

The algorithmic neuroanatomy of action-outcome learning

Richard W Morris¹⁻³, Amir Dezfouli², Kristi R Griffiths^{2,4}, Mike E Le Pelley^{1,3} & Bernard W Balleine^{1,2}

1. School of Psychology, University of New South Wales (UNSW), NSW, Australia
2. Brain & Mind Centre, University of Sydney, Camperdown, NSW, Australia
3. ARC Centre of Excellence in Cognition and its Disorders, Macquarie University, NSW, Australia
4. Brain Dynamics Centre, Westmead Millennium Institute, NSW, Australia

Correspondence:

Bernard Balleine
Decision Neuroscience Laboratory
Level 4 Matthews Building
School of Psychology
University of NSW
Kensington, NSW 2052 Australia

Tel: +61 2 9385 1857

Email: bernard.balleine@unsw.edu.au

Keywords: reinforcement learning, causal, computational neuroscience, fMRI

Word count:	Abstract	143
	Main text	4449
	Display items	7
	References	46

Abstract

Although it is well known that animals can encode the consequences of their actions and can use this information to control action selection and evaluation, it is not known what learning rules control action-outcome (AO) learning. Here we trained participants to encode specific AO associations whilst undergoing functional imaging (fMRI) and used computational modelling to evaluate competing models. This analysis revealed that a Kalman filter, which learned the unique causal effect of each action, best characterized AO learning and found the medial prefrontal cortex differentiated the unique effect of actions from background effects. We subsequently extended these findings to show that mPFC participated in a circuit with parietal cortex and caudate nucleus to segregate distinct contributions to AO learning. The results extend our understanding of goal-directed learning and demonstrate that sensitivity to the causal relationship between actions and outcomes guides goal-directed learning rather than contiguous state-action relations.

1 There is now considerable evidence that animals are capable of encoding the consequences of their actions
2 and that they use that information to select, evaluate and initiate future actions.¹⁻⁵ Although it is clear,
3 therefore, that such learning involves encoding of the action-outcome relationship, the learning rules that
4 govern that learning have yet to be established. Normatively, this relationship has been described by the
5 formalism ΔP (Figure 1A); which captures the effect of the action on outcome delivery over and above any
6 background effects^{6,7}. Nevertheless, how we distinguish these effects is unclear.

7 There are, in fact, very few mechanistic models of how we and other animals detect and encode the AO
8 contingency. Historically, models of associative learning have been proposed as an algorithmic account of
9 causal learning, including learning about actions.^{8,9} These models, and their modern variants,¹ assume that
10 background conditioning plays a key role because the background competes with actions via a summed
11 error term, to uniquely predict the outcome. More recently model-based reinforcement learning models
12 (MBRL), which forgo a summed error-term and instead incorporate a model of the task structure, have
13 been proposed as a general account of goal-directed action.^{4,10,11} MBRL represents the task structure in a
14 covariance matrix that represents the contiguity, rather than the unique causal relationship, between
15 actions and outcomes (i.e., the state-action transitions).

16 Here we evaluated the neural computations of learning about causal actions (AO learning) in the human
17 brain. We found a simple Kalman filter¹²⁻¹⁴ that combined prediction-error learning with the covariance
18 structure of the environment explained the acquisition of AO learning better than models based on
19 covariance or prediction-error alone. This iterative model attributed prediction-errors to different causal
20 variables adjusted by their covariance, in order to uniquely predict the outcome. Critically, model-based
21 fMRI revealed activity in the medial prefrontal cortex (mPFC) and the dorsal anterior cingulate cortex
22 (dACC) tracked changes in the predictive value of actions and the background with respect to specific
23 outcomes, separately. Furthermore, we found the mPFC participated in a network with the striatum and
24 posterior parietal cortex to segregate the influence of different causes via their covariance, a unique
25 prediction of the Kalman algorithm. These findings reveal, for the first time, an integrated corticostriatal
26 network that combines prediction-errors with covariance in order to learn how our actions control our
27 environment.

28 Results

29 We performed two fMRI experiments to probe the computational mechanisms by which the brain learns
30 the AO contingency and distinguishes it from any background effects. In typical neuroscience experiments,
31 actions and their consequences are offered in discrete trials where there is no ambiguity about unique
32 effects. Instead we used a free-operant task without a discrete trial structure, along with noncontingent
33 outcomes, which required participants to infer the causal effects of their actions.

34 AO contingency degradation revealed people learned the unique effects of their actions

35 Experiment 1 involved training hungry participants with two actions for distinct food outcomes, selected
36 before the task by each participant (e.g., button 1 = M&Ms, button 2 = BBQ flavored crackers). Training
37 occurred according to a 'constant probability' schedule.¹⁵ During the fMRI test, at the end of every second
38 the schedule recorded whether a button was pressed and then presented an outcome onscreen with a
39 conditional probability $P(O_i|A_i) = 0.2$, for each action. (Actual food outcomes were provided at the end of
40 all testing). In order to selectively degrade the causal relationship of one AO contingency while equating the
41 reward value of both actions, the schedule also presented one of the outcomes onscreen if neither button
42 had been pressed, i.e., $P(O_1|\sim A_1, \sim A_2) = 0.2$. Under this outcome-specific degradation schedule, delivery of
43 the noncontingent outcome diminishes the reward value of both actions equally (since reward can now be
44 obtained without taking either action). However the noncontingent outcome will selectively reduce the
45 *causal* relationship of only one action (A1) and not the other (A2), because the noncontingent outcome is
46 indistinguishable from the outcomes caused by one action, but easily distinguishable from the outcomes
47 caused by the other action (see methods).¹⁶

48 Figure 2A illustrates that a preference for the non-degraded action (Con) clearly emerged over time as
49 people were exposed to the differences between each AO contingency. Figure 2B (left panel) shows that
50 overall the mean number of Con actions was greater than the degraded actions (Deg). Causal ratings
51 collected at the end of each two minute block also confirmed people judged the Con action more causal
52 than the Deg action, shown in Figure 2C (left panel).

53 The selective impact of noncontingent outcomes on actions and causal judgments reveals our sensitivity to
54 the unique effect of our actions, even when the probability of reinforcement among actions is equal. The
55 noncontingent outcomes made it more difficult for participants to distinguish the outcome they caused
56 from the outcome that would have occurred anyway. As a result, the perceived causal efficacy of that
57 action was reduced. We conducted a follow-up test after the fMRI scan, under the same contingencies,
58 with the addition that each noncontingent outcome was now signaled by a yellow-light cue (Figure 1B & D).
59 The signal reduces the uncertainty about the noncontingent outcomes, which allows participants to once
60 again distinguish the unique effect of their own actions. The results found the addition of the signal
61 restored responding (Figure 2B, Signaled) as well as causal judgments of the degraded action (Deg) to the
62 same level as the Con action (Figure 2C, Signaled). The restoration of actions and causal judgments by the
63 signal implies that learning about the base-rate or ‘background conditioning’ plays a key role in learning the
64 unique causal strength of our actions.

65 **Action-selection reflected causal learning rather than reinforcement learning**

66 Figure 3A shows the correlation between causal judgments and causal actions for each person in
67 Experiment 1 was high and significant, consistent with causal learning guiding action-selection rather than
68 the frequency of reinforcement or immediate temporal effect of reward. To check our experimental control
69 over each contingency in this free-response task, we confirmed the noncontingent outcomes selectively
70 degraded the experienced contingency of the degraded action: post-hoc analysis revealed the mean
71 contingencies experienced for the Con and Deg action were $\Delta P = 0.18$ and $\Delta P = 0.07$ respectively, paired t-
72 test $t_{29} = 12.06$, $p = .8E-12$. The positive contingency (ignoring noncontingent outcomes) for the Deg action
73 was $P(O1|A1) = 0.17$, and very similar to the Con action ($P(O2|A2) = 0.18$, paired t-test $t_{29} = 0.90$, $p = .37$),
74 confirming that serendipitous differences in the positive contingency were not responsible for the results.
75 Importantly, each person received noncontingent outcomes in each block, with the exception of one
76 person who received noncontingent outcomes in only 4 out of 6 blocks. We also checked whether any
77 serendipitous reward contingency existed for the Con action. Figure 3B (blue) shows the correlation
78 between the number of Con actions and the total number of outcomes (contingent + noncontingent)
79 delivered was close to zero across participants, confirming there was no serendipitous reward contingency

80 that may have influenced preference for the Con action. Conversely there was no negative contingency
81 between Deg actions and total outcomes (Figure 3B, red). Furthermore, the distribution of delays between
82 each outcome and the preceding action did not differ within a 10-s interval for Con and Deg actions (Figure
83 3C), confirming that the immediate temporal contiguity with reward was not differentially influencing
84 action-selection. Finally, pre-test preference ratings of the snack food outcomes confirmed both rewards
85 were equally liked. The mean (95% confidence interval) ratings on a 7-point Likert scale were 5.8 CI[5.5,
86 6.2] and 6.3 CI[5.9, 6.6], for BBQ crackers and M&M respectively. Thus, action-selection did not reflect any
87 post hoc or serendipitous differences in reward contingency or contiguity.

88 **Modelling revealed a Bayesian prediction-error best explained AO learning**

89 We simulated and fit three models: a prediction-error model with a summed prediction-error term, a MBRL
90 model with a covariance matrix, and a Kalman algorithm with both, to determine which best explained
91 behavioral performance in Experiment 1. The prediction-error model assumed actions and background cues
92 competed for causal strength via a summed error-term (see methods). Simulation (Supplementary figure 1)
93 confirmed this model resolved the unique effect of the causal action (i.e., converged to ΔP). By contrast,
94 the MBRL used a transition matrix (updated via a state prediction-error) to represent the covariance
95 between each action, outcome and background. Simulation confirmed the covariance learned by this
96 model was insufficient to distinguish the causal action. In particular, the reward value of the free outcomes
97 outcompeted both actions equally in the MBRL, which did not learn or prefer the causal action. Finally, the
98 Kalman algorithm assumed actions compete with background effects via a summed prediction-error term,
99 however the amount learned about each is adjusted by the covariance between them. That is, when causal
100 variables covary (i.e., positive covariance), the model cannot distinguish their separate influence and so
101 adjusts them together. However with negative covariance then the effects can be distinguished and
102 changes in the belief of one cause will inversely affect the other (see Figure 4). By combining the prediction-
103 error term with the covariance, the Kalman filter distinguished causal actions across the widest range of
104 parameter values (Supplementary figure 2).

105 Behavioral fitting indicated action selection was better explained by the Kalman algorithm than a
106 prediction-error model or MBRL. Data from Experiment 1 was used to calculate the posterior probability of

107 the Kalman algorithm, as well as a prediction-error model, a MBRL model (using a transition matrix), and a
108 null model. The null model used the asymptote AO contingencies of each block as action values, thus it was
109 not a learning model but a static model with no temporal dynamics. Comparison with the null model
110 determines whether each learning model can explain how choices depend on the sequence of trial-by-trial
111 feedback. Table 1 shows the negative log likelihoods and relative Bayes factors of each model relative to
112 the null. After fitting each participants' data by maximum-likelihood estimation, the results of a likelihood
113 ratio test (LRT) indicated all learning models predicted significantly more behavior than chance. However,
114 the relative Bayes factor (vs the null model) shows only the optimal Kalman algorithm had a positive value,
115 indicating only this model predicted the acquisition of causal learning over time. This model also explained
116 more choices and the behavior of more individuals than the other models, with a Pseudo-R² of 0.24 and a
117 positive evidence ratio of 2 (20/10 favoring H₁/H₀), which were higher than the respective values for the
118 other learning models. Thus, the majority of subjects and the total evidence favors the simple Kalman
119 algorithm over a static model with no temporal dynamics but perfect asymptote performance. This was not
120 the case for the optimal MBRL, or the prediction-error model (or a causal induction model, see
121 Supplementary Materials), which all explained more variance than chance but had a negative GBF relative
122 to the null and a PER less than 1 (Table 1).

