning strategies task, learning about reward can be supported by an alternative strategy based on episodic memory for trial-unique past outcomes. Healthy adults readily use of both of these strategies in this task.^{31,32} We hypothesized that cerebellar dysfunction would lead specifically to impaired incremental reward learning relative to healthy controls.

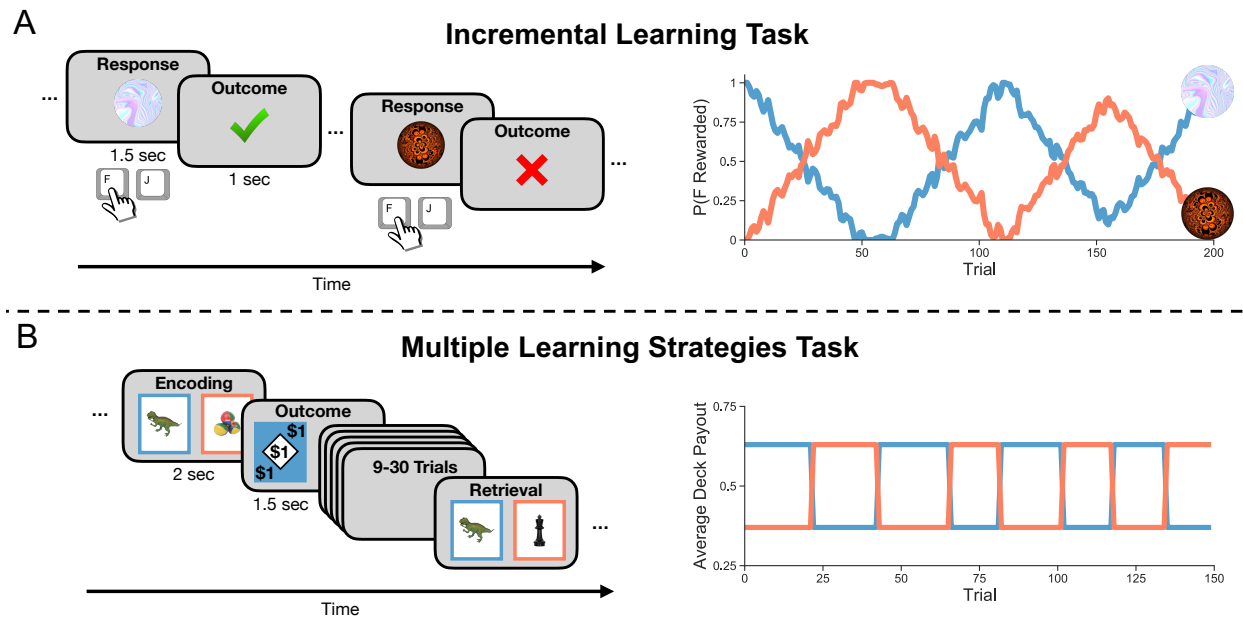


Figure 1 Design of the incremental learning and multiple learning strategies tasks. (A) Left: Trial design for the incremental learning task. Participants saw one of two fractal cues on the screen and were required to press either the F key with their left hand or the J key with their right hand. Following their choice, they received binary probabilistic feedback about whether they were correct or not. Right: Drifting cue-response-reward contingencies over the course of the incremental learning task. The probability that the F key is rewarded is shown for each cue in blue and orange. **(B)** Left: Trial design or the multiple learning strategies task. Participants chose between two decks of cards (one blue and one orange) and received an outcome between \$0-\$1 in intervals of 20 cents. Each card featured a trial-unique object that could repeat once every 9-30 trials. Participants were told that if they saw the same card again, it would be worth the same amount as the first time that it appeared. Right: An example of how average deck value reversed throughout the course of the multiple learning strategies task.

Materials and Methods

Cerebellar Ataxia Participants

Nineteen individuals with cerebellar ataxia (CA) were recruited from the Ataxia Clinic, Columbia University Medical Center and completed both tasks (see **Table 1** for information about basic CA participant demographics and diagnoses). Due to hardware issues, data from one participant on each task was not saved. The first CA participant also completed a shorter pilot version of the incremental learning task, and several changes were made before running this task on the other 18 CA participants. Thus, the final sample for the incremental learning task was 17 CA participants, and the final sample for the multiple learning strategies task was 18 CA participants. Task order was counterbalanced such that 10 CA participants completed the incremental learning task prior to the multiple learning strategies task, and 9 CA participants completed tasks in the opposite order. A neuropsychological battery comprised of the Montreal Cognitive Assessment (MOCA), Beck's Depression Inventory (BDI), MESA digit forward and backward span, trail making test A and B, and the cerebellar cognitive affective syndrome scale (CCAS), lasting approximately 30 minutes was conducted between tasks for each participant. This battery was specifically selected based on the current understanding of the cerebellum's role and association with non-motor symptoms, such as depression,³³ executive function,^{34,35} and attention.³⁶

Healthy Controls

Age- and sex-matched participants were recruited through Amazon Mechanical Turk using the Cloud Research Approved Participants feature.³⁷ To account for potential variability due to online data collection, three matched controls were collected for each CA participant, bringing the total number of controls to 57 (3:1 match). Of these, data from one control was excluded for the multiple learning strategies task due to random responding. Task order was counterbalanced such that the tasks were completed in the identical order to each control's matched CA participant. A modified online neuropsychological battery consisting of 7 measures (see **Supplementary Material**) was completed in between each task for comparison to individuals with CA. Five of these measures (Semantic Fluency, Phonemic Fluency, Category Switching, Similarities and Go No Go) were

Running title: REWARD LEARNING IN THE CEREBELLUM

directly taken from the CCAS, and two others were comprised of the MESA digit forward back backward span. Participant recruitment was restricted to the United States. Before starting each task, all participants were required to score 100% on a quiz that tested their comprehension of the instructions and were made to repeat the instructions until this score was achieved. Informed consent was obtained with approval from the Columbia University Institutional Review Board.

Table 1. Basic CA participant demographics and neuropsychiatric measures

Participant	Age (years)	Sex	Diagnosis	CCAS	MoCA	BDI	QUIP	FDS	BDS	TMTA (sec)	TMTB (sec)
Participant 1	40	M	SCA3	59	21	15	40	8	5	50	150
Participant 2	33	M	SCA3	77	21	9	12	11	5	21	261
Participant 3	20	F	SCA2	92	27	23	10	13	4	41	71
Participant 4	52	F	MSA-C	85	27	4	0	8	3	23	158
Participant 5	61	F	SCA2	92	23	15	38	8	4	49	169
Participant 6	56	M	MSA-C	100	26	7	8	13	4	66	127
Participant 7	52	M	SCA2	62	26	13	55	11	4	105	287
Participant 8	41	F	SCA2	95	29	18	2	10	8	45	97
Participant 9	43	M	SCA1	87	27	0	6	10	6	52	104
Participant 10	62	F	MSA-C	70	21	4	25	6	6	45	121
Participant 11	54	M	SCA2	72	21	0	5	11	5	32	100
Participant 12	67	F	ILOCA	101	29	12	16	12	5	51	88
Participant 13	60	F	SCA3	104	28	1	0	14	9	42	113
Participant 14	51	F	SCA10	74	25	13	8	8	2	35	83
Participant 15	66	M	SCA1	60	23	4	20	7	3	66	127
Participant 16	49	M	IMCA	86	28	17	6	13	7	81	225
Participant 17	54	M	FA	98	26	24	8	10	8	50	80
Participant 18	33	F	FA	84	27	3	26	9	4	38	84
Participant 19	54	F	IMCA	113	28	18	4	11	8	33	62

SCA – Spinocerebellar ataxias, MSA-C – Multiple system atrophy, cerebellar type, ILOCA – Idiopathic late onset cerebellar ataxia, IMCA – Immune-mediated cerebellar ataxia, FA – Friedreich's ataxia, CCAS – Cerebellar cognitive affective/Schmahmann syndrome scale, MoCA – Montreal cognitive assessment, BDI – Beck's depression inventory, QUIP – Questionnaire for impulsive-compulsive disorders in Parkinson's disease, FDS – Forward digit span, BDS – Backward digit span, TMTA – Trail making test part A, TMTB – Trail making test part B

Experiment Tasks

Incremental Learning Task

In the incremental learning task (**Figure 1A**), participants were told that they would be playing a game where they were required to press a key, either F or J, whenever one of two symbols was seen, and that they would receive feedback about whether they had pressed correctly following each trial. They were then informed that it was their job to determine which key they should press for each symbol, and that what key is best will change throughout the experiment. Outcomes were determined by a drifting probability such that each button was correct for each image 50% of the time. Critically, these probabilities differed over time, thus encouraging constant learning throughout the task. Participants were told to press the F key with their left index finger and the J key with their right index finger. The response period during which the symbol remained on the screen lasted 1.5 seconds, with feedback displayed for 1 second immediately following the response period. An intertrial interval featuring a fixation cross was shown for an average of 1 second, but varied between 0.5 and 1.5 seconds. Lastly, to provide a rewarding outcome for correct responses, participants were informed that they could earn bonus money based on their performance. Correct responses were worth an additional cent each.

Multiple Learning Strategies Task

The other task completed by participants was previously developed by our lab³² to measure the relative contribution of incremental learning and episodic memory to decisions (**Figure 1B**). Participants were told that they would be playing a card game where their goal was to win as much money as possible. Each trial consisted of a choice between two decks of cards that differed based on their color (red or blue). Participants had two seconds to decide between the decks. The outcome of each decision was then immediately displayed for 1.5 seconds. Following each decision, participants were shown a fixation cross during the intertrial interval period which varied in length (mean = 1.5 seconds, min = 1 seconds, max = 2 seconds). Decks were equally likely to appear on either side of the screen (left or right) on each trial and screen side was not predictive of outcomes. Participants completed a total of 150 trials.

Running title: REWARD LEARNING IN THE CEREBELLUM

Participants were made aware that there were two ways they could earn bonus money throughout the task, which allowed for the use of incremental learning and episodic memory respectively. First, at any point in the experiment one of the two decks was “lucky”, meaning that the expected value (V) of one deck color was higher than the other ($V_{lucky}=63¢$, $V_{unlucky}=37¢$). Outcomes ranged from \$0 to \$1 in increments of 20¢. Critically, the mapping from V to deck color underwent an unsignaled reversal periodically throughout the experiment, which incentivized participants to utilize each deck’s recent reward history in order to determine the identity of the currently lucky deck. Second, in order to allow us to assess the use of episodic memory throughout the task, each card within a deck featured an image of a trial-unique object that could re-appear once throughout the experiment after initially being chosen. Participants were told that if they encountered a card a second time it would be worth the same amount as when it was first chosen, regardless of whether its deck color was currently lucky or not. On a given trial t , cards chosen once from trials $t - 9$ through $t - 30$ had a 60% chance of reappearing following a sampling procedure designed to prevent each deck’s expected value from becoming skewed by choice, minimize the correlation between the expected value of previously seen cards and deck expected value, and ensure that choosing a previously selected card remained close to 50¢.

Following completion of the multiple learning strategies task, we tested participants’ memory for the trial-unique objects. Participants completed up to 54 three-part memory trials. An object was first displayed on the screen and participants were asked whether or not they had previously seen the object and were given five response options: Definitely New, Probably New, Don’t Know, Probably Old, Definitely Old. If the participant indicated that they had not seen the object before or did not know, they moved on to the next trial. If, however, they indicated that they had seen the object before they were then asked if they had chosen the object or not. Lastly, if they responded that they had chosen the object, they were asked what the value of that object was (with options spanning each of the six possible object values between \$0-1).

