1 **Trace Imbalance in Reinforcement and Punishment Systems Can Mis-reinforce Implicit**

2 **Choices Leading to Anxiety**

3

4 Yuki Sakai[a,b,†], Yutaka Sakai[c,†], Yoshinari Abe[b], Jin Narumoto[b], Saori C. Tanaka[a,*]

5

6 [a]ATR Brain Information Communication Research Laboratory Group, 2-2-2 Hikaridai Seika-

7 Cho, Soraku-Gun, Kyoto 619-0288, Japan

8 [b]Department of Psychiatry, Graduate School of Medical Science, Kyoto Prefectural University of

9 Medicine, 465 Kajii-Cho, Kawaramachi-Hirokoji, Kamigyo-Ku, Kyoto 602-8566, Japan

10 [c]Brain Science Institute, Tamagawa University, 6-1-1, Tamagawa-gakuen, Machida, Tokyo 194-

11 8610, Japan

12

13 [†]These authors contributed equally and should be considered co-first authors.

14

15 *Corresponding author: Saori C. Tanaka

16 Computational Neuroscience Laboratories, Advanced Telecommunications Research Institute

17 International, 2-2-2 Hikaridai, Seika-Cho, Soraku-Gun, Kyoto 619-0288, Japan.

18 Tel: +81-774-95-1250; Fax: +81-774-95-1236.

19 E-mail: xsaori@atr.jp

20

1 **Abstract**

2 Nobody wants to experience anxiety. However, anxiety may be induced by our own implicit

3 choices that are mis-reinforced by some imbalance in reinforcement learning. Here we focused

4 on obsessive-compulsive disorder (OCD) as a candidate for implicitly learned anxiety.

5 Simulations in the reinforcement learning framework showed that agents implicitly learn to

6 become anxious when the memory trace signal for past actions decays differently for positive

7 and negative prediction errors. In empirical data, we confirmed that OCD patients showed

8 extremely imbalanced traces, which were normalized by serotonin enhancers. We also used

9 fMRI to identify the neural signature of OCD and healthy participants with imbalanced traces.

10 Beyond the spectrum of clinical phenotypes, these behavioral and neural characteristics can be

11 generalized to variations in the healthy population.

## Introduction

Humans sometimes experience anxiety. Given that no one likes to feel anxious, we tend to believe that anxiety is passively driven by external factors. However, is it possible that our own choices implicitly induce anxiety? We make mental choices when thinking or mind-wandering, as well as in response to the external world. We are always learning appropriate choices through our experiences, and some choices are learned in subliminal situations[1, 2]. Indeed, some evidence suggests that many of our decisions and actions in everyday life are shaped by implicit learning processes[3]. Hence, our anxiety might be induced by our own choices that are reinforced through implicit learning.

When investigating the possibility of implicit choices leading to anxiety, we must consider its theoretical and biological implementation. Implicit learning requires repetitive experiences[1, 2]. Reinforcement learning theory provides a framework for learning appropriate choices through repetitive experiences[4]. Reinforcement learning reinforces or punishes recent choices based on the difference in the actual outcome from the prediction, called the prediction error. A positive prediction error reinforces recent choices, whereas a negative error punishes them. When examining the biological implementation of implicit learning, it is essential to consider previous work demonstrating that the activity of dopamine-projecting neurons in the midbrain resembles prediction error[5]. Dopaminergic neurons widely project to cortical and subcortical areas. One of the main targets of dopamine is the striatum, which is located in the basal ganglia[6]. In the striatum, distinct types of dopamine receptors—D1 and D2 receptors—are expressed in exclusive neuron groups, which are considered to be involved in distinct circuit patterns, namely, direct and indirect pathways[6]. These distinct neural systems are thought to play different roles in reinforcement learning. Artificial activations of D1- and D2-expressing neurons

1     after a certain behavioral action reinforce and punish the action, respectively[7]. Direct and indirect

2     pathway neurons respond to positive and negative outcomes, respectively[8]. Dopamine-dependent

3     synaptic plasticity in corticostriatal synapses, which underlies reinforcement learning, shows

4     opposing dopamine dependence for D1 and D2 neurons[9, 10]. These results imply that

5     reinforcement and punishment of recent actions may be reflected through different neural

6     systems: the direct and indirect pathways.

7        In the implementation of implicit learning, we should also consider delay as well as

8     reinforcement/punishment. In general, the outcome of a certain choice is available after a certain

9     delay. Therefore, reinforcement and punishment should be assigned to recent choices within a

10     certain time scale. This is called credit assignment, which is implemented as eligibility traces in

11     reinforcement learning[4]. An eligibility trace can be implemented as a memory trace in each

12     synapse[10-12]. In this case, the trace time scales for reinforcement and punishment should be

13     separately controlled in direct and indirect pathways. Reinforcement learning theory requires that

14     both trace time scales be equal. However, the requirement cannot be completely realized in

15     distinct neural systems.

16        Here, we propose a computational model that enables the implementation of imbalanced

17     learning, and we aimed to determine whether the imbalance induces some defects in specific

18     situations. We focused on obsessive-compulsive disorder (OCD) as a candidate for such an

19     imbalanced condition for two reasons. First, abundant evidence suggests that imbalance between

20     direct and indirect pathways is central to the pathophysiology of OCD[13-15]. Second, obsession in

21     OCD can be an implicit choice leading to anxiety. In patients with OCD, an intrusive thought

22     becomes lodged in their mind and drives obsessive anxiety[16]. To neutralize and relieve the

23     anxiety, they tend to perform a compulsive action, despite its cost (e.g., excessively washing

their hands or repeatedly checking their keys), even though the anxiety will spontaneously

diminish before long if no action is taken[17].

This article proposes a computational model to reproduce the abnormal repetition of

obsession and compulsion in OCD symptoms, incorporating two hypotheses: (1) obsession is

induced by patients' implicit choices, and (2) eligibility trace time scales are imbalanced for

reinforcement and punishment in the reinforcement learning framework. We demonstrate that the

imbalance in trace time scales mis-reinforces implicit choices leading to obsession, resulting in a

spiral of repetitive obsession and compulsion. We tested our hypothesis in a behavioral task for

healthy participants and participants with OCD and evaluated the neural substrates of imbalanced

eligibility trace time scales using resting-state functional magnetic resonance imaging (rs-fMRI).


**Results**

**OCD-like behavior in a separate eligibility trace model**

Eligibility traces determine the credit assignment of the outcome prediction error of recent

actions[4]. Each trace represents the recent frequency of an action in a specific state within a time

scale determined by trace decay factor $v$ (see Online Methods). When a prediction error occurs,

the choice probability of each action in each state is updated by the product of the eligibility

trace and the prediction error (Figure 1a). Thus, a positive prediction error reinforces recent

actions within a time scale, whereas a negative error punishes them. We assumed that the

eligibility traces for reinforcement and punishment might be implemented in distinct neural

systems (red and blue in Figure 1a; see Online Methods). Theoretically, the trace factors in the

two systems, $v^{\pm}$, should be balanced: $v^+ = v^-$. However, the separate neural implementation

1    makes it difficult to perfectly maintain the balance. Here, we assumed that the balance was not

2    maintained: $\nu^+ \neq \nu^-$.

3          We modeled mental states involved in anxiety into stochastic transitions between two

4    states: relief and anxiety (Figure 1b). We assumed an action to relieve anxiety as an option in the

5    anxiety state, which could stochastically permit a transition from the anxiety state to the relief

6    state at a certain cost. The action might become compulsive in OCD, and we labeled it

7    "compulsion". We unified any other options in the anxiety state which could relieve the anxiety

8    at lower probability without any cost into the option "other". Our fundamental assumption was

9    that the transition from the relief state to the anxiety state would be caused by the individual's

10   own choice, which might become obsessive in OCD. We labeled this option "obsession". We

11   also unified any other options in the relief state which maintained the relief state into "other". We

12   assumed that every stay in the anxiety state would produce a negative outcome.

13          No positive outcome was introduced in the anxiety-relief transition model. Hence,

14   maintaining "other" in the relief state is clearly optimal. Normal reinforcement learning should

15   reduce the obsession rate. However, if the learning system has some defects, such as an

16   imbalance in trace factors, then the learning system may fail to reduce the obsession and fall into

17   a spiral of obsession and compulsion. To show this possibility, we simulated actor-critic learning

18   (see Online Methods), a typical method of reinforcement learning models[4] in the anxiety-relief

19   state-transition (Figure 1b), incorporating a separate eligibility trace (Figure 1a). When the

20   imbalance in the trace factors was moderate, actor-critic learning was able to reduce the

21   obsession rate, even when the simulation started at a high obsession rate (Figure 1c). In contrast,

22   extreme imbalance induced an increase in the obsession rate that culminated in a spiral of

23   repetitive obsession and compulsion (Figure 1d). Such OCD-like behavior was reduced by

1    preventing the compulsion during exposure to the anxiety state (green bar in Figure 1d). This is

2    the behavioral therapy of exposure and response prevention (ERP) and is one of the first-line

3    treatments of OCD[18].