123 **Model-based fMRI revealed the mPFC distinguishes causal actions from background effects**

124 We evaluated whether the brain learned about causal actions, as described by the Kalman filter, by
125 regressing the model-derived learning estimates against the image data collected in Experiment 1.
126 Independent learning estimates for actions and background (ΔAO and ΔXO ; see Methods) were included as
127 parametric modulators of a stick (delta) function of response and outcome times. We included outcomes in
128 the delta function in order to include times when the action was present as well as absent (the background
129 was assumed to be always present). Whole-brain analysis revealed learning about actions and the
130 background occurred in distinct regions of the mPFC (Figure 5A & B). Action learning (ΔAO) appeared in a
131 medial region of the superior frontal gyrus (BA9, global peak MNI co-ordinates: -15 +47 +40, $Z = 4.71$, $F_{1,29} =$
132 37.12 , $FWE = .031$). At the same time, learning related changes to the background estimates (ΔXO)
133 appeared in the dorsal anterior cingulate cortex (BA32, global peak MNI: -9 +41 +22, $Z = 5.19$, $F_{1,29} = 50.20$,

134 FWE = .004), as well as smaller changes in the left caudate (MNI: -15 +14 +7, $Z = 4.88$, $F_{1,29} = 41.80$, FWE =
135 .017), and cuneus (MNI: -3 -64 +34, $Z = 4.39$, $F_{1,29} = 37.28$, FWE = .04). No other region survived multiple
136 comparison correction in this whole-brain analysis.

137 The results described so far indicate the background expectancy produced by the noncontingent outcome
138 plays a key role in learning the unique effect of our actions. According to prediction-error models (including
139 the Kalman algorithm), this background expectancy will be violated whenever the noncontingent outcome
140 does not follow an action. This implies that a negative AO contingency will be learned under certain
141 conditions (i.e., inhibitory learning). For example, in Experiment 1 we explicitly arranged that O1 sometimes
142 occurs after A1 but never after A2 (in order to equate the reward value of both actions), which results in a
143 negative contingency between A2-O1. We tested a regressor of changes to this negative AO contingency, as
144 learned by the Kalman algorithm, in the whole-brain. BOLD activity in a ventral medial prefrontal region,
145 including the anterior cingulate (BA32) and medial orbitofrontal cortex (BA10), learned the negative AO
146 contingency, global peak MNI: -3 +50 -11, $Z = 3.52$, $F_{1,29} = 18.08$, FWE = .011 (Supplementary figure 3). These
147 imaging results are consistent with recent reports in rodents that the medial orbitofrontal cortex is critical
148 for learning about unobserved outcomes, and in particular a negative AO contingency.¹⁷

149 **Model-based fMRI showed the posterior parietal cortex tracks covariance between causes.**

150 The results so far indicate causal actions are distinguished from the background in the mPFC. A unique
151 feature of the Kalman filter is that the covariance term distinguishes the influence of candidate causes by
152 updating ΔAO and ΔXO together when the covariance is positive and updating them in opposite directions
153 when the covariance is negative (Figure 4). We tested whether any brain regions tracked the covariance
154 between actions and background by entering the covariance values as a parametric modulator. A whole-
155 brain analysis revealed bilateral activity in posterior parietal cortex (BA40) was significantly associated with
156 the covariance term, left global peak MNI: -57 -55 +37, $Z = 5.83$, $F_{1,29} = 72.78$, FWE < .001 and right peak
157 MNI: +42 -67 +43, $Z = 5.34$, $F_{1,29} = 54.67$, FWE = .002 (Figure 5B).

158 **ROI analysis showed prediction-error and covariance converge in caudate**

159 Substantial evidence exists that the ventral striatum tracks or receives reward prediction-errors,¹⁸⁻²⁰ while
160 dorsal striatal regions track action values.¹⁸ In the present example of AO learning, the prediction-error
161 represents the deviations between the observed outcome and the summed total causal expectancy (action
162 + background). We tested whether the striatum tracks this summed error term, by including it as a
163 parametric modulator in an anatomical ROI analysis of the striatum. Figure 5C & 5D shows BOLD responses
164 in a posterior region of the caudate body (green) tracked the summed errors (ROI peak MNI: : +15 +11 +4, Z
165 = 4.17, $t_{29} = 4.94$, FWE = .002 svc) while activity in the anterior caudate (red) was associated with the
166 covariance between actions and background (ROI peak MNI: -15 +23 +7, Z = 3.42, $t_{29} = 3.84$, FWE = .029 svc)
167 These regions were more medial and dorsal than those implicated in reward prediction-error signals but
168 similar to regions implicated in instrumental learning.¹⁸ Thus, the caudate appears to receive sufficient
169 information to segregate the influence of different events and may play an important role in selectively
170 distinguishing causal actions from background effects.

171 **DCM revealed the caudate segregates the effect of the action from background effects**

172 To further determine the caudate's role in distinguishing control, we performed a dynamic causal model
173 (DCM) analysis.²¹ We tested two possibilities shown in Figure 5E & 5F. In Model 1 the caudate is a site of
174 convergence of the updated values from the prefrontal cortex to enable action-selection. In Model 2 the
175 caudate segregates the prediction-error to update the estimates of the action and background separately.
176 Bayesian model selection revealed the relative log-evidence for model 2 was 85.31, which corresponds to
177 strong evidence in favor of segregation. A random-effects analysis (Figure 5G & 5H) revealed a large
178 majority of participants were significantly more consistent with the segregate model than the integration
179 model (the exceedance probability was 99 percent), and it was more likely to be true for any random
180 subject (the posterior probability of the segregate model was 80 percent).

181 **PPI showed cortex and caudate interact when causal actions must be distinguished by their covariance**

182 Experiment 2 replicated the key fMRI results in an independent sample of naive participants, using a design
183 that allowed us to assess the effect of distinguishing the free outcomes on the corticostriatal network we
184 identified above. The experiment used a single AO contingency and varied whether or not the free
185 outcomes were distinguishable from the earned outcomes in a block-by-block fashion, to test the

186 interaction between causality and caudate activity. For half the blocks we used the same outcome for both
187 free and earned outcomes (i.e., as in Experiment 1), while in the other half the earned outcomes were
188 different from the free outcomes. Using distinct outcomes in half the blocks allowed the participant to
189 discern the causal effect of their actions (i.e., equivalent to signaling the free outcomes, as in the follow-up
190 to Experiment 1). Figure 6A shows that degradation reduced total actions and causal ratings, however
191 providing distinct free outcomes restored causal actions and judgments. As before, we fitted the Kalman
192 algorithm to the data using maximum likelihood estimation. The optimal model predicted significantly
193 more choices than chance, the mean group average likelihood per trial was 57 percent (95% CI: 53-60). A
194 functional ROI (fROI) analysis using masks generated from the significant results in the mPFC and caudate of
195 Experiment 1 confirmed that learning about the background (ΔXO) occurred in the same dorsal ACC region
196 (Figure 6C, blue), ROI peak MNI: -3 +36 +38, $Z = 3.94$, $t_{19} = 5.06$, FWE = .02. Meanwhile learning values for
197 the action (ΔAO) occurred in the same region of BA9 (Figure 6C, violet), ROI peak MNI: -2 +47 +46, $Z = 3.04$,
198 $t_{19} = 3.56$, $p = .001$. Covariance between the action and background was tracked in the caudate (Figure 6C,
199 right), ROI peak MNI: -12 +20 +1, $Z = 3.91$, $t_{19} = 5.01$, FWE = .002. We used a whole-brain PPI analysis to
200 determine whether any cortical regions interacted with the caudate when free outcomes were
201 indistinguishable from earned outcomes. A single region in the right parietal junction interacted with the
202 caudate when free outcomes were the same (versus different), shown in Figure 6E, global peak MNI: +54 -
203 58 +30, $Z = 4.68$, $F_{1,19} = 24.90$, FWE = .01. This region overlapped with the right posterior parietal cortex
204 (BA40) identified in Experiment 1 (Figure 6E, red). Together these results implicate the posterior parietal
205 cortex in tracking the covariance with the background, and interacting with the caudate when the unique
206 effect of our actions must be distinguished from the background.

207 Discussion

208 We sought to establish the learning rules that govern AO learning in instrumental conditioning and their
209 neural bases. We found that the medial prefrontal cortex participates in a circuit that detects and
210 segregates the unique causal effects of our actions from other background effects, and more importantly,
211 that this segregation was generated by a Bayesian prediction-error, described by a Kalman filter, which uses

212 a summed prediction-error term along with the covariance between potential causes to distinguish the
213 unique effect of actions from background effects. Furthermore, the caudate appears to be a key point of
214 integration of the covariance term and prediction-error; it segregates the summed prediction-error into
215 separate update values for each causal belief. Thus, this model represents a simple, iterative Bayesian
216 model of change, that unlike other computational-level models,²² provides an algorithmic account of AO
217 learning that can be instantiated in the neural code.²³

218 Many results have emphasized the critical role of the medial prefrontal cortex in AO learning, however the
219 exact nature of this role has been unspecified. Indeed, there is a wealth of evidence that in the rat, the
220 prelimbic region of the medial prefrontal cortex is critical for the acquisition of goal-directed actions.^{2,24-29}
221 Computational models of medial prefrontal cortex function in humans such as the PRO model¹ assume the
222 *dorsal anterior cingulate* signals negative prediction-errors during AO learning, consistent with other
223 prediction-error models.⁸ Such models only offer a partial explanation of the results we observed here: The
224 update to the background effect (ΔXO) represents the learned probability of the noncontingent outcome
225 that must be adjusted against the probability of the contingent outcome. This is consistent with negative
226 prediction-errors when positive changes to the background ($+\Delta XO$) occur in the context of a negative
227 covariance term, and so result in negative changes to the AO contingency ($-\Delta AO$). At other times, positive
228 changes to the background probability ($+\Delta XO$) will occur in line with changes to the AO contingency,
229 consistent with positive prediction-errors. This occurs when the covariance is positive, for instance, in
230 situations in which we cannot distinguish the unique effect of our actions from potential background
231 causes.

232 Our results also distinguished a separate region of the mPFC, near the medial surface of BA9, whose activity
233 represented updates to the unique effect of the action (ΔAO) and replicated the involvement of the mPFC
234 in an independent sample, highlighting the reliability of the current findings (Figure 6C). We further showed
235 using PPI (Figure 6D & E) that the caudate interacts with the parietal cortex when free outcomes are
236 indistinguishable from the earned outcome (i.e., when there was no signal or the free outcome was the
237 same as the earned outcome). The selective interaction between parietal cortex and caudate arises under
238 the additional demands when there is no observable information to distinguish control. When

239 noncontingent outcomes are indistinguishable, the covariance between our actions and the background is
240 the only information that can be used to distinguish them. The right posterior parietal junction identified in
241 the PPI was the same cortical region identified in Experiment 1 as tracking the covariance term (along with
242 the caudate in the subcortex). It also overlaps with the cortical region previously implicated in learning the
243 transition matrix during model-based reinforcement learning,³⁰ suggesting, across studies and laboratories,
244 that this region represents the covariance structure of the environment. While PPI does not indicate the
245 direction of influence, these results are consistent with a network extending from the parietal cortex to
246 caudate to mPFC, which tracks the covariance between actions and other events and then segregates the
247 error-term to learn about and distinguish between the influence of different causes.