Computational Models

In order to best capture subjective estimates of incrementally constructed value on each task, we fit computational models to participants' choices. Below we describe each of these models in detail.

Q Learning Models

We modeled incremental reward learning using a Q Learning model, which is a standard model-free reinforcement learner that assumes a store value (Q) for each deck is updated over time^{26,28}. Q is then referenced on each decision in order to guide choices. After each outcome, r_t , the value for an option 1 Q_1 is updated according to the following rule³ if that option is chosen:

$$Q_{1,t+1} = Q_{1,t} + \alpha(r_t - Q_{1,t}) \quad (1)$$

And is not updated if a different option is chosen:

$$Q_{1,t+1} = Q_{1,t} \quad (2)$$

Likewise, if a different option is chosen, its value is updated equivalently. Large differences between estimated value and outcomes therefore have a larger impact on updates, but the overall degree of updating is controlled by the learning rate, α , which is a free parameter constrained to lie between 0 and 1.

For the incremental learning task, the model learned separate Q values for each cue and button combination, such that four Q values were estimated in total. Decisions on this task were then modeled using the following softmax:

$$P(\text{Choose } F) = \sigma(\beta_{0,1} + \beta_{0,2} + \beta_{1,1}(Q_{F,1} - Q_{J,1}) + \beta_{1,2}(Q_{F,2} - Q_{J,2})) \quad (3)$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad (4)$$

such that four inverse temperatures β were estimated to capture a bias toward choosing a key for each cue ($\beta_{0,1}$ and $\beta_{0,2}$) and sensitivity to incrementally learned value for each cue ($\beta_{1,1}$ and $\beta_{1,2}$).

This model is referred to as the “Q Learner” model throughout the text.

Running title: REWARD LEARNING IN THE CEREBELLUM

For the multiple learning strategies task, the model learned separate Q values for each deck color, such that two Q values were estimated in total. Decisions on this task were then modeled using the following softmax:

$$P(\text{ChooseRed}) = \sigma(\beta_1(Q_R - Q_B) + \beta_2(\text{OldValue}) + \beta_3(\text{Old})) \quad (5)$$

such that three inverse temperatures β were estimated to capture sensitivity to incrementally learned value (β_1), sensitivity to the value of previously seen objects (β_2), and a bias toward choosing the deck featuring a previously seen object regardless of its value (β_3). The predictor *OldValue* was the coded true value of a previously seen object (ranging from 0.5 if the value was \$1 on the red deck or \$0 on the blue deck to -0.5 if the value was \$0 on the red deck and \$1 on the blue deck) and the predictor *Old* was coded as 0.5 if the red deck featured a previously seen object and -0.5 if the blue deck did instead. For both of these predictors, trials that did not feature a previously seen object were coded as 0. This model is referred to as the “Hybrid” model throughout the text.

Biased Responder Model

For both tasks, we compared the performance of the Q Learning models to a model which made choices that were completely independent of reward information. For the incremental learning task, this model was simply:

$$P(\text{ChooseF}) = \sigma(\beta_{0,1} + \beta_{0,2}) \quad (6)$$

such that choices depended only on choosing a button to press for each cue throughout the experiment. For the multiple learning strategies task, this model was:

$$P(\text{ChooseRed}) = \sigma(\beta_0) \quad (7)$$

such that choices depended only on preferring one deck over the other throughout the experiment. Our logic in using this model as a baseline was that responses captured by the Q learning models should, at a minimum, outperform a biased responder that did not consider reward in order for it to make meaningful predictions about participants’ behavior.

Posterior Inference and Model Comparison

Model parameters for each participant were estimated using Bayesian inference. The joint posterior was approximated using No-U-Turn Sampling³⁸ as implemented in stan³⁹. Four chains with 2000 samples (1000 discarded as burn-in) were run for a total of 4000 posterior samples per model per subject. Chain convergence was determined by ensuring that the Gelman-Rubin statistic \hat{R} was close to 1 for all parameters. For the incremental learning task, the Q learner did not converge for one CA participant, and so that individual and their matched controls were removed from further model-based analyses. For the multiple learning strategies task, all models for all participants converged.

Under this approach, the likelihood function for all models can be written as:

$$c_t \sim \text{Bernoulli}(\theta_t) \quad (8)$$

where c_t is 1 if the subject chose F (in the response mapping task) or red (in the multiple learning strategies task). Here, θ_t is the linear combination of inverse temperature parameters and predictors explained above for each model. For the Q learning models, the learning rate, α , had the following weakly informative prior:

$$\alpha \sim \beta(0,1) \quad (9)$$

For all models, every inverse temperature parameter had the following weakly informative prior:

$$\beta \sim \mathcal{N}(0,5) \quad (10)$$

Model fit was assessed using approximate leave-one-out cross validation estimated using Pareto-smoothed importance sampling⁴⁰. The expected log pointwise predictive density (ELPD) was computed and used as a measure of out-of-sample predictive fit for each model.

Bayesian Observers

In order to provide a normative performance benchmark, we also simulated beliefs about incremental value as estimated by Bayesian observers for each task. For the incremental learning task, this learner was a Kalman Filter⁴¹ with observation noise set to the true standard deviation of

Running title: REWARD LEARNING IN THE CEREBELLUM

outcomes and process noise set to 1.5. For the multiple learning strategies task, this learner was a reduced Bayesian change-point detection model⁴² with hazard rate equal to the true proportion of deck reversal trials for each participant in the task. Choices in the incremental learning task were made according to which button the observer believed was the most likely to be rewarded for each cue at each time point. Choices in the multiple learning strategies task were made differently depending on whether a previously seen object was present. For trials in which no previously seen object was shown, the observer responded according to its beliefs about deck value. For trials in which a previously seen object was present, however, the observer compared the value of that object to its belief about deck value for the opposing deck and chose accordingly. In this way, the observer was augmented with “perfect” episodic memory.

Regression Models

Mixed effects Bayesian regressions were used to test effects of group (CA participant or control). Group membership was allowed to vary randomly by CA participant identifier, *pid*, such that CA participants and matched controls were assigned the same ID. In these models, *GroupID* was coded as -0.5 for CA participants and 0.5 for controls. We additionally controlled for working memory ability by including backwards digit span scores, *dsBwd*, as a standardized covariate in these analyses.

For the incremental learning task, we assessed behavioral incremental learning performance using the following logistic regression:

$$\begin{aligned} p(\text{Correct}) = & \sigma(\beta_0 + b_{0,pid[t]} + \text{GroupID}_t(\beta_1 + b_{1,pid[t]}) + \\ & pFReward1_t(\beta_2 + b_{2,pid[t]}) + \text{GroupID}_t \times pFReward1_t(\beta_3 + b_{3,pid[t]}) + \\ & pFReward1_t^2(\beta_4 + b_{4,pid[t]}) + \text{GroupID}_t \times pFReward1_t^2(\beta_5 + b_{5,pid[t]}) + \\ & dsBwd\beta_6 + RT\beta_7) \end{aligned} \quad (11)$$

Here, and in all regressions described in this section, β stands for fixed effects and b stands for random effects of CA participant ID. The predictor *pFReward1* indicates the true underlying difficulty of the task and is the probability that the F key was rewarding for cue 1. A second-order polynomial was included for this predictor as extreme values indicate portions of the task that are easier and middling values indicate portions of the task that were more difficult. Interaction effects

Running title: REWARD LEARNING IN THE CEREBELLUM

of this predictor and group were included to capture differences in sensitivity to the underlying task difficulty between the groups. Lastly, the reaction time, RT , on each decision was included as a standardized covariate in this analysis to account for any differences that may be due to slowed responding by individuals with CA on this task.

For both the incremental learning and multiple learning strategies tasks, we assessed whether there were differences between the groups on Q learning model performance compared to the baseline biased responder model with the following linear regression:

$$ELPDDifference = \beta_0 + b_{0,pid[t]} + GroupID_t(\beta_1 + b_{1,pid[t]}) + dsBwd\beta_2 \quad (12)$$

where *ELPDDifference* was the difference in model performance (Q Learning model ELPD - Biased Responder ELPD; see above) for each subject.

For the multiple learning strategies task, we assessed behavioral incremental learning performance using the following logistic regression:

$$p(ChooseLucky) = \sigma(\beta_0 + b_{0,pid[t]} + T_{-3:3} \times GroupID_t(\beta_{1:7} + b_{1:7,pid[t]}) + dsBwd\beta_8) \quad (13)$$

In this regression, we grouped trials according to their distance from a reversal, up to three trials prior to ($t = -3: -1$), during ($t = 0$), and after ($t = 1: 3$) a reversal occurred. We then dummy coded them to measure their effects on the degree to which the lucky deck was chosen and interacted each dummy coded regressor with group to measure how this was affected by group membership.

We then assessed the degree to which each group used either incrementally learned deck value, the value of previously seen objects, or a bias toward previously seen objects regardless of their value as estimated by the Hybrid Q learning model using a simple linear regression of the following form for each of these inverse temperature parameters and groups:

$$InvTemp_s = \beta_0 + dsBwd\beta_1 \quad (14)$$

Here we interested primarily in the intercept, β_0 , as this determined the degree to which each group's inverse temperatures were above zero.

Running title: REWARD LEARNING IN THE CEREBELLUM

We also assessed the impact of group on subsequent memory performance following the multiple learning strategies task using the following linear regression:

$$Dprime = \beta_0 + b_{0,pid[t]} + GroupID_t(\beta_1 + b_{1,pid[t]}) + dsBwd\beta_2 \quad (15)$$

where *Dprime* is the signal detection measure d' , which is the difference in z scored hit rate and false alarm rate for each participant.

We were also interested in determining whether there were any differences in reaction times between individuals with CA and matched controls due to motor impairment. For both the incremental learning and multiple learning strategies tasks we did this by assessing whether there were any differences in reaction time between groups:

$$RT = \beta_0 + b_{0,pid[t]} + GroupID_t(\beta_1 + b_{1,pid[t]}) \quad (16)$$

where *RT* was the median reaction time across trials in either task. A separate regression of this form was used for each of the two tasks.

Lastly, we assessed whether there were differences on each neuropsychological measure, *Measure*, using a linear regression for each measure that adhered to the following form:

$$Measure = \beta_0 + b_{0,pid[t]} + GroupID_t(\beta_1 + b_{1,pid[t]}) \quad (17)$$

For all regression analyses, fixed effects are reported in the text as the mean of each parameter's marginal posterior distribution alongside 95% credible intervals, which indicate where 95% of the posterior density falls. Parameter values outside of this range are unlikely given the model, data, and priors. Thus, if the range of likely values does not include zero, we conclude that a meaningful effect was observed.

Results

Impaired reward learning in the incremental learning task

Our first goal was to assess CA participants' baseline ability to learn incrementally from reward using the incremental learning task. On this task, CA participants made overall fewer correct choices compared to healthy controls ($\beta_{Group} = -0.88$, 95% *CI* = $[-1.55, -0.144]$; **Figure**

2A). CA participants' choices were less correct throughout the entirety of the task, even during periods of learning where action-outcome contingencies were more deterministic (e.g. close to 100%) compared to more difficult periods of learning ($\beta_{Group \times pFReward1^2} = -5.49$, 95% $CI = [-7.57, -3.52]$; **Figure 2B-C**. Overall, this difference in performance indicates that CA participants did not learn from reward feedback. Although CA participants responded slightly more slowly than healthy controls on this task ($\beta_{Group} = -115.81$, 95% $CI = [-201.26, -33.55]$; **Supplementary Figure 1**), we included reaction times as a covariate in the above regression analysis to ensure that differences in choice accuracy were not attributed to motor slowing in CA participants.