4    We identified the condition for the trace factors $\nu^\pm$ in which OCD-like behavior might

5    emerge by numerical simulations in a representative set of other model parameters (Figure 1e;

6    see Online Methods). In addition to actor-critic learning (leftmost in Figure 1e), we simulated

7    SARSA (State-Action-Reward-State-Action) and Q-learning (right in Figure 1e), other typical

8    reinforcement learning models[4]. The color map represents the fraction of 100 simulation runs in

9    which obsession was reinforced at each pair of $(\nu^+, \nu^-)$. We identified the conditions of OCD-

10   like behavior in the region $\nu^+ > \nu^-$ (Figure 1e). Namely, the trace scale for reinforcement was

11   longer than that for punishment. We also identified the condition in which the behavioral therapy

12   of ERP would not work and localized it to a region with more robust imbalance (Figure 1f).

13   We obtained the theoretical conditions as functions of the other parameters (see the

14   Supplementary Note). Each OCD patient presents a certain compulsive action for specific

15   anxieties among many types of anxiety. Therefore, the corresponding parameters of the anxiety-

16   relief transition model (Figure 1b) should differ by the type of anxiety. Next, we derived the

17   condition for the learning parameters in which OCD-like behavior emerges in some type of

18   anxiety in actor-critic, SARSA, and Q-learning and obtained the common condition that $\nu^+ > \nu^-$

19   (see Section 3.3 in the Supplementary Note). This result suggests that a person with $\nu^+ > \nu^-$ has a
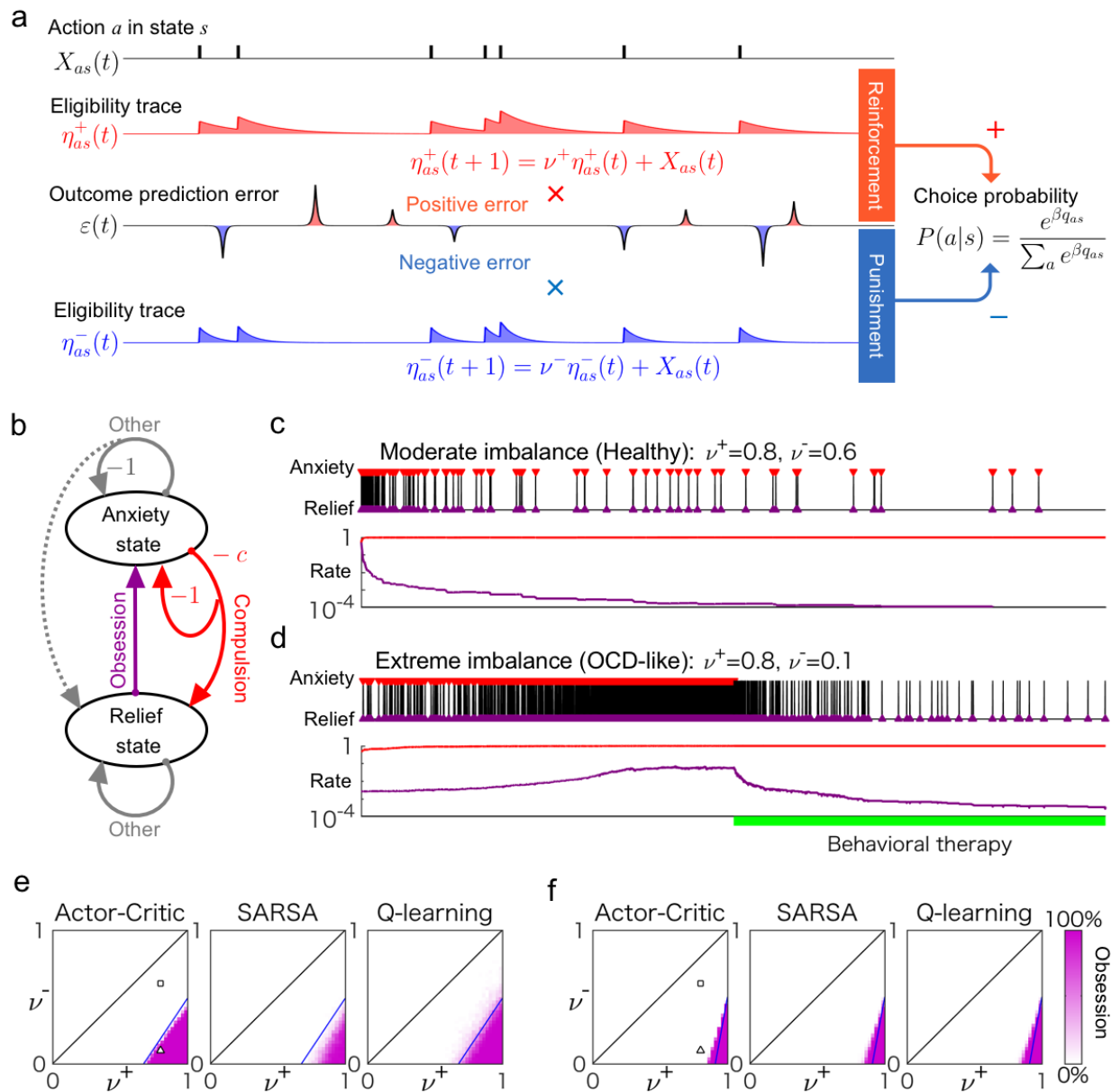
20   risk of OCD.

**Figure 1.** OCD-like behavior in a separate eligibility trace model. (a) Schema of the separate eligibility trace model. (b) State transitions in the anxiety-relief model. (c) Simulation of the separate eligibility trace model in the anxiety-relief state-transition in the case of a moderate imbalance in the trace factors $\nu^{\pm}$. The black polylines represent the temporal pattern of state transitions. The red and magenta triangles represent the occurrences of compulsions and obsessions, respectively. The lower plots show the temporal patterns of the compulsion (red

1  curve) and obsession (magenta curve) rates on a logarithmic scale. (d) As in (c), but in the case

2  of extreme imbalance. In the latter half, we demonstrated behavioral therapy by preventing the

3  compulsion in the anxiety state (green bar). (e) Conditions of OCD-like behavior in the space of

4  the trace factors $\nu^{\pm}$ with different types of learning algorithms: actor-critic, SARSA, and Q-

5  learning. The magenta color scale represents the percentage of the obsession rate increase among

6  100 times simulations around optimal behavior (zero obsession rate). Solid blue lines represent

7  the theoretically derived boundary (see the Supplementary Note). The square and triangle

8  represent the trace factors $\nu^{\pm}$ used in (c) and (d). (f) As in (e), but, in this condition, the

9  behavioral therapy did not work.

10

11  **Participants and behavioral task**

12  To test the suggestion that $\nu^+ > \nu^-$, derived from our computational model of OCD, we applied

13  the delayed feedback task to patients with OCD (n = 33) and healthy controls (HCs) (n = 168).

14  Fifteen HCs were excluded for medical or experimental reasons, and the subsequent analysis was

15  conducted in 33 OCD patients and 153 HCs (Supplementary Table 1). Considering the

16  therapeutic effect of the serotonin reuptake inhibitor (SRI) in OCD[18], SRI might normalize the

17  imbalanced setting of $\nu^+ > \nu^-$. Because a meta-analysis of SRI treatment suggested that higher

18  doses of SRI are more effective in the treatment of OCD versus other psychiatric conditions,

19  such as major depressive disorder[19], we divided the OCD patients into two groups by the dose

20  equivalence of SRIs[20] (Online Methods and Supplementary Table 2): OCD patients with higher

21  SRI doses were grouped into $OCD_{HighSRI}$ (n = 10), and those with lower SRI doses or no

22  psychotropic medications were $OCD_{Low-NoSRI}$ ($OCD_{LowSRI}$, n = 10; $OCD_{NoSRI}$, n = 13). There

23  were no significant differences in obsessive-compulsive and depressive symptoms evaluated by

1    the Yale-Brown Obsessive-Compulsive Scale[21] and the 17-item Hamilton Depression Rating

2    Scale (HDRS)[22] between the $OCD_{HighSRI}$ and $OCD_{Low\text{-}NoSRI}$ groups (Supplementary Table 1). In

3    addition, no patients had current major depressive disorder as a comorbidity (Online Methods).

