

Computationally-informed insights into anhedonia and treatment by κ -opioid receptor antagonism

Bilal A. Bari^{1,2}, Andrew D. Krystal^{3,4}, Diego A. Pizzagalli^{2*}, Samuel J. Gershman^{5,6*}

¹Department of Psychiatry, Massachusetts General Hospital, Boston, MA, USA

²McLean Hospital, Harvard Medical School, Belmont, MA, USA

³Department of Psychiatry, University of California San Francisco, San Francisco, CA, USA

⁴Departments of Psychiatry and Behavioral Sciences and Radiology, Duke University School of Medicine, Durham, NC, USA

⁵Department of Psychology and Center for Brain Science, Harvard University, Cambridge, MA, USA

⁶Center for Brains, Minds, and Machines, Massachusetts Institute of Technology, Cambridge, MA, USA

*Denotes co-senior authorship

Correspondence: Bilal A. Bari, bbari@mgh.harvard.edu

1 Abstract

2 Anhedonia, the loss of pleasure, is prevalent and impairing. Parsing its computational basis
3 promises to explain its transdiagnostic character. We argue that one manifestation of anhedonia—
4 reward insensitivity—may be linked to limited memory capacity. Further, the need to economize
5 on limited capacity engenders a perseverative bias towards frequently chosen actions. Anhedonia
6 may also be linked with deviations from optimal perseveration for a given memory capacity, a
7 pattern that causes *inefficiency* because it results in less reward for the same memory cost. To test
8 these hypotheses, we perform secondary analysis of a randomized controlled trial testing κ -opioid
9 receptor (KOR) antagonism for anhedonia, as well as analyses of three other datasets. We find
10 that anhedonia is associated with deficits in efficiency but not memory, whereas KOR antagonism
11 (which likely elevates tonic dopamine) increases memory and efficiency. KOR antagonism therefore
12 has distinct cognitive effects, only one related to anhedonia.

13 Introduction

14 Anhedonia, the loss of pleasure or lack of reactivity to pleasurable stimuli, is observed in many
15 psychiatric illnesses, including major depressive disorder, bipolar disorder, schizophrenia, anx-
16 iety disorders, post-traumatic stress disorder, substance use disorders, autism, and attention-
17 deficit/hyperactivity disorder [1, 2, 3, 4, 5, 6, 7, 8, 9]. The transdiagnostic character of anhedonia
18 suggests a common mechanism across disorders. The most systematic attempts to formalize this
19 common mechanism have utilized concepts from reinforcement learning [10]. Early models posited
20 that anhedonia corresponds to a reduction in reward sensitivity [11, 12], but the predictions of
21 these models have not been consistently validated, suggesting a more complex picture [13]. Here,
22 we argue that one neglected source of complexity is the interplay between reward sensitivity and
23 cognitive capacity limits.

24 In reinforcement learning theory, states (e.g., stimuli, context) are mapped to actions by a
25 learned policy. The amount of memory needed to store a policy is dictated by the mutual infor-
26 mation between states and actions; any physical system (such as the brain) has a limited memory
27 capacity. One implication of limited capacity is reward insensitivity and, thus, some aspects of
28 anhedonia may arise from cognitive resource limitations.

29 Under capacity limits, policies must be *compressed* by discarding some state information [14,
30 15, 16]. This results in the tendency to reuse frequently chosen actions across multiple states—a
31 form of *perseveration*, the tendency to repeat actions independently of their reinforcement history.
32 The theory of policy compression is normative: it specifies an optimal level of perseveration for
33 a given capacity limit. Empirically, compression strategies may differ, with some policies yielding
34 more reward than others for the same cost. We refer to deviations from optimal perseveration
35 as *inefficiency* because it results in a suboptimal use of finite memory (less reward for the same
36 memory utilization). This phenotype is conceptually distinct from capacity, and can be measured
37 separately. We argue here that capacity and efficiency may be key phenotypes for understanding
38 cognitive disturbances in anhedonia. We show that these can be estimated from behavioral data
39 on a widely used behavioral assay, the Probabilistic Reward Task (PRT), and that they reveal new
40 aspects of anhedonia that would otherwise have been invisible.

41 We also address the underlying neural mechanisms and treatment implications. Our previous
42 work suggested that tonic dopamine should determine the allocation of cognitive resources for
43 task performance based on reward rate [17, 18]. Reduction in tonic dopamine should therefore
44 produce insensitivity of task performance to reward rate [19]. It stands to reason that increasing
45 tonic dopamine should increase reward sensitivity. We demonstrate that this is consistent with the
46 effects of κ opioid receptor (KOR) antagonism, which elevates tonic dopamine [20, 21, 22, 23, 24].
47 We find that efficiency also increases, suggesting that tonic dopamine may not only determine the
48 amount of resources available but also the efficiency of their allocation. Mechanistically, this might
49 be implemented through dopamine-dependent changes in learning rate for perseveration. Finally,
50 we find that anhedonia is associated with changes in efficiency but not memory, highlighting the
51 clinical utility of distinguishing these computational phenotypes.

52 Results

53 Policy complexity and efficiency in anhedonia after κ -opioid receptor antagonism

54 We performed a secondary analysis of an 8-week, multicenter, placebo-controlled, double-blind,
55 randomized trial to test the effects of KOR antagonism for anhedonia (Figure 1A) [25, 26]. Because
56 this trial identified a significant treatment effect of KOR antagonism for anhedonia (as measured by
57 the Snaith-Hamilton Pleasure Scale, SHAPS), we sought to understand the cognitive basis of this
58 improvement. We analyzed a total of 55 participants (KOR antagonist group: $N = 24$; placebo
59 group: $N = 31$) who completed both a baseline and post-treatment Probabilistic Reward Task
60 (PRT). Owing to previously-reported baseline differences in anhedonia between the two groups

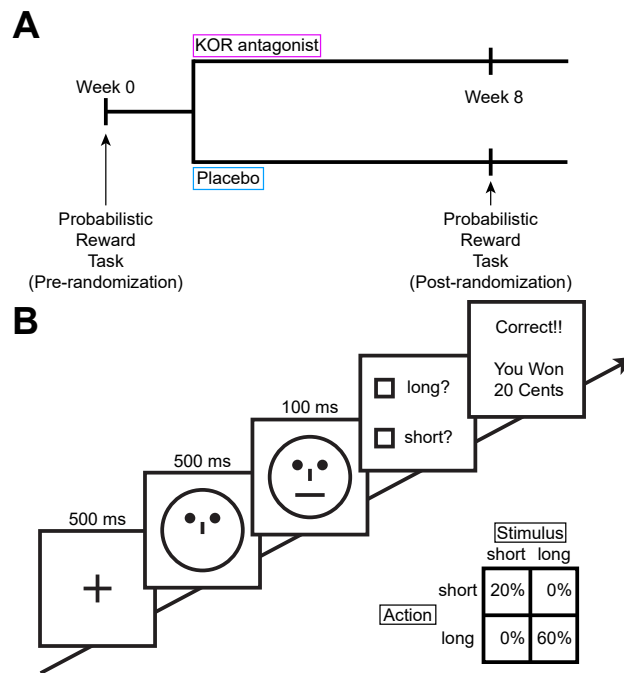


Figure 1: Trial and task design.

A) Participants were randomized to 8 weeks of placebo ($N = 31$) or a KOR antagonist ($N = 24$) and completed the PRT at baseline and at week 8.

B) On each trial of the PRT, participants fixated on a cross, followed by the presentation of a face without a mouth, followed by either a short (11.5mm) or long (13mm) mouth in the face. Participants responded by pressing one of two keyboard keys and completed 200 trials in two blocks of 100 trials. The bottom right shows an example reward schedule where the long stimulus is rewarded more often than the short stimulus. The mapping between response, stimulus, and reward was counterbalanced between participants.

61 (mean SHAPS \pm SD: placebo 33.03 ± 5.54 ; KOR 37.29 ± 8.89 , $p = 0.0338$), we analyzed the
 62 pre-treatment groups separately.

63 The PRT is a reward-based decision making task that asks participants to discriminate two
 64 similar stimuli (Figure 1B) [27, 28]. Unbeknownst to participants, one of the two stimuli yields
 65 reward more often than the other when correctly identified. According to the theory of policy
 66 compression [16], performance in this task (average reward) depends on the amount of information
 67 participants encode about the underlying state (i.e., the stimulus identity), quantified by the mu-
 68 tual information between states and actions—a participant’s *policy complexity*. Each participant is
 69 assumed to have a capacity limit (upper bound on policy complexity), which delimits their achiev-
 70 able performance. If participants maximally utilize their capacity, their average reward should
 71 fall along an optimal reward-complexity frontier, as shown in Figure 2A,B. In the PRT, maximal
 72 reward can be obtained at a policy complexity of 1 bit, corresponding to a policy that perfectly
 73 discriminates the two stimuli. At the other extreme, a subject with no capacity will generate a
 74 policy that ignores the stimuli entirely. Participant policies tend to lie close to the optimal frontier,
 75 indicating that they are utilizing most of their capacity. At the low end of the policy complexity
 76 range, participant policies fall off the optimal frontier (Figure 2F,G), indicating under-utilization
 77 of resources (inefficiency)—a pattern also observed in previous studies [15, 29].

78 At 8 weeks, placebo treatment resulted in a decrease in both policy complexity and reward,
 79 while KOR antagonism yielded an increase in both (Figure 2C). This resulted in significant between-
 80 group differences for both policy complexity (Figure 2D; mean change in policy complexity (post-
 81 treatment minus baseline) \pm SEM: placebo, -0.0245 ± 0.0141 ; KOR, 0.0281 ± 0.0211 , $p = 0.0362$)
 82 and reward (Figure 2E; mean change in reward (post-treatment minus baseline) \pm SEM: placebo,