Next, to more formally assess participants' performance on this task, we fit a standard Q learning model to participants' responses. This model captures the extent to which each participant incorporated trial-by-trial outcomes into running estimates of the value of pressing each button in response to each cue, as well as whether choices are based on these estimates. As a baseline, we compared the performance of this model to a biased responder that merely estimated the extent to which each participant pressed one button over the other, regardless of outcome, in response to each cue. While healthy controls' responses were well described by the Q learning model, this model did no better than the biased responder at predicting CA participants' decisions, thus demonstrating that CA participants engaged in little-to-no incremental learning (**Figure 2D**). On a measure of estimated out-of-sample predictive performance, controls were substantially better fit by the Q learner compared to the biased responder, while this improvement in fit was largely absent for CA participants ($\beta_{Group} = 30.94$, 95% $CI = [16.465, 46.0]$). Thus, while healthy controls incorporated feedback into their estimates about the relationship between cue and action at each timepoint, CA participants generally did not.

Running title: REWARD LEARNING IN THE CEREBELLUM

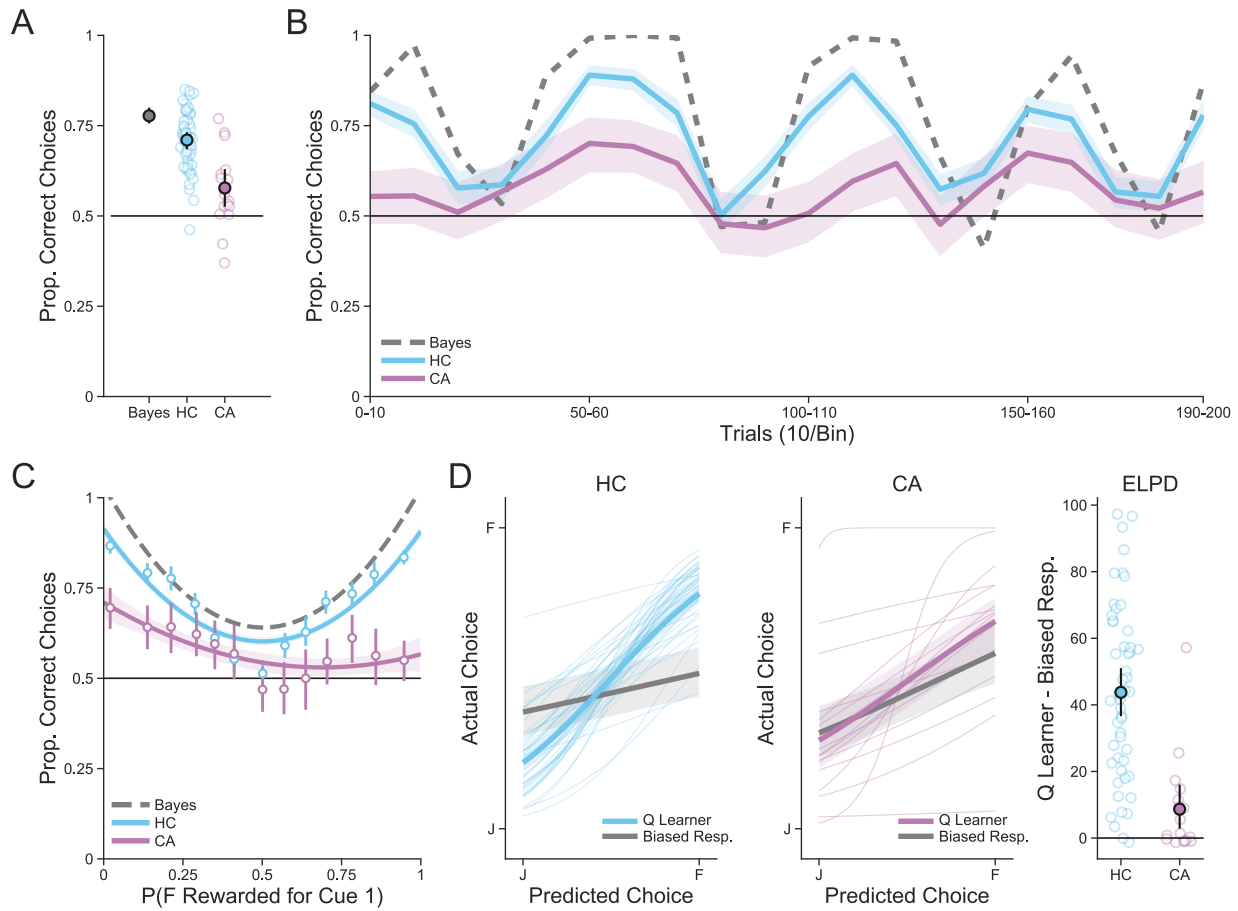


Figure 2 Performance on the incremental learning task. (A) Performance on the incremental learning task averaged across all trials for healthy controls (HC) and CA participants compared to a Bayesian observer in gray, which represents normative performance on the task. Individual points are averages for each subject and filled in points represent group-level averages. Error bars are 95% confidence intervals. (B) Performance on the incremental learning task over time. Each timepoint represents ten trials. Lines are group averages and bands are 95% confidence intervals. For normative comparison, the performance of the Bayesian observer is shown as a dotted gray line. (C) Performance on the incremental learning task as a function of task difficulty, which is indexed by the true underlying probability that pressing the F key was the correct response ($>50\%$) on each trial. Points represent group level averages from 13 bins with an equal number of trials, lines represent the fit of a second-order linear model, and error bars and bands represent 95% confidence intervals. (D) Model performance of the Q Learner and baseline Biased Responder models. Left: Posterior predictive performance. Individual lines represent Q learner fits for each individual, whereas thick lines represent the group-level average fit (with the Q Learner in color and Biased Responder in gray). Bands represent 95% confidence intervals. Right: The difference in estimated out-of-sample predictive performance (as measured by expected log pointwise predictive density; ELPD) between the Q Learner and the Biased Responder model for each group. Individual points are the ELPD difference for each subject and filled in points represent group-level averages. Error bars are 95% confidence intervals.

Together, these results indicate that individuals with CA are impaired at reward learning from trial-and-error.

Impaired incremental reward learning but intact episodic memory in the multiple learning strategies task

After establishing that CA participants were impaired in a task that measured solely incremental reward learning, we wanted to examine both the specificity and generalizability of this impairment by i) providing an alternative means of online reward-based decision making alongside incremental learning and ii) altering the incremental learning task structure to measure responses to reversal events rather than drifting probabilities. The multiple learning strategies task was thus used to accomplish both of these goals.

Consistent with the results of the incremental learning task, CA participants in the multiple learning strategies task were less responsive to reward outcomes compared to controls (**Figure 3A**). Specifically, controls tended to choose the lucky deck more than CA participants immediately prior to a reversal ($\beta_{Group \times t=t-1} = 0.397$, 95% $CI = [0.002, 0.807]$), and this tendency was disrupted by reversals; CA participants did not show this pattern ($\beta_{Group \times t=0} = -0.897$, 95% $CI = [-1.28, -0.535]$), and remained below chance performance after a reversal occurred ($\beta_{Group \times t=t+1} = -0.595$, 95% $CI = [-0.984, -0.21]$). This indicates that CA participants were unable to learn which deck had the higher expected value at any given time throughout the task.

We next assessed the extent to which both incrementally constructed value and episodic value contributed to choice in a combined Hybrid choice model. This model combined a standard Q learning model with three inverse temperature parameters that captured each participants' sensitivity to estimated deck value, the true value of previously seen objects, and a bias toward choosing previously seen objects regardless of their value (**Figure 3B-C**). This model of hybrid choice outperformed a biased responder, which again served as a baseline, for both CA participants and controls as there was no difference between groups in estimated out-of-sample predictive performance ($\beta_{Group} = -1.631$, 95% $CI = [-11.425, 8.444]$; **Supplementary Figure 2**).

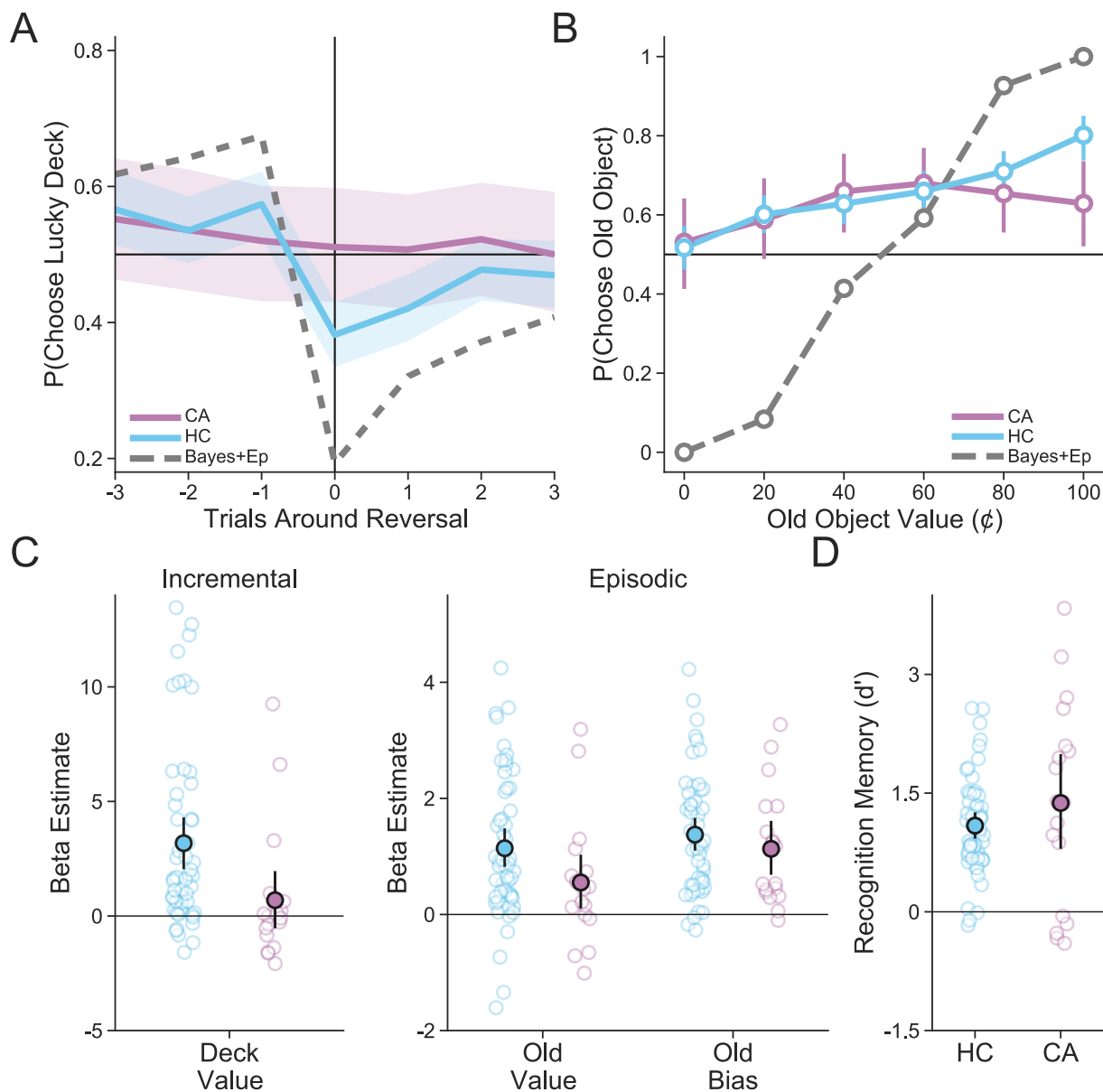


Figure 3 Performance on the multiple learning strategies task. (A) Deck learning performance on the multiple learning strategies task as indicated by the proportion of trials on which the currently lucky deck was chosen as a function of how distant those trials were from a reversal in deck value. Performance for both healthy controls (HC) and CA participants is shown alongside a Bayesian observer with perfect episodic memory for visual comparison. Lines represent group averages and bands represent 95% confidence intervals. (B) Object value usage on trials in which a previously seen object appeared. Points represent group averages and error bars represent 95% confidence intervals. (C) Inverse temperature estimates from the Hybrid model. Individual points represent estimates for each subject, group-level averages are shown as filled in points and error bars represent 95% confidence intervals. (D) Recognition memory performance on the subsequent memory task. Individual points represent each participant's d' prime score, filled in points are group-level averages and error bars are 95% confidence intervals.