4          We used the delayed feedback task to evaluate the trace factors $\nu^{\pm}$ as in our previous

5    research[23] except for the presented stimuli (abstract cues). Briefly, participants chose one of two

6    options displayed on the screen by pressing a left or right button in each trial (Figure 2a). Each

7    stimulus represented different outcomes (+10, +40, −10, or −40 yen) and a delay to the feedback

8    (immediately after a button press or three trials later) (Figure 2b and c). For example, +40(0)

9    represented a gain of 40 yen within the current trial (immediate reward) and −10(3) represented a

10   loss of 10 yen after three trials (delayed punishment). The participants were not told about the

11   stimulus-outcome associations shown in Figure 2b. They received money after the experiment in

12   proportion to the total outcome obtained in the delayed feedback task. To maximize the total

13   outcome, the participants needed to learn the appropriate stimulus choices by correctly assigning

14   the credit of the current feedback to the recent choices that caused the current feedback. The

15   difference in learning effects between immediate and delayed feedbacks reflected the trace
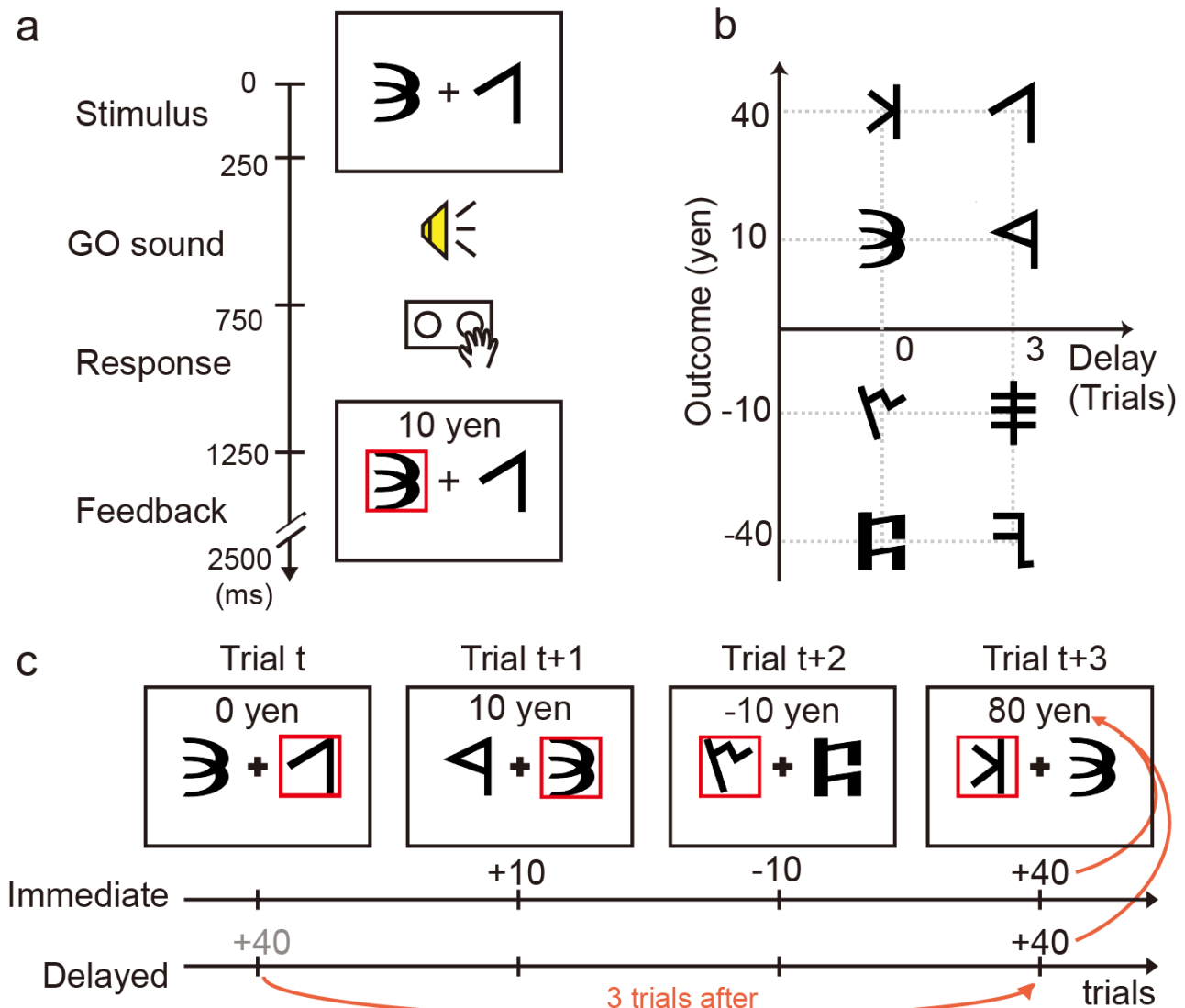
16   factors $\nu^{\pm}$.

Figure 2. Delayed feedback task. (a) Two abstract cues were displayed in each trial. When the participants heard an auditory cue (beep), they chose one of the stimuli within 1 s. A single trial lasted 2.5 s. The sequence of the stimuli pair was pseudorandom. (b) Outcome-delay mapping for each stimulus. The mapping was different among participants. (c) An example of a delayed feedback task. If participants chose the stimuli with a delay [+40(3)] at trial $t$, that outcome was not displayed immediately. If they chose the stimuli with no delay [+40(0)] at trial $t+3$, the sum of the delayed and immediate outcomes was displayed at trial $t+3$ (in this case +80 yen).

11

**Behavioral results**

Based on our simulation of OCD-like behavior with trace factors $\nu^+ > \nu^-$ (Figure 1e), patients with OCD would show impaired learning in the stimuli with delayed feedback. Therefore, the optimal choice rates between the representative four pairs of stimuli with the same delay and different magnitudes [pairs with no delays: +40(0) vs. +10(0) and −10(0) vs. −40(0); pairs with delays: +40(3) vs. +10(3) and −10(3) vs. −40(3)] were compared among groups. Consistent with our trace factor hypothesis, the OCD$_{Low-NoSRI}$ group showed impaired learning of stimuli with delays and performed at chance levels (Figure 3). Notably, OCD$_{HighSRI}$ patients exhibited intact learning similar to HCs in all pairs of interests.

A mixed-design two-way repeated-measures ANOVA with a within-participant factor of sessions (1–6 sessions) and a between-participant factor of groups (OCD$_{Low-NoSRI}$ patients, OCD$_{HighSRI}$ patients, and HCs) was conducted to clarify the between-group differences. There were significant interactions in the learning of stimulus pairs with delays. In the 10(3) vs. 40(3) pair, the interaction between session and group ($F(7.05, 645.34) = 2.32$, $p = 0.024$, $\eta_p^2 = 0.025$) and the simple main effects of group were significant in sessions 4 ($F(2, 183) = 3.97$, $p = 0.021$, $\eta_p^2 = 0.042$) and 5 ($F(2, 183) = 6.27$, $p = 0.0023$, $\eta_p^2 = 0.064$). Bonferroni-Holm–corrected post-hoc comparisons confirmed OCD$_{Low-NoSRI}$ patients < HCs ($t(183) = 2.81$, $p_{adjusted} = 0.016$) in session 4 and OCD$_{Low-NoSRI}$ patients < HCs ($t(183) = 3.54$, $p_{adjusted} = 0.0015$) in session 5 (Figure 3).