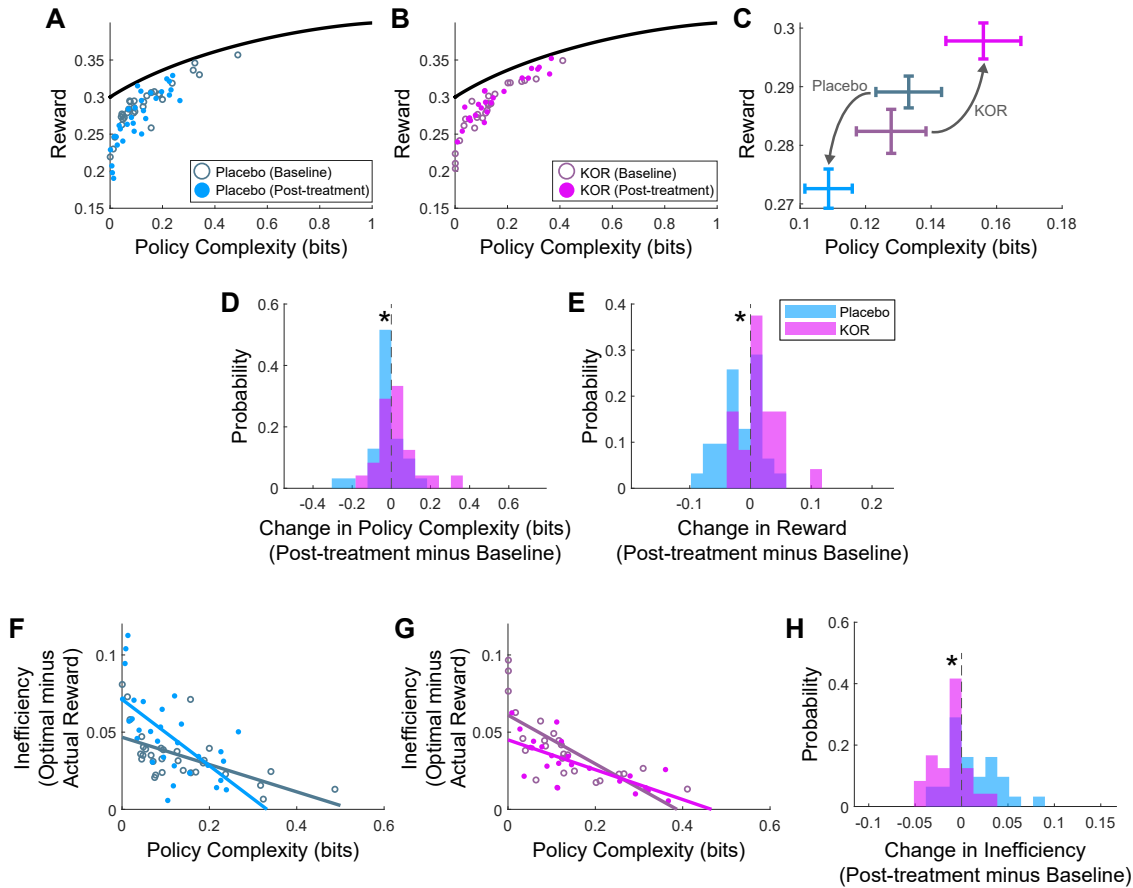


Figure 2: Changes in policy complexity and efficiency as a function of KOR antagonism.

A,B) Reward-complexity relationship for the placebo and KOR groups at baseline and post-treatment. The black line shows the reward-complexity frontier, which indicates the optimal reward as a function of policy complexity.

C) Mean \pm SEM reward-complexity relationship as a function of treatment (placebo or KOR antagonism) and time (baseline or post-treatment).

D) Change in policy complexity (post-treatment minus baseline) as a function of treatment.

E) Change in reward (post-treatment minus baseline) as a function of treatment.

F,G) Relationship between inefficiency and complexity for the placebo and KOR groups. Overlaid lines are from a linear mixed-effects model fitting inefficiency as a function of policy complexity, treatment, and time.

H) Change in inefficiency (post-treatment minus baseline) as a function of treatment.

83 $-0.0165 \pm 5.61 \times 10^{-3}$; KOR, $0.0154 \pm 6.53 \times 10^{-3}$, $p = 4.81 \times 10^{-4}$). Following treatment, the
84 KOR group also became significantly more efficient compared to the placebo group (Figure 2H;
85 mean change in inefficiency (post-treatment minus baseline) \pm SEM: placebo, $0.0130 \pm 4.80 \times 10^{-3}$;
86 KOR, $-0.0109 \pm 4.04 \times 10^{-3}$, $p = 5.68 \times 10^{-4}$). Thus, KOR antagonism increases average reward
87 through a combination of increasing both policy complexity and efficiency.

88 Policy compression makes the additional prediction that more complex policies should result
89 in slower response times, since the brain must inspect more bits to find a coded state [16, 18, 30].
90 Indeed, we found that KOR antagonism, relative to placebo, slowed participants down (mean
91 change in response times (post-treatment minus baseline) \pm SEM: placebo, $-59.3\text{ms} \pm 23.4$; KOR,
92 $13.6\text{ms} \pm 20.4$, $p = 0.0274$).

93 To better understand how KOR treatment changed the relationship between inefficiency and
94 policy complexity, we fit a linear mixed effects model predicting inefficiency as a function of policy
95 complexity, treatment, and time. We identified two relevant effects: a significant treatment \times
96 time interaction (coefficient = -0.0405 , $p = 4.23 \times 10^{-5}$), which has the effect of lowering the
97 intercept, and a significant policy complexity \times treatment \times time interaction (coefficient = 0.187 ,
98 $p = 1.76 \times 10^{-3}$), which has the effect of increasing the slope. The combination of the change in
99 intercept and slope has the net effect of increasing efficiency as a function of policy complexity,
100 revealing that KOR treatment increases efficiency independent of its changes to complexity. We
101 will develop this insight further with our reinforcement learning modeling. Overall, these results
102 suggest two orthogonal effects of KOR treatment: increases in complexity and increases in efficiency.
103 Stated another way, participants gain increased cognitive resources *and* make better use of those
104 resources.

105 Reinforcement learning model of KOR antagonism

106 We developed a cost-sensitive reinforcement learning model to gain insight into how KOR antag-
107 onism affects decision making. We adapted a Q -learning model, ubiquitous in the reinforcement
108 learning literature [31]. This model estimates the expected reward associated with each action for
109 each stimulus (called Q -values) and updates these estimates by learning from the outcome (pres-
110 ence or absence of reward). Since the optimal policy under policy compression contains a marginal
111 action probability term to engender perseveration (state-independent actions), we augmented our
112 model with a marginal action probability term that was similarly estimated on a trial-by-trial ba-
113 sis. Our model contained a reward learning rate, α_{learn} , to govern the learning of action values, a
114 perseveration learning rate, α_{persev} , to govern the learning of the marginal action probability, and a
115 reward sensitivity parameter, β , that determines the balance between action values and persevera-
116 tion in driving behavior. The β parameter is linked to capacity, where higher capacity is associated
117 with higher values of β . Given the structure of our model, β is equivalent to a parameter scaling
118 reward magnitude, as has been posited in anhedonia [12].

119 To model the effects of treatment, we allowed KOR and placebo to scale these parameters. Based
120 on formal model comparison (Extended Data Table 1), we selected a model that separately scaled
121 the perseveration learning rate, α_{persev} , and the reward sensitivity, β , as a function of treatment.
122 We confirmed that our model could recover α_{persev} and β , the parameters of interest (Extended
123 Data Table 2). To provide confidence in the ability of our model to capture key characteristics
124 of the data, we first fit the model to participant data and then had the model perform the PRT
125 (using the parameter estimates for each participant) to generate a synthetic dataset (Extended
126 Data Figure 1). This simulated dataset captured all key features of our data (see Supplementary
127 information).

128 Having confirmed that our model could generate realistic data and recover parameters of inter-
129 est, we turned our attention to parameter estimates to better understand how treatment affected
130 decision making. We found that placebo and KOR antagonism scaled the perseveration learning
131 rate, α_{persev} in opposite directions (Figure 3A; posterior 95% credible interval; placebo, -2.96 to
132 -0.82 ; KOR, 0.61 to 1.96). The difference between KOR antagonism and placebo corresponds to

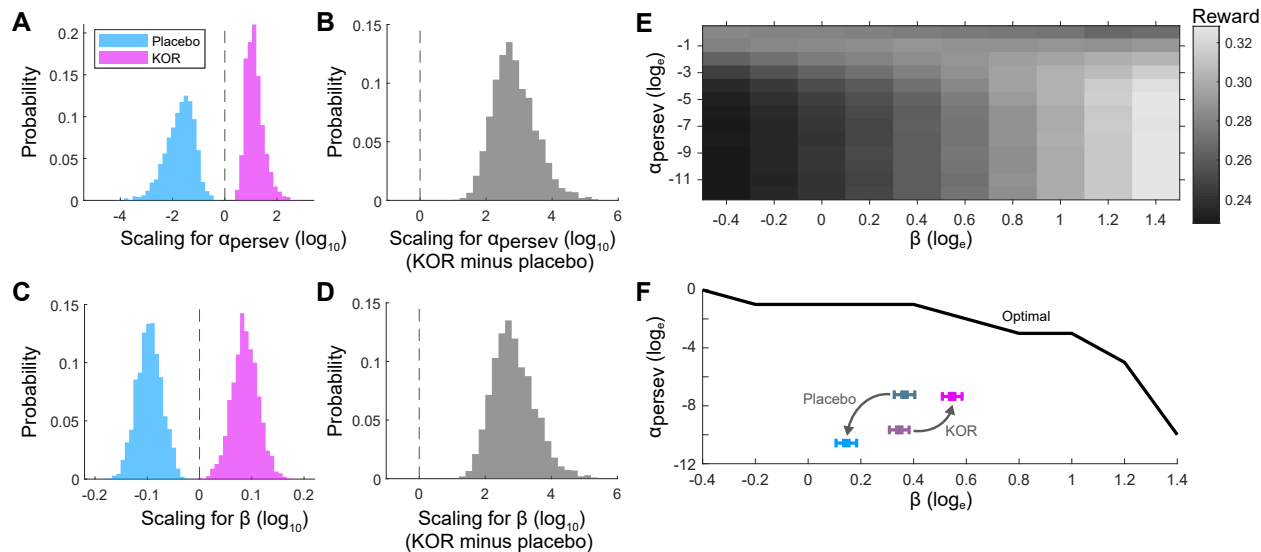


Figure 3: Scaling of reinforcement learning parameters as a function of KOR antagonism.

- A) Posterior distribution of parameter values for scaling of α_{persev} as a function of treatment. Scaling is multiplicative, where values greater than 0 indicate that treatment increases the parameter value, whereas values less than 0 indicate that treatment decreases the parameter value.
- B) Posterior distribution of treatment effect for scaling α_{persev} , estimated as the difference in scaling between KOR and placebo.
- C) Posterior distribution of parameter values for scaling of β as a function of treatment.
- D) Posterior distribution of treatment effect for scaling β .
- E) Heatmap showing mean reward obtained as a function of α_{persev} and β .
- F) Effect of treatment in parameter space. Black line shows the optimal α_{persev} for each value of β .

133 the net effect of treatment on α_{persev} , which was positive and excluded 0, showing that treat-
 134 ment increases perseveration (Figure 3B; difference in posterior 95% credible interval (KOR minus
 135 placebo), 1.77 to 4.32). We similarly found that placebo and KOR antagonism scaled the reward
 136 sensitivity, β , in opposite directions (Figure 3C; posterior 95% credible interval; placebo, -0.143 to
 137 -0.050; KOR, 0.037 to 0.138), with a treatment effect that was positive and excluded 0 (Figure 3D;
 138 difference in posterior 95% credible interval (KOR minus placebo), 0.114 to 0.254).

