

1 **Effects of depression on prefrontal–striatal goal-directed and habitual control**

2
3 Suyeon Heo^{1,2} and Sang Wan Lee^{*1,2,3,4}

4
5 ¹ Department of Bio and Brain Engineering,

6 ² Brain and Cognitive Engineering Program,

7 ³ KAIST Institute for Health Science Technology,

8 ⁴ KAIST Institute for Artificial Intelligence,

9 Korea Advanced Institute of Science Technology (KAIST),

10 Daejeon 34141, Republic of Korea

11
12 *Corresponding author (sangwan@kaist.ac.kr)

13 14 **ABSTRACT**

15 Depression is characterized by deficits in the reinforcement learning (RL) process. Although
16 many computational and neural studies have extended our knowledge of the impact of
17 depression on RL, most focus on habitual control (model-free RL), yielding a relatively poor
18 understanding of goal-directed control (model-based RL) and arbitration control to find a
19 balance between the two. We investigate the effects of depression on goal-directed and
20 habitual control in the prefrontal–striatal circuitry. We find that depression is associated with
21 attenuated state and reward prediction error representation in the insula and caudate, a
22 disruption of arbitration control in the predominantly inferior lateral prefrontal cortex and
23 frontopolar cortex, and suboptimal value–action conversion. These findings fully characterize
24 how depression influences different levels of RL, challenging previous conflicting views that
25 depression simply influences either habitual or goal-directed control. Our study creates
26 possibilities for various clinical applications, such as early diagnosis and behavioral therapy
27 design.

28 29 **INTRODUCTION**

30 Major depressive disorder (MDD) has received considerable attention, as the lifetime
31 prevalence of the disorder is higher than 10% worldwide¹. MDD is characterized by deficits in

1 decision-making^{2,3} and its underlying reward learning processes⁴. Recently, with the
2 development of computational models, several studies have explored how depression
3 influences the reward learning system (for a review, see (Chen et al., 2015)).

4 Reinforcement learning (RL), the process of learning to develop a behavioral policy to
5 maximize reward⁵, has been known to be guided by the two distinct RL strategies: model-
6 based (MB) RL and model-free (MF) RL, each of which guides goal-directed and habitual RL,
7 respectively⁶. Model-based RL guides context-sensitive and goal-directed behaviors through
8 a sophisticated process in which the learning agent makes decisions by simulating an internal
9 environmental model, whereas model-free RL is associated with habitual responses to reward-
10 predicting stimuli based on learned associations between stimuli and rewards^{6,7}. Mounting
11 evidence suggests that depression is characterized by impairments in either model-based or
12 model-free RL. For example, behaviors in depressive individuals can be accounted for by
13 impaired model-based RL⁸⁻¹⁰ or a transition from model-based to model-free RL¹¹. However,
14 most studies have only explored the effect of depression on model-free RL. For instance,
15 depressive people exhibit an impaired ability to learn stimulus–reward associations
16 accompanying inaccurate representations of reward prediction error¹²⁻¹⁷ or abnormal learning
17 rate control^{11,18}.

18 Impairment in RL is associated with not only the onset of depressive symptoms, but also the
19 development of depression. For instance, stress, one of the major risk factors for
20 depression^{19,20}, can induce deficits in RL. Previous findings have shown that people exhibit a
21 reduced ability to engage in model-based RL under conditions of chronic^{21,22} and acute^{23,24}
22 stress. These findings suggest a gradual impairment of RL from the very early stages of
23 depression.

24 Although these studies have contributed to our understanding of depression in the context of
25 RL, it is still unclear whether depression is best characterized by model-free RL, model-based
26 RL, or an interaction between the two, or how depression influences the neural circuits guiding
27 goal-directed and habitual behavioral control. Moreover, little is known about how these cases
28 extend to early or mild depression.

29 Here, we aim to provide a computational and neural account of how depression affects goal-
30 directed and habitual control in the prefrontal–striatal circuitry. First, to investigate the effects
31 of depression on model-based and model-free RL, we ran a model comparison analysis to
32 identify a version of arbitration control intended to account for various behavioral traits of
33 depression. In particular, our computational models consider sub-optimality of RL, allowing us
34 to explain choice behavior patterns across a wide spectrum of depression. We combine this

1 with model-based functional magnetic resonance imaging (fMRI) to identify the parametric
2 effects of depression on neural systems associated with model-based and model-free RL. In
3 the subsequent analysis, we attempt to fully characterize how depression disrupts the
4 arbitration between model-based and model-free RL by combining the results from the
5 computational modeling and model-based fMRI analyses and the multi-voxel pattern analysis
6 (MVPA).

7

8 **RESULTS**

9 **Effects of Depression on Behavior Performance of Goal-directed and Habitual Learning**

10

11 [Figure 1] Markov decision task structure

12

13 Sixty-three participants conducted a sequential two-stage Markov decision task²⁵. Of the
14 subjects, 28 were scanned with an fMRI while performing the task. The task manipulated
15 both prediction uncertainty and task goals to dissociate goal-directed and habitual behavior
16 control (Figure 1; for more detail, see Methods). The task consisted of four types of blocks
17 (high/low state–action–state transition uncertainty x specific/flexible goal condition). Before
18 each experiment, participants completed the Center for Epidemiologic Studies Depression
19 (CES-D) questionnaire²⁶ (For the distribution of participant’s depression severity, see
20 Supplementary Figure S2).

21

22 [Figure 2] Behavioral results

23

24 Overall task performance is negatively correlated with the self-reported depression score.
25 The accumulated reward decreases significantly as individual depression score (CES-D)
26 increases (correlation coefficient estimate=-0.584 [$p=4.94e-07$]; Figure 2a). The proportion of
27 optimal choices, the measure that quantifies the extent to which a subject’s choice reflects an
28 optimal policy, is also inversely proportional to the CES-D score (correlation coefficient
29 estimate=-0.567 [$p=9.32e-05$]; Figure 2b). Finally, choice consistency, the proportion of
30 making the same choice as in previous trials, decreases as the CES-D score increases
31 (correlation coefficient estimate=-0.472 [$p=1.24e-06$]; Figure 2c). These results demonstrate

1 that depression has a damaging effect on both goal-directed and habitual control performance,
2 leading to suboptimal choices.

3

4 **Computational Model of Arbitration Control Allowing for Suboptimal Decision-Making**

5

6 [Figure 3] Computational model of dynamic arbitration control

7

8 We adopted the previous dynamic arbitration control hypothesis that respective prediction
9 uncertainty—specifically, the amount of uncertainty in the state and reward prediction error of
10 model-based and model-free RL—mediates the trial-by-trial value integration of the model-
11 based and model-free systems²⁵. To fully explore the effects of depression on arbitration
12 control, however, a model should be flexible enough to account for any individual variability
13 arising from suboptimal learning and decision-making.

14 To consider this, we redesigned the arbitration control scheme to allow for sub-optimality of
15 RL in both learning values and converting learned values into choice behavior. We included
16 the former by introducing separate learning rates for model-based and model-free RL and the
17 latter by defining an exploitation sensitivity parameter as a function of model preference for
18 either model-based or model-free RL (for more detail, see Methods). This model setting
19 reflects the hypothesis that prediction uncertainty mediates not only value integration, but also
20 value–action conversion (Figure 3a).

21 We compared prediction performance for the five different versions of arbitration control,
22 including the original arbitration model and four other versions implementing our hypothesis
23 in different ways. We used the Bayesian information criteria (BIC) as a performance measure
24 to preclude overfitting. We found that the version implementing our hypothesis, in which the
25 degree of exploitation is determined by the weighted sum of the model-based and model-free
26 exploitation parameters with the model choice probability, explains the subjects' behavior
27 ($t(62)=2.46$, $[p=0.017]$; paired t-test comparison with the second-best model, Arb_α),
28 significantly better than the original model²⁵ (Figure 3b; for a model comparison, see
29 Supplementary Table S1; for estimated model parameters, see Supplementary Table S2). Our
30 model comparison result not only corroborates the previous finding that prediction uncertainty
31 mediates the arbitration between model-based and model-free RL, but also demonstrates the
32 effect of prediction uncertainty on both value integration and value–action conversion.