248 We found the competitive allocation of causal belief to actions relative to the background is a form of
249 selective learning closely related to cue-competition models in associative learning.⁸ An important
250 difference is the covariance matrix of the Kalman filter, in which the off-diagonal terms track the
251 covariation between events. In our model, the covariance allowed the learner to distinguish or segregate
252 the effects of an action from the background in the absence of that action, as shown in Figure 4. This
253 allowed the model to reason counterfactually about what would have happened if an action had not
254 occurred. In this manner, the covariance is analogous to heuristically motivated formalizations of within-
255 event learning (e.g., negative alpha) which allows learning about absent events in recent versions of cue-
256 competition models.^{31,32} Our simple Kalman filter thus combines key features of contemporary associative
257 learning and model-based reinforcement learning.

258 These results also make an important contribution to the common claim that goal-directed learning is
259 analogous to MBRL. In general, MBRL is concerned with building a model of the environment, given the
260 state caused by each action (i.e., the covariance or transition matrix). In such models, “state prediction-
261 errors” and the covariance matrix they update³⁰ only describe the contiguity between states and, as a
262 result, the covariance matrix cannot learn or accurately represent a causal relationship. By contrast, causal
263 learning is not concerned with the transition probabilities between different states, but rather the trade-
264 offs between competing contingencies to determine whether that state was caused by an action or not. This
265 is a primary difference between causal models and MBRL. The question of which action caused which state

266 is arguably more fundamental to goal-directed learning. For example, the prediction-errors we observed
267 here are critically different from “state prediction-errors” because they are adjusted for the probability that
268 another state (the background) may also cause the outcome. This adjustment leads to a representation of
269 the unique causal strength of our actions. Hence, an important implication of this proposal is that MBRL *per*
270 *se* may be sufficient for maximizing reward but it does not provide a complete account of goal-directed
271 learning since it is unable to calculate, and so is insensitive to, the causal relationship between actions and
272 states.

273 Unlike some other computational models of causal learning, the causal estimates learned by our Kalman
274 filter converge to ΔP , a normative measure of causal strength. Other researchers have argued that ΔP does
275 not provide the best approximation of human causal inference because changing the base-rate probability
276 of an outcome while holding ΔP constant modulates causal judgements (e.g., “the base-rate illusion”).
277 However the base-rate illusion is considerably weaker in free-response, instrumental learning where trials
278 are not explicitly segmented.³³⁻³⁶ Furthermore, when learning about causal effects, active intervention is a
279 more reliable guide to causal relations than is sheer observation, largely because actions constitute one
280 basic way to control for possible alternative causes.³⁷⁻⁴⁰ Humans are able to reason suppositionally or
281 counterfactually about what would be expected to happen if an intervention is made or not made, and
282 midbrain dopamine neuron firing⁴¹ along with our Bayesian model reflects these counterfactual action
283 values. For these reasons, AO learning may not suffer the same biases as other forms of causal learning that
284 are based on passive observation.

285 In conclusion, learning about the causal effects of our actions, as required for goal-directed learning and as
286 investigated here, appears to reflect features of traditional associative models such as competition for
287 predictive value, as well as modern conventions such as environmental structure (covariance). In our hands,
288 these features were combined in a highly simplified, iterative Kalman filter that learned a probability
289 distribution over action-outcome contingencies to provide a novel account of AO learning. In our results
290 there was impressive agreement across experiments and replications that distinct regions of a
291 corticostriatal network distinguished the unique causal effect of actions from those of the background.
292 More generally, the results revealed how our neuroanatomy performs Bayesian computations, consistent

293 with growing evidence that the brain learns and make decisions on the basis of probability
294 distributions.^{23,42-44}

295

296 **Methods**

297 **Participants.**

298 In Experiment 1, 31 right-handed English speaking volunteers, aged between 19 and 51 (mean age 30.5)
299 were scanned. One participant was removed due to excessive head movement (> 2 mm), thus $n = 30$ (18
300 females). On the basis of a power analysis (Supplementary Material), scans from 23 right-handed, English
301 speaking volunteers, aged between 17 and 32 (mean age 26) were considered for Experiment 2. Three
302 participants were removed due to excessive head movement (> 2 mm), thus $n = 20$ (11 females). All
303 participants were free of food allergies, neurological or psychiatric disease, and psychotropic drugs, and
304 reported strong liking of the snack foods we provided as reward. Informed consent to participate was
305 obtained and the study was approved by the Human Research Ethics Committee at the University of Sydney
306 (HREC no. 12812). Participants were reimbursed \$45 AUD in shopping vouchers, in addition to the snack
307 foods they earned during the test session.

308 **AO contingency degradation task.**

309 In each experiment, participants were instructed not to eat three hours prior the appointment. Pre-testing
310 involved obtaining preference ratings on a 7-point scale for each of three snacks (M&Ms, BBQ flavored
311 crackers, chocolate cookies), from which the two most similarly preferred snacks were selected for the
312 experiment.

313 Experiment 1 involved learning two AO contingencies concurrently. Participants were instructed they could
314 liberate snack foods (BBQ flavored crackers and M&Ms) from a vending machine by tilting it to the left or
315 right (by pressing either a left or right button), and that sometimes the vending machine would also release
316 a snack for free. They were also instructed to find the best action for releasing snacks. Outcomes were
317 indicated by the presentation of a visual stimulus depicting the snack for 1-s duration (a particular snack
318 food, e.g., M&M or BBQ cracker), during which further outcomes could not be earned. The relationship

319 between actions and outcomes were constant across blocks for each participant (e.g., left = M&M and right
320 = BBQ crackers for all blocks). Each block lasted 120-s, and the software controlling the task PsychoPy2
321 v1.8,^{45,46} divided each block into 120 one-second intervals to determine the outcome rate. Participants
322 were unaware of the 1-s intervals, and they responded freely using the index finger on their right hand to
323 press the left or right button on a Lumina MRI-compatible response pad (LU-400, Cedrus Corporation, CA).
324 An action (tilt left or tilt right) earned a particular outcome with a probability $P = 0.2$ if that action had
325 occurred in the preceding 1-s interval. If both actions occurred in the preceding interval then only the most
326 recent action was considered for reinforcement. A free outcome was delivered with $P = 0.2$ if neither action
327 had been made. This schedule ensured two important features: 1) that there was no serendipitous
328 contingency between an action and a free outcome, which would result in a higher reward contingency for
329 the contingent action¹⁶, and 2) the earned outcome appeared at a varying interval up to one second after a
330 successful action, which is sufficient to introduce ambiguity into the perceived AO contingency.⁴⁷
331 Participants completed six blocks; the outcome (BBQ cracker or M&M) that was subject to contingency
332 degradation was counterbalanced across blocks (ABBAAB). At the end of each block, participants rated how
333 causal each action was with respect to each outcome on a Likert scale from 1 (not at all) to 7 (very causal).
334 A follow-up test was conducted after the scan. The test setting, duration and programmed AO
335 contingencies in the follow-up test were exactly the same as in the scanner, with the addition of a 1-s
336 yellow light cue displayed on the front of the virtual vending machine immediately prior to the delivery of
337 each free reward. At the end of all testing, participants received all snacks that had been delivered
338 onscreen during test.

339 Experiment 2 involved learning a single AO contingency. The session was arranged in 12 blocks of 60-s
340 duration, and in each block the participant responded freely for a single snack food reward
341 (counterbalanced between BBQ crackers and M&Ms). As before, in each block the positive contingency was
342 $P(O|A) = 0.2$ in every second a response was made. The probability of a free outcome in every second when
343 no response was made, i.e., $P(O|\sim A)$, varied between 0, 0.1 and 0.2 across blocks in a counterbalanced
344 order. Conversely, ΔP varied from 0.2, 0.1 to 0. In addition, Experiment 2 varied whether the free outcome
345 was the same snack or a different snack as the earned outcome in each block, in an ABBA order. In half the

346 blocks the earned and free outcomes were different which effectively signaled the free outcomes and
347 allowed the participant to discern the causal effect of their actions. For the other half of blocks the free
348 outcomes were the same snack food as the earned outcome, thus making it difficult to discern unique
349 causal effects.

350 **Behavior data analysis**

351 The behavioral data consisted of the rate of responding during each block and the causal ratings obtained
352 at the end each block. Experiment 1 tested for differences between Con and Deg actions in the proportion
353 of total responses, as well as mean causal ratings. In each case, a Shapiro-Wilk test confirmed the data did
354 not violate the assumption of normality and differences were assessed by paired t-tests (two-tailed).
355 Experiment 2 tested the main effect of the outcome condition (Same versus Different free outcomes) and
356 its linear interaction with the contingency condition ($\Delta P = 0.2, 0.1, \text{ and } 0.0$), using a 2 x 3 repeated
357 measures ANOVA (two-tailed). Mauchly's test was used to detect violations of sphericity, in which case the
358 Greenhouse-Geisser correction was applied.

359 **Model-based analysis.**

360 For each of the learning models described below, the real-time occurrence of outcomes was modelled with
361 a logistic function $f(x) = 1 / (1 + e^{\infty(D-k)})$ to produce a binary result (0,1) determining whether the outcome
362 will or will not be associated with the prior action. D is the delay between the previous action and outcome,
363 and k is the temporal threshold included as a free parameter in the model-fitting described below.

364 Prediction-error learning. The prediction-error model adopted a standard delta rule exemplified in the
365 Rescorla-Wagner model⁸ and adapted for vector-valued predictions in modern reincarnations.¹ This allows
366 multiple action outcomes to be predicted simultaneously, each with its own summed error-term. In this
367 model, the predicted outcome \hat{o} is a weighted sum of actions and background cues ($\hat{o} = Va$). Updates to the
368 weights (V) occur by the prediction error:

$$369 \quad \Delta V = \alpha(o - Va) \quad (1)$$

370 where α is a free parameter controlling the learning rate. In this way, the model replicates the prediction-
371 error learning of the Kalman algorithm below, but without adjustment by a covariance matrix so all changes
372 are restricted to actions on the current trial.

373 Model-Based Reinforcement Learning. The MBRL model was adapted from the FORWARD model described
374 in Glascher et al (2010), which uses experience with state transitions to update an estimated matrix of
375 transition probabilities. The transition matrix (T) held the current estimate of the probability of
376 transitioning from action \mathbf{a} (a binary vector indicating one of three possibilities: make action A1, make
377 action A2, or Wait) to an outcome state \mathbf{o} (a binary vector indicating one of three possibilities: outcome O1
378 delivered, outcome O2 delivered, or no outcome delivered [background state]). Wait actions occurred at
379 the end of every second in which no other action occurred. In the T matrix, the different actions were
380 represented in different columns, while the different outcomes were represented in different rows. The
381 transitions were initialized to uniform distributions connecting each action and outcome. Upon each step,
382 having taken action \mathbf{a} and arrived in outcome state \mathbf{o} , the FORWARD learner estimates the expected
383 outcomes on the basis of the current transition matrix ($\hat{\mathbf{o}} = T\mathbf{a}$), and computes a state prediction-error ΔT :

$$384 \quad \Delta T = \alpha(\mathbf{o} - \hat{\mathbf{o}}) \quad (2)$$

385 Updates to the transition matrix T of the observed transition occur via ΔT :

$$386 \quad T(i) = T(i) + \Delta T \quad (3)$$

387 where i is the column corresponding to the taken action.