Importantly, this indicates that the behavior of both CA participants and controls was well described by the hybrid choice model, which is expected if CA participants are unimpaired at episodic value learning. For each group, we then assessed whether sensitivity differed from zero and, if so, concluded that participants in that group made choices that were affected by each possible predictor. While healthy controls incorporated deck value into their decisions ($\beta_{HC} = 3.173$, 95% $CI = [2.181, 4.189]$), CA participants generally did not ($\beta_{CA} = 0.681$, 95% $CI = [-0.668, 2.066]$). This reward learning deficit was specific to value acquired incrementally, however, because CA participants and controls were both sensitive to episodic value ($\beta_{HC} = 1.373$, 95% $CI = [1.095, 1.654]$; $\beta_{CA} = 1.13$, 95% $CI = [0.59, 1.643]$) and were both similarly biased by previously seen objects regardless of their value ($\beta_{HC} = 1.142$, 95% $CI = [0.798, 1.477]$; $\beta_{CA} = 0.551$, 95% $CI = [0.028, 1.056]$).

We additionally had each participant complete a subsequent memory test for a subset of objects shown during the multiple learning strategies task. There was no difference in recognition memory performance between groups ($\beta_{Group} = -0.487$, 95% $CI = [-1.144, 0.157]$). This result provides further evidence that CA participants were unimpaired at using episodic memory throughout the task relative to their stark impairments in incremental learning. Lastly, CA participants and healthy controls demonstrated no differences in reaction time on this task, suggesting that the behavioral differences reported here cannot be attributed to motor slowing in CA ($\beta_{Group} = -86.40$, 95% $CI = [-222.28, 44.76]$; **Supplementary Figure 1**).

Controlling for effects of non-motor deficits and disease subtype

We next sought to ensure that the differences in the tasks reported here were specific to deficits in reward learning rather than general cognitive impairment. Controlling for cognitive impairment is particularly important because recent work^{43,44} has suggested that incremental learning experiments tax higher level functions, like executive control and working memory, in addition to learning from reward prediction error. To address this issue and assess possible cognitive impairment, we conducted a battery of neuropsychological measures on CA participants (see **Methods**). Of these, a subset of measures were also completed by healthy controls (**Supplementary Table 1; Supplementary Figure 3**). We found no differences in performance

between groups on all measures except for the backwards digit span task, which indexes working memory ability, and on which healthy controls scored higher than CA participants ($\beta_{Group} = -2.57$, 95% $CI = [-4.17, -0.92]$). Backwards digit span scores were thus included as covariates in all regression analyses (see **Methods**) in order to control for impacts of this performance difference on impairments in incremental learning, and all effects described in the above sections thus control for differences in working memory performance. To further ensure that CA participants' deficient incremental learning was not due to broad cognitive impairment, we also repeated all analyses excluding seven CA participants (and their matched controls) with mild cognitive impairment (MCI), as indicated by scoring lower than 26 on the MoCA (**Table 1**). While CA participants with MCI consisted of some of the lowest performing participants in our sample (**Supplementary Figures 4 and 5**), we found no differences in the results across both tasks when they were excluded. It is therefore unlikely that CA participants' impaired reward learning ability is due to either working memory deficits or cognitive decline more broadly.

We next sought to further characterize the nature of CA participants' reward learning impairment by looking at the relationship between incremental learning sensitivity, as measured by the Q learning models in each task, and performance on our neuropsychological battery. The extent to which CA participants learned about cues in the incremental learning task related only to total CCAS score ($r = 0.84$, $p < 0.001$, *Bonferroni corrected*; **Supplementary Table 2**), suggesting that the specific contributions of the cerebellum to cognition may impact performance in this task. The CCAS scale was recently developed to measure the exact types of cognitive impairment that result from damage to the cerebellum.⁴⁵ Because more focal cerebellar lesions tend to lead to lower total CCAS scores,⁴⁶ this provides further evidence of the necessity for the cerebellum to successfully perform the incremental learning task. The relationship between total CCAS score and performance was driven by the timed portions of the CCAS scale (e.g. the Semantic Fluency and Category Switching measures; **Supplementary Table 3**), suggesting a potential effect of slowed responses in the incremental learning task. While CA participants did indeed respond more slowly than healthy controls on this task (see above), we controlled for this difference in our behavioral analysis. There was also no relationship between any of the timed CCAS measures and reaction time on the incremental learning task (**Supplementary Figure 1**), further suggesting that these measures captured independent timing-related impairment. Finally,

there was no relationship between any measure and incremental learning ability in the multiple learning strategies task (**Supplementary Table 2**).

Finally, we addressed the possibility that the subset of our sample of CA participants consisting of diagnoses that were less restricted to the cerebellum, namely the three individuals with multiple system atrophy (MSA), could be responsible for the deficits reported here. We repeated all analyses with these three participants excluded and found no differences in the results (**Supplementary Figures 6 and 7**).

Discussion

The results of the present work demonstrate that individuals with cerebellar dysfunction, represented by CA cases in our cohort, are impaired at reward learning. While the cerebellum and basal ganglia have traditionally been treated as making separate contributions to learning,^{5,21} recent findings have called this dichotomy into question.^{8–19} This work has suggested that, alongside its role in motor learning, the cerebellum likely operates in concert with the basal ganglia to support reinforcement learning from reward. Our study corroborates these findings from animal models,^{10–19} providing evidence that the human cerebellum is necessary for learning associations from reward. In comparison to age- and sex-matched healthy controls, CA participants were impaired at reward-based learning from trial-and-error. Further, CA participants retained the ability to employ an alternative strategy based in episodic memory to guide their decisions, demonstrating that this impairment is specific to incremental learning. These results challenge the idea that the cerebellum is used primarily for motor learning and shed light on how multiple neural systems may interact with one another to support learning in the non-motor domain.

Our findings join a litany of recent research suggesting that the cerebellum plays a broad role in human cognition.^{22,47–49} Indeed, individuals with damage to the cerebellum demonstrate impairment in a wide range of cognitive functions including cognitive control⁵⁰ and impulsivity.⁵¹ Human functional neuroimaging studies have also revealed cerebellar activity in a variety of different non-motor tasks.^{22,48} Many of these functions are likely supported by the robust bidirectional connections the cerebellum shares with the prefrontal cortex.^{52,53} In particular, recent findings have indicated that individuals with CA have heightened domain-specific impulsive and

Running title: REWARD LEARNING IN THE CEREBELLUM

compulsive behaviors, which is a common symptom of underlying reward system dysfunction.^{54,55} Our study adds to this work by suggesting that the cerebellum is additionally necessary for reward learning in humans.

While there is growing evidence validating the implication of the cerebellum in reward-based learning in animals, there is only limited work on this topic in humans. Early imaging studies, for example, demonstrated cerebellar BOLD activity in patients with substance use disorder who performed reward-based learning tasks²³ and experienced cravings,²⁴ and also in response to unexpected reward.²⁵ However, it remains unknown how cerebellar damage impacts reward learning, as investigations of reward learning in the cerebellum are rare. While two previous studies employed reward-based experimental tasks in individuals with isolated ischemic lesions of the cerebellum,^{56,57} results until this point have remained far from conclusive. Thoma *et al.*⁵⁶ used a reward-based learning task consisting of an initial acquisition phase in which eight participants with cerebellar damage were rewarded for learning associations between colors and symbols followed by a reversal portion in which they had to disremember previously acquired knowledge and learn new associations for each cue. While participants with cerebellar damage demonstrated no impairment at acquiring new, reward-based knowledge, they were selectively impaired at learning from a single reversal. While this study complements our findings, we found evidence for more global impairment: CA participants in both of our tasks were unable to learn associations from reward on a trial-by-trial basis. Rustemeier *et al.*⁵⁷ took a different approach by asking twelve individuals with cerebellar damage to learn a simple acquisition task from probabilistic feedback and subsequently transfer this knowledge to re-arranged stimuli. While participants were unimpaired behaviorally at this task, electroencephalographic (EEG) results revealed that they may process reward-based feedback differently from controls. Our findings support this interpretation and further suggest that processing of trial-by-trial feedback is not just different, but impaired, in individuals with cerebellar damage. Finally, while another related study showed impairment in learning from reinforcement in twelve participants with cerebellar damage,²⁰ this study was focused specifically on the motor domain.

While our findings suggest that the cerebellum is necessary for incremental reward learning, they cannot speak to the neural circuitry underlying this role. One intriguing possibility is that the

cerebellum may operate in tandem with the basal ganglia—canonically seen as the seat of reinforcement learning in the brain^{5,21}—to learn about reward incrementally. Reward prediction error signals in midbrain dopamine neurons that provide input to the basal ganglia^{27,29} have also been found to be encoded by cerebellar neurons.^{9,15,17,19} Further, through excitatory projections to the ventral tegmental area, the cerebellum has widespread reciprocal connections with the basal ganglia and has recently been shown to influence reward-driven behavior through these projections.^{8,58} While reinforcement learning via the basal ganglia and supervised learning via the cerebellum have typically been treated as fulfilling entirely separate roles,^{5,21} these systems appear to be more interdependent than previously thought. Future investigations of the relationship between the basal ganglia and cerebellum are needed to clarify the exact mechanisms underlying reinforcement learning in the brain.

There are several potential limitations to the findings of our study. Our study participants included a large variety of different conditions causing cerebellar dysfunction, including some with MSA and spinocerebellar ataxia (SCA). While some of these pathologies are predominantly restricted to the cerebellum, non-cerebellar brain areas and circuits, such as the dopaminergic system in MSA, could also be involved. However, there was no change in the reported reward-based learning deficits with subgroup analyses comparing CA participants whose conditions are known to be multisystemic and those whose conditions show more isolated cerebellar pathology. Second, while cognitive impairment due to neurodegenerative disease could potentially contribute to some of the deficits measured here, we accounted for this possibility by establishing that the incremental reward learning deficits reported here persist regardless of MCI status. We also collected basic neuropsychological measures from all participants, and CA participants were not different from controls on the vast majority of measures. We focused particularly on possible contributions of working memory given recent work suggesting that working memory plays an important role in incremental reward learning.^{43,44} While CA participants and controls performed similarly on the forward digit span task, CA participants were somewhat impaired at backwards digit span. We controlled for this difference by including backwards digit span scores as covariates in all regression analyses.