Regarding the −10(3) vs. −40(3) pair, the interaction between session and group ($F(7.42, 678.8) = 2.09$, $p = 0.039$, $\eta_p^2 = 0.022$) and the simple main effects of group were significant in sessions 5 ($F(2, 183) = 7.59$, $p = 0.0007$, $\eta_p^2 = 0.077$) and 6 ($F(2, 183) = 5.83$, $p = 0.0035$, $\eta_p^2 = 0.060$). Post-hoc comparisons confirmed significantly impaired learning in the OCD$_{Low-NoSRI}$

12

1　group $OCD_{Low-NoSRI}$ patients < HCs (t(183) = 3.87, $p_{adjusted}$ = 0.0005) and $OCD_{Low-NoSRI}$ patients <

2　$OCD_{HighSRI}$ patients(t(183) = 2.38, $p_{adjusted}$ = 0.036) in session 5 and $OCD_{Low-NoSRI}$ patients < HCs

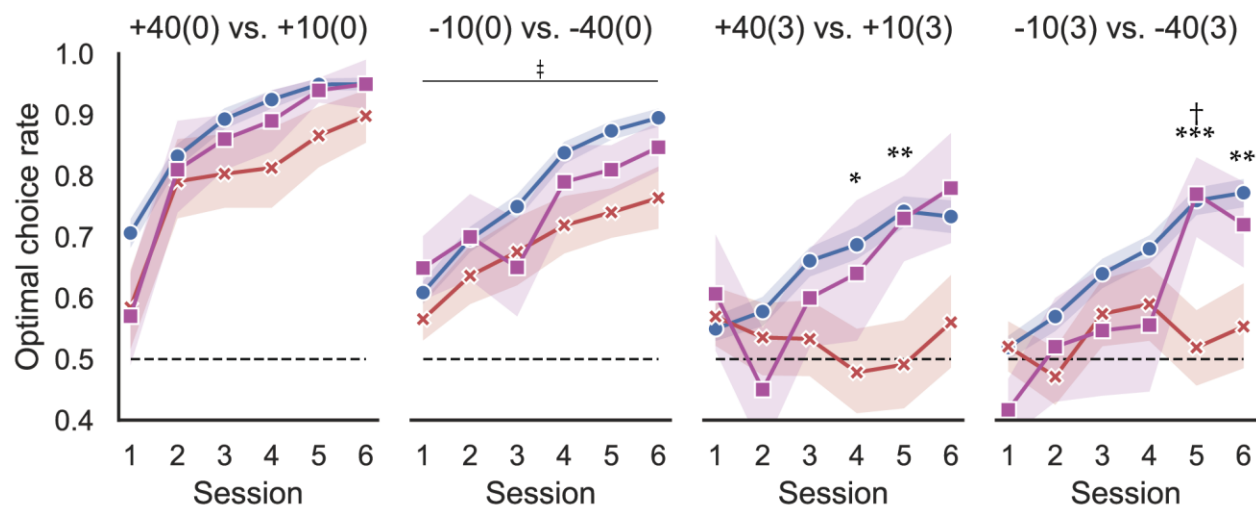3　(t(183) = 3.41, $p_{adjusted}$ = 0.0024) in session 6 (Figure 3).



4

5　Figure 3. Optimal choice rates of the delayed feedback task. Each panel represents the result for

6　each pair of stimuli [from left to right: pairs with no delays, +40(0) vs. +10(0) and −10(0) vs.

7　−40(0); pairs with delays, +40(3) vs. +10(3) and −10(3) vs. −40(3)]. The line and colored area

8　represented the mean and standard error of the optimal choice rates of groups (red cross, $OCD_{Low-}$

9　$_{NoSRI}$ patients; magenta square, $OCD_{HighSRI}$ patients; blue circle, HCs). The horizontal dashed line

10　at 0.5 is the chance level. *, **, and ***$p_{adjusted}$ < 0.05, 0.01, and 0.001 with Bonferroni-Holm–

11　corrected post-hoc comparisons of the simple main effects of group ($OCD_{Low-NoSRI}$ patients <

12　HCs); †$p_{adjusted}$ < 0.05, Bonferroni-Holm–corrected post-hoc comparisons of the simple main

13　effects of group ($OCD_{Low-NoSRI}$ patients < $OCD_{HighSRI}$ patients); ‡$p_{adjusted}$ < 0.05, Bonferroni-Holm–

14　corrected post-hoc comparisons of the main effect of group ($OCD_{Low-NoSRI}$ patients < HCs).

15

16　**Computational model-based behavioral analysis and results**

1    We fitted the behavioral data with the actor-critic learning model using the learning rate ($\alpha$),

2    exploration-exploitation degree ($\beta$), and separation of the eligibility trace for reinforcement and

3    punishment ($\nu^{\pm}$) (see Online Methods). Because we were specifically interested in the individual

4    variance represented by reinforcement learning parameters, we fitted parameters in each

5    participants' data independently using maximum a posteriori estimation rather than pooled data

6    as a group[24].

7        To clarify the $\nu^{+}/\nu^{-}$ distribution of each group, we projected each participants' estimated

8    parameter to $\nu^{+}/\nu^{-}$ space and visualized it using kernel density estimation (Figure 4a). Consistent

9    with our computational model simulation of OCD, $OCD_{Low-NoSRI}$ participants were distributed to

10   the $\nu^{+} > \nu^{-}$ imbalanced area, whereas the distribution was balanced in HCs and $OCD_{HighSRI}$

11   participants (Figure 4a). To compare the $\nu^{+}/\nu^{-}$ distribution among groups, we confirmed the

12   multivariate homogeneity of group dispersions ($F(2, 183) = 1.06$, $p = 0.35$) and applied

13   permutational multivariate analysis of variance (PERMANOVA). The $\nu^{+}/\nu^{-}$ distribution was

14   significantly different (PERMANOVA: $F(1, 184) = 6.41$, $p = 0.0039$), with Bonferroni-Holm–

15   corrected post-hoc pairwise comparisons revealing a significant difference between $OCD_{Low-}$

16   $_{NoSRI}$ participants and HCs ($p_{adjusted} = 0.014$) (Figure 4a). Parameters $\alpha$ and $\beta$ were not

17   significantly different among groups (Kruskal-Wallis test, $p > 0.05$).

18       We conducted clustering analysis using hierarchical density-based spatial clustering of

19   applications with noise (HDBSCAN)[25] to evaluate the diversity in HCs. HDBSCAN revealed

20   two clusters in HCs: a balanced $\nu$ cluster (Figure 4a, blue circle; $n = 83$) and an imbalanced $\nu$

21   cluster (Figure 4a, white circle; $n = 59$); the remaining 11 HCs were not clustered. Obsessive-

22   compulsive trends and a propensity to adhere to fine-grained details were compared between

23   clusters using five subscales of the Padua Inventory (PI; "Checking", "Dirt", "Doubt",

14

1    "Impulse", "Precision")[26] and the Attention to Detail subscale of the Autism Quotient (AQ)[27],

2    respectively. The imbalanced HC cluster showed significantly higher scores in the PI Checking

3    score (Brunner-Munzel test, statistic = −2.11, p = 0.019) and the AQ Attention to Detail score

4    (Brunner-Munzel test, statistic = −2.79, p = 0.0030) than the balanced HC cluster (Figure 4b).

5    Significant differences were still found after false discovery rate correction for multiple

6    comparisons among scales (PI Checking, $p_{adjusted}$ = 0.046; AQ Attention to Detail, $p_{adjusted}$ =

7    0.015). In OCD participants taking SRIs (n = 20), the equivalent dose of SRIs was significantly

8    correlated with the imbalanced settings of v ($v^{+}$–$v^{-}$) [Spearman's rank correlation (20); r =

9    −0.61, p = 0.0045; Figure 4c]. That is, higher doses of SRIs normalized the imbalanced settings
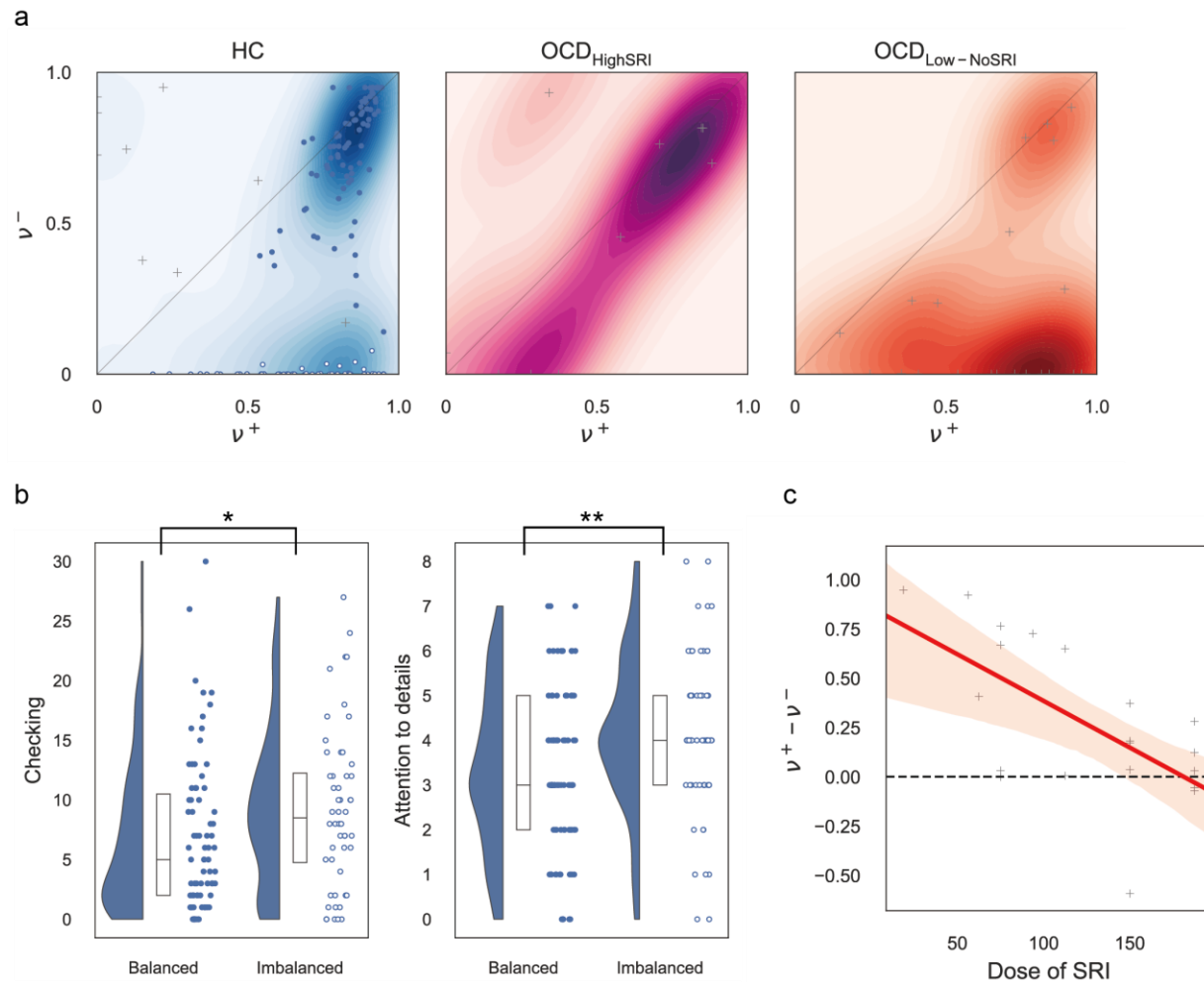
10   of v in participants with OCD.