388 Kalman algorithm. The aim of the Kalman algorithm was to learn the unique causal weight of each action
389 over and above the background (i.e., ΔP). The algorithm builds a probabilistic representation of the causal
390 weights (\mathbf{w}) of each input (actions and background cues) predicting each outcome (\mathbf{o}), representing causal
391 beliefs. The causal beliefs are represented by a multivariate normal density $N(\mathbf{w}|\boldsymbol{\mu}, C)$ with a prior Gaussian
392 distribution, and after observing each outcome the causal weights are updated by changes in the mean and
393 variance of each distribution (see below). The mean $\boldsymbol{\mu}$ represents the belief in the unique causal strength
394 while the variance C captures the uncertainty around that belief. When the variance is large, there is large

395 uncertainty regarding the true causal strength. The updating equations for the mean and variance of the
 396 causal weight have the following form:

$$397 \quad \Delta\boldsymbol{\mu} = (\nu + \mathbf{a}^T C \mathbf{a})^{-1} (\mathbf{o} - \boldsymbol{\mu}^T \mathbf{a}) C \mathbf{a} \quad (4)$$

$$398 \quad \Delta C = -(\nu + \mathbf{a}^T C \mathbf{a})^{-1} C \mathbf{a} \mathbf{a}^T C \quad (5)$$

399 where ν is a free parameter capturing outcome variance. For each second, the execution of an action and
 400 the constant background context are represented as a binary input vector (\mathbf{a}) and an outcome vector (\mathbf{o})
 401 represents the delivery of the outcomes (O1 and O2). The first term $(\nu + \mathbf{a}^T C \mathbf{a})^{-1}$ represents the total
 402 certainty (inverse sum of outcome uncertainty and belief uncertainty) and it governs the learning rate. $\boldsymbol{\mu}^T \mathbf{a}$
 403 is a vector representing the learned causal weights on the basis of the current inputs (\mathbf{a}). The difference
 404 between the observed outcome and learned causal weights ($\mathbf{o} - \boldsymbol{\mu}^T \mathbf{a}$) is a vector of outcome-specific
 405 prediction errors, each of which represents a summed error-term for O1 and O2. The rightmost term $C \mathbf{a}$ is
 406 the product of the covariance matrix and input vector, and it allows for changes to the mean belief about
 407 actions otherwise correlated with the background context but absent on the present trial. In this manner,
 408 the Kalman filter is able to distinguish the unique influence of actions from the context during
 409 noncontingent outcomes. Importantly, the covariation between each action and the background are
 410 tracked in the off-diagonal elements of C , which allowed us to test a unique prediction of this model. Thus
 411 changes to the mean beliefs ($\Delta\boldsymbol{\mu}$) depend on the prediction-error as well as the covariance matrix C .

412 Null model. For comparison, we also described a null model without any temporal dynamics but ideal
 413 asymptote performance. The null model assumed that the probability of taking each action was
 414 proportional to the final ΔP obtained in each block, so $Q_{\text{right}} = \Delta P_{\text{right}}$ and $Q_{\text{left}} = \Delta P_{\text{left}}$.

415 Policy. The policy of each model was the same. In each learning model, each action had a unique causal
 416 relationship with two outcomes representing the belief regarding that particular AO contingency (e.g., $\boldsymbol{\mu}_{a,o}$).
 417 So for each action, we selected the highest causal belief associated with that action $Q_a = \arg \max_a \boldsymbol{\mu}_{a,o}$ and
 418 then used it to probabilistically explain the action choices of each participant using the softmax rule:

$$419 \quad \pi_{\text{right}} = e^{\tau Q_{\text{right}}} / (e^{\tau Q_{\text{right}}} + e^{\tau Q_{\text{left}}}) \quad (6)$$

420 Bayesian model comparison. We generated observation models based on the three learning models
421 described above as well as the null model, and fit them to each subject's behaviour separately using
422 maximum-likelihood estimation. A non-linear optimization was achieved using the fmincon function in
423 MATLAB R2014B (The Mathworks Inc., MA, USA) over 100 random starting values for each subject. We
424 measured the overall goodness of fit using the average likelihood per trial of the best fit model for each
425 subject. The average likelihood per trial was calculated as the exponent of the sum of log likelihoods
426 divided by the number of trials (responses) for each subject. We compared models by aggregating the
427 probability of the data given the model over each subject's fit (single-subject Bayesian Information Criterion
428 [BIC] score) to estimate the model evidence for the full dataset. The aggregate for each model was then
429 compared to compute a Group Bayes Factor (GBF). We also report the number of subjects for whom
430 individual model comparison gave the same answer as the GBF and the positive evidence ratio (PER), where
431 "positive" evidence for one model versus another exists if the log Bayes factor is larger than three (Kass &
432 Raftery, 1995).

433 **Image data analysis**

434 MRI data were acquired on a 3-Tesla GE Discovery using a 32-channel head coil. A T1-weighted high-
435 resolution structural scan was acquired for each subject for screening and registration with a 1-mm³ voxel
436 resolution (TR: 7200 ms, TE: 2.7 ms, 176 sagittal slices, 1-mm thick, no gap, 256 x 256 x 256 matrix). For
437 BOLD acquisition, we acquired echo planar image (EPI) volumes comprising 52 axial slices in an ascending
438 interleaved acquisition order (TR: 2910 ms, TE: 20 ms, FA: 90 degrees, FOV: 240 mm, matrix: 128 x 128,
439 acceleration factor: 2, slice gap: 0.2 mm) with a voxel resolution of 1.88 x 1.88 x 2.0 mm. Slices were angled
440 15 degrees from AC-PC to reduce signal loss in the OFC. In Experiment 1, 343 EPIs were acquired while in
441 Experiment 2, 260 EPIs were acquired.

442 Data were analysed using SPM8 (www.fil.ion.ucl.ac.uk/spm). Preprocessing and statistical analysis were
443 conducted separately for each experiment. The first four images were automatically discarded to allow for
444 T1 equilibrium effects, then images were slice time corrected to the middle slice and realigned with the first
445 volume. The structural image was coregistered to the mean functional image, segmented and warped to
446 MNI space. The warp parameters were then used to normalise the resampled functional images (2 mm³).

447 Images were then smoothed with a Gaussian kernel of 8-mm full-width half maximum to improve
448 sensitivity for group analysis.

449 **Model-based fMRI analysis**

450 For each first-level GLM analyses, we constructed a vector of delta values for action causal beliefs (ΔAO)
451 and background causal beliefs (ΔXO), generated with the parameters provided by the group maximum
452 likelihood estimation (MLE).⁴⁸ For ΔAO , the delta values were taken for the current action contingency (i.e.,
453 A1-O1 or A2-O2), while for ΔXO the delta values were taken for the current outcome (B-O1 or B-O2). To
454 test for brain activity tracking the unique changes in each vector, we entered ΔAO and ΔXO as parametric
455 modulators of a stick function that included both response and reward times in an event-related design.
456 While these update signals will fluctuate independently, there will be some collinearity when the
457 covariance is zero. Collinearity is a problem when trying to determine unique effects associated with each
458 regressor. However, the variance inflation factor can be used to indicate if a collinearity problem is present.
459 The variance inflation factor was 1.23, which is within the bounds of a conservative threshold < 5 .⁴⁹
460 Nevertheless, to remove any residual collinearity between these regressors, each regressor was entered as
461 the second modulator to ensure it was adjusted for the prior regressor using the default orthogonalize
462 routine in SPM.⁴⁹ Each GLM also contained rating periods and six movement parameters. Betas were
463 estimated with a 128-s high-pass filter and AR1 correction for auto-correlation. The resulting beta images
464 were included in a group-level random effects analysis in SPM one-sample t-tests. SPM F-contrasts (two-
465 tailed) were used to create whole-brain statistical parametric maps, corrected for multiple comparisons
466 using a voxel-level FWE- $p < .05$. SPM t-contrasts (one-tailed) were used in each ROI analysis, corrected for
467 multiple comparisons using FWE (svc) in the case of anatomical ROIs (Experiment 1) and uncorrected at $p <$
468 $.001$ (svc) in the case of independent functional ROIs (Experiment 2).

469 DCM analysis. Each of the volumes-of-interest (VOI) for the DCM analysis was spatially defined according to
470 the group results of the relevant GLM analysis. The BA9 VOI was defined by the significant cluster from the
471 analysis of ΔAO in the group results. This significant group cluster was used to construct a binary mask and
472 this mask was then used to define the VOI and extract the first eigenvector for all individuals, adjusted for
473 the ΔAO and ΔXO regressors. All subthreshold voxels within the mask were included which were $p < .5$

474 (uncorrected), which roughly corresponds to all voxel activity positively related to ΔAO . The mPFC VOI was
475 extracted in the same manner but using the significant cluster from the analysis of ΔXO in the group
476 analysis. The caudate VOI was defined by a group ROI analysis of press rate restricted to the striatum ($p <$
477 $.05$, small volume corrected), with a single cluster of 48 voxels in the right caudate (peak MNI: +15 +10 +6),
478 and otherwise extracted in the same manner as other VOIs.

479 PPI analysis. The psychological term was the block condition, whether or not the free rewards were
480 distinguishable within that block. The physiological term was the timeseries from the anterior caudate in
481 each participant ($n = 20$) using a group fROI mask from the GLM analysis of the covariance. We constructed
482 the interaction term in SPM8 (per defaults) and included all three terms in the first-level GLM. Finally we
483 tested for regions of interaction in the whole-brain, corrected for multiple comparisons FDR $-q < .05$.

484 **Data availability**

485 The analyses in this report were conducted by RWM (unblinded). Data is available upon request.

486 Unthresholded statistical maps are available for viewing and download at

487 <http://neurovault.org/collections/VXWZKTWE/>. Experimental programs and Matlab code to generate

488 simulations can be downloaded from <http://balleinelab.com>

Acknowledgments

This study was supported by a Laureate Fellowship from the Australian Research Council (ARC; #FL0992409) awarded to BWB. RWM was supported by National Health and Medical Research Council (NHMRC) Project Grant #1069487, and the ARC Centre of Excellence in Cognition and its Disorders (Macquarie University). MJG was supported by the NHMRC R.D. Wright Biomedical Career Development Fellowship (APP1061875). MLP was supported by an ARC Future Fellowship (FT100100260) and BWB by a Senior Principal Research Fellowship from the NHMRC. The programs used in the behavioral tasks are available for download at [http:// balleinelab.com](http://balleinelab.com).

The authors declare no biomedical financial interests or potential conflicts of interest.