Running title: REWARD LEARNING IN THE CEREBELLUM

Taken together, our findings suggest that the human cerebellum is necessary for reward learning. These results provide new constraints on models of non-motor learning and suggest that the cerebellum and basal ganglia work in tandem to support learning from reinforcement.

Funding

J.N. was supported by the NSF Graduate Research Fellowship (1644869). S.H.K. was supported by NINDS R01NS104423, NINDS R01 NS118179, NINDS R01 NS124854, and National Ataxia Foundation. D.S. was supported by an NSF CRCNS award (1822619), NIMH R01 MH121093 and the Kavli Foundation.

Competing interests

The authors report no competing interests.

References

1. Raymond JL, Lisberger SG, Mauk MD. The Cerebellum: A Neuronal Learning Machine? *Science*. 1996;272(5265):1126-1131. doi:10.1126/science.272.5265.1126
2. Llinás R, Welsh JP. On the cerebellum and motor learning. *Curr Opin Neurobiol*. 1993;3(6):958-965. doi:10.1016/0959-4388(93)90168-X
3. Ito M, Itō M. *The Cerebellum and Neural Control*. Raven Press; 1984.
4. Marr D. A theory of cerebellar cortex. *J Physiol*. 1969;202(2):437-470. doi:10.1113/jphysiol.1969.sp008820
5. Doya K. What are the computations of the cerebellum, the basal ganglia and the cerebral cortex? *Neural Netw*. 1999;12(7):961-974. doi:10.1016/S0893-6080(99)00046-5
6. Wolpert DM, Miall RC, Kawato M. Internal models in the cerebellum. *Trends Cogn Sci*. 1998;2(9):338-347. doi:10.1016/S1364-6613(98)01221-2
7. Raymond JL, Medina JF. Computational Principles of Supervised Learning in the Cerebellum. *Annu Rev Neurosci*. 2018;41:233-253. doi:10.1146/annurev-neuro-080317-061948

Running title: REWARD LEARNING IN THE CEREBELLUM

8. Caligiore D, Arbib MA, Miall RC, Baldassarre G. The super-learning hypothesis: Integrating learning processes across cortex, cerebellum and basal ganglia. *Neurosci Biobehav Rev.* 2019;100:19-34. doi:10.1016/j.neubiorev.2019.02.008
9. Hull C. Prediction signals in the cerebellum: Beyond supervised motor learning. Ivry RB, ed. *eLife.* 2020;9:e54073. doi:10.7554/eLife.54073
10. Sendhilnathan N, Goldberg ME. The Mid-Lateral Cerebellum Is Necessary for Reinforcement Learning. *BioRxiv.* 2020. doi:10.1101/2020.03.20.000190
11. Sendhilnathan N, Semework M, Goldberg ME, Ipata AE. Neural Correlates of Reinforcement Learning in Mid-lateral Cerebellum. *Neuron.* 2020;106(1):188-198.e5. doi:10.1016/j.neuron.2019.12.032
12. Sendhilnathan N, Ipata A, Goldberg ME. Mid-lateral cerebellar complex spikes encode multiple independent reward-related signals during reinforcement learning. *Nat Commun.* 2021;12(1):6475. doi:10.1038/s41467-021-26338-0
13. Larry N, Yarkoni M, Lixenberg A, Joshua M. Cerebellar climbing fibers encode expected reward size. Raymond JL, Calabrese RL, Edgley SA, eds. *eLife.* 2019;8:e46870. doi:10.7554/eLife.46870
14. Carta I, Chen CH, Schott AL, Dorizan S, Khodakhah K. Cerebellar modulation of the reward circuitry and social behavior. *Science.* 2019;363(6424):eaav0581. doi:10.1126/science.aav0581
15. Heffley W, Hull C. Classical conditioning drives learned reward prediction signals in climbing fibers across the lateral cerebellum. Calabrese RL, Raymond JL, Lerner T, Llano I, Khodakhah K, eds. *eLife.* 2019;8:e46764. doi:10.7554/eLife.46764
16. Wagner MJ, Kim TH, Savall J, Schnitzer MJ, Luo L. Cerebellar granule cells encode the expectation of reward. *Nature.* 2017;544(7648):96-100. doi:10.1038/nature21726
17. Heffley W, Song EY, Xu Z, et al. Coordinated cerebellar climbing fiber activity signals learned sensorimotor predictions. *Nat Neurosci.* 2018;21(10):1431-1441. doi:10.1038/s41593-018-0228-8
18. Kostadinov D, Beau M, Blanco-Pozo M, Häusser M. Predictive and reactive reward signals conveyed by climbing fiber inputs to cerebellar Purkinje cells. *Nat Neurosci.* 2019;22(6):950-962. doi:10.1038/s41593-019-0381-8

Running title: REWARD LEARNING IN THE CEREBELLUM

19. Ohmae S, Medina JF. Climbing fibers encode a temporal-difference prediction error during cerebellar learning in mice. *Nat Neurosci.* 2015;18(12):1798-1803. doi:10.1038/nn.4167
20. Therrien AS, Wolpert DM, Bastian AJ. Effective reinforcement learning following cerebellar damage requires a balance between exploration and motor noise. *Brain J Neurol.* 2016;139(Pt 1):101-114. doi:10.1093/brain/awv329
21. Doya K. Complementary roles of basal ganglia and cerebellum in learning and motor control. *Curr Opin Neurobiol.* 2000;10(6):732-739. doi:10.1016/s0959-4388(00)00153-7
22. King M, Hernandez-Castillo CR, Poldrack RA, Ivry RB, Diedrichsen J. Functional boundaries in the human cerebellum revealed by a multi-domain task battery. *Nat Neurosci.* 2019;22(8):1371-1378. doi:10.1038/s41593-019-0436-x
23. Volkow ND, Wang GJ, Ma Y, et al. Expectation enhances the regional brain metabolic and the reinforcing effects of stimulants in cocaine abusers. *J Neurosci Off J Soc Neurosci.* 2003;23(36):11461-11468.
24. Grant S, London ED, Newlin DB, et al. Activation of memory circuits during cue-elicited cocaine craving. *Proc Natl Acad Sci U S A.* 1996;93(21):12040-12045.
25. Ramnani N, Elliott R, Athwal BS, Passingham RE. Prediction error for free monetary reward in the human prefrontal cortex. *NeuroImage.* 2004;23(3):777-786. doi:10.1016/j.neuroimage.2004.07.028
26. Sutton RS, Barto AG. Reinforcement Learning: An Introduction. :352.
27. Houk JC, Adams JL, Barto AG. A model of how the basal ganglia generate and use neural signals that predict reinforcement. In: *Models of Information Processing in the Basal Ganglia.* Computational neuroscience. The MIT Press; 1995:249-270.
28. Rescorla RA, Wagner AR. 3 A Theory of Pavlovian Conditioning : Variations in the Effectiveness of Reinforcement and Nonreinforcement. In: ; 1972.
29. Schultz W, Dayan P, Montague PR. A Neural Substrate of Prediction and Reward. *Science.* 1997;275(5306):1593-1599. doi:10.1126/science.275.5306.1593

Running title: REWARD LEARNING IN THE CEREBELLUM

30. Kuo SH. Ataxia. *Contin Minneap Minn*. 2019;25(4):1036-1054.
doi:10.1212/CON.0000000000000753
31. Nicholas J, Daw ND, Shohamy D. Uncertainty alters the balance between incremental learning and episodic memory. *BioRxiv*. 2022. doi:10.1101/2022.07.05.498877
32. Duncan K, Semmler A, Shohamy D. Modulating the Use of Multiple Memory Systems in Value-based Decisions with Contextual Novelty. *J Cogn Neurosci*. Published online July 19, 2019:1-13.
doi:10.1162/jocn_a_01447
33. Hariri AR. The Emerging Importance of the Cerebellum in Broad Risk for Psychopathology. *Neuron*. 2019;102(1):17-20. doi:10.1016/j.neuron.2019.02.031
34. Bellebaum C, Daum I. Cerebellar involvement in executive control. *The Cerebellum*. 2007;6(3):184-192. doi:10.1080/14734220601169707
35. Beuriat PA, Cohen-Zimmerman S, Smith GNL, Krueger F, Gordon B, Grafman J. A New Insight on the Role of the Cerebellum for Executive Functions and Emotion Processing in Adults. *Front Neurol*. 2020;11. Accessed August 1, 2022. <https://www.frontiersin.org/articles/10.3389/fneur.2020.593490>
36. Mannarelli D, Pauletti C, Currà A, et al. The Cerebellum Modulates Attention Network Functioning: Evidence from a Cerebellar Transcranial Direct Current Stimulation and Attention Network Test Study. *The Cerebellum*. 2019;18(3):457-468. doi:10.1007/s12311-019-01014-8
37. Litman L, Robinson J, Abberbock T. TurkPrime.com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behav Res Methods*. 2017;49(2):433-442. doi:10.3758/s13428-016-0727-z
38. Hoffman MD, Gelman A. The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. :31.
39. Team SD. *Stan Reference Manual*. Accessed October 12, 2021. https://mc-stan.org/docs/2_28/reference-manual/index.html
40. Vehtari A, Gelman A, Gabry J. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Stat Comput*. 2017;27(5):1413-1432. doi:10.1007/s11222-016-9696-4

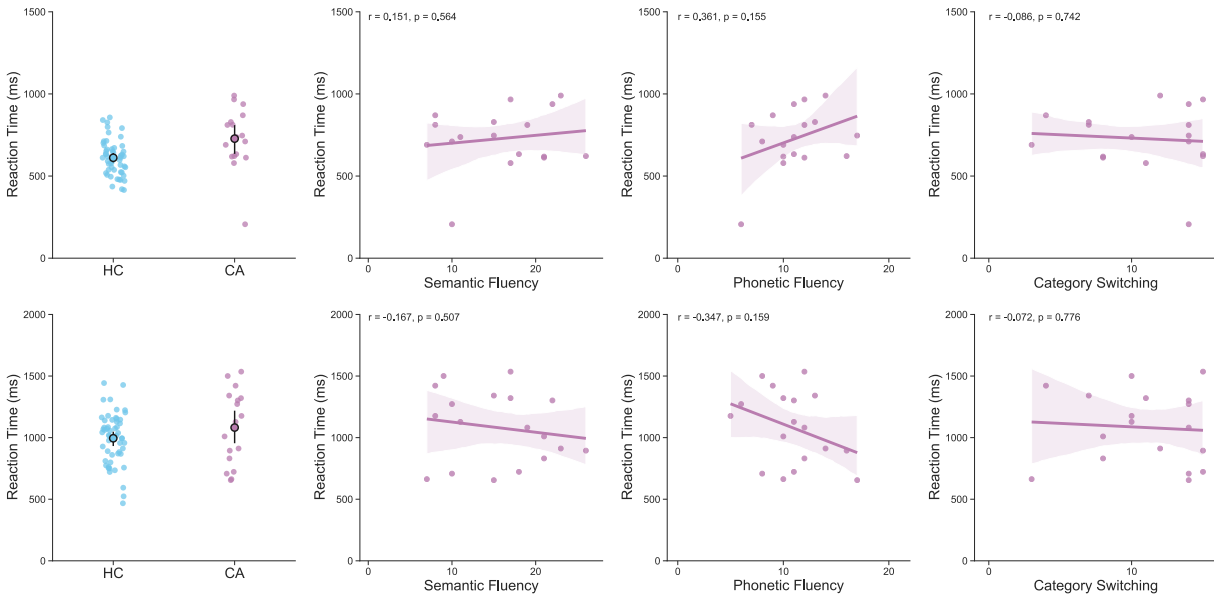
Running title: REWARD LEARNING IN THE CEREBELLUM