1

2 **Goal-directed and Habitual Control in the Prefrontal–Striatal Circuitry**

3

4 [Figure 4] Neural correlates of dynamic arbitration control

5

6 To further examine whether our model explains the neural activity patterns of brain areas
7 previously implicated in model-based and model-free RL, we ran a model-based fMRI analysis
8 in which each of the key signals of our computational model were regressed against the fMRI
9 data.

10 First, we replicated previous findings concerning the neural correlates of prediction error for
11 the model-based and model-free systems. The state prediction error (SPE) was found
12 bilaterally in the insula and the dorsolateral prefrontal cortex (dlPFC) (all $p < 0.05$ in cluster-
13 level corrected). The reward prediction error (RPE) was correlated with neural activity in both
14 sides of the ventral striatum ($p < 0.05$ in family-wise error [FWE] corrected). These results are
15 fully consistent with previous findings^{25,27,28} (Figure 4, Supplementary Table S3).

16 We also successfully replicated previous findings supporting the neural hypothesis of
17 arbitration control. We found the max reliability signal, the key signal used to mediate
18 arbitration between model-based and model-free RL, in the bilateral inferior lateral PFC (ilPFC,
19 left: $p < 0.05$ in cluster-level corrected; right: $p < 0.05$ in small-volume corrected [SVC] at [8,6,-
20 2]) and the frontopolar cortex (FPC, $p < 0.05$ in cluster-level corrected), fully consistent with
21 previous results^{25,29} (Figure 4, Supplementary Table S3).

22 Next, we tested the brain areas implicated in value computation. The chosen value of the
23 model-based system (Q_{MB}) was found to be encoded in the precentral gyrus ($p < 0.05$ in cluster-
24 level corrected) and the orbital and medial PFC ($p < 0.05$ in SVC at [-14,36,-8])^{28,30}. The chosen
25 value of the MF system (Q_{MF}) was found in the dorsal ACC, supplementary motor area,
26 premotor cortex, dorsolateral PFC ($p < 0.05$ in cluster-level corrected), and dorsomedial PFC
27 ($p < 0.05$ in SVC at [9,35,40])^{31,32}. Notably, this model-free value signal was also found in the
28 posterior putamen, the brain area known to be involved in valuation for habitual learning
29 ($p < 0.05$ in SVC at [-27,-4,1])³²⁻³⁴. We also tested for the integrated value signal, expressed as
30 a sum of the value estimates of the model-based and model-free systems weighted by the
31 arbitration control signal (P_{MB}). The ventromedial PFC was positively correlated with the
32 difference between the integrated value signals for the chosen and unchosen actions ($p < 0.05$
33 in cluster-level corrected), fully consistent with previous reports on choice values³⁵⁻³⁷ (Figure

1 4, Supplementary Table S4).

2 Unlike the previous arbitration hypothesis²⁵, our computational model also predicted that
3 arbitration control influences how integrated values are converted into actual choices. Finally,
4 we attempted to identify the brain regions involved in value–action conversion. We found that
5 the inferior parietal lobe, insula ($p < 0.05$ in FWE corrected), middle frontal gyrus, globus
6 pallidus, FPC, supplementary motor area, and thalamus ($p < 0.05$ in cluster-level corrected) are
7 positively correlated with the probability value of the chosen action, referred to as the output
8 value of the softmax function³⁸. This finding is consistent with previous findings indicating
9 stochastic action selection³⁹. Other brain areas, such as the orbitofrontal cortex, superior
10 temporal gyrus, middle frontal gyrus, supramarginal gyrus ($p < 0.05$ in FWE corrected), medial
11 PFC, superior frontal gyrus ($p < 0.05$ in cluster-level corrected), posterior medial cortex, and
12 lateral PFC, are negatively correlated with the chosen action probability. The negative
13 encoding of stochastic action selection in the posterior medial cortex and the lateral PFC also
14 replicates previous findings, as these regions have been implicated in the valuation of
15 counterfactual choices^{35,39,40} (Figure 4, Supplementary Table S5).

16

17 **Effects of Depression on Goal-directed and Habitual Learning**

18

19 [Figure 5] Parametric effect of depression on goal-directed and habitual learning

20

21 To fully explore how depression affects the neural computations underlying model-based
22 and model-free RL, we examined the relationship between the individual depression score
23 and neural representations in each brain region implicated in model-based and model-free
24 RL. We found evidence indicating the effect of depression on RL in multiple brain areas
25 encoding prediction errors. The correlation coefficient between left insula activation and the
26 SPE ($[-36, 20, -4]$, $z = 4.49$), which represents the efficiency of neural encodings of the SPE in
27 the left insula, was inversely proportional to the depression score (estimated correlation
28 coefficient = -0.396 , $p = 0.037$; Figure 5a). We also found a significant negative correlation
29 between the neural efficiency for encoding the RPE in the bilateral caudate (left: $[-4, 6, -4]$,
30 $z = 4.50$, right: $[4, 8, -4]$, $z = 5.30$) and the individual depression score (estimated correlation
31 coefficient = $-0.412/-0.376$, $p = 0.029/0.049$ for the left and right, respectively; Figure 5b).

32 These findings directly demonstrate how depression affects value updates for goal-directed
33 and habitual learning.

1

2 **Effects of Depression on Prefrontal Arbitration Control**

3

4 [Figure 6] Parametric effect of depression on prefrontal arbitration control

5

6 We also tested for the neural effects of depression on arbitration control. First, we found
7 that the individual depression score is significantly correlated with the learning rate for MF
8 reliability estimation, the key variable required to quantify the reliability of predictions made
9 by the model-free RL strategy based on the RPE (estimated correlation coefficient=0.335,
10 $p=0.007$; Figure 6a, left; Supplementary Figure S1). This offers a theoretical prediction that
11 depression entails over-sensitivity to the RPE, making arbitration control more sensitive to
12 the prediction of the model-free system (Figure 6a, right). This implies that the brain regions
13 implicated in mediating arbitration control focus on information about the reliability of the
14 model-free system, rather than encoding the reliability information of the RL system
15 controlling behavior at the moment (i.e. max reliability), leading to the disruption of neural
16 computation underlying normal arbitration control.

17 To test the prediction that depression disrupts neural processing pertaining to arbitration
18 control, we conducted a GLM analysis with the max reliability signal. We found that the effect
19 size (parameter estimates from the GLM analysis) of the max reliability signal for FPC
20 ([8,44,40], $z=3.83$) was negatively correlated with the depression score (estimated
21 correlation coefficient=-0.441, $p=0.019$; Figure 6b). Moreover, the effect sizes of the max
22 reliability signal for the bilateral iLPFC ([-52,26,16], z -score=4.48; [42,20,-8], z -score=3.61)
23 and FPC ([8,44,40], z -score=3.83), the brain areas previously implicated in arbitration
24 control^{25,29,41}, tended to be lower in the depressive group (CES-D score \geq 16) than in the
25 control group (CES-D score $<$ 16) (Figure 6c, left; one-way ANOVA; $F_{1,26}=2.76$ [$p=0.108$],
26 $F_{1,26}=2.93$ [$p=0.099$], $F_{1,26}=5.18$ [$p=0.031$] for the left and right iLPFC and FPC, respectively).

27 To further evaluate the prediction that depression makes neural processing for arbitration
28 control more sensitive to the model-free system, we ran an MVPA for the bilateral iLPFC and
29 FPC. This analysis quantifies the amount of information concerning the reliability of the
30 model-free system embedded in these brain areas. We used a support vector machine, an
31 optimal neural network for prediction and generalization, to conduct a binary classification of
32 model-free reliability (high vs. low; upper/lower 33rd percentile threshold) and compared the
33 prediction performance of the control and depression groups.

1 We found that the prediction performance of model-free reliability in the bilateral ilPFC and
2 FPC was significantly higher in the depression group than in the control group (Figure 6c,
3 right; one-way ANOVA; $F_{1,26}=4.57$ [$p=0.042$], $F_{1,26}=4.33$ [$p=0.047$], $F_{1,26}=4.89$ [$p=0.036$] for
4 the left and right ilPFC and FPC, respectively). On the other hand, the same analysis found
5 no significant inter-group differences in model-based reliability signal or max reliability signal
6 (Supplementary Table S6). Taken together, these results strongly support our arbitration
7 control hypothesis that depression is associated with increased sensitivity to RPE, leading to
8 instable arbitration in which the reliability of predictions of the model-free system becomes
9 predominant and the reliability of predictions of the model-based system becomes less
10 influential.