15

Figure 4. Estimated parameters and their relationships with the clinical characteristics. (a) Distribution of participants in each group projected in the $v^+/v^-$ space using the kernel density estimation (blue, HCs; magenta, $OCD_{HighSRI}$ patients; red, $OCD_{Low-NoSRI}$ patients). The horizontal axis represents $v^+$ and the vertical axis represents $v^-$. The diagonal line represents the balanced line in the $v^+/v^-$ space. Blue and white circles in the leftmost figure represent the balanced and imbalanced clusters in HCs, respectively. Other participants were depicted using +. (b) PI Checking and AQ Attention to Detail scores in the balanced and imbalanced clusters in HCs. The half-violin plot, box plot, and strip plot represent the probability density, interquartile range and median, and raw data, respectively. Blue and white circles represent the balanced and

16

1  imbalanced clusters. * and ** represent significant differences (Brunner-Munzel test, * statistic =

2  −2.11, p = 0.019; ** statistic = −2.79, p = 0.0030). (c) Significant correlation between the SRI

3  dose and $v^+-v^-$ [Spearman's rank correlation (20); r = −0.61, p = 0.0045]. The line and colored

4  areas are the regression line and the 95% confidence interval (+: each OCD patient taking SRIs).

5  The dashed line represents the balanced settings of $v^{\pm}$.

6

7  **Neural signatures of OCD and the imbalanced $v^+ > v^-$ HC cluster**

8  The above behavioral results revealed the presence of imbalanced clusters in both the OCD

9  patients and the HCs. We next explored such imbalanced conditions of $v^+ > v^-$ in neural circuits.

10  Recently, various behavioral and demographic characteristics and psychiatric conditions have

11  been thought to be represented in a brain network, even in the resting state[14, 15, 28-30]. Here, we

12  explored the neural signature of nonmedicated patients with OCD hypothetically related to $v^+ >$

13  $v^-$ and extended it into an imbalanced ($v^+ > v^-$) HC cluster using rs-fMRI. First, we constructed a

14  whole-brain functional connectivity (FC) matrix using Cole-Anticevic Brain-wide Network

15  Partition (CAB-NP)[31] and compared 49 $OCD_{NoSRI}$ patients and 53 HCs (dataset A; Online

16  Methods and Supplementary Table 4) using network-based statistics (NBS)[32] (initial threshold, t

17  = 3.87; 10,000 random permutations; see Online Methods). We detected a significant network

18  component related to OCD (OCD network) with significantly increased connectivity in

19  $OCD_{NoSRI}$ patients compared with HCs ($p_{adjusted}$ = 0.022, Figure 5a). The OCD network

20  comprises nodes in the dorsolateral frontal cortex (DLPFC), parietal cortex, retrosplenial cortex,

21  and the hippocampal formation and parahippocampal region. These brain regions mainly belong

22  to the default mode network (DMN) or frontoparietal network (FPN) (Supplementary Figure 1).

23  To confirm the robustness of the OCD network, we compared the mean FC of the OCD network

17

1    between 10 $OCD_{NoSRI}$ patients and 18 HCs in the entirely independent dataset (dataset B; Online

2    Methods and Supplementary Table 5). Similarly, the mean FC within the OCD network was

3    significantly higher in $OCD_{NoSRI}$ patients than in HCs (Figure 5b, Brunner-Munzel test, statistic

4    = −3.11, p = 0.0027). Finally, we explored whether the imbalanced ($v^+ > v^-$) HC cluster showed

5    OCD-like characteristics also regarding a functional network. We compared all FCs of the OCD

6    network between 10 HCs in the imbalanced and balanced clusters (dataset C: independent from

7    dataset A and B; Online Methods and Supplementary Table 6) detected in our delayed feedback

8    task. We found that the FC between the DLPFC and presubiculum was significantly increased in

9    the imbalanced ($v^+ > v^-$) HC cluster, similar to OCD (Figure 5c, Brunner-Munzel test, statistic =

10   −2.94, p = 0.0051). The FC still showed a significant trend after false discovery rate correction

11   for multiple comparisons among the FCs of the OCD network ($p_{adjusted}$ = 0.066).
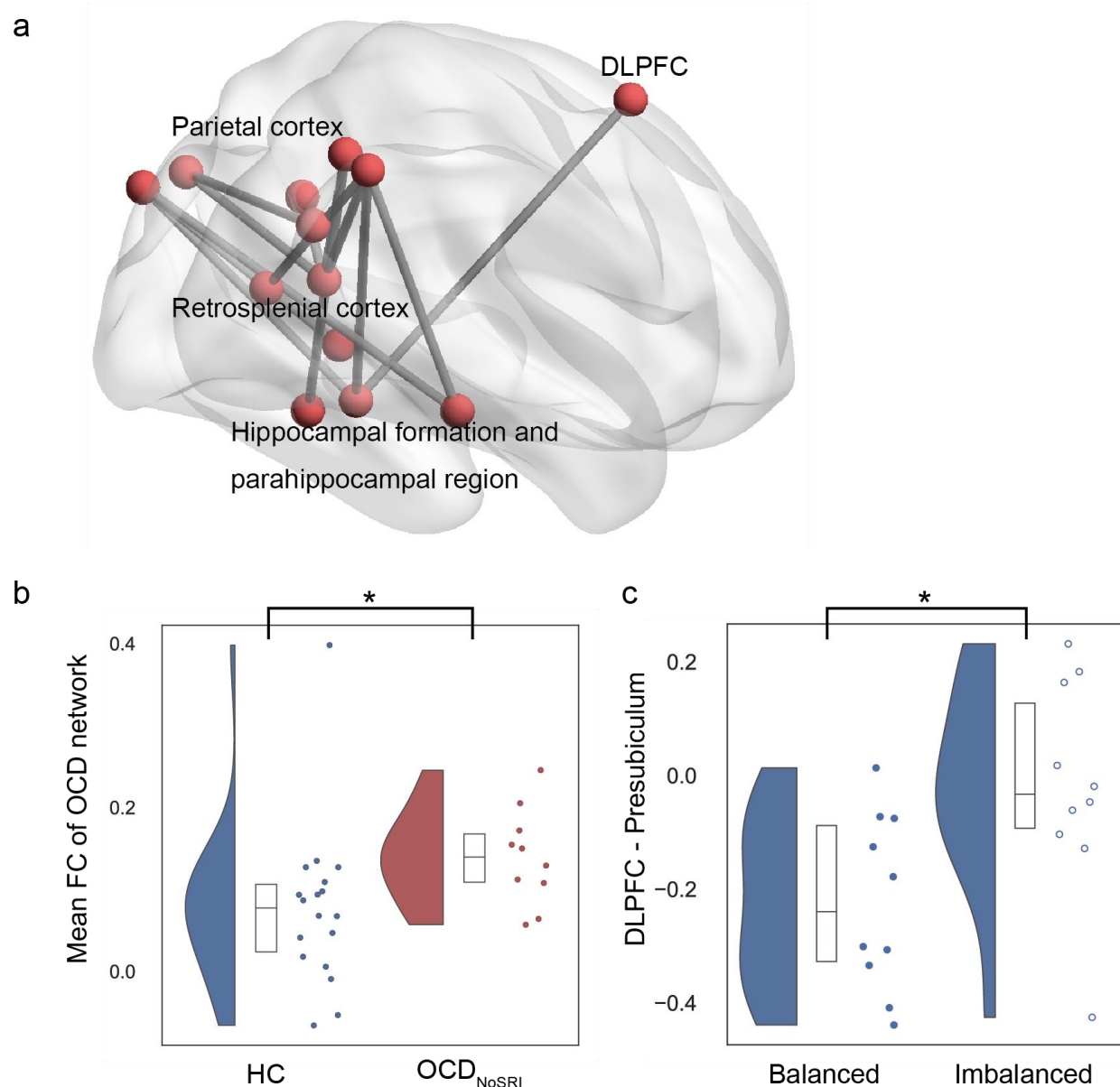
Figure 5. The OCD network and its extension to the imbalanced ($v^+ > v^-$) HC cluster. There was