41. Kalman RE. A New Approach to Linear Filtering and Prediction Problems. *J Basic Eng.* 1960;82(1):35-45. doi:10.1115/1.3662552
42. Nassar MR, Wilson RC, Heasley B, Gold JJ. An Approximately Bayesian Delta-Rule Model Explains the Dynamics of Belief Updating in a Changing Environment. *J Neurosci.* 2010;30(37):12366-12378. doi:10.1523/JNEUROSCI.0822-10.2010
43. Collins AGE, Frank MJ. How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci.* 2012;35(7):1024-1035. doi:10.1111/j.1460-9568.2011.07980.x
44. Yoo AH, Collins AGE. How Working Memory and Reinforcement Learning Are Intertwined: A Cognitive, Neural, and Computational Perspective. *J Cogn Neurosci.* 2022;34(4):551-568. doi:10.1162/jocn_a_01808
45. Hoche F, Guell X, Vangel MG, Sherman JC, Schmahmann JD. The cerebellar cognitive affective/Schmahmann syndrome scale. *Brain.* 2018;141(1):248-270. doi:10.1093/brain/awx317
46. Chirino-Pérez A, Marrufo-Meléndez OR, Muñoz-López JL, et al. Mapping the Cerebellar Cognitive Affective Syndrome in Patients with Chronic Cerebellar Strokes. *The Cerebellum.* 2022;21(2):208-218. doi:10.1007/s12311-021-01290-3
47. McDougale SD, Tsay JS, Pitt B, et al. Continuous manipulation of mental representations is compromised in cerebellar degeneration. *Brain J Neurol.* Published online February 24, 2022:awac072. doi:10.1093/brain/awac072
48. Buckner RL. The cerebellum and cognitive function: 25 years of insight from anatomy and neuroimaging. *Neuron.* 2013;80(3):807-815. doi:10.1016/j.neuron.2013.10.044
49. Koziol LF, Budding D, Andreasen N, et al. Consensus Paper: The Cerebellum's Role in Movement and Cognition. *Cerebellum Lond Engl.* 2014;13(1):151-177. doi:10.1007/s12311-013-0511-x
50. Alexander MP, Gillingham S, Schweizer T, Stuss DT. Cognitive impairments due to focal cerebellar injuries in adults. *Cortex J Devoted Study Nerv Syst Behav.* 2012;48(8):980-990. doi:10.1016/j.cortex.2011.03.012

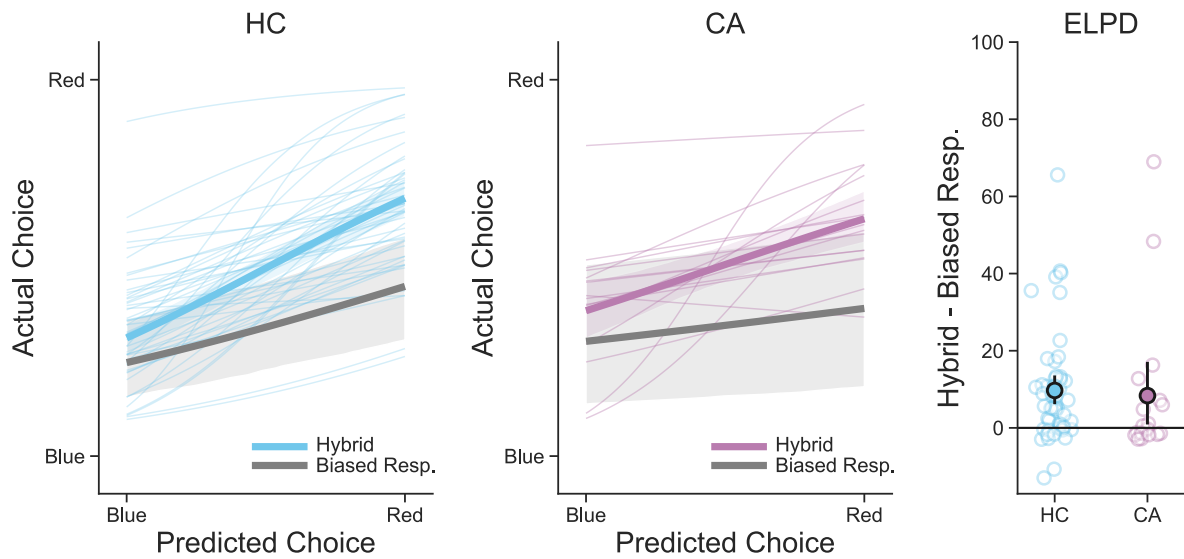
Running title: REWARD LEARNING IN THE CEREBELLUM

51. Amokrane N, Lin CYR, Desai NA, Kuo SH. The Impact of Compulsivity and Impulsivity in Cerebellar Ataxia: A Case Series. *Tremor Hyperkinetic Mov.* 10:43. doi:10.5334/tohm.550
52. Buckner RL, Krienen FM, Castellanos A, Diaz JC, Yeo BTT. The organization of the human cerebellum estimated by intrinsic functional connectivity. *J Neurophysiol.* 2011;106(5):2322-2345. doi:10.1152/jn.00339.2011
53. Middleton FA, Strick PL. Cerebellar Projections to the Prefrontal Cortex of the Primate. *J Neurosci.* 2001;21(2):700-712. doi:10.1523/JNEUROSCI.21-02-00700.2001
54. Amokrane N, Viswanathan A, Freedman S, et al. Impulsivity in Cerebellar Ataxias: Testing the Cerebellar Reward Hypothesis in Humans. *Mov Disord.* 2020;35(8):1491-1493. doi:10.1002/mds.28121
55. Chen TX, Lin CYR, Aumann MA, et al. Impulsivity Trait Profiles in Patients With Cerebellar Ataxia and Parkinson Disease. *Neurology.* 2022;99(2):e176-e186. doi:10.1212/WNL.0000000000200349
56. Thoma P, Bellebaum C, Koch B, Schwarz M, Daum I. The Cerebellum Is Involved in Reward-based Reversal Learning. *The Cerebellum.* 2008;7(3):433. doi:10.1007/s12311-008-0046-8
57. Rustemeier M, Koch B, Schwarz M, Bellebaum C. Processing of Positive and Negative Feedback in Patients with Cerebellar Lesions. *Cerebellum Lond Engl.* 2016;15(4):425-438. doi:10.1007/s12311-015-0702-8
58. Caligiore D, Pezzulo G, Baldassarre G, et al. Consensus Paper: Towards a Systems-Level View of Cerebellar Function: the Interplay Between Cerebellum, Basal Ganglia, and Cortex. *The Cerebellum.* 2017;16(1):203-229. doi:10.1007/s12311-016-0763-3