12 **Effects of Depression on Value–Action Conversion**

14 [Figure 7] Parametric effect of depression on value-action conversion

16 Our computational model also explains how the degrees of control allocated to the model-
17 based and model-free systems influence how value is converted into actual choice
18 (exploitation sensitivity). For example, exploitative and explorative choices are associated with
19 high and low exploitation sensitivity, respectively. We found that the exploitation parameter for
20 model-based RL is negatively correlated with the individual depression score (correlation
21 coefficient= -0.412 , $p=0.001$; Figure 7a, left), indicating that subjects with higher depression
22 scores exhibit more exploratory choices when their choices are guided by model-based RL
23 (Figure 7a, right).

24 In the subsequent neural analysis, we explored the relationships between the depression
25 score and the neural representations of value–action conversion. The parameter estimates of
26 the probability value of taking the chosen action are significantly correlated with the depression
27 score in two brain regions. The parameter estimates in the two seed regions—the right
28 superior temporal gyrus (STG; $[-46,-20,-8]$, $z=4.21$) and left middle temporal gyrus (MTG; $[56,-$
29 $32,-8]$, $z=3.55$)—are positively and negatively correlated with the CES-D score, respectively
30 (correlation coefficient= 0.430 , $p=0.025$ for STG; correlation coefficient= -0.430 , $p=0.022$ for
31 MTG; Figure 7b).

1 **DISCUSSION**

2 By combining a computational model allowing for sub-optimality in learning and decision-
3 making, a model-based fMRI analysis, and an MVPA, the present study fully characterizes
4 how depression influences the different levels and stages of RL: value computation, prefrontal
5 arbitration control for value integration, and value–action conversion. We found that
6 depression has a parametric effect on neural representations of prediction error for model-
7 based and model-free systems, respectively, explaining how depression hampers value
8 computation ability. Another intriguing finding is that the brain areas implicated in arbitration
9 control, bilateral iPFC and FPC, become more sensitive to the predictions of the model-free
10 system in people with depression, indicating that depression disrupts the balance between
11 goal-directed and habitual control. We also found that depression increases the tendency to
12 make exploratory choices during model-based control, but not during model-free control.

13

14 **Computational Theory of Prefrontal Goal-directed and Habitual Control**

15 The present study's computational model of dynamic arbitration control of model-based and
16 model-free RL allows us to explore the full parametric effects of depression on prefrontal goal-
17 directed and habitual control. Although mounting evidence suggests that prediction uncertainty
18 might be a key variable for prefrontal arbitration control^{6,25,27}, little is known about the
19 computational reasons people with depression tend to exhibit behavioral biases towards either
20 goal-directed or habitual behavior. Addressing this issue involves a few challenges. First,
21 simply evaluating the two separate hypotheses contradicts the prevailing view that the brain
22 circuitries guiding goal-directed and habitual behavior interact with each other. Second, there
23 is no guarantee that a rational arbitration control model is flexible enough to explain the
24 individual variability associated with depression. Third, exploring a depression-specific model
25 based on the assumption that depression follows a computational regime that substantially
26 deviates from rational decision-making may enable us to explain severe depression, but
27 cannot explain a continuum extending from a normal to a severely depressed state.

28 To fully address these issues, we considered a computational model of dynamic arbitration
29 control allowing for individual variability in suboptimal learning and decision making.
30 Intriguingly, we found that the influence of prediction uncertainty is not confined to value
31 integration²⁵, but extends as far as value–action conversion. This also allowed us to test the
32 full effect of depression on decision-making at different computation levels: model-based and
33 model-free reinforcement learning, arbitration control for value integration and value–action
34 conversion.

1

2 **Effects of Depression on Neural Representations of Prediction Error**

3 Our study found that depression has a parametric effect on the neural representations of the
4 two distinct types of prediction errors associated with model-based and model-free RL: SPE
5 and RPE. The neural analysis revealed that depression scores were correlated with an
6 attenuation of the SPE signal in the left insula and the RPE in the bilateral caudate.

7 Dopamine is crucial for both model-free and model-based RL. Numerous previous studies
8 have reported dopamine's role in guiding the RPE^{42,43}, and a recent finding discussed the
9 essential role of dopamine in stimulus-stimulus associative learning⁴⁴, implicating the
10 involvement of dopamine in SPE representation. Depression is characterized by decreased
11 dopamine levels^{45,46}, which may impair learning in both model-free and model-based systems.
12 In fact, RPE signals in depression have reportedly been reduced in various experimental
13 conditions (both Pavlovian learning¹² and instrumental learning¹³⁻¹⁷). Our study not only
14 corroborates previous findings concerning RPE deficits in depression, but also further
15 suggests that depression may impact neural representations of SPE.

16

17 **Effects of Depression on Prefrontal Arbitration Control**

18 One interesting prediction of the model is that CES-D is positively correlated with the
19 parameter value for controlling the learning rate for updating model-free reliability based on
20 the RPE, indicating that reliability estimation for the model-free strategy is very sensitive to
21 RPE changes. This suggests that, in people with high CES-D scores, arbitration control might
22 be predominantly driven by the model-free reliability signal, rather than by a fair comparison
23 of the model-free and model-based reliability signals. We explored this possibility through a
24 combination of a general linear model (GLM) and MVPA.

25 Our GLM analysis showed the negative effect of the depression on the neural representations
26 of arbitration control in the prefrontal cortex. In bilateral inferior lateral PFC and frontopolar
27 cortex, the brain areas reportedly encoding the key variable for arbitration control²⁵, neural
28 representations tend to be weaker in the high CES-D score group. Critically, the subsequent
29 MVPA shows that the amount of reliability information of the model-free system is significantly
30 higher in the high CES-D score group. Taken together, these findings theoretically implicate
31 that depression may hinder the PFC's ability to estimate the reliability of each learning strategy
32 from the corresponding prediction error.

33

1 **Effects of Depression on Valuation–Action Conversion**

2 The present study also provides a computational and neural account of how depression
3 causes sub-optimal action selection. Our computational model predicts that depression
4 increases the tendency to make exploratory choices during model-based control, rather than
5 model-free control.

6 This finding could also clarify the two conflicting views of choice consistency behaviors in
7 depression^{47,48}. Beever et al. (2013) found no significant difference in exploration pattern in a
8 reward-maximizing task between a normal and a depressed group. Blanco et al. (2013), on
9 the other hand, found that a depression group tended to explore more. This conflict might be
10 attributable to differences in task structure. Beever et al.'s study used a task with a relatively
11 stable environmental structure, such that people performed tasks relying on the model-free
12 system. This is consistent with our view that depression has a relatively weak influence on
13 exploration during model-free control. However, the reward structure used in Blanco's study
14 encouraged more frequent policy changes, accommodating the need for model-based control.
15 This is also consistent with our view that exploratory choice behavior becomes more
16 pronounced during model-based control.

17 The neural results of the present study, which show that the STG response is higher in people
18 with depression, are fully consistent with previous finding that STG response increases when
19 people switch to other options rather staying⁴⁹. Our results address not only the implication for
20 the role of STG in exploration, but also how depression influences exploration at the neural
21 level.

22 Our finding that the degree of exploitation decreases as CES-D score increases (shown in
23 Figure 7) explains why reward sensitivity is reduced in people with depression. A decreasing
24 degree of exploitation decreases the tendency to convert a learned policy into an actual choice,
25 reducing the efficiency of translating changes in the reward structure into changes in actual
26 choices. This is also consistent with the view that our model's degree of exploitation parameter
27 can be interpreted as reward sensitivity^{16,50}. In addition to supporting existing evidence of
28 declined reward sensitivity in depression^{51–53}, the present study advances the view by
29 proposing that this tendency becomes stronger during model-based control.

30

31 **Potential clinical applications**

32 The present findings suggest how depression influences goal-directed and habitual control
33 in the prefrontal–striatal circuitry. The full characterization of the effects of depression on

1 different stages of learning and decision-making creates possibilities for various clinical
2 applications. First, our neural results explain why such brain stimulus techniques as repetitive
3 transcranial magnetic stimulation (rTMS) and deep brain stimulus (DBS) to the frontopolar
4 cortex⁵⁴ or ventral striatum⁵⁵ are effective in alleviating depressive symptoms. Second, our
5 theoretical idea suggests that behavioral therapy to reduce sensitivity to reward prediction
6 errors might help people with depression regain a balance between goal-directed and habitual
7 control. Intriguingly, our findings also indicate a parametric effect of depression on learning
8 and decision-making in a relatively young age group (average=22.8 yrs). Considering the
9 onset of MDD is approximately 25 to 45 yrs⁵⁶, our study offers possibilities for not only
10 investigating how mild depression transitions to MDD, but also developing clinical applications
11 for the early diagnosis of MDD.