139 To gain insight into how scaling these parameters affects decision making, we simulated datasets
 140 where we only changed parameters of interest (Extended Data Figure 1; Extended Data Table 3).
 141 Increasing only α_{persev} produces an increase in efficiency and a small decrease in policy complexity.
 142 The increase in efficiency manifests as a change in the intercept, but not the slope, of the relationship
 143 between inefficiency and policy complexity. Increasing only β produces a relatively large increase
 144 in policy complexity, which is consistent with the theoretical link between larger β and increased
 145 capacity. It also produces an increase in efficiency for low-complexity policies. Increasing both
 146 α_{persev} and β , like we find for KOR antagonism, produces both an increase in policy complexity
 147 and an increase in efficiency. The increase in efficiency manifests as a change in both the intercept
 148 (decrease) and the slope (increase) of the relationship between inefficiency and policy complexity,
 149 like our empirical findings.

150 We gained insight into the relationship between KOR antagonism and optimal behavior by
 151 visualizing the relationship between α_{persev} , β , and reward, while holding α_{learn} fixed (Figure 3E).
 152 As β increases, for the optimal α_{persev} , the net reward obtainable also increases, consistent with
 153 our theory linking higher β to higher capacity and higher capacity to greater reward. We also find
 154 that increasing perseverative learning is most beneficial at lower values of β (i.e., lower capacity),
 155 consistent with the idea that perseveration is increasingly optimal as subjects become more resource

156 limited. In Figure 3F, we can see that the effect of KOR antagonism is to shift both α_{persev} and β
157 closer to an optimal regime. A notable finding is the increased α_{persev} at baseline for the placebo
158 group relative to the KOR group. This is consistent with the baseline difference in SHAPS between
159 these groups, with the placebo group having lower SHAPS: the larger α_{persev} estimates for this
160 group is closer to the optimal regime and is consistent with less severe anhedonia.

161 Policy complexity and efficiency as a function of hedonic tone

162 Because the original study identified a significant improvement in the SHAPS following KOR an-
163 tagonism [25], we sought to identify which mechanism—increased policy complexity, increased
164 efficiency, or both—is associated with anhedonia. We first examined the relationship between he-
165 donic tone and reward learning in a non-clinical population. We recruited 100 participants from
166 Amazon Mechanical Turk and implemented a version of the PRT suitable for online delivery [32].
167 Participants completed the SHAPS and reported a wide range of scores (mean SHAPS \pm SD: 11.45
168 \pm 6.54, range 0 to 36). We show the reward-complexity relationship in Figure 4A. For visualization
169 purposes only, we perform a median split of participants on the basis of SHAPS.

170 Unlike the effects of KOR antagonism, we found that SHAPS did not predict policy complexity
171 (coefficient = -5.24×10^{-3} , $p = 0.241$). We did, however, identify a significant relationship
172 with inefficiency. We fit a linear regression predicting inefficiency as a function of SHAPS and
173 policy complexity and identified a significant intercept change (coefficient for effect of SHAPS =
174 9.55×10^{-3} , $p = 6.54 \times 10^{-3}$) but not a significant slope change (coefficient for SHAPS \times policy
175 complexity interaction = -0.0182 , $p = 0.394$). Given our simulations exploring the effects of
176 changing parameters (Extended Data Figure 1), a change of intercept without a change of slope is
177 consistent with hedonic tone affecting perseveration (α_{persev}) and not capacity (β).

178 We reanalyzed two prior PRT datasets and found similar effects on the relationship between
179 inefficiency and policy complexity (Extended Data Figure 2). The first was a transdiagnostic
180 sample of patients (control group: $N = 25$; clinical group: $N = 41$, 18 with bipolar disorder, 23
181 with major depressive disorder) [33, 34]. These groups differed significantly in baseline anhedonia
182 (mean anhedonic Beck Depression Inventory-II subscore \pm SD: control, 0.72 ± 1.02 ; clinical, 5.22
183 ± 3.78 , $p = 2.22 \times 10^{-7}$; mean Mood and Anxiety Symptom Questionnaire-Anhedonic Depression
184 subscale \pm SD: control, 51.5 ± 12.6 ; clinical, 77.1 ± 19.3 , $p = 1.45 \times 10^{-7}$). Consistent with
185 differences in anhedonia, when we analyzed inefficiency as a function of policy complexity and
186 group, we identified a significant intercept difference (coefficient for clinical group = 5.29×10^{-3} ,
187 $p = 0.022$) without a concurrent slope difference (coefficient for policy complexity \times clinical group
188 interaction = -1.71×10^{-3} , $p = 0.443$). We additionally found no difference in policy complexity
189 between the two groups (mean policy complexity \pm SEM: control, 0.371 ± 0.043 ; clinical, $0.333 \pm$
190 0.024 , $p = 0.412$).

191 The second dataset we analyzed was a test of a longstanding hypothesis relating reduced
192 dopamine to anhedonia [35, 36]. In this double-blinded study, participants received either placebo
193 or low-dose pramipexole—thought to reduce phasic dopamine release—and performed the PRT
194 (placebo group: 13; pramipexole group: 11) [37]. When we analyzed inefficiency as a function of
195 policy complexity and treatment, we identified a significant intercept effect (coefficient for treat-
196 ment = 7.82×10^{-3} , $p = 0.043$) without a significant slope effect (coefficient for policy complexity
197 \times treatment = -1.90×10^{-3} , $p = 0.615$). We also found no difference in policy complexity as
198 a function of treatment (mean policy complexity: placebo, 0.297 ± 0.043 ; pramipexole, $0.319 \pm$
199 0.057 , $p = 0.757$).

200 Reinforcement learning model of hedonic tone

201 We next fit a reinforcement learning model similar to the one we used for the KOR dataset, except
202 now we allowed α_{persev} and β to scale as a function of SHAPS. We found that increases in SHAPS
203 were associated with less perseveration (Figure 4C; posterior 95% credible interval: -0.739 to -

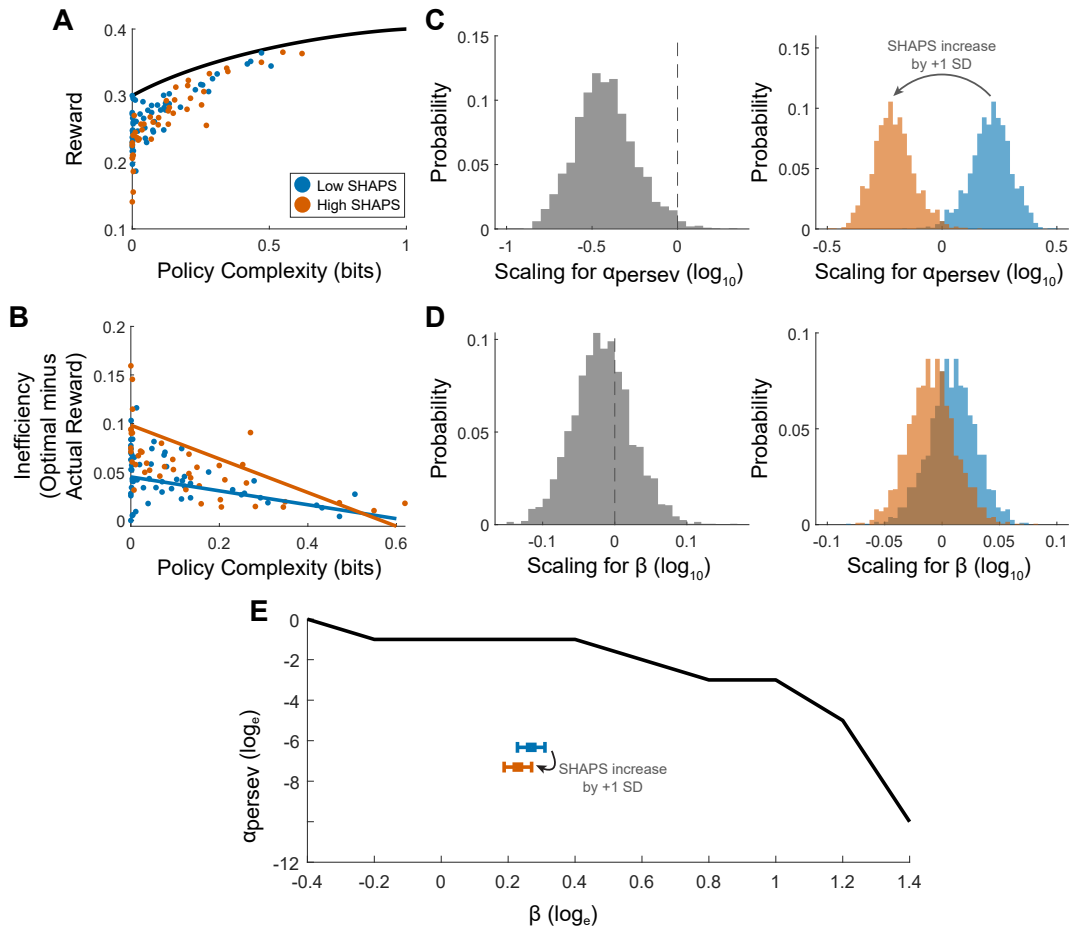


Figure 4: Changes in complexity, efficiency, and reinforcement learning parameters as a function of hedonic tone.

A) Reward-complexity tradeoff as a function of hedonic tone. For illustration only, participants are median split on the basis of SHAPS scores into ‘Low SHAPS’ (low anhedonia) and ‘High SHAPS’ (high anhedonia).

B) Inefficiency-complexity relationship as a function of hedonic tone. For illustration only, the color lines are regression fits denoting extremes of SHAPS in our dataset (blue is lowest SHAPS = 0, orange is highest SHAPS = 36).

C) Left: Posterior distribution of parameter values for scaling of α_{persev} as a function of SHAPS. Right: An example demonstrating scaling for an increase in SHAPS of 1 SD (from 0.5 SD below the mean (blue) to 0.5 SD above the mean (orange)).

D) Left: Posterior distribution of parameter values for scaling of β as a function of SHAPS. Right: Scaling for the same increase in SHAPS.

E) Effect of variation in SHAPS in parameter space. Black line shows the optimal α_{persev} for each value of β .

204 0.046). In contrast, anhedonia had no effect on modulating β , in contrast to KOR antagonism
205 (Figure 4D; posterior 95% credible interval: -0.097 to 0.066). In parameter space, the net effect
206 of an increase in SHAPS is to move participants away from an optimal regime (Figure 4E). Taken
207 together, these data support the notion that hedonic tone spans the axis of efficiency, not capacity.

208 Discussion

209 We leveraged a theory of resource-limited reinforcement learning to shed light on the cognitive
210 structure of anhedonia. Building on prior work demonstrating impairments in reward sensitivity,
211 we decomposed these impairments into separate effects of policy complexity (state-dependence of
212 an action policy) and efficiency (utilization of cognitive resources). We found that KOR antagonism
213 affected both of these measures, whereas anhedonia is associated only with reduced efficiency.