References

- 1 Alexander WH & Brown JW (2011). Medial prefrontal cortex as an action-outcome predictor. *Nat Neurosci* 14, 1338-1344.
- 2 Balleine BW & Dickinson A (1998). Goal-directed instrumental action: contingency and incentive learning and their cortical substrates. *Neuropharmacology* 37, 407-419.
- 3 Balleine BW & O'Doherty JP (2010). Human and rodent homologies in action control: corticostriatal determinants of goal-directed and habitual action. *Neuropsychopharm* 35, 48-69.
- 4 Doll BB, Simon DA & Daw ND (2012). The ubiquity of model-based reinforcement learning. *Curr Opin Neurobiol* 22, 1075-1081.
- 5 Tanaka SC, Balleine BW & O'Doherty JP (2008). Calculating consequences: brain systems that encode the causal effects of actions. *J Neurosci* 28, 6750-6755.
- 6 Allan LG (1980). A Note on Measurement of Contingency between 2 Binary Variables in Judgment Tasks. *Bulletin of the Psychonomic Society* 15, 147-149.
- 7 Maier SF & Seligman MEP (1976). Learned Helplessness - Theory and Evidence. *Journal of Experimental Psychology-General* 105, 3-46.
- 8 Rescorla RA & Wagner AR. in *Classical Conditioning II* (eds A. H. Black & W. F. Prokasy) 64-99 (Appleton-Century-Crofts, 1972).
- 9 Dickinson A (2001). Causal learning: an associative analysis. *Quarterly Journal of Experimental Psychology* 54B, 3-25.
- 10 Daw ND, Niv Y & Dayan P (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nat Neurosci* 8, 1704-1711.
- 11 Solway A & Botvinick MM (2012). Goal-directed decision making as probabilistic inference: a computational framework and potential neural correlates. *Psychol Rev* 119, 120-154.
- 12 Dayan P, Kakade S & Montague P (2000). Learning and selective attention. *Nature neuroscience* 3 Suppl, 1218-1223.
- 13 Sutton RS (1992). Gain adaptation beats least squares? *Proceedings of the Seventh Yale Workshop on Adaptive and Learning Systems*, 161-166.
- 14 Meinhold RJ & Singpurwalla ND (1983). Understanding the Kalman Filter. *American Statistician* 37, 123-127.
- 15 Hammond LJ (1980). The effect of contingency upon the appetitive conditioning of free-operant behavior. *J Exp Anal Behav* 34, 297-304.
- 16 Dickinson A & Mulatero CW (1989). Reinforcer Specificity of the Suppression of Instrumental Performance on a Non-Contingent Schedule. *Behavioural Processes* 19, 167-180.
- 17 Bradfield LA, Dezfouli A, van Holstein M, Chieng B & Balleine BW (2015). Medial Orbitofrontal Cortex Mediates Outcome Retrieval in Partially Observable Task Situations. *Neuron* 88, 1268-1280.
- 18 O'Doherty J, Dayan P, Schultz J, Deichmann R, Friston K *et al.* (2004). Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science* 304, 452-454.
- 19 Pessiglione M, Seymour B, Flandin G, Dolan RJ & Frith CD (2006). Dopamine-dependent prediction errors underpin reward-seeking behaviour in humans. *Nature* 442, 1042-1045.
- 20 Tobler PN, O'Doherty JP, Dolan RJ & Schultz W (2006). Human neural learning depends on reward prediction errors in the blocking paradigm. *J Neurophysiol* 95, 301-310.
- 21 Daunizeau J, David O & Stephan KE (2011). Dynamic causal modelling: a critical review of the biophysical and statistical foundations. *Neuroimage* 58, 312-322.
- 22 Jacobs RA & Kruschke JK (2011). Bayesian learning theory applied to human cognition. *Wiley Interdiscip Rev Cogn Sci* 2, 8-21.
- 23 Guest O & Love BC (2016). What the Success of Brain Imaging Implies about the Neural Code. *bioRxiv*
- 24 Yin HH, Knowlton BJ & Balleine BW (2005). Blockade of NMDA receptors in the dorsomedial striatum prevents action-outcome learning in instrumental conditioning. *Eur J Neurosci* 22, 505-512.
- 25 Ostlund SB & Balleine BW (2007). The contribution of orbitofrontal cortex to action selection. *Ann N Y Acad Sci* 1121, 174-192.

- 26 Ostlund SB & Balleine BW (2007). Orbitofrontal cortex mediates outcome encoding in Pavlovian but
not instrumental conditioning. *J Neurosci* 27, 4819-4825.
- 27 Corbit LH & Balleine BW (2003). The role of prelimbic cortex in instrumental conditioning. *Behav*
Brain Res 146, 145-157.
- 28 Dias-Ferreira E, Sousa JC, Melo I, Morgado P, Mesquita AR *et al.* (2009). Chronic stress causes
frontostriatal reorganization and affects decision-making. *Science* 325, 621-625.
- 29 Naneix F, Marchand AR, Di Scala G, Pape JR & Coutureau E (2009). A role for medial prefrontal
dopaminergic innervation in instrumental conditioning. *J Neurosci* 29, 6599-6606.
- 30 Glascher J, Daw N, Dayan P & O'Doherty JP (2010). States versus rewards: dissociable neural
prediction error signals underlying model-based and model-free reinforcement learning. *Neuron*
66, 585-595.
- 31 Dickinson A & Burke J (1996). Within-compound associations mediate the retrospective reevaluation
of causality judgements. *Quarterly Journal of Experimental Psychology Section B-Comparative and*
Physiological Psychology 49, 60-80.
- 32 Van Hamme LJ & Wasserman EA (1994). Cue Competition in Causality Judgments - the Role of
Nonpresentation of Compound Stimulus Elements. *Learning and Motivation* 25, 127-151.
- 33 Vallee-Tourangeau F, Murphy RA & Baker AG (2005). Contiguity and the outcome density bias in
action-outcome contingency judgements. *Quarterly Journal of Experimental Psychology Section B-*
Comparative and Physiological Psychology 58, 177-192.
- 34 Wasserman EA, Chatlosh DL & Neunaber DJ (1983). Perception of Causal Relations in Humans -
Factors Affecting Judgments of Response Outcome Contingencies under Free-Operant Procedures.
Learning and Motivation 14, 406-432.
- 35 Wasserman EA, Elek SM, Chatlosh DL & Baker AG (1993). Rating Causal Relations - Role of
Probability in Judgments of Response Outcome Contingency. *Journal of Experimental Psychology-*
Learning Memory and Cognition 19, 174-188.
- 36 Vallee-Tourangeau F & Murphy RA (1999). Action-effect contingency judgment tasks foster
normative causal reasoning. *Proceedings of the Twenty First Annual Conference of the Cognitive*
Science Society, 821-821.
- 37 Cheng PW & Buehner M. in *Oxford Handbook of Thinking and Reasoning* (eds K. J. Holyoak & R. G.
Morrison) Ch. 12, 210-233 (Oxford University Press, 2012).
- 38 Holyoak KJ & Cheng PW (2011). Causal learning and inference as a rational process: the new
synthesis. *Annu Rev Psychol* 62, 135-163.
- 39 Pearl J (2009). *Causality: Models, Reasoning, and Inference*. 2nd edn, Cambridge University Press.
- 40 Steyvers M, Tenenbaum JB, Wagenmakers EJ & Blum B (2003). Inferring causal networks from
observations and interventions. *Cognitive Science* 27, 453-489.
- 41 Kishida KT, Saez I, Lohrenz T, Witcher MR, Laxton AW *et al.* (2016). Subsecond dopamine
fluctuations in human striatum encode superposed error signals about actual and counterfactual
reward. *Proc Natl Acad Sci U S A* 113, 200-205.
- 42 Goussev V (2004). Does the brain implement the Kalman filter? *Behavioral and Brain Sciences* 27,
404-405.
- 43 Morris RW, Dezfouli A, Griffiths KR & Balleine BW (2014). Action-value comparisons in the
dorsolateral prefrontal cortex control choice between goal-directed actions. *Nat Commun* 5, 4390.
- 44 Pouget A, Beck JM, Ma WJ & Latham PE (2013). Probabilistic brains: knowns and unknowns. *Nat*
Neurosci 16, 1170-1178.
- 45 Peirce JW (2007). PsychoPy—Psychophysics software in Python. *Journal of Neuroscience Methods*
162, 8-13.
- 46 Peirce JW (2009). Generating stimuli for neuroscience using PsychoPy. *Frontiers in*
Neuroinformatics 2
- 47 Shanks DR & Dickinson A (1991). Instrumental judgment and performance under variations in
action-outcome contingency and contiguity. *Mem Cognit* 19, 353-360.
- 48 Daw ND. in *Decision making, affect, and learning: Attention and performance XXIII* Vol. 23 (eds M.
R. Delgado, E. A. Phelps, & T. W. Robbins) Ch. 1, 3-38 (Oxford University Press, 2011).
- 49 Mumford JA, Poline JB & Poldrack RA (2015). Orthogonalization of regressors in fMRI models. *PLoS*
One 10, e0126255.

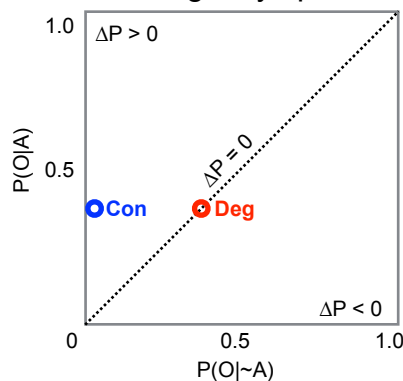
Table 1. Model evidence and comparison scores

	Free parameters	Negative log likelihood	LRT X^2	LRT p	Pseudo- R^2	Relative Bayes Factor ($H_1 - H_0$)	No. favoring H_1/H_0 model
Kalman filter	ν , k , and τ	14,056	$X^2_3 = 7,790$	$< 1E-30$	0.24	1573	22/6
Prediction-error	a , k , and τ	16,343	$X^2_1 = 5,124$	$< 1E-30$	0.13	-506	18/11
MBRL	a , k , and τ	18,853	$X^2_1 = 1,04$	$.7E-15$	0.01	-2959	1/25
Null model (H_0)	τ	15,992	$X^2_1 = 5,826$	$< 1E-30$	0.13		

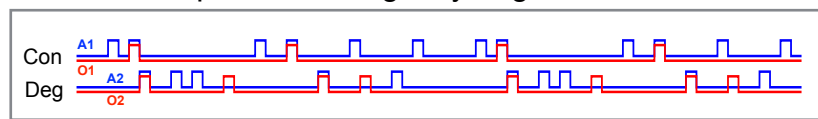
Legend: For each model, optimal model evidence (aggregate negative log likelihood scores), model significant differences from chance (Likelihood ratio test, LRT), model fit (Pseudo- R^2), as well as Bayesian model comparisons among each alternate model relative to the informed model (H_0). ν is prediction uncertainty, α is learning rate, k is the AO delay temporal threshold, and τ is inverse temperature (exploitation/exploration).

Figure 1

A AO contingency space



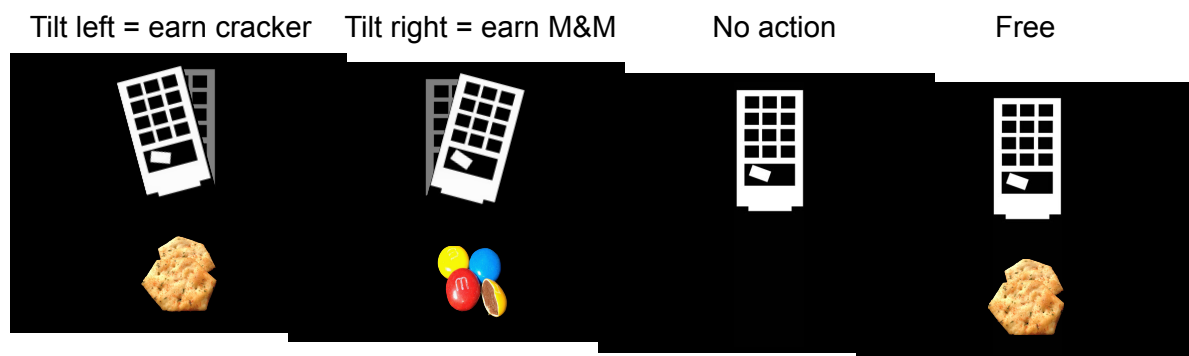
B Outcome-specific contingency degradation