12

13 **METHODS**

14 **Participants**

15 Sixty-five right-handed Koreans (28 females; mean age of 22.8 ± 3.8) participated in the study.
16 Participants were recruited from the local society through the online announcement. Only 28
17 subjects were scanned with fMRI during the task. Two subjects whose total accumulated
18 reward are below the chance-level (mean amount of rewards with 10,000 random simulations)
19 were excluded from the analysis. Thus, a total of sixty-three behavioral data and twenty-eight
20 fMRI data were left for the analysis. To acquire the depressive level of individuals, people were
21 instructed to complete the Center for Epidemiologic Studies Depression (CES-D)
22 questionnaires²⁶ before the experiment (For the distribution of participant's depression severity,
23 see Supplementary Figure S2).

24 No subjects had the history of neurological or psychiatric disease. Every subject provided
25 written consent to the experimental protocols which were approved by the Institutional Review
26 Board (IRB) of the Korea Advanced Institute of Science and Technology (KAIST).

27

28 **Task**

29 We used the sequential two-stage Markov decision task proposed to dissociate model-based
30 and model-free learning strategy²⁵. In this task, subjects make a binary choice (either left or
31 right) and proceed to the next state with a certain probability. When the next state is appeared
32 in the screen, participants make another choice. The two consecutive choices is followed by
33 a transition to an outcome state. Subjects perform 100 trials in the pre-learning session to

1 learn the structure of the task. Four main sessions with 80 trials on average follow the pre-
2 learning session. Participants are instructed to collect as many coins as possible in the main
3 sessions.

4 The task consists of two conditions: a specific-goal condition and flexible-goal condition. The
5 goal condition is indicated by a color of a box at the beginning of each trial. In the specific-goal
6 condition, participants are presented with a box with a specific color (red, blue, yellow). A
7 monetary reward is given only when the coin color matches with the color of the given box. In
8 the flexible-goal condition, on the other hand, participants are given a white coin box with
9 which all types of coins become redeemable. Two types of state-transition probability are used
10 to control the uncertainty of the environment. The state-transition probability (0.5, 0.5) and (0.9,
11 0.1) is intended to implement the highly-uncertain and relatively less uncertain environment,
12 respectively. The four types of block (2 goal-conditions x 2 uncertainty conditions) are
13 presented in pseudo-random order. Each block consists of 4-6 trials.

14

15 **Computational Models**

16 The computational model of this study is motivated by the previous arbitration control
17 hypothesis that prediction uncertainty of the model-based and model-free RL is a key variable
18 to guide value integration of the two corresponding systems²⁵. The model consists of the three
19 processes: value learning, arbitration, and action selection (Figure 3a).

20 In the value learning stage, both a model-based and model-free system learn action values
21 for each state. A model-based system uses state prediction error (SPE = 1-expected transition
22 probability) to update the state-action-state transition probability, by using a FORWARD
23 learning²⁷ and learns action values by combining the learned state-action-state transition
24 probability and reward in the outcome state. For a model-free system, on the other hand, the
25 state-action value learning is based on RPE (RPE = actual value-expected value). It is
26 implemented with a SARSA algorithm⁵.

27 In the arbitration process, the reliability estimation of the model-based system was
28 implemented with a hierarchical empirical Bayes method using the history of the SPE, the
29 reliability of the model-free system was implemented with the Pearce-hall associability rule
30 using an unsigned RPE. These estimated reliability signal were then used to guide the
31 competition between the two systems, which is implemented with a dynamic two-state
32 transition model. The output of this model is a model choice probability (P_{MB}), used as the
33 control weight for value integration of the two systems. Finally, in the action selection stage,

1 the model selects the action stochastically using softmax rule³⁸. For more details, refer Lee et
2 al (2014).

3 In this study, we suggested two variants of arbitration control, allowing for sub-optimality in
4 value learning, arbitration, and action selection: one version with separate model-based and
5 model-free learning and another version with a dynamic exploitation. The former type of the
6 model assumes the different learning rates of a model-based and a model-free system. The
7 latter class of models is based on the former model, with the further assumption that the
8 degree of exploitation, an indicator of optimality of the RL agent's policy, is a function of the
9 model choice probability, P_{MB} . We tested three different types of exploitation as follows:
10 logistics, linear, weighted linear. Note that in all cases, we set the parameters of the model in
11 a way that is reduced to the original RL with a single exploitation parameter.

12

13 **fMRI Data Acquisition**

14 Functional imaging was conducted on a 3T Siemens (Magnetom) Verio scanner located in
15 the KAIST brain imaging center (Daejeon). Forty-two axial slices were acquired with
16 interleaved-ascending order at the resolution of 3 mm x 3 mm x 3 mm, covering the whole
17 brain. A one-shot echo-planner imaging pulse sequence was utilized (TR = 2800 ms; TE = 30
18 ms; FOV = 192 mm; flip angle = 90°). The high resolution structural image was also acquired
19 for each subject to the resolution of 0.7 mm X 0.7 mm x 0.7 mm.

20

21 **fMRI Data Pre-processing**

22 Images were processed and analyzed using the SPM12 software (Wellcome Department of
23 Imaging Neuroscience, London, UK). The first two volumes were removed to reduce T1
24 equilibrium effects. The EPI images were corrected for slice timing, motion movement and
25 spatially normalized to the standard template imaging provided by SPM software.

26 For the general linear model analysis (GLM), normalized images were smoothed with 6mm
27 FWHM Gaussian Kernel and a high-pass filter (128s cut-off) was applied to remove the noise.

28 For the multivoxel pattern analysis (MVPA), unsmoothed EPI image data was used. De-
29 trending and z-scoring were processed to reduce the linear trends and to match the range of
30 the signal.

31

1 **General Linear Model Analysis (GLM)**

2 Subject-specific value-related signals and arbitration control signals were computed from the
3 arbitration model, and the signals were regressed against voxel-wise signals from the EPI
4 image set. The order of the regressors is as follows: prediction error from the model-based
5 system (SPE) and the model-free system (RPE), reliability comparison signal which is a key
6 variable for arbitration control ($=\max(\text{Rel}_{\text{MB}}, \text{Rel}_{\text{MF}})$; max reliability), the chosen value of model-
7 based system (Q_{fwd}), the chosen value of model-free system (Q_{sarsa}), the difference between
8 chosen and unchosen integrated values (Q_{arb}) and the probability of selecting chosen action
9 ($P_{\text{chosen action}}$). The regressors were serially non-orthogonalized in the GLM analysis to prevent
10 the effect of regressor orders in the interpretation of the results. MARSBAR software
11 (<http://marsbar.sourceforge.net>) was used to extract parameter estimates from the region of
12 interest⁵⁷.

13

14 **Multivoxel Pattern Analysis (MVPA)**

15 The MVPA analysis was conducted to quantify types and amounts of information encoded in
16 specific the region of interest (ROI). The classification performance is regarded as the amount
17 of information pertaining to the variable of interest. Three ROIs, left/right inferior lateral
18 prefrontal cortex (ilPFC) and Frontopolar prefrontal cortex (FPC), were selected, which were
19 known to engage in the arbitration control process. Masks of each brain region were
20 functionally defined from the GLM analysis. We used the clusters whose response to the Max
21 reliability signal survived after the whole-brain correction ($p < 0.001$, uncorrelated) as a mask
22 for each ROI. We set the BOLD response time 4-6s.

23 A binary Support Vector Machine (SVM) classifier was applied to learn voxels patterns with
24 each ROI. For each subjects data, the SVM was trained to best match its output to a binarized
25 reliability-related signal (MB reliability, MF reliability, or Max reliability); the 33th and 67th
26 percentile threshold were used to define the two classes, 'high value group' and 'low value
27 group', respectively. All voxels in the mask were used for learning. The input dimension is 196,
28 79, 294 for left ilPFC, right ilPFC and FPC, respectively. The average number of data from
29 each subject are 350 for MB reliability, 549 for MF reliability and 543 for Max Reliability. Thirty-
30 fold cross validation was conducted for evaluation. All processes were implemented based on
31 the Princeton Multi-Voxel Pattern Analysis toolbox⁵⁸. Finally, an ANOVA analysis was
32 conducted to compare signal prediction accuracy between the normal and depression group;
33 the two subject groups were defined by using the standard cutoff criteria of CES-D score,
34 16^{59,60}.