214 The finding that anhedonia is not associated with reduced complexity is surprising, in part, be-
215 cause complexity determines reward sensitivity, and reward insensitivity appears to be the cardinal
216 feature of anhedonia (though see [13] for more nuance). There are a number of explanations for this
217 apparent disconnect. One is that anhedonia may be more psychologically related to the concept
218 of ‘liking,’ the pleasure associated with reward, rather than ‘wanting,’ the motivation furnished by
219 reward learning [38], both of which are relevant for anhedonia. In our paradigm, reward sensitivity
220 is related to ‘wanting,’ which would render the PRT an inappropriate assay to measure deficits
221 in ‘liking.’ Further, the SHAPS is not designed to disambiguate these different aspects of reward
222 processing, but newer scales such as the Dimensional Anhedonia Rating Scale [39], the Temporal
223 Experience of Pleasure Scale [40], and the Positive Valence Systems Scale [41] provide insight into
224 the multidimensional nature of anhedonia.

225 It is also plausible that anhedonia might be a consequence of reduced reward *learning*, not
226 reduced reward sensitivity [42]. A limitation of our study is that our model could not recover
227 the reward learning rate (Extended Data Figure 1). It is worth noting that our findings seem at
228 odds with an influential reinforcement learning account of anhedonia implicating decreased reward
229 sensitivity as the key causal variable [12]. Interestingly, the parameterization of that model links
230 increased reward sensitivity with increased perseveration. Our model orthogonalizes reward sensi-
231 tivity from perseveration, suggesting what was previously identified as blunted reward sensitivity
232 may have been impaired perseverative learning (see Supplementary materials).

233 Under our computational framework, perseveration is closely related to habits, since habits can
234 be similarly thought of as state-independent actions within a particular context [43]. A prediction
235 of our findings is that anhedonia may not only manifest as a deficit in perseveration, but may
236 also manifest as a deficit in habit formation. Intriguingly, recent work on the origin of habits has
237 revealed that they are largely divorced from reward [44]. If true, this would highlight a cognitive
238 deficit in anhedonia unrelated to reward processing. Altogether, our findings motivate a future
239 research program studying habit formation in anhedonia, both important for better understanding
240 this symptom and because it may form the basis of clinically relevant behavioral interventions.

241 The aspect of KOR antagonism which appears to be unrelated to anhedonia (increased policy
242 complexity) suggests relevant clinical utility outside of anhedonia. As one example, we hypothe-
243 size KOR antagonism may prove beneficial in treating cognitive deficits in chronic schizophrenia,
244 a clinically-relevant domain with pressing needs for psychopharmacological treatment. Cognitive
245 deficits in schizophrenia are well-established [45] and cognitive deficits are among the strongest
246 predictors of functional outcomes [46]. Despite decades of effort, there are no first-line pharma-
247 cotherapies for cognitive symptoms in schizophrenia [45] (though recently-developed muscarinic
248 acetylcholine receptor agonists show promise [47, 48, 49]). We recently demonstrated that patients
249 with chronic schizophrenia have *reduced* policy complexity relative to healthy control participants
250 [29]. It stands to reason that increasing complexity in chronic schizophrenia, perhaps via KOR an-
251 tagonism, might treat a subset of cognitive deficits and improve functional outcomes. Although it
252 may seem counterproductive to administer dopaminergic drugs in schizophrenia, numerous studies

253 have shown that dopamine-releasing agents can be safe to administer in this population [50, 51, 52].
254 Neurobiologically, our finding that KOR antagonism increases complexity is similar to our
255 previous results following administration of dopaminergic medications in Parkinson’s disease [18].
256 A new subtlety of our findings here is that tonic dopamine may control the efficiency of resource
257 allocation, a finding that is perhaps related to the role of dopamine in habit formation [53, 54, 55].
258 Further, anhedonia may be related to more subtle disruptions in the dopaminergic system than
259 had been previously thought, as more global disruptions would likely reduce complexity as well.

260 Conclusion

261 We leveraged computational principles to identify two mechanisms of action of KOR antagonism—
262 one related to anhedonia (increase in efficiency), and one unrelated to anhedonia (increase in
263 policy complexity). We hypothesize that the increase in complexity can be leveraged for other
264 indications, including possibly cognitive deficits in psychosis. Our results provide a clear example
265 of the potential for computational psychiatry to provide transdiagnostic insights that integrate
266 across levels of analysis.

267 **Methods**

268 **KOR antagonism: randomized control trial design and participants**

269 We conducted a secondary analysis of a phase 2a clinical trial designed to test the efficacy of a novel
270 κ -opioid receptor (KOR) antagonist for the treatment of anhedonia [25, 26, 56]. The trial was an 8-
271 week, multicenter, placebo-controlled, double-blind, randomized study in a transdiagnostic sample
272 of participants with anhedonia. Active drug was JNJ-67953964 (Aticaprant, previously CERC-
273 501 and LY2456302), a selective KOR antagonist dosed at 10mg daily. Since this trial used a
274 biomarker-based proof-of-mechanism approach, the preregistered primary outcome was a change in
275 functional magnetic resonance imaging of the ventral striatum during reward anticipation, measured
276 at baseline and 8 weeks. Preregistered secondary outcomes were a change in the mean Snaith-
277 Hamilton Pleasure Scale (SHAPS), a clinically-validated measure of anhedonia [57], assessed every
278 2 weeks, and a change in the response bias - a measure of reward learning - on the Probabilistic
279 Reward Task. The trial was preregistered at [NCT02218736](#). We report here a secondary analysis
280 of the Probabilistic Reward Task, which was not part of the preregistered protocol.

281 Participants were aged 21 to 65, recruited from six US centers, had a SHAPS of at least 20
282 (assessed using dimensional scoring guidelines [58]), and had a DSM-IV TR diagnosis of major
283 depressive disorder, bipolar I or II depression, generalized anxiety disorder, social phobia, panic
284 disorder, or post-traumatic stress disorder. Participants were enrolled after providing informed
285 consent to a protocol approved by each local institutional review board. Our dataset for secondary
286 analysis consisted of 55 patients (KOR antagonist group: $N = 24$ [44%]; mean age \pm SD, $39.2 \pm$
287 13.9 years; 10 males [42%]; placebo group: $N = 31$ [56%]; mean age \pm SD, 40.8 ± 13.7 years; 12
288 males [39%]) [26]).

289 **Non-clinical population: study design and participants**

290 We conducted an online-based study to assess how variation in hedonic tone affects reward learning
291 in a non-clinical population. We recruited 100 participants (mean age \pm SD, 41.9 ± 11.5 ; 62 males
292 [62%]) from Amazon Mechanical Turk. We selected our sample size based on an effect size we
293 assumed would be half of what we identified for the KOR dataset ($f^2 = 0.1297$) with a desired
294 power of 90% to maximize the probability of identifying an effect. These participants completed the
295 Probabilistic Reward Task followed by a demographic survey and the SHAPS. Participants gave
296 informed consent, and the Harvard University Committee on the Use of Human Subjects approved
297 the experiment.

298 **Clinical population: study design and participants**

299 We reanalyzed data from patient populations performing the PRT [33, 34]. The dataset consisted
300 of 66 total participants (control group: $N = 25$ [38%]; mean age \pm SD, 38.4 ± 10.8 ; 14 males
301 [56%]; clinical group: $N = 41$ [62%]; mean age \pm SD, 41.9 ± 10.3 ; 24 males [59%]; 18 with
302 bipolar disorder [44%], 23 with major depressive disorder [56%]). The control participants had
303 no psychiatric diagnosis and were taking no psychoactive medications. In addition to the PRT,
304 participants completed the Beck Depression Inventory-II and the Mood and Anxiety Symptom
305 Questionnaire.

306 **Pramipexole administration: study design and participants**

307 We reanalyzed data from a double-blind, randomized trial assessing the effect of pramipexole, a
308 D2/D3 receptor agonist, on reward learning in the PRT [37]. Participants (placebo group: 13
309 [54%]; mean age \pm SD, 24.8 ± 3.2 ; 8 males [62%]; pramipexole group: 11 [46%]; mean age \pm SD,
310 26.0 ± 5.8 ; 6 males [56%]) were randomized to placebo or pramipexole. In the pramipexole group,
311 participants received a single 0.5mg dose, a low dose thought to act as a dopamine antagonist and

312 reduce phasic dopamine release. Participants completed the PRT 2 hours after receiving placebo
313 or pramipexole.

314 **Snaith-Hamilton Pleasure Scale (SHAPS)**

315 The SHAPS is a 14-item questionnaire used to assess anhedonia across four domains: inter-
316 est/pastimes, social interaction, sensory experience, and food/drink. Participants are asked to
317 respond to pleasurable situations (e.g., I would enjoy being with my family or close friends) with
318 one of the following responses on the basis of the last few days: strongly disagree, disagree, agree,
319 strongly agree. According to dimensional scoring guidelines [58], scores range from 1 for strongly
320 agree to 4 for strongly disagree, yielding a range of 14 to 56, with higher scores corresponding to
321 greater anhedonia. The SHAPS is the only clinical measure of anhedonia that significantly changes
322 with treatment in clinical trials [25, 59, 60].

323 **Probabilistic Reward Task (PRT)**

324 The PRT is a computerized decision making task designed to elicit learning in response to reward
325 [27, 28]. On each trial, participants observe one of two difficult-to-discriminate stimuli and are asked
326 to report which stimulus they observed. In the clinical trial, stimuli consisted of cartoon faces with
327 either a short mouth (11.5 mm) or a long mouth (13 mm) presented for 100 ms and participants
328 responded by pressing one of two keyboard keys ('z' or '/'). Participants completed 200 trials in
329 two 100 trial blocks, instead of 300 trials as usual, owing to time constraints imposed by the clinical
330 trial [25]. In the online-based task, stimuli consisted of images of either 10 squares/7 circles or 7
331 squares/10 circles (with 8 variations of each) and participants reported whether they observed more
332 squares or circles with one of two keyboard keys ('A' or 'L') [32]. Participants completed 300 trials
333 in three 100 trial blocks. Importantly, and unbeknownst to participants, correctly responding to
334 one stimulus yielded reward on 60% of trials ('rich' stimulus) while correctly responding to the other
335 stimulus yielded reward on 20% of trials ('poor' stimulus). They were instructed that not all correct
336 responses would yield a reward. The rich/poor stimuli and responses were counterbalanced across
337 participants in both studies. For our analyses, we excluded the first 25 trials to allow behavior to
338 stabilize. Our findings were qualitatively similar if we changed this trial exclusion threshold.