C With signaled non-contingent outcomes



D Experiment 1 Outcome-specific degradation



E Experiment 1 Signaled free outcomes

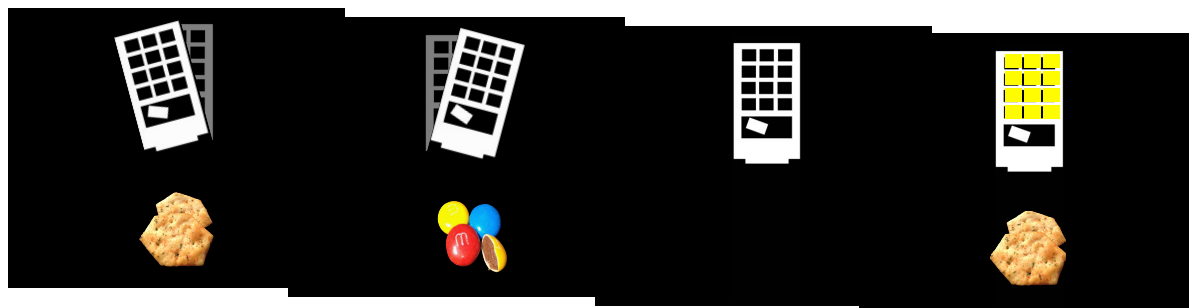


Figure 1. The action-outcome relationship in **A**) contingency space, where $\Delta P = P(O|A) - P(O|\sim A)$, that is, a positive ΔP exists when the conditional probability of an outcome given an action is greater than the probability of the outcome in the absence of the action (Con) while ΔP approaches zero as these conditional probabilities become equal (e.g., Deg); **B**) Outcome-specific degradation schedule where $P(O1|A1) = P(O2|A2)$ while the addition of noncontingent outcomes (O2) produces differences in ΔP (i.e., $\Delta P = 0.4$ and 0 for Con and Deg, respectively); **C**) Signaled schedule where the cue (pink) marks the noncontingent outcomes **D**) The action-outcome relationships presented onscreen (counterbalanced) in the degradation test in the MRI and **E**) the signaled follow-up test after the MRI, where a 1-s visual cue (yellow) indicated the delivery of each noncontingent outcome

Figure 2

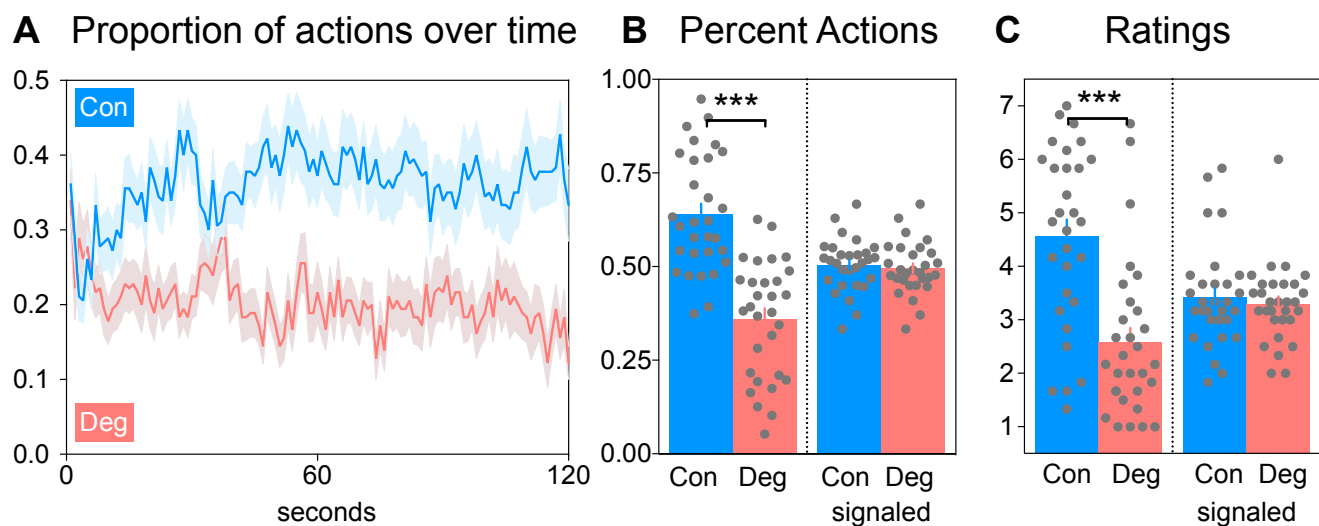


Figure 2. Experiment 1 behavioural results ($N = 30$) **A)** Mean (shaded = SEM) probability of each action over time shows that on average the contingent action was gradually selected over the degraded action **B)** The mean (errorbars SEM) percent of contingent actions was significantly greater than degraded actions when free outcomes were unsignalled, paired t-test $t_{29} = 4.15$, $***p = .0002$. However when free outcomes were signalled then degraded actions were restored, paired t-test $t_{29} = 0.75$, $p = .46$. **C)** Mean (errorbars SEM) causal judgments of the contingent action were greater than the degraded action when free outcomes were unsignalled, paired t-test $t_{29} = 3.94$, $***p < .0004$, and this difference was removed when the free outcomes were signalled, paired t-test $t_{29} = 0.88$, $p = .39$.

Figure 3

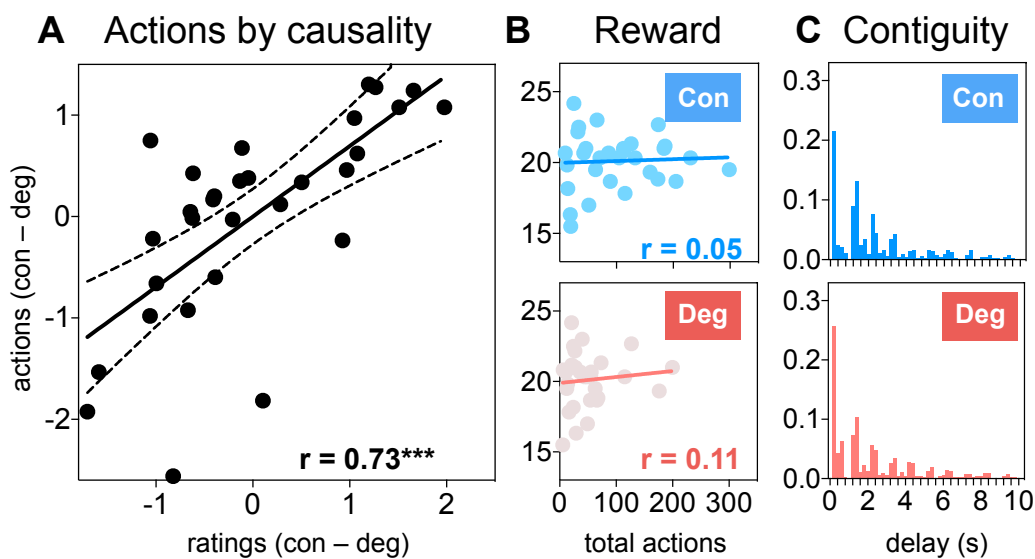


Figure 3. Causality explains action selection better than reward in Experiment 1 ($N = 30$) **A**) The difference (in z-score units) in contingent over degraded actions and causal judgments among participants was correlated, $F_{1,28} = 26.20$, $^{***}p = .00002$ (dotted lines 95% CI). **B**) No correlation existed between the number of contingent actions (blue, $F_{1,28} = 0.07$, $p = .79$), or degraded actions (red, $F_{1,28} = 0.30$, $p = .59$) and the total number of outcomes among participants. **C**) Frequency histogram of the experienced delays between actions and reward shows the distribution was similar for both actions, Kolmogorov-Smirnov $D_{78} = 0.09$, $p = 0.99$

Figure 4

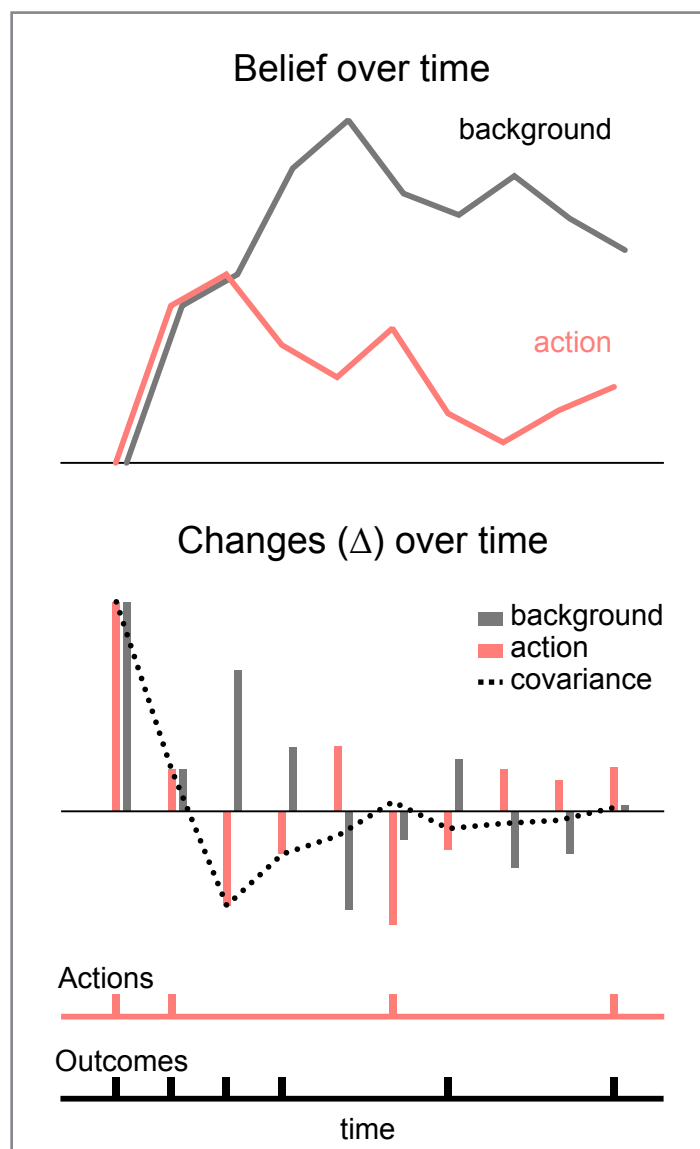


Figure 4. The Kalman algorithm changes beliefs over time according to a prediction-error adjusted for covariance **A)** Under a degradation schedule, belief in the action and background initially increase together as actions co-occur with outcomes. However with free outcomes, the belief in the background diverges from the action. **B)** Changes in the background and action occur in the same direction when covariance between the action and background is positive, but when the covariance is negative then beliefs move in opposite directions. **C)** Action and outcome events over time in this example