1

2 **Data Availability**

3 All data analyzed in this study will be available upon request.

4

5 **ACKNOWLEDGEMENTS**

6 This research was supported by the Brain Research Program through the National Research
7 Foundation of Korea (NRF) funded by the Ministry of Science, ICT & Future Planning (NRF-
8 2016M3C7A1914448), NRF funded by the Korea government (MSIT) (No. NRF-2017R1C 1B
9 2008972), the research fund of the KAIST (Korea Advanced Institute of Science and
10 Technology) (Grant code: G04150045), and Institute for Information & Communications
11 Technology Promotion (IITP) grant funded by the Korea government (No. 2017-0-00451).

12

13 **Author Contributions**

14 S.H. and S.W.L. designed the study, analyzed the data and wrote the manuscript. S.H.
15 conducted the experiments.

16

17 **Competing interests**

18 No conflict of interests.

19

20 **References**

- 21 1. Lim, G. Y. *et al.* Prevalence of Depression in the Community from 30 Countries
22 between 1994 and 2014. *Sci. Rep.* **8**, 2861 (2018).
- 23 2. Leykin, Y., Roberts, C. S. & DeRubeis, R. J. Decision-Making and Depressive
24 Symptomatology. *Cognit. Ther. Res.* **35**, 333–341 (2011).
- 25 3. Must, A., Horvath, S., Nemeth, V. L. & Janka, Z. The Iowa Gambling Task in
26 depression - what have we learned about sub-optimal decision-making strategies?
27 *Front. Psychol.* **4**, 732 (2013).
- 28 4. Admon, R. & Pizzagalli, D. A. Dysfunctional Reward Processing in Depression. *Curr.*
29 *Opin. Psychol.* **4**, 114–118 (2015).

- 1 5. Sutton, R. S. & Barto, A. G. *Introduction to Reinforcement Learning*. (MIT Press,
2 1998).
- 3 6. Daw, N. D., Niv, Y. & Dayan, P. Uncertainty-based competition between prefrontal
4 and dorsolateral striatal systems for behavioral control. *Nat. Neurosci.* **8**, 1704–1711
5 (2005).
- 6 7. Doya, K., Samejima, K., Katagiri, K. & Kawato, M. Multiple Model-Based
7 Reinforcement Learning. *Neural Comput.* **14**, 1347–1369 (2002).
- 8 8. Markman, K. & Miller, A. *Depression, Control, and Counterfactual Thinking:
9 Functional for Whom? Journal of Social and Clinical Psychology* **25**, (2006).
- 10 9. Quelhas, A. C., Power, M. J., Juhos, C. & Senos, J. Counterfactual thinking and
11 functional differences in depression. *Clin. Psychol. Psychother.* **15**, 352–365 (2008).
- 12 10. Huys, Q. J. M., Daw, N. D. & Dayan, P. Depression: a decision-theoretic analysis.
13 *Annu. Rev. Neurosci.* **38**, 1–23 (2015).
- 14 11. Maddox, W. T. *et al.* Elevated Depressive Symptoms Enhance Reflexive but not
15 Reflective Auditory Category Learning. *Cortex.* **58**, 186–198 (2014).
- 16 12. Kumar, P. *et al.* Abnormal temporal difference reward-learning signals in major
17 depression. *Brain* **131**, 2084–2093 (2008).
- 18 13. Gradin, V. B. *et al.* Expected value and prediction error abnormalities in depression
19 and schizophrenia. *Brain* **134**, 1751–1764 (2011).
- 20 14. Dombrovski, A. Y., Szanto, K., Clark, L., Reynolds, C. F. & Siegle, G. J. Reward
21 signals, attempted suicide, and impulsivity in late-life depression. *JAMA psychiatry* **70**,
22 1 (2013).
- 23 15. Ubl, B. *et al.* Altered neural reward and loss processing and prediction error signalling
24 in depression. *Soc. Cogn. Affect. Neurosci.* **10**, 1102–1112 (2015).
- 25 16. Rothkirch, M., Tonn, J., Kohler, S. & Sterzer, P. Neural mechanisms of reinforcement
26 learning in unmedicated patients with major depressive disorder. *Brain* **140**, 1147–
27 1157 (2017).
- 28 17. Kumar, P. *et al.* Impaired reward prediction error encoding and striatal-midbrain
29 connectivity in depression. *Neuropsychopharmacology* **43**, 1581–1588 (2018).
- 30 18. Chase, H. W. *et al.* Approach and avoidance learning in patients with major

- 1 depression and healthy controls: relation to anhedonia. *Psychol. Med.* **40**, 433–440
2 (2010).
- 3 19. Plieger, T., Melchers, M., Montag, C., Meermann, R. & Reuter, M. Life stress as
4 potential risk factor for depression and burnout. *Burn. Res.* **2**, 19–24 (2015).
- 5 20. Radley, J. J. *et al.* STRESS RISK FACTORS AND STRESS-RELATED
6 PATHOLOGY: NEUROPLASTICITY, EPIGENETICS AND ENDOPHENOTYPES.
7 *Stress* **14**, 481–497 (2011).
- 8 21. Dias-Ferreira, E. *et al.* Chronic stress causes frontostriatal reorganization and affects
9 decision-making. *Science* **325**, 621–625 (2009).
- 10 22. Radenbach, C. *et al.* The interaction of acute and chronic stress impairs model-based
11 behavioral control. *Psychoneuroendocrinology* **53**, 268–280 (2015).
- 12 23. Schwabe, L. & Wolf, O. T. Stress-induced modulation of instrumental behavior: from
13 goal-directed to habitual control of action. *Behav. Brain Res.* **219**, 321–328 (2011).
- 14 24. Otto, A. R., Raio, C. M., Chiang, A., Phelps, E. A. & Daw, N. D. Working-memory
15 capacity protects model-based learning from stress. *Proc. Natl. Acad. Sci.* **110**, 20941
16 LP-20946 (2013).
- 17 25. Lee, S. W., Shimojo, S. & O’Doherty, J. P. Neural computations underlying arbitration
18 between model-based and model-free learning. *Neuron* **81**, 687–699 (2014).
- 19 26. Radloff, L. S. The CES-D Scale: A Self-Report Depression Scale for Research in the
20 General Population. *Appl. Psychol. Meas.* **1**, 385–401 (1977).
- 21 27. Glascher, J., Daw, N., Dayan, P. & O’Doherty, J. P. States versus rewards:
22 dissociable neural prediction error signals underlying model-based and model-free
23 reinforcement learning. *Neuron* **66**, 585–595 (2010).
- 24 28. O’Doherty, J. P., Cockburn, J. & Pauli, W. M. Learning, Reward, and Decision
25 Making. *Annu. Rev. Psychol.* **68**, 73–100 (2017).
- 26 29. Bogdanov, M., Timmermann, J. E., Gläscher, J., Hummel, F. C. & Schwabe, L.
27 Causal role of the inferolateral prefrontal cortex in balancing goal-directed and
28 habitual control of behavior. *Sci. Rep.* **8**, 9382 (2018).
- 29 30. Wunderlich, K., Dayan, P. & Dolan, R. J. Mapping value based planning and
30 extensively trained choice in the human brain. *Nat. Neurosci.* **15**, 786–791 (2012).