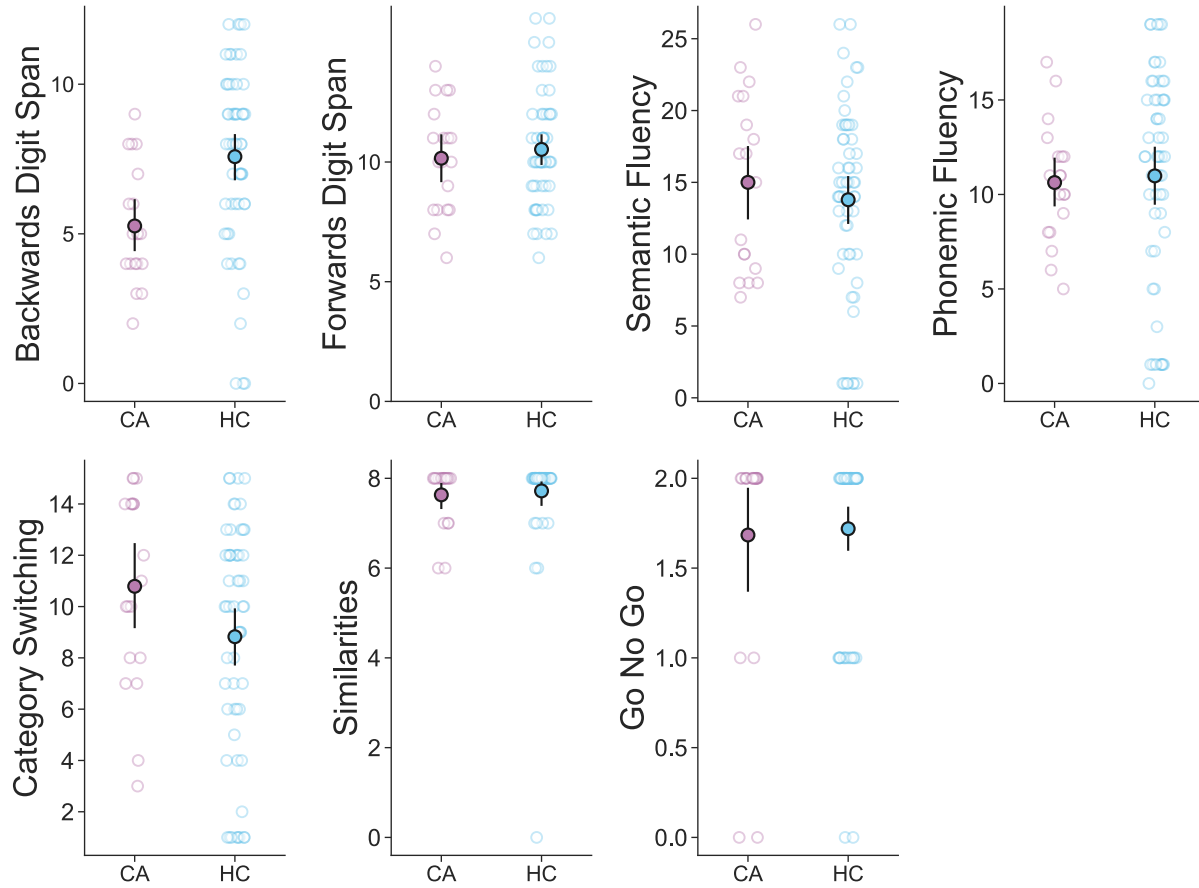
Supplementary Material



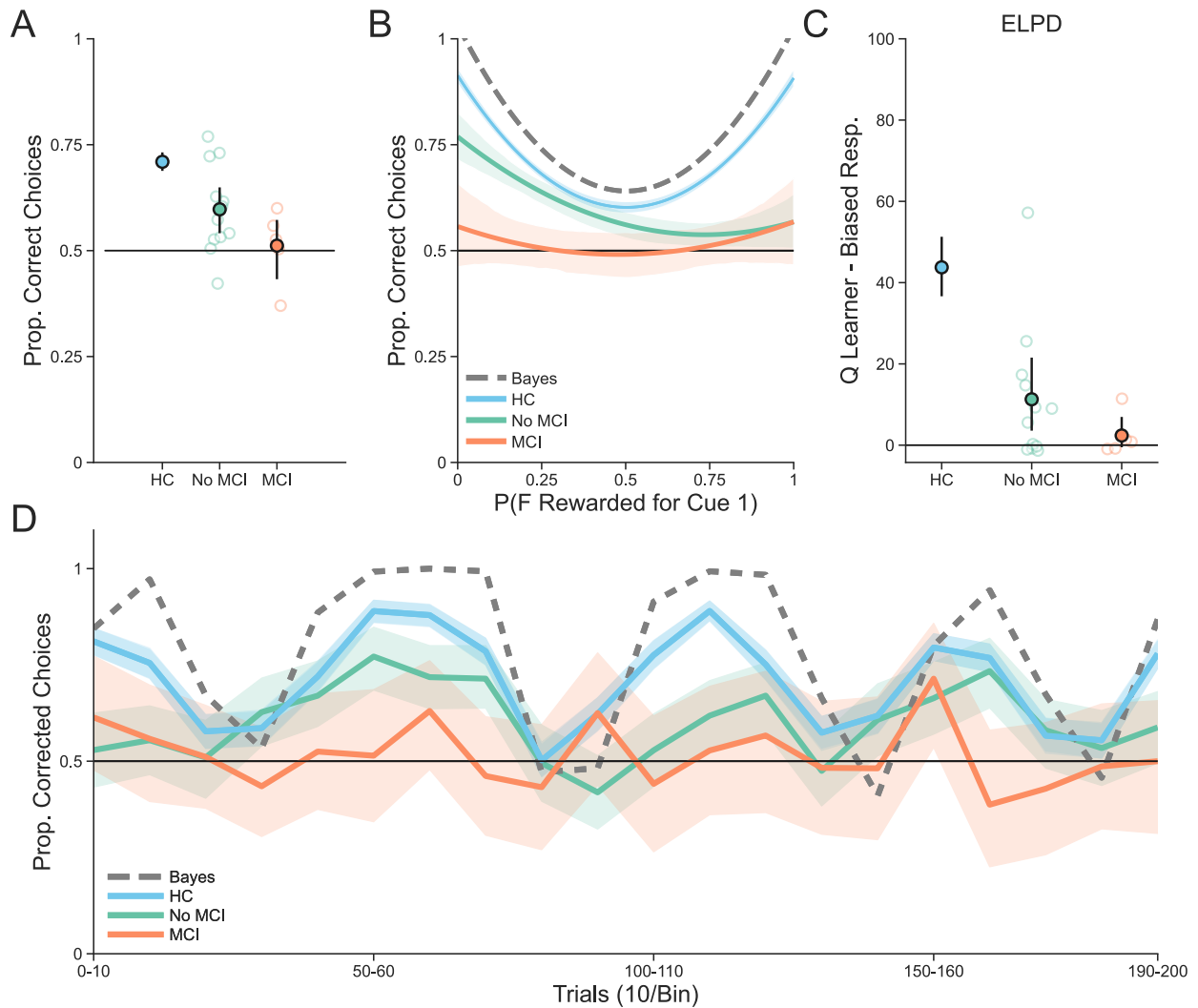
Supplementary Figure 1 Reaction time and relationship to timing-dependent CCAS measures. Reaction time in the incremental learning (**Top row**) and multiple learning strategies (**Bottom row**) tasks. Lefthand panels show median RTs for CA participants compared to healthy controls (HC). Righthand panels show the relationship between CCAS measures which yielded a significant relationship with incremental value sensitivity in the incremental learning task and were conducted under time pressure, and reaction time, for CA participants.



Supplementary Figure 2 Model performance on the multiple learning strategies task for the Hybrid Learner and Biased Responder models. Left: Posterior predictive performance. Individual lines represent Hybrid Learner fits for each individual, whereas thick lines represent the group-level average fit (with the Hybrid Learner in color and Biased Responder in gray). Bands represent 95% confidence intervals. **Right:** The difference in estimated out-of-sample predictive performance (as measured by expected log pointwise predictive density) between the Hybrid Learner and the Biased Responder model for each group. Individual points are the ELPD difference for each subject and filled in points represent group-level averages. Error bars are 95% confidence intervals.



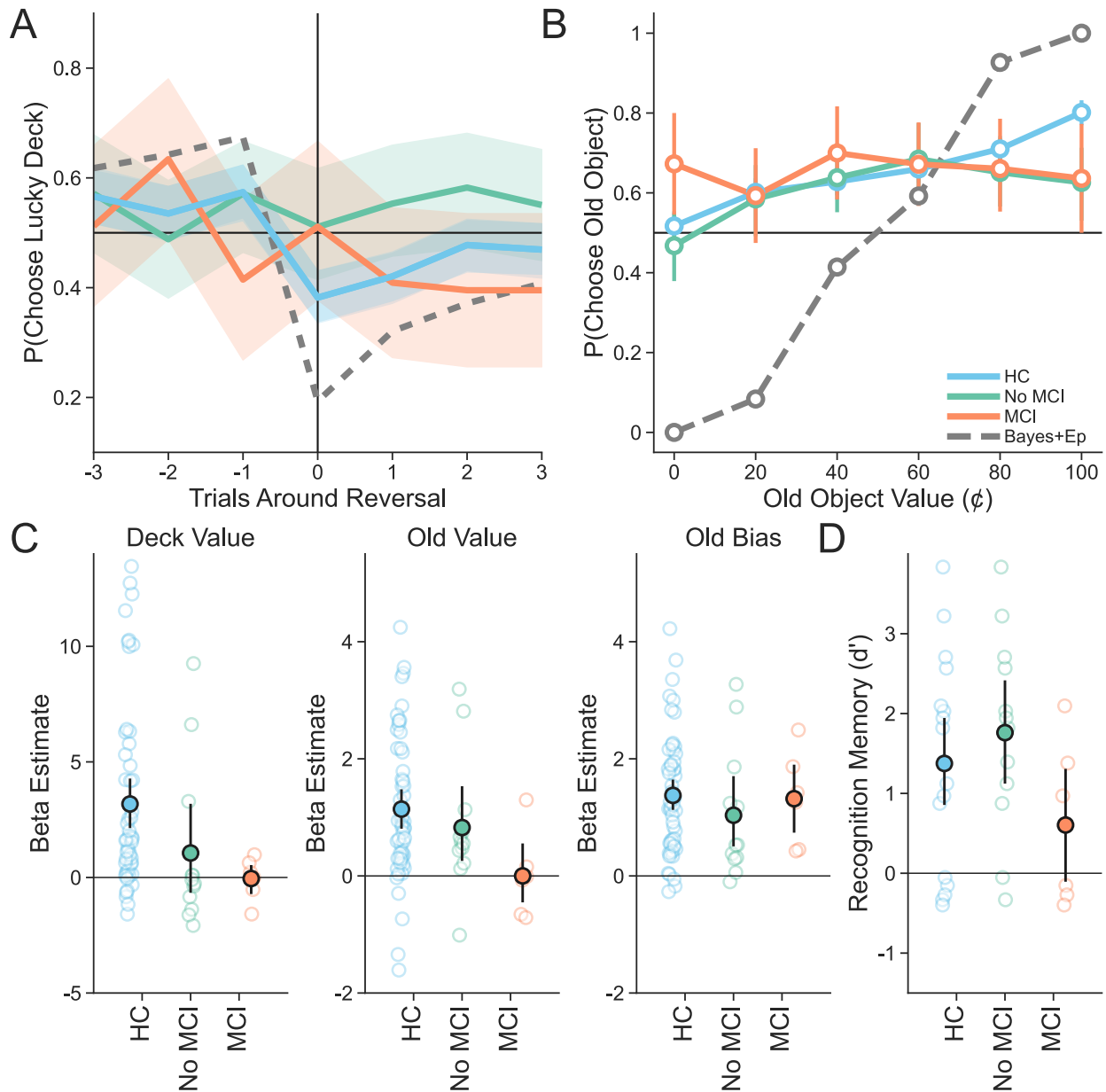
Supplementary Figure 3 Neuropsychological test performance for the subset of neuropsychological measures that were completed by both CA participants and healthy controls (HC). Scores from individual participants are plotted as empty circles behind group-level means plotted as filled circles with uncertainty represented by 95% confidence intervals.



Supplementary Figure 4 CA participant performance on the incremental learning task separated by mild cognitive impairment (MCI versus all others) and compared to healthy controls (HC). (A) Performance on the incremental learning task averaged across all trials. Individual points are averages for each subject and filled in points represent group-level averages. Error bars are 95% confidence intervals. **(B)** Performance on the incremental learning task as a function of task difficulty, which is indexed by the true underlying probability that pressing the F key was the correct response ($>50\%$) on each trial. Points represent group level averages from 13 bins with an equal number of trials, lines represent the fit of a second-order linear model, and error bars and bands represent 95% confidence intervals. **(C)** The difference in estimated out-of-sample predictive performance (as measured by expected log pointwise predictive density; ELPD) between the Q Learner and the Biased Responder model for each group. Individual points are the ELPD difference for each subject and filled in points represent group-level averages. Error bars are

Running title: REWARD LEARNING IN THE CEREBELLUM

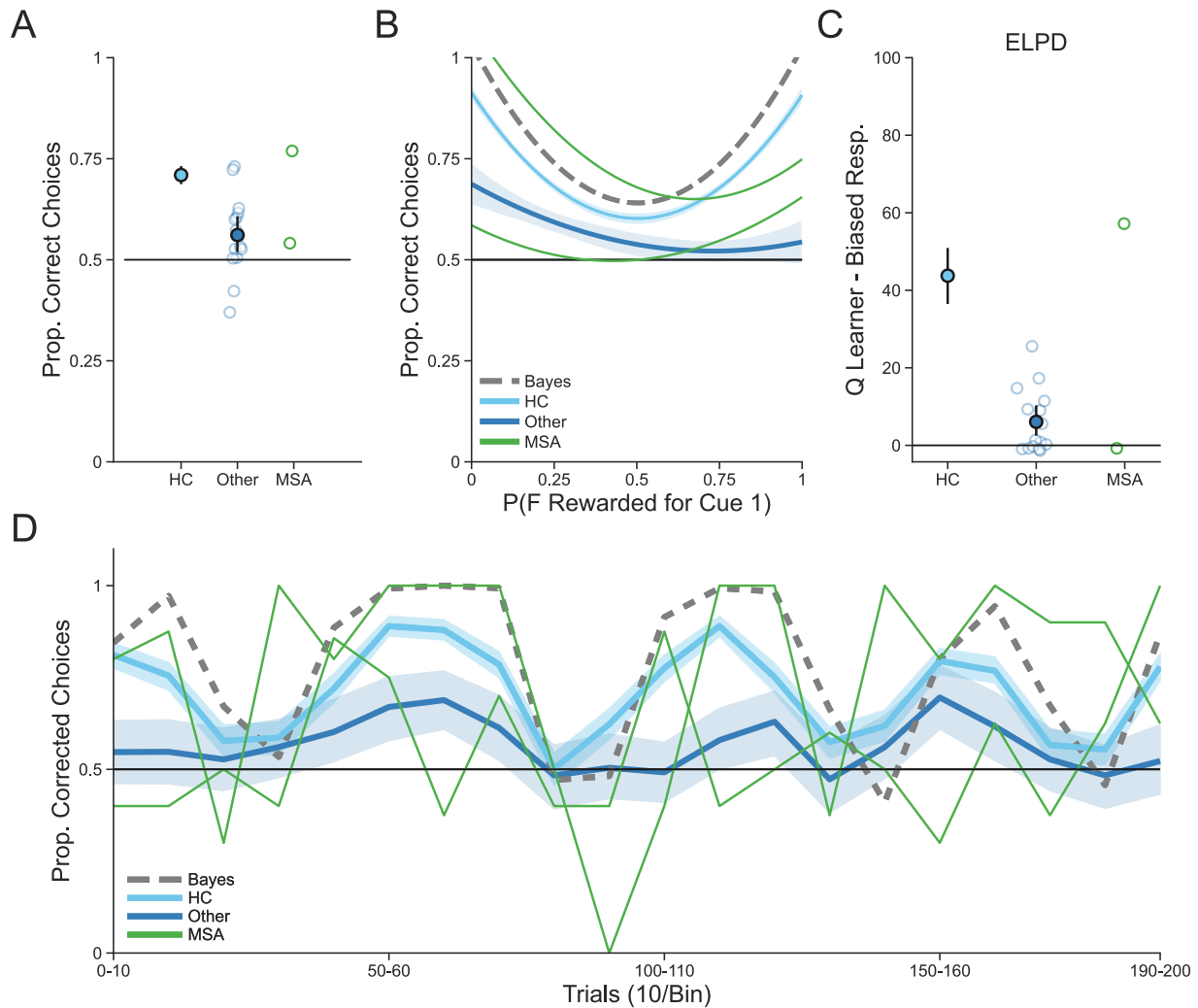
95% confidence intervals. **(D)** Performance on the incremental learning task over time. Each timepoint represents ten trials. Lines are group averages and bands are 95% confidence intervals. For normative comparison, the performance of the Bayesian observer is shown as a dotted gray line.