339 **Policy compression: a capacity limit applied to decisions**

340 All information processing systems—the human brain included—must contend with resource lim-
341 itations when making decisions. These constraints take on many forms, including computational
342 costs [61], metabolic costs [62], interference costs [63], and others [64]. Under policy compression,
343 we formalize the cognitive cost as the mutual information between states and actions, the policy
344 *complexity*:

$$345 I^\pi(S; A) = \sum_s P(s) \sum_a \pi(a|s) \log \frac{\pi(a|s)}{P(a)} \quad (1)$$

346 where $P(a) = \sum_s P(s)\pi(a|s)$ is the marginal probability of choosing action a under the policy.
347 In general, we assume that policies are subject to a capacity constraint, an upper bound, C , on
348 policy complexity. Shannon's noisy channel theorem states that the minimum expected number of
349 bits to transmit a signal across a noisy information channel without error is equal to the mutual
350 information. Therefore, if the optimal policy requires more memory than the subject possesses,
351 then it must *compress* the policy, or render it less state-dependent. We define the optimal policy,
 π^* , as:

$$352 \pi^* = \underset{\pi}{\operatorname{argmax}} V^\pi, \text{ subject to } I^\pi(S; A) \leq C \quad (2)$$

352 where V^π is the expected reward under policy π :

$$V^\pi = \sum_s P(s) \sum_a \pi(a|s) Q(s, a) \quad (3)$$

353 and $Q(s, a)$ is the expected reward for taking action a in state s .

354 We can express our constrained optimization problem in the following unconstrained Lagrangian
355 form:

$$\pi^* = \operatorname{argmax}_\pi \beta V^\pi - I^\pi(S; A) - \sum_s \lambda(s) \left(\sum_a \pi(a|s) - 1 \right) \quad (4)$$

356 where $\beta \geq 0, \lambda(s) \geq 0$ are Lagrangian multipliers. Solving this equation reveals that the optimal
357 policy takes on the following form:

$$\pi^*(a|s) \propto \exp[\beta Q(s, a) + \log P^*(a)] \quad (5)$$

358 where $P^*(a)$ is the optimal marginal action distribution, which can be interpreted as a form of
359 perseveration.

360 The optimal policy takes the form of the familiar softmax distribution, common in the reinforce-
361 ment learning literature. Here, the Lagrange multiplier, β , plays the role of the inverse temperature
362 parameter. Note that although β typically takes on the role of balancing exploration/exploitation
363 in reinforcement learning, we made no such appeals in deriving this policy. Moreover, β is a function
364 of the policy complexity:

$$\beta^{-1} = \frac{dV^\pi}{dI^\pi(S; A)} \quad (6)$$

365 At high policy complexity, when $\frac{dV^\pi}{dI^\pi(S; A)}$ is shallow, the optimal β is large and the policy is domi-
366 nated by Q -values. At low policy complexity, the optimal β is close to 0, and Q -values have minimal
367 impact on the policy. Moreover, when β is small, the perseveration term, $\log P^*(a)$, dominates,
368 and the policy is largely state-independent.

369 To construct the empirical reward–complexity curves, in both datasets, we computed the average
370 reward according to equation 3, where $P(s) = [0.5, 0.5]$ and $Q(s, a) = \begin{bmatrix} 0.2 & 0.6 \\ 0.6 & 0.2 \end{bmatrix}$, by construction,
371 and $\pi(a|s)$ was calculated from empirical action frequencies. We estimated mutual information by
372 computing the empirical action frequencies for each state for each session.

373 Reinforcement learning modeling

We constructed a cost-sensitive Q -learning model which contains three parameters ($\alpha_{\text{learn}}, \alpha_{\text{persev}}$, and β) and estimates action values, $Q(s, a)$, and marginal action probability, $P(a)$, to generate actions according to the following policy, mimicking the optimal policy under policy compression:

$$\begin{aligned} \Delta Q(s, a) &= \alpha_{\text{learn}} [r - Q(s, a)] \\ \Delta P(a) &= \alpha_{\text{persev}} [\pi(a|s) - P(a)] \\ \pi(a|s) &\propto \exp[(\beta Q(s, a) + \log(P(a)))] \end{aligned}$$

where $r = 1$ if the current trial is rewarded and 0 otherwise. The key feature of our model is a mechanism that allows treatment to multiplicatively scale α_{persev} and β (obtained after model comparison, see below). The model scales parameters in the following manner:

$$\begin{aligned} \alpha_{\text{persev}} &= \alpha_{\text{persev, baseline}} \cdot 10^{s_{\text{persev, treatment}}} \\ \beta &= \beta_{\text{baseline}} \cdot 10^{s_{\text{beta, treatment}}} \end{aligned}$$

A scaling value of 0 results in no scaling, > 0 results in an increase, and < 0 results in a decrease. For our online study, we scaled parameters as a function of the z -scored SHAPS in the following

manner:

$$\begin{aligned}\alpha_{\text{persev}} &= \alpha_{\text{persev,baseline}} \cdot 10^{s_{\text{persev}} \cdot \text{SHAPS}} \\ \beta &= \beta_{\text{baseline}} \cdot 10^{s_{\text{beta}} \cdot \text{SHAPS}}\end{aligned}$$

We initialized $Q(s, a)$ at 0 and $P(a)$ at 0.5 and we assumed scaling terms equaled 0 on sessions without treatment. We included all trials for analysis. Learning rates were constrained not to exceed 1. We constructed hierarchical models to obtain estimates of each parameter. Parameters were drawn from the following distributions:

$$\begin{aligned}\alpha_{\text{learn}} &\sim \text{Beta}(a_{\text{learn}}, b_{\text{learn}}) \\ \alpha_{\text{persev,baseline}} &\sim \text{Beta}(a_{\text{persev}}, b_{\text{persev}}) \\ \beta_{\text{baseline}} &\sim \text{Gamma}(a_{\text{beta}}, b_{\text{beta}})\end{aligned}$$

where we used $\text{Cauchy}^+(0,5)$ as a weakly-informative prior for each parameter. The gamma distribution was parameterized with a shape (a_{beta}) and scale (b_{beta}) parameter. Finally, the scaling terms were drawn according to

$$\begin{aligned}s_{\text{persev,treatment}} &\sim N(0, 1) \\ s_{\text{beta,treatment}} &\sim N(0, 1)\end{aligned}$$

374 α_{learn} , $\alpha_{\text{persev,baseline}}$, and β_{baseline} were constrained at the group level (one parameter per partici-
375 pant) and scaling terms were constrained at the treatment level (one parameter per treatment).

376 We initially fit a model that scaled all parameters ($s_{\text{learn,treatment}}$, $s_{\text{persev,treatment}}$, $s_{\text{beta,treatment}}$),
377 which produced an estimate of $s_{\text{learn,treatment}}$ that did not differ from 0, suggesting that treatment
378 did not effect α_{learn} . We therefore compared this ‘full’ model to the ‘reduced’ model we present above
379 (which does not scale α_{learn}). We performed model comparison using Pareto-smoothed importance
380 sampling leave-one-out cross-validation to estimate the expected log predictive density, a validated
381 measure of Bayesian model evaluation [65]. We found that our reduced model produced a similar
382 fit. We next compared our reduced model to three simpler variants: one that only scaled α_{persev} ,
383 one that only scaled β , and one with no scaling of any parameters. Model comparison favored the
384 model we present above which scales α_{persev} and β (Extended Data Table 1).

385 We next performed posterior predictive checks. We used the mean of each parameter as a point
386 estimate and simulated 200 trials of the PRT for each participant to mimic the dataset that was
387 used to fit the model. We analyzed this simulated dataset in the exact manner we analyzed the
388 ground-truth dataset.

389 To provide confidence in our interpretation of parameter changes, we tested the ability of our
390 model to recover known parameters. Using the same fictive, simulated dataset as above, we fit our
391 reinforcement learning model and obtained recovered parameter estimates. We computed Pearson’s
392 correlation between the known and recovered parameters (Extended Data Table 2).

393 For our heatmap of reward obtained with different α_{persev} and β combinations, we ran 3,000
394 independent simulations of the PRT for each combination of parameters. We fixed α_{learn} at 0.1423,
395 the mean posterior estimate across all participants. We performed a grid search across thirteen
396 logarithmically-spaced α_{persev} values from e^{-12} to e^0 , and ten β values from $e^{-0.4}$ to $e^{1.4}$.

397 Models were fit using R 4.2.2 (accessed with RStudio 2022.12.0+353) using the Rstan package
398 (version 2.26.13). We performed model comparison using the loo package (version 2.5.1).

399 Statistical analyses

400 For all group-level differences, we computed two-sided t -tests. In the KOR dataset, owing to
401 repeated measures, we fit linear mixed effects models to predict 1) inefficiency and 2) the probability
402 of choosing the richer option. Independent variables were policy complexity, treatment (placebo

403 or KOR), and time (baseline or post treatment), with a random intercept per participant. For the
404 other datasets (online non-clinical, clinical, and pramipexole), we fit a linear regression to predict
405 inefficiency. For the online non-clinical dataset, the dependent variables were policy complexity
406 and z -scored SHAPS. For the clinical dataset, they were policy complexity and group (control or
407 clinical). For the pramipexole dataset, they were policy complexity and treatment (placebo or
408 pramipexole). All analyses were 2-sided with an α of 0.05.

409 **Competing interests**

410 A.D.K. has been a consultant for Eisai, Axsome, Big Health, Harmony, Idorsia, Jazz, Janssen,
411 Takeda, Millenium Merck, Neurocrine, Neurawell, Otsuka, Evecxia and Sage Research and re-
412 ceived support from the NIH, the Ray and Dagmar Dolby Family Fund, Janssen, Jazz. Neurocrine,
413 Attune, Harmony, and Axsome. Over the past 3 years, D.A.P. has received consulting fees from
414 Boehringer Ingelheim, Compass Pathways, Engrail Therapeutics, Karla Therapeutics, Neumora
415 Therapeutics, Neurocrine Biosciences, Neuroscience Software, Otsuka, Sage Therapeutics, Sama
416 Therapeutics, Sunovion Therapeutics, and Takeda; he has received honoraria from the Ameri-
417 can Psychological Association, Psychonomic Society and Springer (for editorial work) as well as
418 Alkermes; he has received research funding from the Brain and Behavior Research Foundation,
419 Dana Foundation, Wellcome Leap, Millennium Pharmaceuticals, and NIMH; he has received stock
420 options from Compass Pathways, Engrail Therapeutics, Neumora Therapeutics, and Neuroscience
421 Software. D.A.P. has a financial interest in Neumora Therapeutics, which has licensed the copyright
422 to the PRT through Harvard University. The interests of D.A.P. were reviewed and are managed by
423 McLean Hospital and Mass General Brigham in accordance with their conflict-of-interest policies.
424 No funding from these entities was used to support the current work, and all views expressed are
425 solely those of the authors. B.A.B. and S.J.G. declare no competing interests.