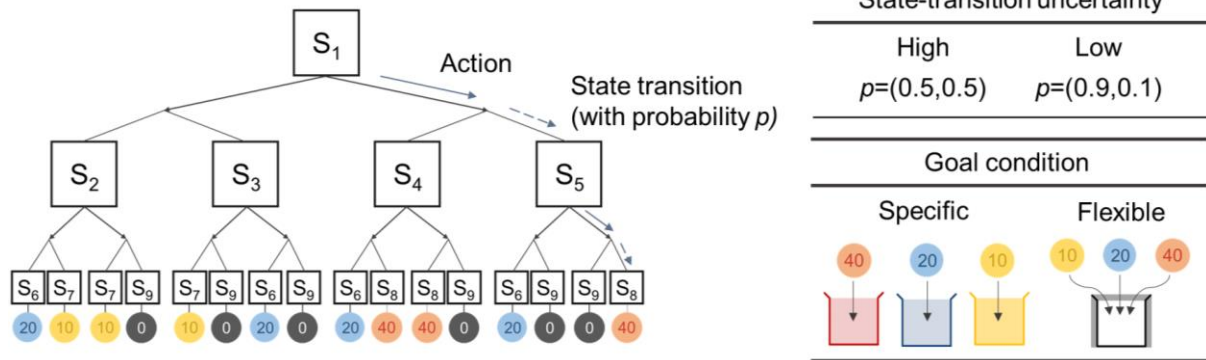
- 1 31. Rowe, J. B., Hughes, L. & Nimmo-Smith, I. Action selection: A race model for selected
2 and non-selected actions distinguishes the contribution of premotor and prefrontal
3 areas. *Neuroimage* **51**, 888–896 (2010).
- 4 32. Hare, T. A., Schultz, W., Camerer, C. F., O’Doherty, J. P. & Rangel, A.
5 Transformation of stimulus value signals into motor commands during simple choice.
6 *Proc. Natl. Acad. Sci. U. S. A.* **108**, 18120–18125 (2011).
- 7 33. Tricomi, E., Balleine, B. W. & O’Doherty, J. P. A specific role for posterior dorsolateral
8 striatum in human habit learning. *Eur. J. Neurosci.* **29**, 2225–2232 (2009).
- 9 34. Horga, G. *et al.* Changes in corticostriatal connectivity during reinforcement learning
10 in humans. *Hum. Brain Mapp.* **36**, 793–803 (2015).
- 11 35. Boorman, E. D., Behrens, T. E. J., Woolrich, M. W. & Rushworth, M. F. S. How green
12 is the grass on the other side? Frontopolar cortex and the evidence in favor of
13 alternative courses of action. *Neuron* **62**, 733–743 (2009).
- 14 36. Jocham, G., Hunt, L. T., Near, J. & Behrens, T. E. J. A mechanism for value-guided
15 choice based on the excitation-inhibition balance in prefrontal cortex. *Nat. Neurosci.*
16 **15**, 960 (2012).
- 17 37. Jocham, G. *et al.* Dissociable contributions of ventromedial prefrontal and posterior
18 parietal cortex to value-guided choice. *Neuroimage* **100**, 498–506 (2014).
- 19 38. R Luce, D. *Individual Choice Behavior: A Theoretical Analysis.* New York **115**, (2005).
- 20 39. Daw, N. D., O’Doherty, J. P., Dayan, P., Seymour, B. & Dolan, R. J. Cortical
21 substrates for exploratory decisions in humans. *Nature* **441**, 876–879 (2006).
- 22 40. Boorman, E. D., Behrens, T. E. & Rushworth, M. F. Counterfactual choice and
23 learning in a Neural Network centered on human lateral frontopolar cortex. *PLoS Biol.*
24 **9**, (2011).
- 25 41. Cole, M. W., Repovs, G. & Anticevic, A. The frontoparietal control system: a central
26 role in mental health. *Neuroscientist* **20**, 652–664 (2014).
- 27 42. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and
28 reward. *Science* **275**, 1593–1599 (1997).
- 29 43. Hollerman, J. R. & Schultz, W. Dopamine neurons report an error in the temporal
30 prediction of reward during learning. *Nat. Neurosci.* **1**, 304–309 (1998).

- 1 44. Sharpe, M. J. *et al.* Dopamine transients are sufficient and necessary for acquisition
2 of model-based associations. *Nat. Neurosci.* **20**, 735–742 (2017).
- 3 45. Dunlop, B. W. & Nemeroff, C. B. The role of dopamine in the pathophysiology of
4 depression. *Arch. Gen. Psychiatry* **64**, 327–337 (2007).
- 5 46. Malhi, G. S. & Berk, M. Does dopamine dysfunction drive depression? *Acta Psychiatr.*
6 *Scand.* **115**, 116–124 (2007).
- 7 47. Blanco, N. J., Otto, A. R., Maddox, W. T., Beevers, C. G. & Love, B. C. The influence
8 of depression symptoms on exploratory decision-making. *Cognition* **129**, 563–568
9 (2013).
- 10 48. Beevers, C. G. *et al.* Influence of depression symptoms on history-independent
11 reward and punishment processing. *Psychiatry Res.* **207**, 53–60 (2013).
- 12 49. Paulus, M. P., Feinstein, J. S., Leland, D. & Simmons, A. N. Superior temporal gyrus
13 and insula provide response and outcome-dependent information during assessment
14 and action selection in a decision-making situation. *Neuroimage* **25**, 607–615 (2005).
- 15 50. Huys, Q. J. M., Pizzagalli, D. A., Bogdan, R. & Dayan, P. Mapping anhedonia onto
16 reinforcement learning: a behavioural meta-analysis. *Biol. Mood Anxiety Disord.* **3**, 12
17 (2013).
- 18 51. Steele, J. D., Kumar, P. & Ebmeier, K. P. Blunted response to feedback information in
19 depressive illness. *Brain* **130**, 2367–2374 (2007).
- 20 52. Chen, C., Takahashi, T., Nakagawa, S., Inoue, T. & Kusumi, I. Reinforcement
21 learning in depression: A review of computational research. *Neurosci. Biobehav. Rev.*
22 **55**, 247–267 (2015).
- 23 53. Alloy, L. B., Olino, T., Freed, R. D. & Nusslock, R. Role of Reward Sensitivity and
24 Processing in Major Depressive and Bipolar Spectrum Disorders. *Behav. Ther.* **47**,
25 600–621 (2016).
- 26 54. Downar, J. & Daskalakis, Z. J. New Targets for rTMS in Depression: A Review of
27 Convergent Evidence. *Brain Stimul. Basic, Transl. Clin. Res. Neuromodulation* **6**,
28 231–240 (2013).
- 29 55. Delaloye, S. & Holtzheimer, P. E. Deep brain stimulation in the treatment of
30 depression. *Dialogues Clin. Neurosci.* **16**, 83–91 (2014).
- 31 56. Kessler, R. C. *et al.* Age of onset of mental disorders: a review of recent literature.

- 1 *Curr. Opin. Psychiatry* **20**, 359–364 (2007).
- 2 57. Brett, M., Anton, J.-L. L., Valabregue, R. & Poline, J.-B. Region of interest analysis
3 using an SPM toolbox - Abstract Presented at the 8th International Conference on
4 Functional Mapping of the Human Brain, June 2-6, 2002, Sendai, Japan. *Neuroimage*
5 **16**, Abstract 497 (2002).
- 6 58. Detre, G. J. *et al.* The Multi-Voxel Pattern Analysis (MVPA) toolbox. *Ohbm* (2006).
- 7 59. Comstock, G. W. & Helsing, K. J. Symptoms of depression in two communities.
8 *Psychol. Med.* **6**, 551–563 (1976).
- 9 60. Weissman, M. M., Sholomskas, D., Pottenger, M., Prusoff, B. A. & Locke, B. Z.
10 Assessing depressive symptoms in five psychiatric populations: a validation study.
11 *Am. J. Epidemiol.* **106**, 203–214 (1977).