no overlap of participants among the results in (a), (b), and (c). L and R represent left and right,

respectively. (a) Significantly increased functional network component in $OCD_{NoSRI}$ participants

compared with HCs ($p_{adjusted} = 0.022$). (b) Significantly increased mean FC within the OCD

network of $OCD_{NoSRI}$ participants compared with HCs in the completely independent dataset (*

Brunner-Munzel test, statistic $= -3.11$, p $= 0.0027$). (c) The imbalanced HC cluster showed

1    OCD-like increased FC between the DLPFC and presubiculum (* Brunner-Munzel test, statistic

2    = −2.94, p = 0.0051). In (b) and (c), the half-violin plot, box plot, and strip plot represent the

3    probability density, interquartile range and median, and raw data, respectively.

4

5    **Discussion**

6    In this study, we showed that implicit choices in our mind can induce anxiety, which is related to

7    imbalance in eligibility trace time scales for reinforcement and punishment. Specifically, we

8    constructed a computational model of OCD using a separate eligibility trace model (Figure 1a–d)

9    and found extremely imbalanced trace factors $v^+ > v^-$ in OCD (Figure 4a) and its neural substrate

10   (Figure 5a) in our empirical data. In addition, behavioral therapy (ERP) and psychotropic

11   medication (SRIs), which are the first-line treatments for OCD, were reflected in our

12   computational model (Figure 1d) and behavioral results (Figure 4c), respectively.

13           While the theoretical framework of the eligibility trace has long been conceptualized in

14   the field of reinforcement learning[4], its empirical evidence has been reported relatively

15   recently[10, 33-35]. The latest theories regarding synaptic plasticity have proposed that the co-

16   activation of pre- and postsynaptic neurons sets a flag at the synapse (eligibility trace) that leads

17   to a weight change only if additional factors (i.e., reinforcement or punishment) are present while

18   the flag is set[12]. These additional factors could be implemented by the phasic activity of some

19   neuromodulators, such as dopamine, which is supposed to represent the prediction error[5].

20   Although the detailed mechanisms concerning how the trace time scales for reinforcement and

21   punishment are modulated remain unclear, our theoretical consideration of the anxiety-relief

22   transition model showed that the imbalanced trace factors $v^+ > v^-$ could lead to the condition of

23   repetitive choices of anxiety and its relief in the model, similar to obsession and compulsion in

1    OCD (Figure 1a–d). Our experimental data from OCD patients and HCs strongly support the

2    predictions from our computational model, namely, the apparent impairment in the learning with

3    delayed feedback (Figure 3) and the extremely imbalanced trace factors $\nu^+ > \nu^-$ in OCD (Figure

4    4a). It is noteworthy that the imbalanced trace factors $\nu^+ > \nu^-$ are quite convincing because the

5    conventional pathophysiological model of OCD suggests excess tone in the direct pathway over

6    the indirect pathway[13-15], which are supposed to be related to $\nu^+$ and $\nu^-$, respectively.

7            There have been several computational models of OCD, all of which primarily focused

8    on compulsive behaviors[36-38]. Compulsive actions are thought to be reinforced by the rewarding

9    effect of relief from anxiety[17]. Excessive anxiety will induce excessive reinforcement of

10   compulsion leading to habit formation[38]. This excessively habitual compulsion is considered a

11   cause of OCD[38, 39]. However, obsessive thoughts that drive anxiety also increase with the

12   severity of OCD symptom[17]. Although the essence of OCD symptoms is the abnormal repetition

13   of obsession and compulsion[17], there is so far no unified model to explain the growth of both

14   obsession and compulsion. Our unified model can represent the vicious circle of obsession and

15   compulsion, which are the phenomenological characteristics of OCD[17] (Figure 1a–d). Moreover,

16   our model extends our understanding of the therapeutic effects of first-line treatments of OCD.

17   While ERP seems to promote the appropriate choices to prevent obsession, even under

18   imbalanced trace factors (Figure 1d), SRIs resolve the imbalance itself (Figure 4a and c). In

19   practice guidelines for the treatment of OCD[18], the combination therapy of ERP and SRI is

20   recommended if patients do not respond to ERP monotherapy. We identified such a condition

21   where ERP would fail in our computational model (Figure 1f). SRI add-on therapy in this ERP-

22   refractory condition can be viewed as assisting in the normalization of trace factors. Regarding

23   the relationships between the neuromodulator serotonin and the time scales of the eligibility

1    trace, serotonin seems to modulate the synaptic plasticity in many brain regions directly or

2    indirectly through its regulatory effects on other neuromodulatory systems[12, 33, 40, 41]. In our

3    previous study, we also found similar modulatory effects of the trace factor time scales, that is,

4    depletion of the serotonin precursor tryptophan increased the $v^+ > v^-$ imbalance compared with

5    the control condition[23]. There is still abundant scope for further research aimed at elucidating

6    how serotonin modulates the time scale of trace factors.

7        Using rs-fMRI, we found that patients with OCD exhibited significantly increased

8    connectivity of the OCD network component, which consists of the DLPFC, parietal cortex,

9    retrosplenial cortex, hippocampal formation, and parahippocampal region in two independent

10   datasets (Figure 5a and b). Specifically, the OCD network is mainly composed of FCs within the

11   DMN (7 of 13 FCs) and FCs between the DMN and FPN (4 of 13 FCs) (Supplementary Figure

12   1). These results are quite consistent with a recent meta-analysis of rs-fMRI findings in OCD, in

13   which the authors reported FC alterations within and between the DMN and FPN[42]. Moreover,

14   the DLPFC, parietal cortex, and hippocampus have been implicated in encoding the eligibility

15   traces related to reward-based decision-making for different time scales[12, 43-45]. These brain

16   regions might help to deal with the multiple eligibility trace time scales required for

17   reinforcement learning as a functional network, together with the retrosplenial cortex, which has

18   dense anatomical connections with all of the other detected brain regions[46].

19       Beyond the range of clinical phenotypes seen in patients with OCD, a broader continuum

20   of the obsessive-compulsive trait is also observed in behavioral and neural characteristics.

21   Specifically, the imbalanced ($v^+ > v^-$) HC cluster showed a significantly greater propensity for

22   Checking and Attention to Detail (Figure 4b) and increased FC between the DLPFC and

23   presubiculum, which belongs to the OCD network (Figure 5c). These results not only increase

1    the reliability of our clinical findings, but also support the generalizability of our findings to a

2    broader population on the obsessive-compulsive continuum. Further study with greater focus on

3    OCD patients and their unaffected healthy first-degree relatives should be performed to evaluate

4    the potential of our findings as an endophenotype of OCD[47].

5        In this study, we provided evidence that one's own implicit choices can induce anxiety in

6    OCD. However, anxiety can itself manifest in many different forms. In some situations, anxiety

7    may motivate people to take action or strive to meet a goal. In other situations, people may

8    experience anxiety as a symptom of an anxiety disorder such as social anxiety disorder and

9    generalized anxiety disorder. We previously demonstrated the commonality of different types of

10    anxiety, including obsessive-compulsive anxiety, using fMRI[29]. Because we cannot cover all

11    forms of anxiety in this research, further study with greater focus on the relationships between

12    implicit choices and various representations of anxiety from computational aspects should be

13    performed.

14

**Conclusion**

16    Anxiety can feel like a calamity that befalls us. However, our research shows that our own

17    implicit choices can trigger anxiety. Psychiatric symptoms are often thought of as alterations in

18    the mind that are not directly quantifiable, but they can be directly assessed through the creation

19    of appropriate computational models. Although it is currently difficult to identify treatment-

20    resistant patients from their clinical symptoms, our computational model suggests that patients

21    with extremely imbalanced trace scale parameters may not respond to behavioral therapy alone.