426 References

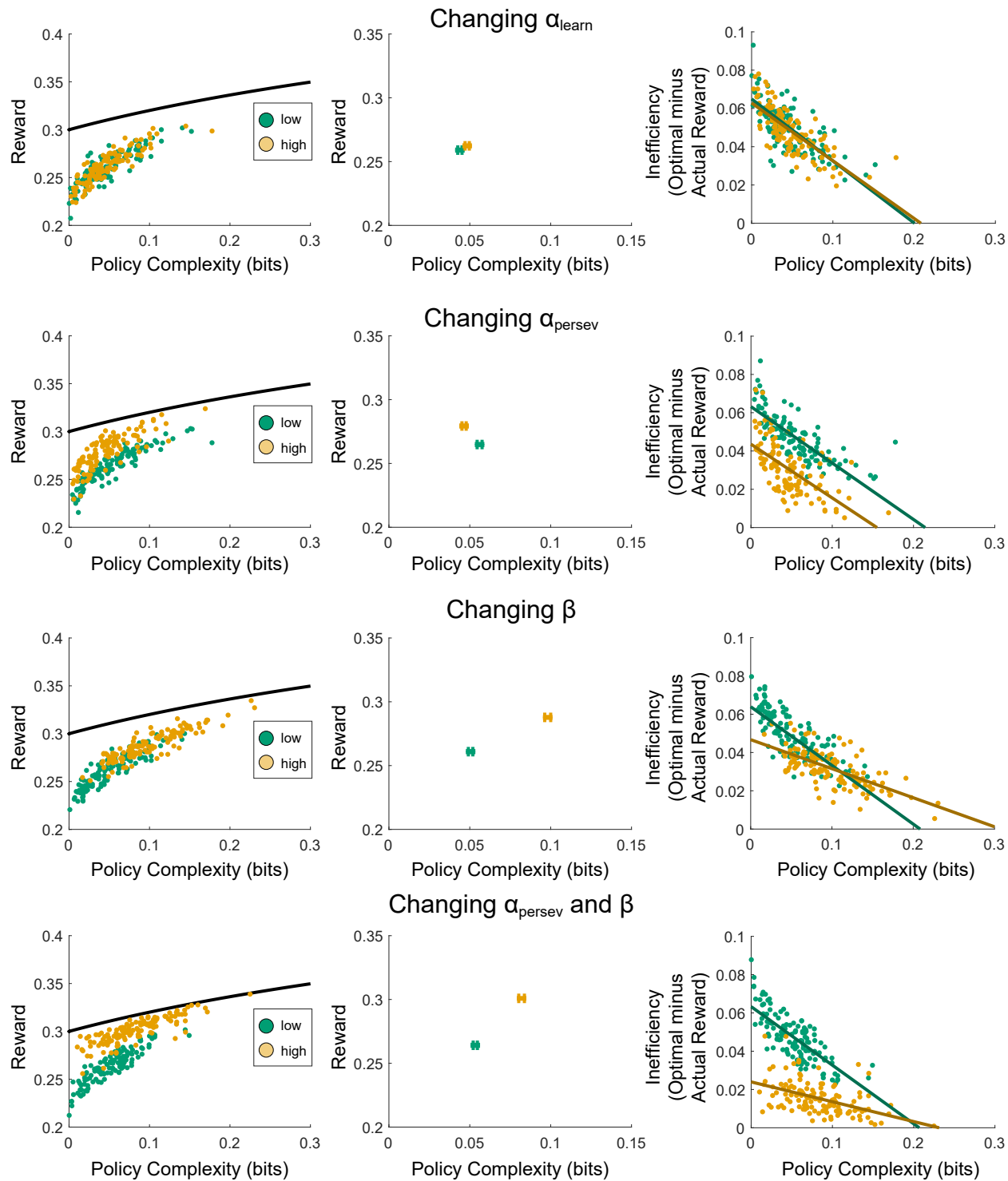
- 427 [1] Gard, D. E., Kring, A. M., Gard, M. G., Horan, W. P. & Green, M. F. Anhedonia in
428 schizophrenia: distinctions between anticipatory and consummatory pleasure. *Schizophrenia*
429 *research* **93**, 253–260 (2007).
- 430 [2] Kashdan, T. B., Zvolensky, M. J. & McLeish, A. C. Anxiety sensitivity and affect regulatory
431 strategies: Individual and interactive risk factors for anxiety-related symptoms. *Journal of*
432 *Anxiety disorders* **22**, 429–440 (2008).
- 433 [3] Hatzigiakoumis, D. S., Martinotti, G., Giannantonio, M. D. & Janiri, L. Anhedonia and
434 substance dependence: clinical correlates and treatment options. *Frontiers in psychiatry* **2**, 10
435 (2011).
- 436 [4] Chevallier, C., Grezes, J., Molesworth, C., Berthoz, S. & Happé, F. Brief report: Selective
437 social anhedonia in high functioning autism. *Journal of autism and developmental disorders*
438 **42**, 1504–1509 (2012).
- 439 [5] Meinzer, M. C., Pettit, J. W., Leventhal, A. M. & Hill, R. M. Explaining the covariance
440 between attention-deficit hyperactivity disorder symptoms and depressive symptoms: The
441 role of hedonic responsivity. *Journal of clinical psychology* **68**, 1111–1121 (2012).
- 442 [6] Nawijn, L. *et al.* Reward functioning in ptsd: a systematic review exploring the mechanisms
443 underlying anhedonia. *Neuroscience & Biobehavioral Reviews* **51**, 189–204 (2015).
- 444 [7] Husain, M. & Roiser, J. P. Neuroscience of apathy and anhedonia: a transdiagnostic approach.
445 *Nature Reviews Neuroscience* **19**, 470–484 (2018).
- 446 [8] Pizzagalli, D. A. *Anhedonia: preclinical, translational, and clinical integration*, vol. 58
447 (Springer Nature, 2022).
- 448 [9] Guineau, M. G. *et al.* Anhedonia as a transdiagnostic symptom across psychological disorders:
449 A network approach. *Psychological Medicine* **53**, 3908–3919 (2023).
- 450 [10] Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction* (MIT Press, 2018).
- 451 [11] Chase, H. *et al.* Approach and avoidance learning in patients with major depression and
452 healthy controls: relation to anhedonia. *Psychological Medicine* **40**, 433–440 (2010).
- 453 [12] Huys, Q. J., Pizzagalli, D. A., Bogdan, R. & Dayan, P. Mapping anhedonia onto reinforcement
454 learning: a behavioural meta-analysis. *Biology of mood & anxiety disorders* **3**, 1–16 (2013).
- 455 [13] Huys, Q. J. & Browning, M. A computational view on the nature of reward and value in anhe-
456 donia. In *Anhedonia: Preclinical, Translational, and Clinical Integration*, 421–441 (Springer,
457 2021).
- 458 [14] Parush, N., Tishby, N. & Bergman, H. Dopaminergic balance between reward maximization
459 and policy complexity. *Frontiers in Systems Neuroscience* **5**, 22 (2011).
- 460 [15] Gershman, S. J. Origin of perseveration in the trade-off between reward and complexity.
461 *Cognition* **204**, 104394 (2020).
- 462 [16] Lai, L. & Gershman, S. J. Policy compression: An information bottleneck in action selection.
463 In *Psychology of Learning and Motivation*, vol. 74, 195–232 (Elsevier, 2021).
- 464 [17] Mikhael, J. G., Lai, L. & Gershman, S. J. Rational inattention and tonic dopamine. *PLoS*
465 *computational biology* **17**, e1008659 (2021).

- 466 [18] Bari, B. A. & Gershman, S. J. Undermatching is a consequence of policy compression. *Journal*
467 *of Neuroscience* **43**, 447–457 (2023).
- 468 [19] Manohar, S. G. *et al.* Reward pays the cost of noise reduction in motor and cognitive control.
469 *Current Biology* **25**, 1707–1716 (2015).
- 470 [20] Carlezon, W. A. *et al.* Depressive-like effects of the κ -opioid receptor agonist salvinorin a on
471 behavior and neurochemistry in rats. *Journal of Pharmacology and Experimental Therapeutics*
472 **316**, 440–447 (2006).
- 473 [21] Bruijnzeel, A. W. kappa-opioid receptor signaling and brain reward function. *Brain research*
474 *reviews* **62**, 127–146 (2009).
- 475 [22] Ebner, S. R., Roitman, M. F., Potter, D. N., Rachlin, A. B. & Chartoff, E. H. Depressive-like
476 effects of the kappa opioid receptor agonist salvinorin a are associated with decreased phasic
477 dopamine release in the nucleus accumbens. *Psychopharmacology* **210**, 241–252 (2010).
- 478 [23] Muschamp, J. W. *et al.* Activation of creb in the nucleus accumbens shell produces anhedonia
479 and resistance to extinction of fear in rats. *Journal of Neuroscience* **31**, 3095–3103 (2011).
- 480 [24] Wallace, C. W., Holleran, K. M., Slinkard, C. Y., Centanni, S. W. & Jones, S. R. Kappa opioid
481 receptors negatively regulate real time spontaneous dopamine signals by reducing release and
482 increasing uptake. *bioRxiv* 2024–02 (2024).
- 483 [25] Krystal, A. D. *et al.* A randomized proof-of-mechanism trial applying the ‘fast-fail’ approach
484 to evaluating κ -opioid antagonism as a treatment for anhedonia. *Nature medicine* **26**, 760–768
485 (2020).
- 486 [26] Pizzagalli, D. A. *et al.* Selective kappa-opioid antagonism ameliorates anhedonic behavior:
487 Evidence from the fast-fail trial in mood and anxiety spectrum disorders (fast-mas). *Neu-*
488 *ropsychopharmacology* **45**, 1656–1663 (2020).
- 489 [27] Tripp, G. & Alsop, B. Sensitivity to reward frequency in boys with attention deficit hyperac-
490 tivity disorder. *Journal of clinical child psychology* **28**, 366–375 (1999).
- 491 [28] Pizzagalli, D. A., Jahn, A. L. & O’Shea, J. P. Toward an objective characterization of an
492 anhedonic phenotype: a signal-detection approach. *Biological psychiatry* **57**, 319–327 (2005).
- 493 [29] Gershman, S. J. & Lai, L. The reward-complexity trade-off in schizophrenia. *Computational*
494 *Psychiatry* **5** (2021).
- 495 [30] Lai, L. & Gershman, S. J. Human decision making balances reward maximization and policy
496 compression. *PsyArXiv* (2023).
- 497 [31] Watkins, C. J. & Dayan, P. Q-learning. *Machine learning* **8**, 279–292 (1992).
- 498 [32] de Leeuw, J. R. prt-test-my-brain. <https://github.com/jodeleeuw/prt-test-my-brain>
499 (2021).
- 500 [33] Pizzagalli, D. A., Iosifescu, D., Hallett, L. A., Ratner, K. G. & Fava, M. Reduced hedonic
501 capacity in major depressive disorder: evidence from a probabilistic reward task. *Journal of*
502 *psychiatric research* **43**, 76–87 (2008).
- 503 [34] Pizzagalli, D. A., Goetz, E., Ostacher, M., Iosifescu, D. V. & Perlis, R. H. Euthymic patients
504 with bipolar disorder show decreased reward learning in a probabilistic reward task. *Biological*
505 *psychiatry* **64**, 162–168 (2008).