12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

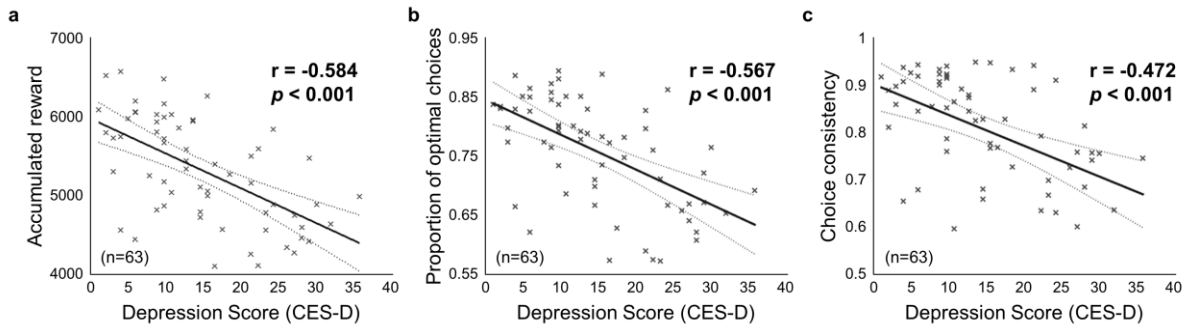
1 **Figures**



2
3 **Figure 1. Markov decision task structure**

4 We use the two-stage Markov decision task proposed by Lee et al. (2014). In each stage,
 5 participants make a binary choice (left or right). After the first choice in the initial state (S_1),
 6 they were moved forward to one of four states in the second stage ($S_2 - S_5$) with certain state-
 7 action-state transition probability p . The transition probability (0.5, 0.5) and (0.9, 0.1)
 8 corresponds to a high-uncertainty and a low-uncertainty environment, respectively. The task
 9 consists of the two goal conditions: a specific-goal and a flexible-goal condition. In the specific
 10 goal condition, subjects can collect coins (redeemable for monetary reward) only if the coin
 11 color matches with the color of the token box (red, blue, yellow). In the flexible goal condition
 12 indicated by the white token box, all types of coins are redeemable.

13
14
15
16
17
18
19
20
21
22
23



1

2

Figure 2. Behavioral results

3

4

5

6

7

8

(a) Relationship between the individual depression score (CES-D) and accumulated reward (n=63). The task performance decreases as the depression score increases. (b) Relationship between depression score and the proportion of optimal choices (n=63). The proportion of optimal choices is inversely proportional to the depression score. (c) Relationship between the depression score and choice consistency in the first state (S_1) (n=63). The choice consistency index is negatively correlated with the depression score.

9

10

11

12

13

14

15

16

17

18

19

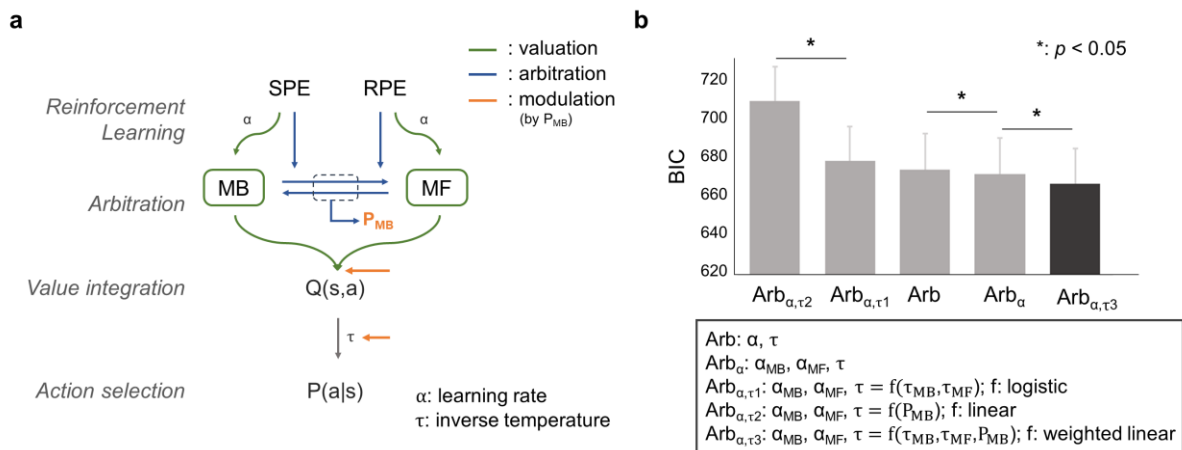
20

21

22

23

24



1

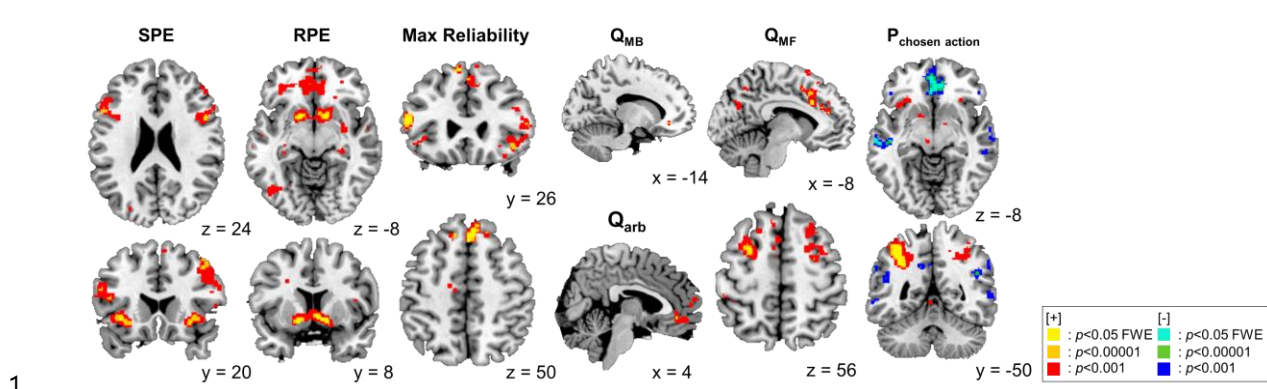
2 **Figure 3. Computational model of dynamic arbitration control**

3 (a) Computational model to investigate dynamic control mechanisms to arbitrate between
 4 model-based (MB) and model-free (MF) RL. Two separate RL systems update an action value
 5 and reliability of its prediction by using prediction errors (SPE and RPE, respectively; green).
 6 The reliability values of the two systems were then used to compute the model choice
 7 probability (P_{MB}) (blue). The model choice probability guides both the value integration and
 8 value-action conversion process; both effects are indicated by the orange arrow. (b) Model
 9 comparison analysis. We used Bayesian Information Criteria (BIC) for comparing the
 10 goodness of fit while penalizing for the model complexity. Arb refers to the computational
 11 model proposed in (Lee et al., 2014). Additional versions of arbitration control consider
 12 separate learning rates and dynamic exploitation. Arb_{α} assigns separate learning rates to two
 13 different systems (i.e. use α_{MB} , α_{MF} instead of α). $Arb_{\alpha, \tau}$ is the same as Arb_{α} , except for the
 14 assumption that the degree of exploitation is a function of the model choice probability. The
 15 best version of model uses the degree of exploitation (τ) as a weighted sum of the model-
 16 based and the model-free exploitation parameter (τ_{MB} and τ_{MF} , respectively) with the model
 17 choice probability P_{MB} (i.e. $\tau = P_{MB} * \tau_{MB} + (1 - P_{MB}) * \tau_{MF}$). The error bar stands for the
 18 standard error mean.

19

20

21



2 **Figure 4. Neural correlates of dynamic arbitration control**

3 Prediction error, value, reliability, action choice probability signals from the proposed model
4 are shown as colored blobs. SPE and RPE refers to state prediction error and reward
5 prediction error, respectively. Max reliability refers to the reliability of whichever system had
6 the highest reliability index on each trial ($=\max(\text{Rel}_{\text{MB}}, \text{Rel}_{\text{MF}})$). Q_{MB} and Q_{MF} indicate the
7 chosen value from the model-based and model-free system. Q_{arb} refers to the difference
8 between integrated chosen action value and unchosen action value. $P_{\text{chosen action}}$ refers to the
9 probability assigned to the chosen action. See more detailed information in Supplementary
10 Table S3-S5.