22    These results suggest that our findings could one day be applied to the appropriate selection of

23    OCD treatment. In addition, psychiatric symptoms have been regarded in recent years as a

1  symptom dimension common to various mental diseases, rather than being specific to a disease.

2  In this study, we focused on patients with OCD and healthy participants, but our approach could

3  be applied to the assessment of the anxiety dimension in various populations[48]. Future research is

4  needed to address these hypotheses in prospective longitudinal cohorts or in larger cohorts with

5  various psychiatric symptoms.

6

7  **Online Methods**

8  **Separate eligibility trace model**

9  Eligibility traces determine the weight assignment of outcome prediction errors for recent actions

10  (Figure 1a). If the eligibility traces for positive and negative prediction errors are implemented in

11  distinct neural systems, the respective eligibility traces $\eta_{as}^{+}(t)$ and $\eta_{as}^{-}(t)$ at time step $t$ for action

12  $a$ ($a = 0, 1, \ldots$) in state $s$ ($s = 0, 1, \ldots$) obey the following equation[4],

13  $$\eta_{as}^{\pm}(t + 1) = v^{\pm}\, \eta_{as}^{\pm}(t) + X_{as}(t),$$

14  where $X_{as}(t) = 1$ when action $a$ is chosen in state $s$ at time step $t$ and $X_{as}(t) = 0$ otherwise. The

15  factors $v^{\pm}$ ($0 < v^{\pm} < 1$) determine the decaying time scales of the eligibility traces. The policy

16  parameters $q_{as}(t)$ ($a = 0, 1, \ldots$) determine choice probability in state $s$ as a soft-max function,

17  $$P(a|s) = e^{\beta q_{as}(t)} \Big/ \sum_{\acute{a}} e^{\beta q_{\acute{a}s}(t)},$$

18  where $\beta$ represents the degree of the exploration-exploitation balance. Each policy parameter is

19  updated as below,

20  $$q_{as}(t + 1) = q_{as}(t) + \alpha[\eta_{as}^{+}(t)\max\{0, \varepsilon(t)\} + \eta_{as}^{-}(t)\min\{0, \varepsilon(t)\}\,],$$

21  where $\varepsilon(t)$ denotes the outcome prediction error and $\max\{0, \varepsilon(t)\}$ and $\min\{0, \varepsilon(t)\}$ represent

22  positive and negative components of the prediction error, respectively (Figure 1a). Theoretically,

23  the trace factors $v^{\pm}$ should be balanced as $v^{+} = v^{-}$. However, the separate neural implementation

24

1    makes it difficult to perfectly maintain the balance. Here, we assumed that the balance could be

2    broken: $v^+ \neq v^-$.

3        The outcome prediction error $\varepsilon(t)$ is determined as a function of the current outcome

4    $r(t)$ and the policy parameters $\{q_{as}\}$. The form depends on learning algorithms. In actor-critic

5    learning, the outcome prediction is based on the state value $v_s = \sum_a q_{as}$, and the prediction error

6    $\varepsilon(t) = r(t) + \gamma v_{s(t+1)} - v_{s(t)}$, where $r(t)$ denotes the current outcome and $s(t)$ and $s(t+1)$

7    denote the current and next states. The parameter $\gamma$ denotes the discount factor that determines

8    the weight of future prediction. In SARSA and Q-learning, the outcome prediction is based on

9    the action value estimated as each policy parameter $q_{as}$, and the prediction error $\varepsilon(t) = r(t) +$

10   $\gamma V_{s(t+1)} - q_{a(t)s(t)}$. The difference is in the term of the next state value: $V_{s(t+1)} = q_{a(t+1)s(t+1)}$

11   in SARSA and $V_{s(t+1)} = \max_a q_{as(t+1)}$ in Q-learning.

12

**Anxiety-relief transition model**

14   We modeled the mental states involved in anxiety into stochastic transitions between two states:

15   relief ($s$=0) and anxiety ($s$=1) (Figure 1b). We assumed two options in each state: "compulsion"

16   ($a$=1) and "other" ($a$=0) in the anxiety state, and "obsession" ($a$=1) and "other" ($a$=0) in the relief

17   state. We defined a matrix $b$ to determine the state-transition probabilities,

18   $$P(s(t+1) = 1 | a(t) = i, s(t) = j) = b_{ij},$$

19   $$P(s(t+1) = 0 | a(t) = i, s(t) = j) = 1 - b_{ij}.$$

20   We assumed that every stay in the anxiety state produced a negative outcome normalized to $-1$.

21   The relative cost of compulsion was denoted as $c$. Although $b_{00}$>0 allows passive anxiety, we

22   focused on $b_{00}$=0 in the main article. More general cases, including passive anxiety, are

23   considered in the Supplementary Note.

1

**Demonstration of OCD-like behavior**

For the example in Figure 1, we fixed the parameters of the anxiety-relief transition as $b_{00}$=0, $b_{10}$=1, $b_{01}$=0.9, $b_{11}$=0.5, and $c$=0.01 and the learning parameters as $\alpha$=0.1, $\beta$=1, $\gamma = 0.5$, and $\nu^+$=0.8, except for $\nu^-$. For Figure 1c, the trace factor for punishment was set as $\nu^-$=0.6 and the initial values of variables as $\eta_{as}^\pm(0) = q_{as}(0) = s(0) = 0$ . The simulation consisted of 100,000 time steps. For Figure 1d, we set the trace factors as $\nu^-$=0.1 and the initial values as $\eta_{as}^\pm(0) = q_{a1}(0) = s(0) = 0$, $q_{00}(0) = 3$, and $q_{10}(0) = -3$, to show that reinforcement of obsession started from even a low obsession rate. After 50,000 time steps, compulsion ($a$=1 in $s$=1) was always prevented, and other ($a$=0 in $s$=1) was forced regardless of the choice probability, demonstrating the behavioral therapy of ERP.

For Figure 1e and f, we evaluated the fraction of 100 simulation runs in which the obsession rate was reinforced from a low obsession rate on $40 \times 40$ grids in ($\nu^+$, $\nu^-$) space. In each simulation, 200 instances of forced obsessions were intermittently caused in the relief state because spontaneous obsession scarcely occurred at a low obsession rate. The simulation was judged to be obsession-reinforced if the choice probability of obsession became larger than the initial value. The initial values of variables were set as $\eta_{as}^\pm(0) = q_{a1}(0) = s(0) = 0$, $q_{00}(0) = 5$, and $q_{10}(0) = -5$. Each forced obsession was caused when $\eta_{as}^\pm(t)$ sufficiently approached zero in the relief state ($\eta_{as}^\pm(t) < 0.001$).

20

**Participants**

In total, 33 patients with OCD and 168 HCs participated in the behavioral task. Fifteen HCs were excluded for medical or experimental reasons, and the subsequent analysis was conducted in 33

1   OCD patients and 153 HCs. There were no significant differences in age, sex, or handedness

2   (Supplementary Table 1). All OCD patients and 13 of the 153 HCs were recruited at Kyoto

3   Prefectural University of Medicine (KPUM), whereas 140 of the 153 HCs were recruited at

4   Advanced Telecommunications Research Institute International (ATR). The Medical Committee

5   on Human Studies at KPUM and the Ethics Committee at ATR approved all procedures in this

6   study. All participants gave written informed consent after receiving a complete description of

7   the study. All methods were carried out following the approved guidelines and regulations.