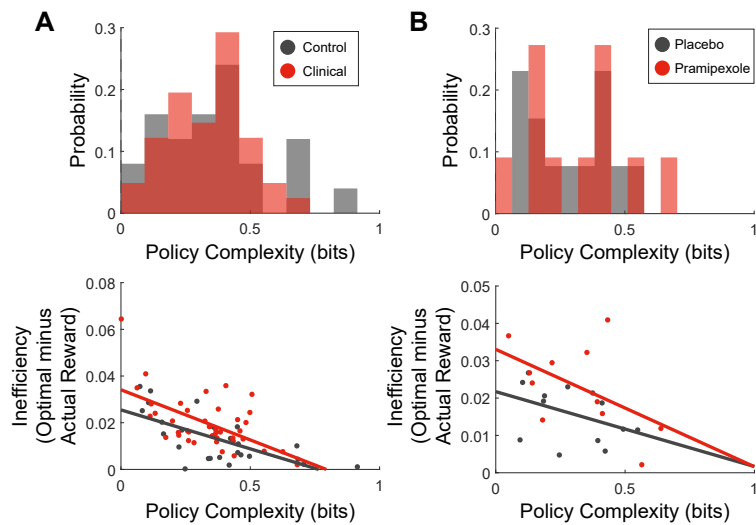
- 506 [35] Wise, R. A. Dopamine and reward: the anhedonia hypothesis 30 years on. *Neurotoxicity*
507 *research* **14**, 169–183 (2008).
- 508 [36] Argyropoulos, S. V. & Nutt, D. J. Anhedonia revisited: is there a role for dopamine-targeting
509 drugs for depression? *Journal of psychopharmacology* **27**, 869–877 (2013).
- 510 [37] Pizzagalli, D. A. *et al.* Single dose of a dopamine agonist impairs reinforcement learning
511 in humans: behavioral evidence from a laboratory-based measure of reward responsiveness.
512 *Psychopharmacology* **196**, 221–232 (2008).
- 513 [38] Berridge, K. C. & Robinson, T. E. What is the role of dopamine in reward: hedonic impact,
514 reward learning, or incentive salience? *Brain research reviews* **28**, 309–369 (1998).
- 515 [39] Rizvi, S. J. *et al.* Development and validation of the dimensional anhedonia rating scale (dars)
516 in a community sample and individuals with major depression. *Psychiatry research* **229**,
517 109–119 (2015).
- 518 [40] Gard, D. E., Gard, M. G., Kring, A. M. & John, O. P. Anticipatory and consummatory
519 components of the experience of pleasure: a scale development study. *Journal of research in*
520 *personality* **40**, 1086–1102 (2006).
- 521 [41] Khazanov, G. K., Ruscio, A. M. & Forbes, C. N. The positive valence systems scale: develop-
522 ment and validation. *Assessment* **27**, 1045–1069 (2020).
- 523 [42] Pizzagalli, D. A. Toward a better understanding of the mechanisms and pathophysiology
524 of anhedonia: are we ready for translation? *American Journal of Psychiatry* **179**, 458–469
525 (2022).
- 526 [43] Robbins, T. & Costa, R. M. Habits. *Current biology* **27**, R1200–R1206 (2017).
- 527 [44] Nebe, S., Kretschmar, A., Brandt, M. C. & Tobler, P. N. Characterizing human habits in the
528 lab. *Collabra: Psychology* **10** (2024).
- 529 [45] Green, M. F. & Harvey, P. D. Cognition in schizophrenia: Past, present, and future.
530 *Schizophrenia Research: Cognition* **1**, e1–e9 (2014).
- 531 [46] Green, M. F. Cognitive impairment and functional outcome in schizophrenia and bipolar
532 disorder. *Journal of Clinical Psychiatry* **67**, 3 (2006).
- 533 [47] Brannan, S. K. *et al.* Muscarinic cholinergic receptor agonist and peripheral antagonist for
534 schizophrenia. *New England Journal of Medicine* **384**, 717–726 (2021).
- 535 [48] Krystal, J. H. *et al.* Emraclidine, a novel positive allosteric modulator of cholinergic m4
536 receptors, for the treatment of schizophrenia: a two-part, randomised, double-blind, placebo-
537 controlled, phase 1b trial. *The Lancet* **400**, 2210–2220 (2022).
- 538 [49] Kaul, I. *et al.* Efficacy and safety of the muscarinic receptor agonist karxt (xanomeline-
539 trospium) in schizophrenia (emergent-2) in the usa: results from a randomised, double-blind,
540 placebo-controlled, flexible-dose phase 3 trial. *The Lancet* **403**, 160–170 (2024).
- 541 [50] Tsoi, D. T.-y., Porwal, M. & Webster, A. C. Efficacy and safety of bupropion for smoking
542 cessation and reduction in schizophrenia: systematic review and meta-analysis. *The British*
543 *Journal of Psychiatry* **196**, 346–353 (2010).
- 544 [51] Englisch, S., Morgen, K., Meyer-Lindenberg, A. & Zink, M. Risks and benefits of bupropion
545 treatment in schizophrenia: a systematic review of the current literature. *Clinical neurophar-*
546 *macology* **36**, 203–215 (2013).

- 547 [52] Lindenmayer, J.-P., Nasrallah, H., Pucci, M., James, S. & Citrome, L. A systematic review of
548 psychostimulant treatment of negative symptoms of schizophrenia: challenges and therapeutic
549 opportunities. *Schizophrenia research* **147**, 241–252 (2013).
- 550 [53] Faure, A., Haberland, U., Condé, F. & El Massioui, N. Lesion to the nigrostriatal dopamine
551 system disrupts stimulus-response habit formation. *Journal of Neuroscience* **25**, 2771–2780
552 (2005).
- 553 [54] Wickens, J. R., Horvitz, J. C., Costa, R. M. & Killcross, S. Dopaminergic mechanisms in
554 actions and habits. *Journal of Neuroscience* **27**, 8181–8183 (2007).
- 555 [55] Wang, L. P. *et al.* Nmda receptors in dopaminergic neurons are crucial for habit learning.
556 *Neuron* **72**, 1055–1066 (2011).
- 557 [56] Fast-Fail Trials in Mood and Anxiety Spectrum Disorders; Kappa Opioid Receptor Phase
558 2a. Clinicaltrials.gov identifier: Nct02218736. [https://clinicaltrials.gov/study/
559 NCT02218736](https://clinicaltrials.gov/study/NCT02218736) (Updated: January 8, 2019. Accessed: February 27, 2024.).
- 560 [57] Snaith, R. P. *et al.* A scale for the assessment of hedonic tone the snaith–hamilton pleasure
561 scale. *The British Journal of Psychiatry* **167**, 99–103 (1995).
- 562 [58] Franken, I. H., Rassin, E. & Muris, P. The assessment of anhedonia in clinical and non-
563 clinical populations: further validation of the snaith–hamilton pleasure scale (shaps). *Journal
564 of affective disorders* **99**, 83–89 (2007).
- 565 [59] Martinotti, G. *et al.* Acetyl-l-carnitine in the treatment of anhedonia, melancholic and neg-
566 ative symptoms in alcohol dependent subjects. *Progress in Neuro-Psychopharmacology and
567 Biological Psychiatry* **35**, 953–958 (2011).
- 568 [60] Di Giannantonio, M. & Martinotti, G. Anhedonia and major depression: the role of agomela-
569 tine. *European Neuropsychopharmacology* **22**, S505–S510 (2012).
- 570 [61] Bossaerts, P., Yadav, N. & Murawski, C. Uncertainty and computational complexity. *Philo-
571 sophical Transactions of the Royal Society B* **374**, 20180138 (2019).
- 572 [62] Gailliot, M. T. & Baumeister, R. F. The physiology of willpower: Linking blood glucose to
573 self-control. *Personality and social psychology review* **11**, 303–327 (2007).
- 574 [63] Musslick, S. *et al.* Parallel processing capability versus efficiency of representation in neural
575 networks. *Network* **8** (2016).
- 576 [64] Shenhav, A. *et al.* Toward a rational and mechanistic account of mental effort. *Annual review
577 of neuroscience* **40**, 99–124 (2017).
- 578 [65] Vehtari, A., Gelman, A. & Gabry, J. Practical bayesian model evaluation using leave-one-out
579 cross-validation and waic. *Statistics and computing* **27**, 1413–1432 (2017).

580 **Extended data**



Extended Data Figure 1: Effect of changing reinforcement learning model parameters on reward-complexity relationship and inefficiency. Parameter values for simulation are given in Extended Data Table 3.



Extended Data Figure 2: Policy complexity and inefficiency for reanalyzed PRT datasets.

- A) Clinical dataset from [33] and [34].
- B) Pramipexole dataset from [37].

Model			Expected log predictive density difference \pm SE	Effective number of parameters (p_loo) \pm SE
Scale α_{learn}	Scale α_{persev}	Scale β		
	x	x	0.0 ± 0.0	100.5 ± 1.2
x	x	x	-1.8 ± 2.5	107.6 ± 1.3
	x		-14.7 ± 6.1	112.9 ± 1.4
		x	-29.2 ± 7.8	100.8 ± 1.1
			-43.1 ± 10.1	98.5 ± 1.1

Extended Data Table 1: Model comparison using Pareto-smoothed importance sampling leave-one out cross validation. A difference in the expected log predictive density of 4 points provides evidence in favor of a model. The first model, which scales α_{persev} and β , is favored over the second, which scales all parameters, since it provides similar expected predictive accuracy with fewer parameters.

Parameter	Pearson Correlation Between Actual and Recovered Parameter (95% CI)
α_{learn}	0.412 (0.183 - 0.607)
α_{persev}	0.941 (0.862 - 0.983)
β	0.926 (0.884 - 0.952)

Extended Data Table 2: Parameter recovery.