Figure 5

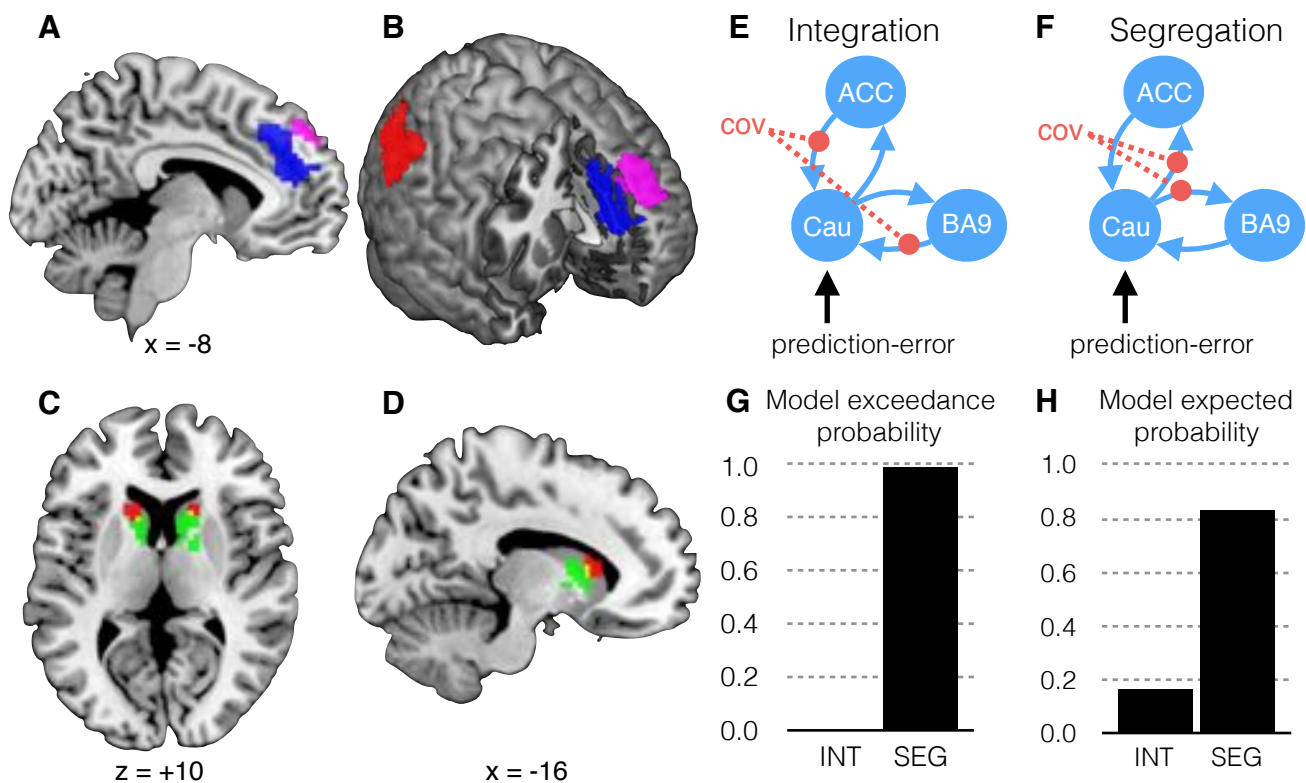


Figure 5 Corticostriatal network for causal learning in Experiment 1, $N = 30$. **A**) Model-derived learning variables were tracked in the medial prefrontal cortex: Model updates to actions (ΔAO) occurred in the medial prefrontal cortex (BA9) (violet voxels FWE cluster $p = .007$). At the same time, model updates to the background (ΔXO) occurred in the dorsal anterior cingulate cortex (ACC), (blue voxels FWE cluster $< .001$); **B**) Cut-away representation showing the spatial relationship of the corticostriatal network, including model covariance in the right posterior parietal cortex (BA40) (red voxels FWE cluster $p = .$). **C**) ROI analysis in the striatum: Red voxels (image threshold $p < .05$ svc) in the anterior caudate tracked the covariance, while green voxels (image threshold $p < .05$ svc) in the caudate body tracked summed prediction-errors. Overlap is indicated in yellow. **D**) Sagittal view of ROI results. **E**) DCM showing the caudate integrates information from separate regions in the mPFC (dACC and BA9), modulated by the covariance between potential causes; **F**) An alternate DCM showing the caudate segregating information in the mPFC, modulated by the covariance between potential causes **G**) The probability the data supports the Segregate model (SEG) is more likely than the Integrate model (INT). **H**) Posterior probability of each model (INT vs SEG) generating the observed data

Figure 6

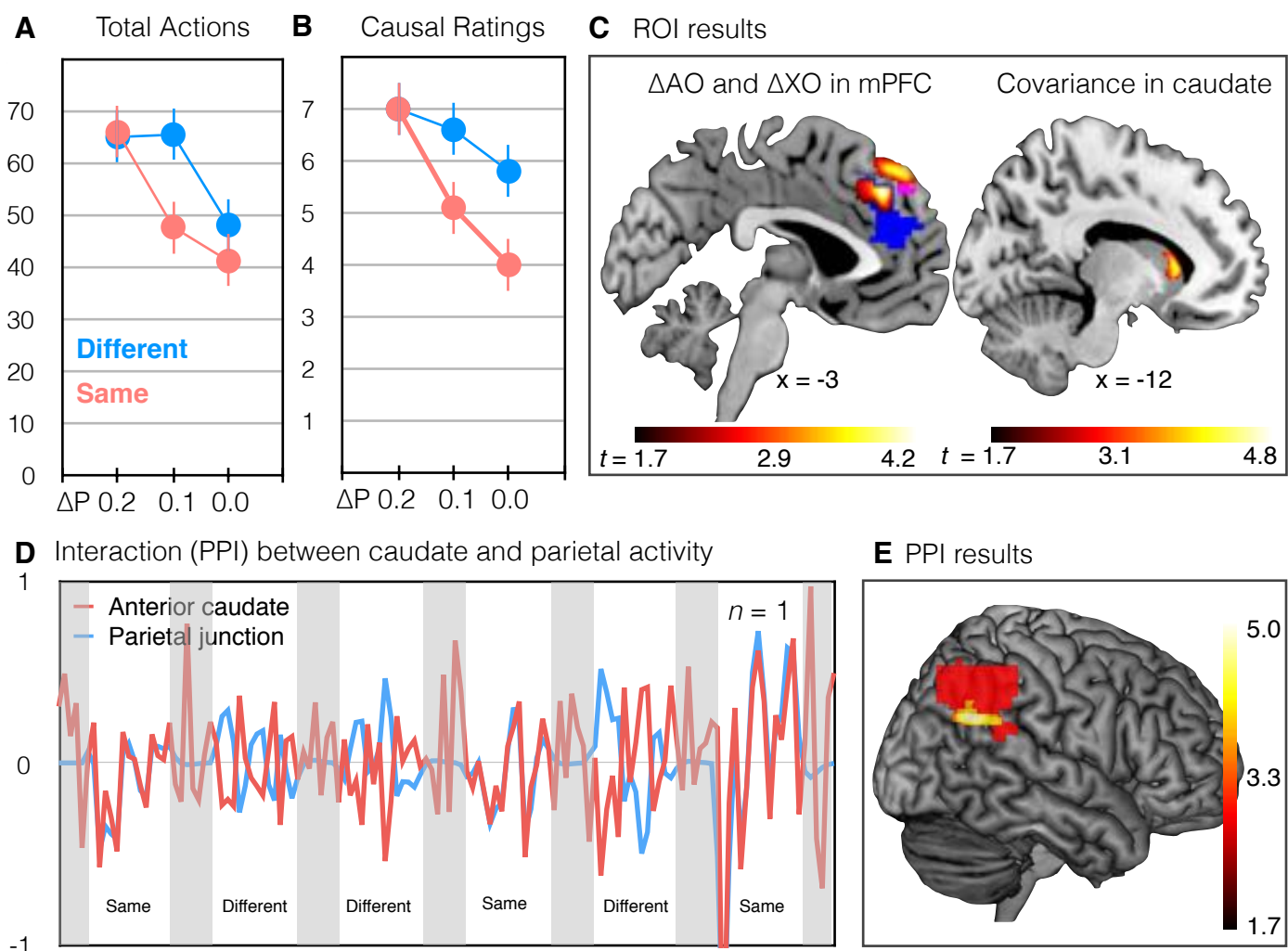
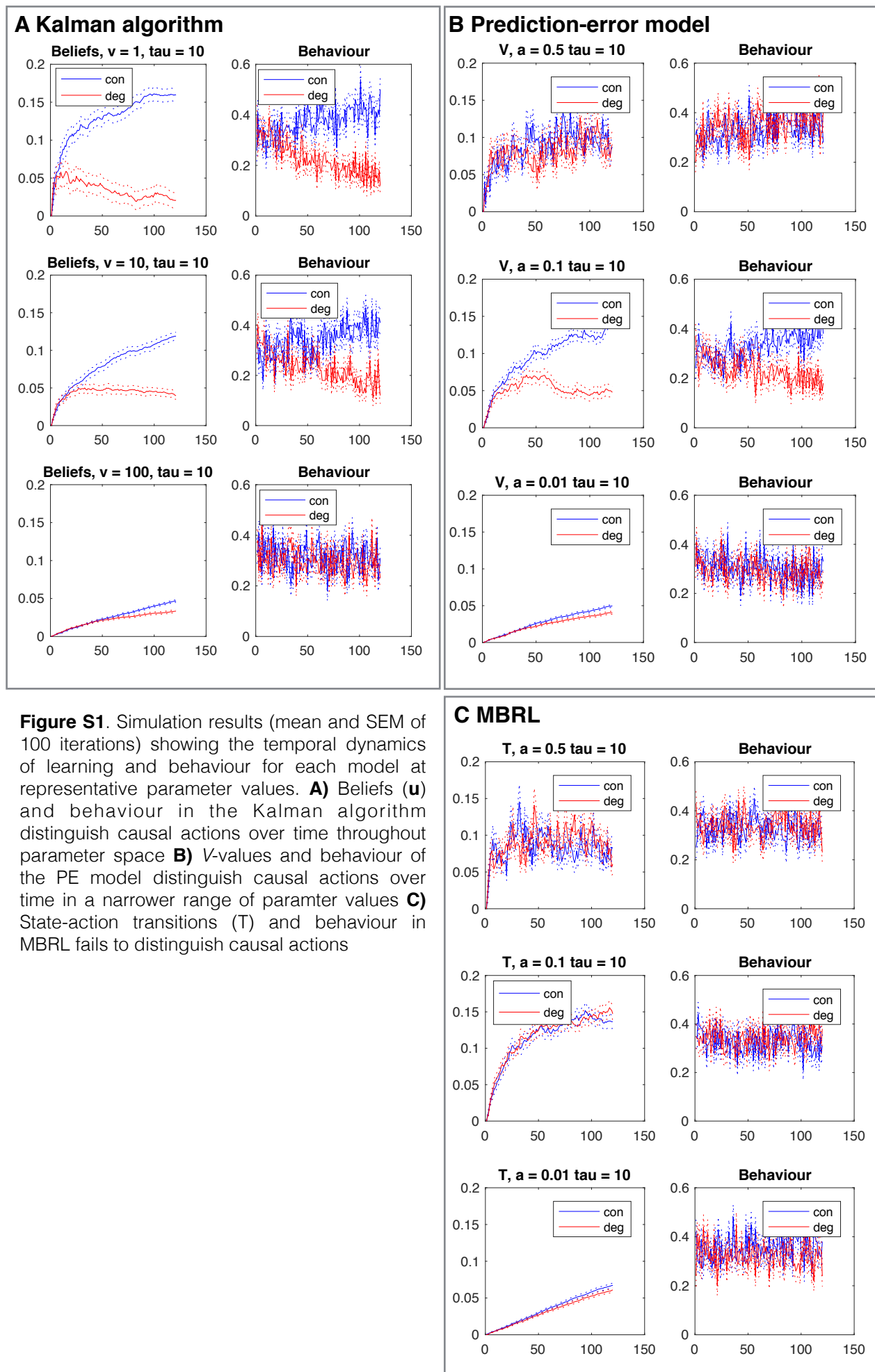


Figure 6. Experiment 2 results ($N = 20$) **A**) Mean (errorbar SEM) total responses were significantly higher when the free outcomes were different from earned outcomes, outcome main effect $F_{1,19} = 12.57$, $p = .02$. This difference decreased as ΔP increased, outcome by contingency interaction $F_{1,5,38} = 7.24$, $p = .005$. **B**) Mean (errorbar SEM) causal judgments were higher when free outcomes were different from earned outcomes ($F_{1,19} = 33.82$, $p = .6E-4$) and this difference decreased as ΔP increased, interaction $F_{2,38} = 10.07$, $p = .002$. **C**) Updates to the action ΔAO occurred in the BA9 fROI (violet, image threshold $p < .001$ svc) while updates to the background ΔXO occurred in the dorsal ACC fROI (blue, image threshold $p < .001$ svc) and the covariance was tracked in the caudate fROI (green, image threshold $p < .001$ svc). **D**) Illustrative results from a single subject showing the caudate and posterior parietal cortex interacted with the causal condition **E**) Right parietal junction activity interacted with caudate activity when noncontingent outcomes were indistinguishable from contingent outcomes (covariance from Experiment 1 shown in red for comparison), image threshold $p < .001$ unc.