11

12

13

14

15

16

17

18

19

20

21

22

23

24

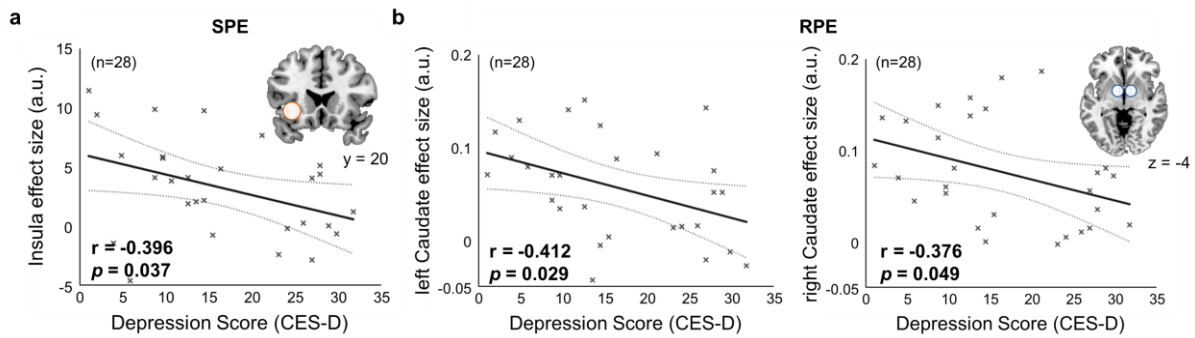


Figure 5. Parametric effect of depression on goal-directed and habitual learning

(a) Depression impacts on neural encoding of SPE information (n=28). The shaded circles represent seed regions for which parameter estimates of the GLM analysis were extracted. The seed region is the left insula, and the parameter estimates were extracted from our GLM analysis which regressed the SPE signal against the BOLD response ($[-36,20,-4]$, z -score=4.49; Figure 4). The estimated effect size for the left insula is negatively correlated with the depression score. (b) Depression impacts on RPE response (n=28). The estimated effect size of RPE for bilateral caudate (left: $[-4,6,-4]$, z -score=4.50, right: $[4,8,-4]$, z -score=5.30) is negatively correlated with the depression score. a.u. stands for arbitrary units.

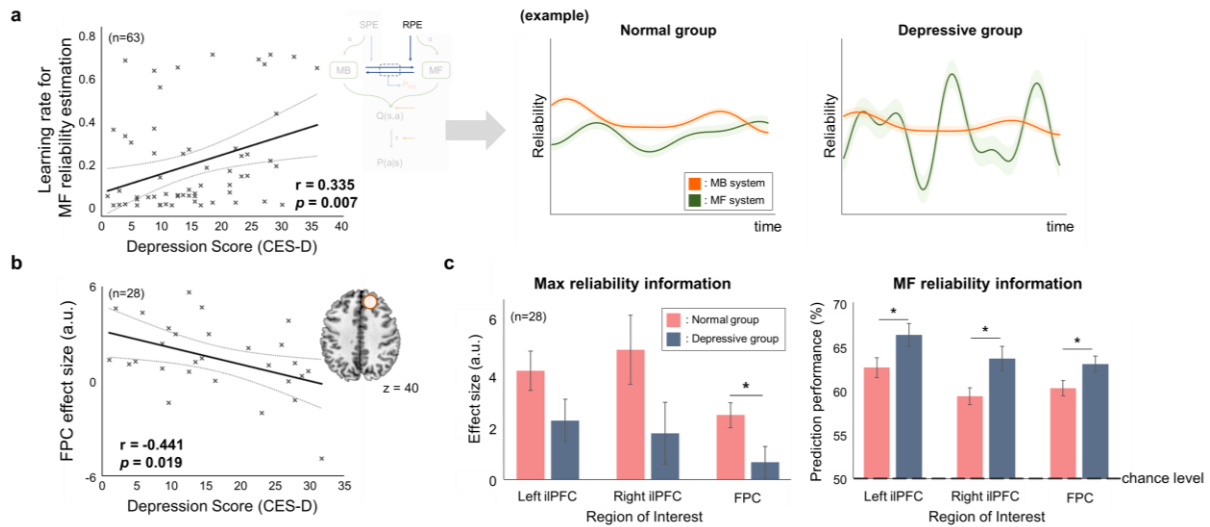
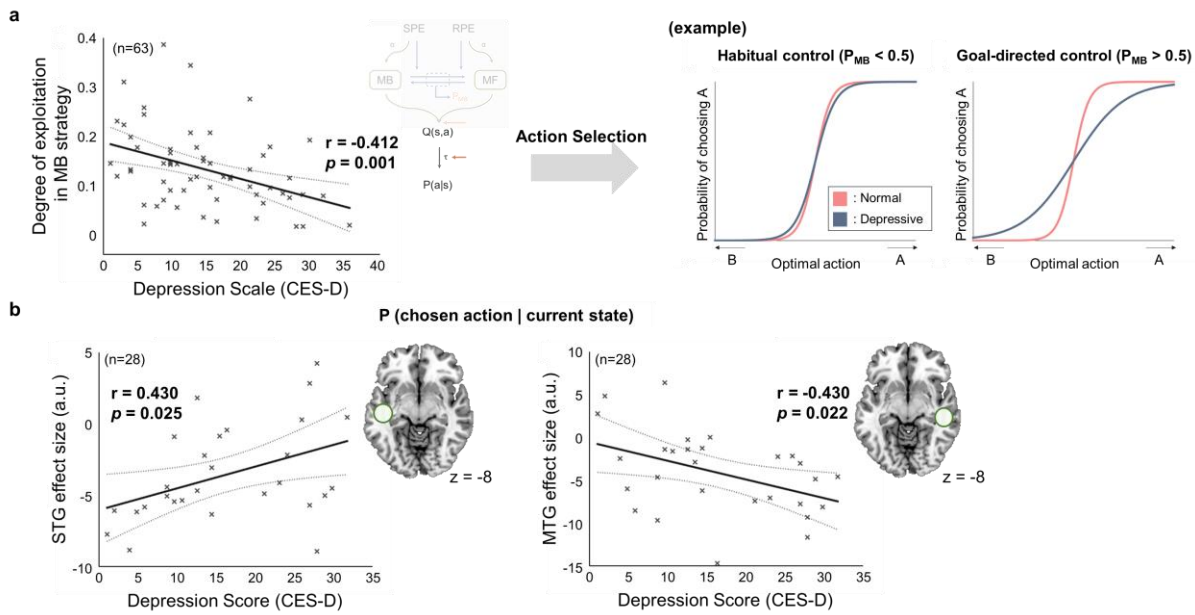


Figure 6. Parametric effect of depression on prefrontal arbitration control

(a) Depression effects on the model-free system’s sensitivity to RPE (n=63; the total number of subjects, including 40 who participated behavioral experiment only and 23 who were also scanned with the fMRI). (Left) Relationship between the depression score and the learning rate parameter for reliability estimation of the model-free system. Individuals with a higher depression score tend to exhibit a higher learning rate, indicating that their reliability estimation for the model-free system is more sensitive to RPE. (Right) Illustrative examples of reliability changes of people with a low (“normal group”) and a high depression score (“depressive group”), each of which is associated with a low and a high learning rate for reliability estimation, respectively. The depressive group shows rapid changes in MF reliability due to higher learning rate. (b) Relationship between the depression score and the estimated effect size of Max reliability for frontopolar cortex (FPC), the brain area implicated in arbitration control (n=28; the total number of subjects scanned with fMRI). The estimated effect size of the Max reliability signal for FPC (coordinates [8,44,40]; z-score=3.83) is negatively correlated with the depression score. (c) Comparison of reliability signal representation performance in normal (n=15) and depressive (n=13) group (GLM and MVPA analysis). The parameter estimates of the Max reliability for the bilateral iIPFC ([-52,26,16], z-score=4.48; [42,20,-8], z-score=3.61) and FPC ([8,44,40], z-score=3.83) tend to decrease in the depressive group. The MVPA analysis with these three seed regions reveals that the amount of information about the model-free reliability was significantly higher in the depressive group in all three regions (one-way ANOVA). Asterisk (*) indicates significant difference at the 0.05 level.



1

2 **Figure 7. Parametric effect of depression on value-action conversion**

3 (a) Depression effects on model parameters (n=63; the total number of subjects, including 40
4 who participated behavioral experiment only and 23 who were also scanned with the fMRI).
5 (Left) Relationship between the depression score and the degree of exploitation in MB strategy
6 (τ_{MB}). The MB exploitation parameter decreases as the depression score increases. (Right)
7 Examples illustrating exploitation parameter effects on action selection. Shown are the
8 softmax functions that convert an action value into a choice probability value, between the
9 normal and the depressive group for goal-directed and habitual control. Compared to the
10 normal group (pink), the depressive group makes more exploratory choices especially as they
11 rely more on the model-based system (blue). (b) Relationship between the depression score
12 and the parameter estimate of the probability of selecting chosen action for the two seed
13 regions, right Superior Temporal Gyrus (STG) ([-46,-20,-8]; z-score=4.21) and the left
14 Middle Temporal Gyrus (MTG) ([56,-32,-8]; z-score=3.55). The effect size of right STG and left
15 MTG increases and decreases with the individual severity of depression, respectively.