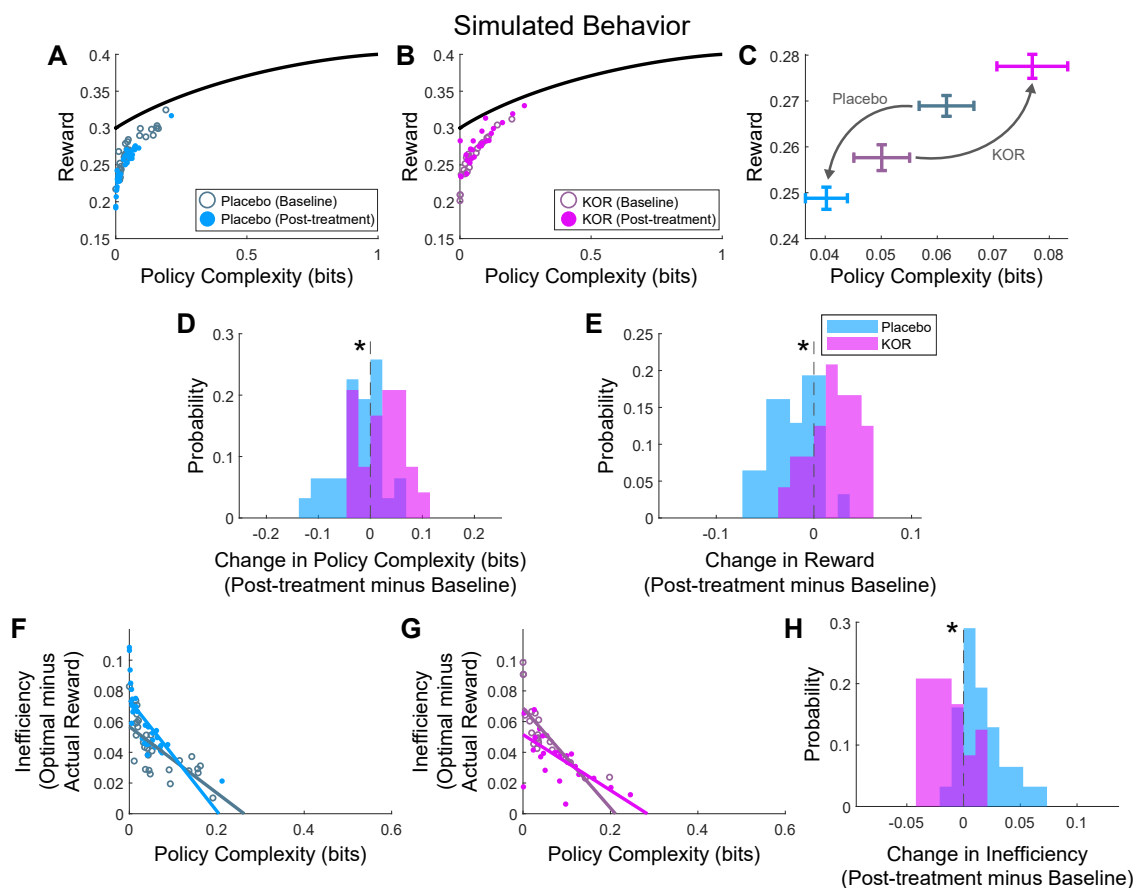
Simulation	Low Parameter	High Parameter	Fixed Parameters
Changing α_{learn}	$\alpha_{\text{learn}} = 0.09$	$\alpha_{\text{learn}} = 0.18$	$\alpha_{\text{persev}} = 2.5 \cdot 10^{-4}$ $\beta = 1.4$
Changing α_{persev}	$\alpha_{\text{persev}} = 2.5 \cdot 10^{-4}$	$\alpha_{\text{persev}} = 2.5 \cdot 10^{-2}$	$\alpha_{\text{learn}} = 0.14$ $\beta = 1.4$
Changing β	$\beta = 1.4$	$\beta = 2.2$	$\alpha_{\text{learn}} = 0.14$ $\alpha_{\text{persev}} = 2.5 \cdot 10^{-4}$
Changing α_{persev} & β	$\alpha_{\text{persev}} = 2.5 \cdot 10^{-4}$ $\beta = 1.4$	$\alpha_{\text{persev}} = 2.5 \cdot 10^{-2}$ $\beta = 2.2$	$\alpha_{\text{learn}} = 0.14$

Extended Data Table 3: Parameters used for Extended Data Figure 1 simulations.

581 **Supplementary information**

582 **Reinforcement learning model of KOR antagonism: behavioral simulations on**
 583 **the Probabilistic Reward Task**

584 Treatment increased policy complexity (Supplementary Figure 1D; mean change in policy com-
 585 plexity (post-treatment minus baseline) \pm SEM: placebo, $-0.0214 \pm 7.69 \times 10^{-3}$; KOR, $0.0269 \pm$
 586 8.32×10^{-3} , $p = 9.07 \times 10^{-5}$) and reward (Supplementary Figure 1E; mean change in reward (post-
 587 treatment minus baseline) \pm SEM: placebo, $-0.0201 \pm 4.36 \times 10^{-3}$; KOR, $0.0199 \pm 4.62 \times 10^{-3}$,
 588 $p = 7.21 \times 10^{-8}$). There was a significant decrease in inefficiency (Supplementary Figure 1H; mean
 589 change in inefficiency (post-treatment minus baseline) \pm SEM: placebo, $0.0161 \pm 3.50 \times 10^{-3}$; KOR,
 590 $-0.0149 \pm 3.53 \times 10^{-3}$, $p = 1.08 \times 10^{-7}$). Using the same linear mixed effects model to predict in-
 591 efficiency as a function of policy complexity, treatment, and time, we found a significant treatment
 592 \times time interaction (coefficient = -0.0346 , $p = 1.94 \times 10^{-7}$) and a significant policy complexity \times
 593 treatment \times time interaction (coefficient = 0.290 , $p = 8.66 \times 10^{-4}$).



Supplementary Figure 1: Simulation: changes in complexity and efficiency as a function of KOR antagonism.

- A,B) Reward-complexity relationship for placebo and KOR groups, at baseline and post-treatment.
- C) Mean \pm SEM reward-complexity relationship as a function of treatment and time.
- D) Change in policy complexity as a function of treatment.
- E) Change in reward as a function of treatment.
- F-G) Inefficiency-complexity relationship for placebo and KOR groups.
- H) Change in inefficiency as a function of treatment.

594 Anhedonia model from Huys et al (2013)

595 Huys et al (2013) developed a reinforcement learning model, fit to PRT data, describing anhedonia as a reduction in reward sensitivity [12]. We will show that the parameterization of reward
596 nia as a reduction in reward sensitivity [12]. We will show that the parameterization of reward
597 sensitivity in this model produces a similar effect as our perseveration term.

In their model, reward prediction errors are computed by scaling binary reward, r , by a reward sensitivity parameter ρ . These reward prediction errors are multiplied by ϵ , the learning rate, to iteratively update Q -values.

$$\begin{aligned}\delta &= \rho r - Q(s, a) \\ \Delta Q(s, a) &= \epsilon \delta\end{aligned}$$

Given the reward structure in the PRT, this has the effect of scaling Q -values by ρ as $\begin{bmatrix} \rho^{0.2} & 0 \\ 0 & \rho^{0.6} \end{bmatrix}$. These Q -values are used to update choice weights, which are fed through a standard softmax decision rule to generate a policy:

$$\begin{aligned}W(s, a) &= \gamma I(s, a) + \zeta Q(s, a) + (1 - \zeta)Q(\bar{s}, a) \\ \pi(a|s) &\propto \exp(W(s, a))\end{aligned}$$

598 The choice weights of this model contain two noteworthy components. The first is an instruction
599 variable, $I(s, a)$, where $I(s, a) = 1$ for the instructed action for a given stimulus, and 0 otherwise.
600 Instructions are scaled by γ to capture how strongly instructions influence choice. The second
601 component describes sensory ambiguity and allows Q -values for the non-presented stimulus - $Q(\bar{s}, a)$
602 - to ‘leak’ into the policy. This is done by the ζ parameter, where $\zeta \in [0.5, 1]$; $\zeta = 1$ describes
603 no sensory ambiguity (only $Q(s, a)$ contributes) and $\zeta = 0.5$ describes complete sensory ambiguity
604 ($Q(s, a)$ and $Q(\bar{s}, a)$ contribute equally).

To see how this sensory ambiguity rule leads to perseveration, we can define $\zeta = \theta + 0.5$, where $\theta \in [0, 0.5]$ and replace ζ in the choice weights:

$$W(s, a) = \gamma I(s, a) + (\theta + 0.5)Q(s, a) + (1 - (\theta + 0.5))Q(\bar{s}, a)$$

which we can rearrange as

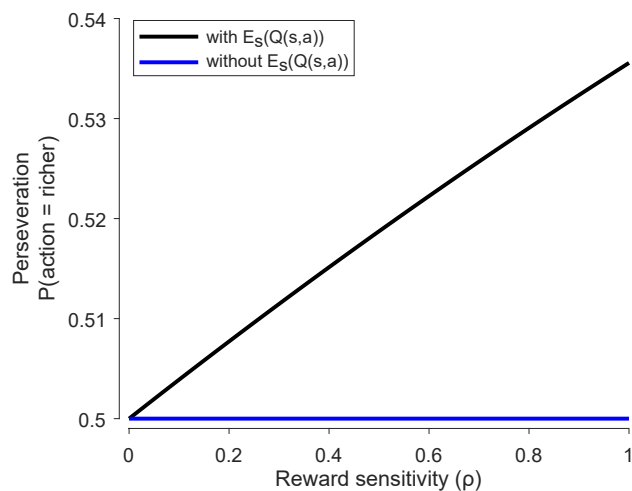
$$W(s, a) = \gamma I(s, a) + \theta(Q(s, a) - Q(\bar{s}, a)) + 0.5Q(s, a) + 0.5Q(\bar{s}, a)$$

Since the states are equiprobable ($p(s) = 0.5$), this latter set of terms, $0.5Q(s, a) + 0.5Q(\bar{s}, a)$, can be written as $\mathbb{E}_s(Q(s, a)) = \sum_s p(s)Q(s, a)$, the expected Q -value of taking action a . We can therefore write the weights as

$$W(s, a) = \gamma I(s, a) + \theta(Q(s, a) - Q(\bar{s}, a)) + \mathbb{E}_s(Q(s, a))$$

605 Written this way, weights are a function of three variables: 1) $I(s, a)$, the instructions, 2) $Q(s, a) -$
606 $Q(\bar{s}, a)$, the difference in Q -values between the observed and non-observed states, to account for
607 sensory ambiguity, and 3) $\mathbb{E}_s(Q(s, a))$, a state-independent value term which can be thought of as
608 a kind of perseveration since it will generate an action bias.

609 We ran a simulation to gain an intuition into how $\mathbb{E}_s(Q(s, a))$ engenders perseveration (Sup-
610 plementary Figure 2). In this simulation, $Q(s, a) = \begin{bmatrix} \rho^{0.2} & 0 \\ 0 & \rho^{0.6} \end{bmatrix}$, meaning $\mathbb{E}_s(Q(s, a)) \propto \rho^{0.6}$ for the
611 richer option and $\propto \rho^{0.2}$ for the leaner option. Intuitively, $\mathbb{E}_s(Q(s, a))$ will proportionally favor the
612 richer option as reward sensitivity grows, leading to an action bias.



Supplementary Figure 2: Increasing reward sensitivity (ρ) in the Huys et al (2013) model leads to perseveration.

Increasing reward sensitivity (ρ) leads to increased perseveration (black). To demonstrate that the $E_s(Q(s,a))$ term is responsible for perseveration, we ran the same simulation with the $E_s(Q(s,a))$ removed (blue). For this simulation, we used $\gamma = 1$ and $\epsilon = 0.25$. Findings were insensitive to choice of γ and ϵ .