Supplementary Figure 1



Supplementary Figure 2

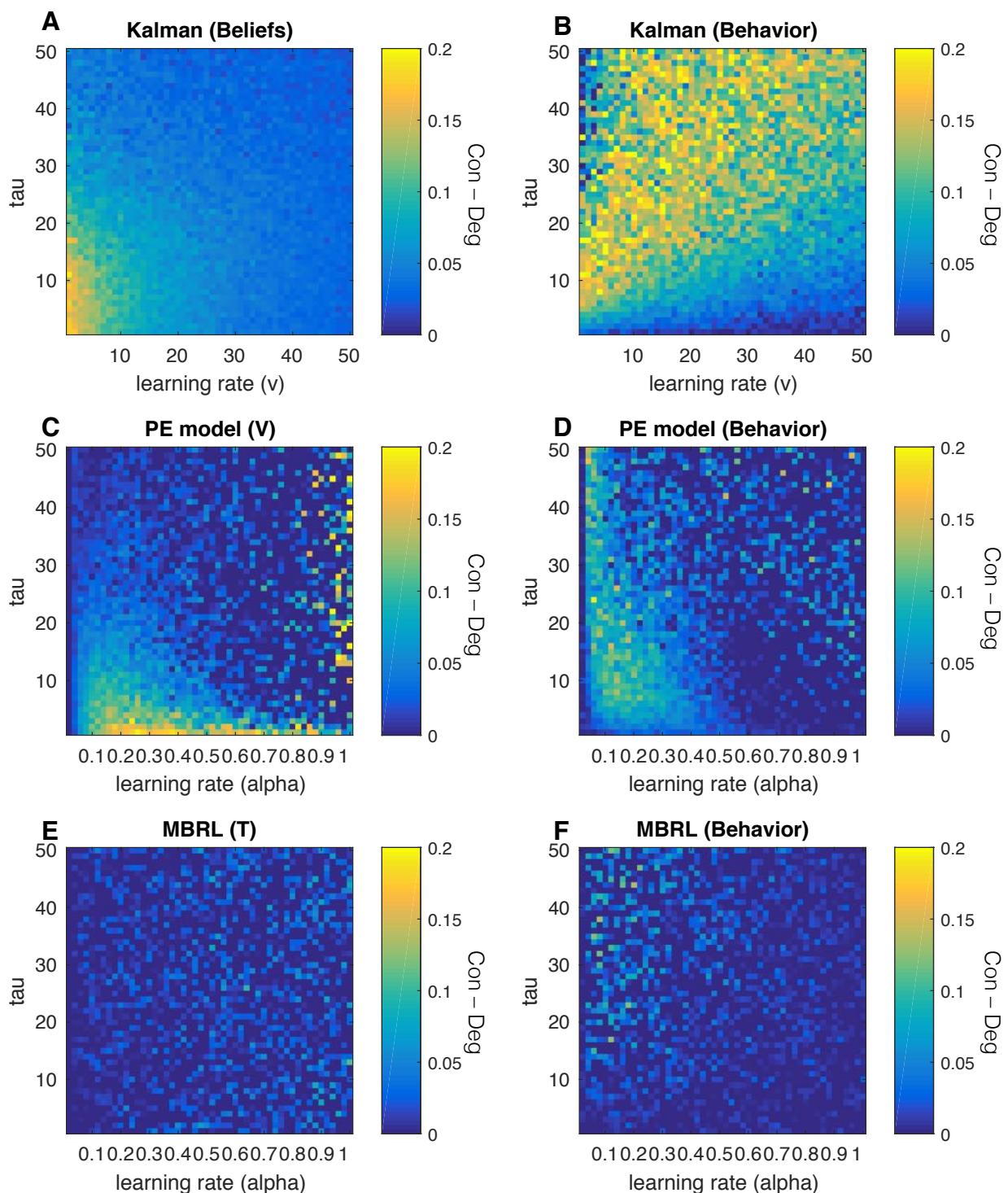


Figure S2. Simulation results showing the parameter space (of learning rate & tau) over which each model distinguishes the causal action during degradation. **A**) Beliefs (\mathbf{u}) in the Kalman algorithm distinguish causal actions throughout parameter space, with the best distinction at low parameter values **B**) Behavior of the Kalman algorithm distinguishes causal actions throughout parameter space **C**) V-values in the PE model distinguish causal actions at low values of tau **D**) Behavior of the PE model partially distinguishes causal actions at low learning rates **E**) State-action transitions (T) in MBRL fail to distinguish causal actions **F**) Behavior of the MBRL fails to distinguish causal actions

Supplementary Figure 3

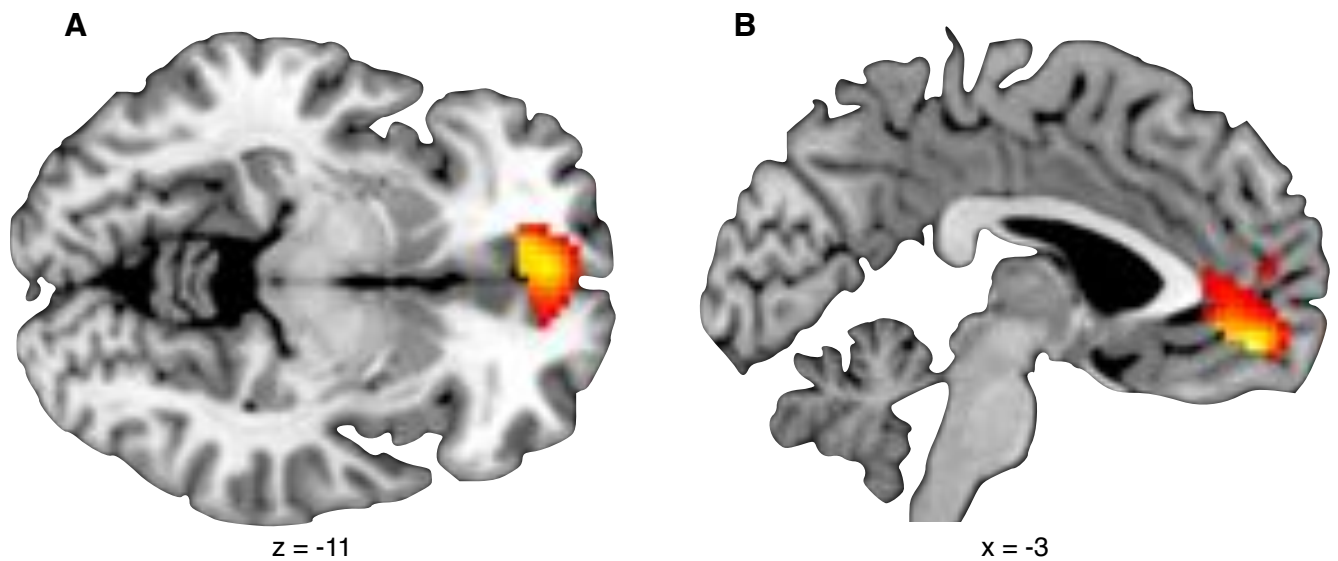


Figure S3 Responses to the alternate AO contingency in the ventromedial prefrontal cortex ($N = 30$), $F_{1,29} = 18.08$, FWE = .011, including **A**) the medial orbitofrontal cortex and **B**) the anterior cingulate. Image thresholded at $p < .001$, uncorrected.

Supplementary Figure 4

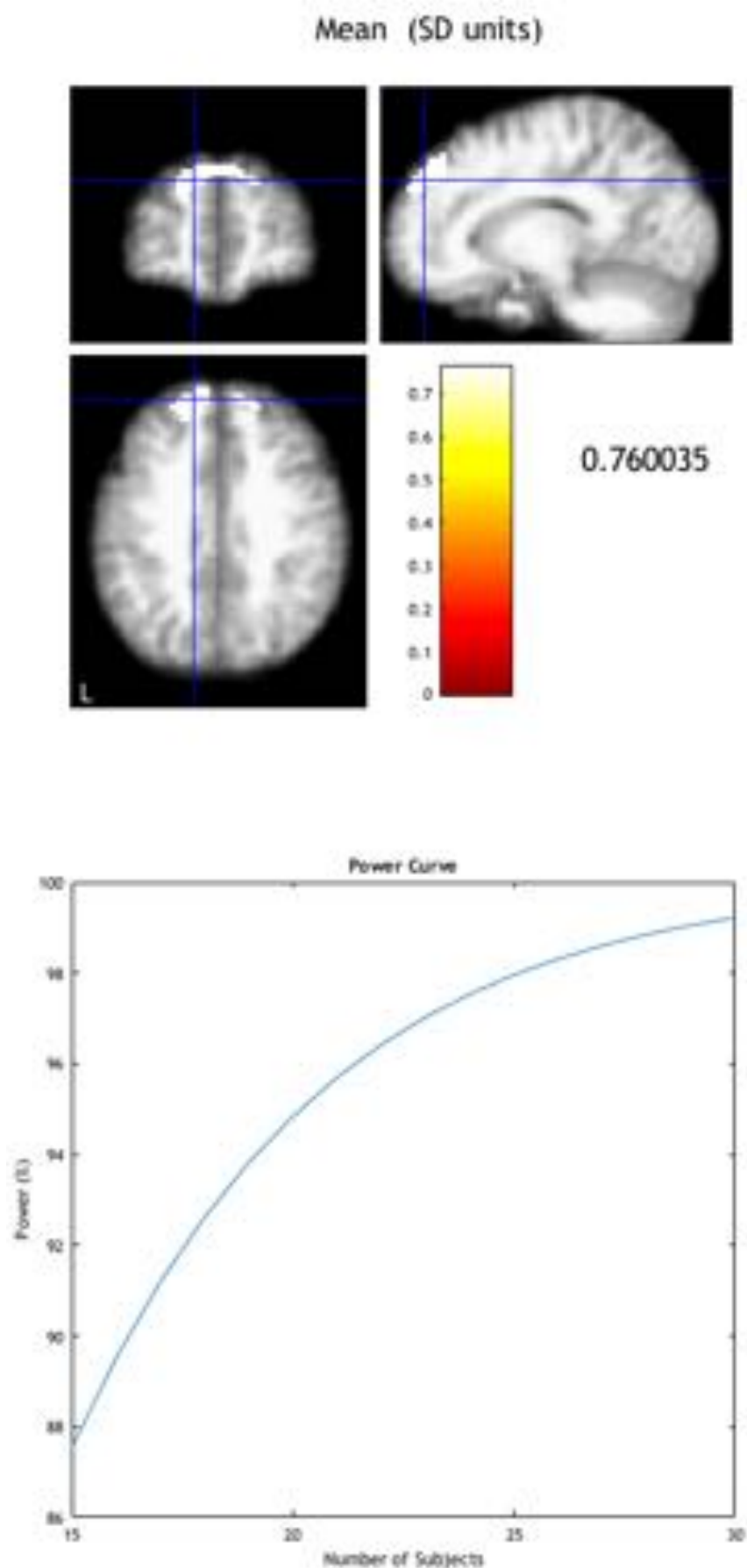


Figure S4. Power analysis of Experiment 1, using the mPFC response during Δ AO, indicated $N > 20$ would achieve 95 percent power