Supplementary Figure 5 CA participant performance on the multiple learning strategies task separated by mild cognitive impairment (MCI versus all others) and compared to healthy controls (HC). (A) Deck learning performance on the multiple learning strategies task as indicated by the proportion of trials on which the currently lucky deck was chosen as a function of how distant those trials were from a reversal in deck value. Lines represent group averages and bands represent 95% confidence intervals. **(B)** Object value usage on trials in which a previously seen object appeared. Points represent group averages and error bars represent 95% confidence intervals. **(C)** Inverse temperature estimates from the Hybrid model. Individual points represent estimates for each subject, group-level averages are shown as filled in

Running title: REWARD LEARNING IN THE CEREBELLUM

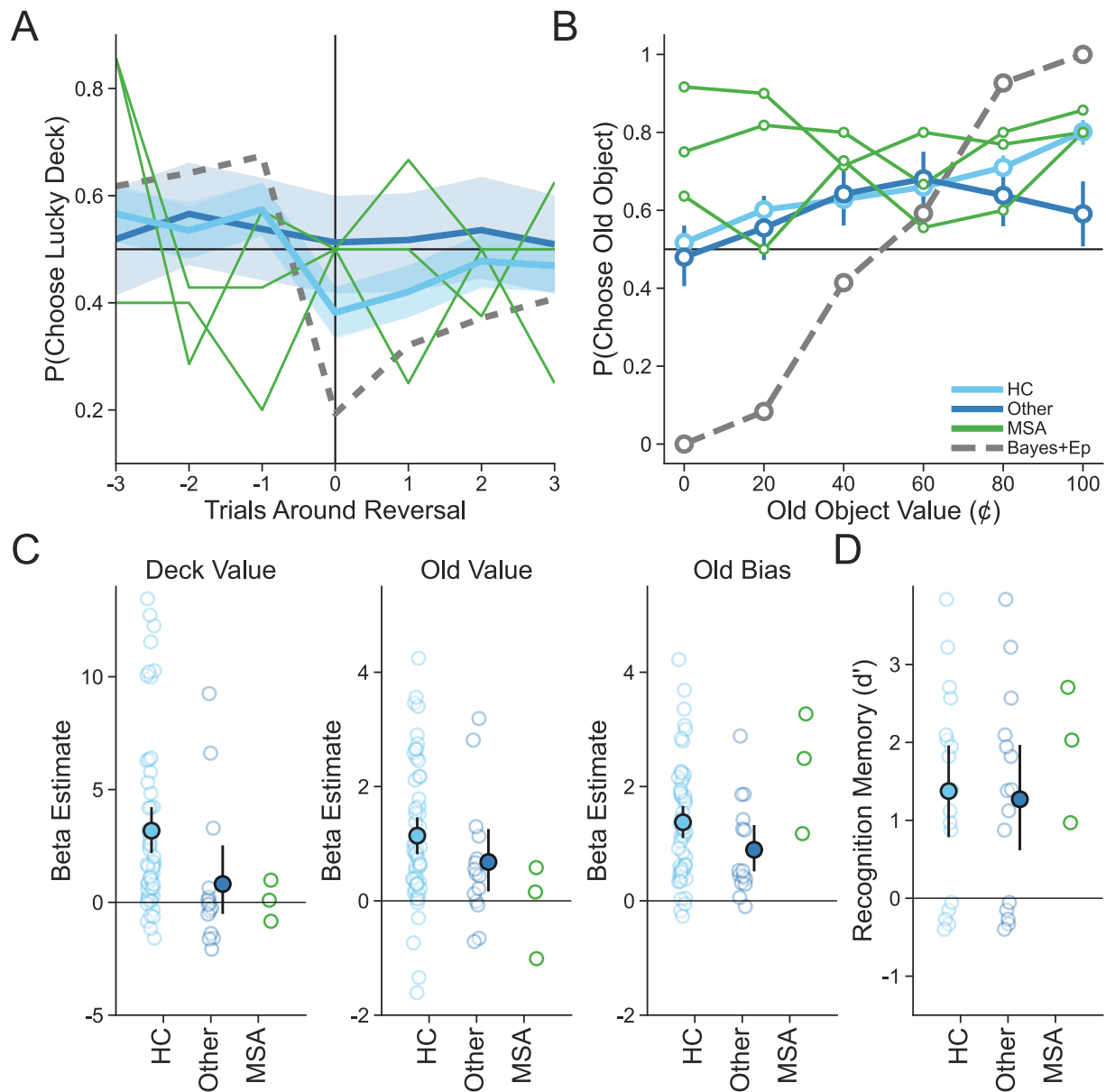
points and error bars represent 95% confidence intervals. **(D)** Recognition memory performance on the subsequent memory task. Individual points represent each participant's d' prime score, filled in points are group-level averages and error bars are 95% confidence intervals.



Supplementary Figure 6 CA participant performance on the incremental learning task separated by diagnosis (MSA individuals versus all others) and compared to healthy controls (HC). (A) Performance on the incremental learning task averaged across all trials. Individual points are averages for each subject and filled in points represent group-level averages. Error bars are 95% confidence intervals. **(B)** Performance on the incremental learning task as a function of task difficulty, which is indexed by the true underlying probability that pressing the F key was the correct response ($>50\%$) on each trial. Points represent group level averages from 13 bins with an equal number of trials, lines represent the fit of a second-order linear model, and error bars and bands represent 95% confidence intervals. **(C)** The difference in estimated out-of-sample predictive performance (as measured by expected log pointwise predictive density; ELPD) between the Q Learner and the Biased Responder model for each group. Individual points are the ELPD difference for each subject and filled in points represent group-level averages. Error bars are

Running title: REWARD LEARNING IN THE CEREBELLUM

95% confidence intervals. **(D)** Performance on the incremental learning task over time. Each timepoint represents ten trials. Lines are group averages and bands are 95% confidence intervals. For normative comparison, the performance of the Bayesian observer is shown as a dotted gray line.



Supplementary Figure 7 CA participant performance on the multiple learning strategies task separated by diagnosis (MSA individuals versus all others) and compared to healthy controls (HC)

(A) Deck learning performance on the multiple learning strategies task as indicated by the proportion of trials on which the currently lucky deck was chosen as a function of how distant those trials were from a reversal in deck value. Lines represent group averages and bands represent 95% confidence intervals. **(B)** Object value usage on trials in which a previously seen object appeared. Points represent group averages and error bars represent 95% confidence intervals. **(C)** Inverse temperature estimates from the Hybrid model. Individual points represent estimates for each subject, group-level averages are shown as filled in

Running title: REWARD LEARNING IN THE CEREBELLUM

points and error bars represent 95% confidence intervals. **(D)** Recognition memory performance on the subsequent memory task. Individual points represent each participant's d' prime score, filled in points are group-level averages and error bars are 95% confidence intervals.

Running title: REWARD LEARNING IN THE CEREBELLUM

Supplementary Table 1. Results of regression analyses assessing differences in neuropsychological test performance between CA participants and healthy controls.

Measure	β Estimate	95% Credible Interval
Backwards	-2.57	[-4.18, -0.92]
Forwards	-0.09	[-1.35, 1.15]
Semantic	1.54	[-2.12, 5.14]
Phonemic	0.04	[-2.75, 2.88]
Category	1.68	[-0.76, 4.10]
Similarities	-0.20	[-0.50, 0.12]
Go No Go	0.0001	[-0.33, 0.33]

Running title: REWARD LEARNING IN THE CEREBELLUM

Supplementary Table 2. Correlations between CA participant-level measures (Total neuropsychological scores and symptom duration) and incremental value sensitivity (as estimated by the Q Learning model in the incremental learning task and as by the Hybrid Q Learning model in the multiple learning strategies task).

Measure	Pearson's R	P Value (Bonferroni Corrected)	Task
Symptom Duration	-0.0661	1	Incremental Learning
SARA (Total)	0.1067	1	Incremental Learning
MoCA (Total)	0.5225	0.1885	Incremental Learning
CCAS (Total)	0.8419	0.0001*	Incremental Learning
BDI (Total)	0.2977	1	Incremental Learning
QUIP (Total)	-0.3892	0.7356	Incremental Learning
Symptom Duration	-0.3699	0.7848	Multiple Learning Strategies
SARA (Total)	-0.2141	1	Multiple Learning Strategies
MOCA (Total)	0.1438	1	Multiple Learning Strategies
CCAS (Total)	0.3849	0.6887	Multiple Learning Strategies
BDI (Total)	0.4641	0.3141	Multiple Learning Strategies
QUIP (Total)	-0.074	1	Multiple Learning Strategies

CCAS – Cerebellar cognitive affective/Schmahmann syndrome scale, SARA – Scale for the Assessment and Rating of Ataxia, MoCA – Montreal cognitive assessment, BDI – Beck's depression inventory, QUIP – Questionnaire for impulsive-compulsive disorders in Parkinson's disease

***p < 0.001

Running title: REWARD LEARNING IN THE CEREBELLUM

Supplementary Table 3. Correlations between CA participant CCAS subscale measures and incremental value sensitivity (as estimated by the Q Learning model) in the incremental learning task.

CCAS Measure	Pearson's R	P Value (Bonferroni Corrected)
Semantic Fluency	0.6693	0.033*
Phonetic Fluency	0.6235	0.0749
Category Switching	0.7168	0.012*
Digit Span Fwd (CCAS)	0.295	I
Digit Span Bwd (CCAS)	0.3146	I
Cube Drawing	0.5009	0.4055
Verbal Recall	0.4224	0.9118
Similarities	0.5881	0.1302
Go No Go	-0.2654	I
Affect	0.2944	I

*p < 0.05