8   Trained, experienced clinical psychiatrists and psychologists assessed all participants.

9        All patients were primarily diagnosed using the Structured Clinical Interview for DSM-

10  IV Axis I Disorders-Patient Edition (SCID)[49]. Exclusion criteria were 1) cardiac pacemaker or

11  other metallic implants or artifacts; 2) significant disease, including neurological diseases,

12  disorders of the pulmonary, cardiac, renal, hepatic, or endocrine systems, or metabolic disorders;

13  3) prior psychosurgery; 4) psychotropic medication except for SRIs; 5) DSM-IV diagnosis of

14  mental retardation and pervasive developmental disorders based on a clinical interview and

15  psychosocial history; and 6) pregnancy. We excluded patients with current DSM-IV Axis I

16  diagnosis of any significant psychiatric illness except OCD as much as possible, and only two

17  patients with trichotillomania, one patient with a tic disorder, one patient with panic disorder,

18  and one patient with bulimia nervosa were included as patients with comorbidities. The

19  experienced clinical psychiatrists or psychologists applied the Yale-Brown Obsessive-

20  Compulsive Scale (Y-BOCS)[21] for clinical evaluation of obsessive-compulsive symptoms in

21  patients with OCD. Handedness was classified based on a modified 25-item version of the

22  Edinburgh Handedness Inventory.

1    We divided OCD into two groups by medication status. Thirteen patients with OCD were

2    drug-free for all types of psychotropic medication, and the remaining 20 patients were taking

3    SRIs only. The SRI and imipramine equivalent doses[20] are summarized in Supplementary Table

4    2. OCD patients with higher SRI doses (higher than or equal to imipramine equivalent dose 150

5    mg) were grouped into the OCD$_{HighSRI}$ group (n = 10) and those with lower SRI doses (less than

6    imipramine equivalent dose 150 mg) or no psychotropic medications were grouped into OCD$_{Low-}$

7    $_{NoSRI}$ (OCD$_{LowSRI}$, n = 10; OCD$_{NoSRI}$, n = 13). The threshold was determined by considering a

8    meta-analysis of the SRI treatment[19] and common clinical doses used in Japan. To evaluate the

9    antidepressant effects of SRI, we compared the depressive symptoms evaluated by the HDRS[22]

10    between OCD$_{HighSRI}$ and OCD$_{Low-NoSRI}$ patients. There was one missing value of the HDRS in the

11    OCD$_{HighSRI}$ group.

12

**Behavioral task and statistical analysis**

14    The delayed feedback task was similar to the one in our previous research[23], except for

15    the presented stimuli. Participants chose one of the two options (abstract cues) displayed on the

16    screen by pressing a left or right button in each trial within 1 s after an auditory cue (Figure 2a).

17    Depending on the selected stimulus, monetary feedback with different outcomes (+10, +40, −10,

18    or −40 yen) was displayed either immediately after the button press or three trials later (Figure

19    2b and c). We did not offer monetary feedback for the first five trials because participants can

20    learn the relationships between stimuli and outcomes or delays quickly. If participants pressed a

21    button before the auditory cue or more than 1 s passed without any button press, −50 yen was

22    displayed as a punishment. Such trials were considered error trials, and delayed outcomes were

23    not considered.

1    At each trial, two abstract cues were displayed side by side on the screen (Figure 2a). We

2    prepared 16 pairs (Figure 2b), counterbalancing the number of appearances of each stimulus. The

3    16 pairs of stimuli were presented in pseudorandom order. Each pair was presented as the

4    scheduled number of trials in each session: each of six pairs [+10(0) vs. +40(0); +10(3) vs.

5    +40(3); −10(0) vs. −40(0); −10(3) vs. −40(3); +10(0) vs. +40(3); −10(0) vs. −40(3)] was

6    presented in 10 trials during a single session, and each of 10 pairs [+40(0) vs. +40(3); +10(0) vs.

7    +10(3); −10(0) vs. −10(3); −40(0) vs. −40(3); +10(3) vs. +40(0); −10(3) vs. −40(0); +10(0) vs.

8    −10(0); +40(0) vs. −40(0); +10(3) vs. −10(3); +40(3) vs. −40(3)] was presented in five trials

9    during a single session. Each participant performed 110 trials during a single session and six

10   sessions in each experiment. About 28 min was required for participants to complete six

11   sessions. At the beginning of each session, the session number was displayed on the screen for

12   2.5 s. Before the task, each participant practiced the test session under the same task settings

13   except for stimuli, and we confirmed that all participants understood the task set.

14   The total outcome (except for punishments related to button press errors), reaction time,

15   and the number of error trials were compared between the OCD and HC groups. Patients with

16   OCD showed a significantly lower total monetary outcome compared with HCs [median

17   (interquartile range): OCD patients, 2890 (1590–3600) yen; HCs, 4550 (3020–5950) yen;

18   Brunner-Munzel test, statistic = 4.08, p = 0.0002], whereas there were no group differences in

19   reaction time [median (interquartile range): OCD patients, 524.6 (487.1–565.4) ms; HCs, 522.6

20   (475.3–575.9) ms; Brunner-Munzel test, p > 0.05)] and number of button press errors [median

21   (interquartile range): OCD patients, 4 (2–9); HC, 4 (2–8); Brunner-Munzel test, p > 0.05). These

22   results suggested impaired learning in OCD individuals compared with HCs. Based on our

23   hypothesis of trace factors $v^+/v^-$, patients with OCD would show impaired learning in the stimuli

1  with delayed feedback. Therefore, the optimal choice rate of the representative four pairs of

2  stimuli with the same delay and different magnitude [pairs with no delays: +40(0) vs. +10(0) and

3  −10(0) vs. −40(0); pairs with delays: +40(3) vs. +10(3) and −10(3) vs. −40(3)] were compared

4  among groups. A mixed-design two-way repeated-measures ANOVA with a within-participants

5  factor of sessions (1–6 sessions) and a between-participants factor of groups ($OCD_{Low-NoSRI}$

6  patients, $OCD_{HighSRI}$ patients, and HCs) was conducted to clarify the between-group differences.

7  Degrees of freedom were corrected using Chi-Muller's epsilon because Mendoza's multisample

8  sphericity test indicated that the assumption of sphericity had been violated. Bonferroni-Holm–

9  corrected post-hoc comparisons were conducted to clarify the between-group differences.

10

**Model comparison and parameter estimation for the delayed feedback task**

12  We fitted the behavioral data with actor-critic learning with separate eligibility traces for positive

13  and negative prediction errors. We defined 16 states for possible pairs of stimuli presented at

14  each trial. Available actions at each state involved the choosing of alternative stimuli. To

15  facilitate model-fitting in the face of limited experimental data from each participant, we used

16  regularizing priors that favored realistic values and maximum a posteriori estimation rather than

17  maximum likelihood estimation[24]. The learning rate α and trace factor ν were constrained to the

18  range of $0 \leq \alpha \leq 0.95$ and $0 \leq \nu \leq 0.95$ with a uniform prior. The exploration-exploitation degree

19  β was constrained to the range of $0 \leq \beta \leq 100$ with a gamma (2,3) prior distribution that favored

20  relatively lower values. We fixed the discount factor $\gamma = 0$, because the term of the next state

21  value was just noise in our delayed feedback task in which the state of each trial was randomly

22  selected. We optimized parameters by minimizing the negative log posterior of the data with

23  different parameter settings using the hyperopt package[50]. Likelihood ratio tests to assess the

1    contribution of the additional parameter in our model (four parameters: $\alpha$, $\beta$, $\nu^+$, $\nu^-$) compared

2    with the standard actor-critic learning model (three parameters: $\alpha$, $\beta$, $\nu$) showed that the

3    additional parameter was justified in 128 of 186 participants ($X^2$ test with one degree of freedom,

4    p < 0.05). Because our target model was validated in behavioral data, we applied the separate

5    eligibility trace model in the substantive analysis.

6         We evaluated the performance of parameter estimation in actor-critic learning with

7    separate eligibility traces. We created 200 simulation data points using a model with the

8    following parameters: $\alpha$, 0.1 ± 0.05 (mean ± SD); $\beta$, 1 ± 0.2; $\nu^+$ and $\nu^-$, randomly selected within

9    0.01–0.95. The parameter estimation was quite accurate[51] (Supplementary Figure 2).

10   Specifically, Pearson's r and the mean absolute error between the true and estimated $\nu^+$ or $\nu^-$

11   were 0.99 and 0.03, respectively (Supplementary Figure 2).

12        We compared the $\nu^+/\nu^-$ distribution in our model among three groups (OCD$_{\text{Low-NoSRI}}$

13   patients, OCD$_{\text{HighSRI}}$ patients, and HCs) using PERMANOVA with the ADNOIS function and

14   10,000 permutations using the Euclidean distance implemented in the statistical package R[52].

15   The multivariate homogeneity of group dispersions was confirmed with the *betadisper* function

16   with 10,000 permutations in R[52]. The learning rate $\alpha$ and inverse temperature $\beta$ were compared

17   using the Kruskal-Wallis test. To further investigate relationships between clinical characteristics

18   and estimated parameters in the HC group, we conducted clustering analysis using HDBSCAN[25].

19   We detected two clusters (balanced cluster, n = 83; imbalanced cluster, n = 59; the remaining 11

20   HCs were not clustered) and evaluated their obsessive-compulsive trait using the five PI

21   subscales: "Checking", "Dirt", "Doubt", "Impulse", "Precision"[26]. In addition, the propensity to

22   adhere to fine-grained details was evaluated using the Attention to Detail subscale of AQ[27].

23   There were 15 missing values in PI (n = 127) and 1 missing value in AQ (n = 141). The

1 therapeutic effects of SRIs were evaluated by the Spearman's rank correlation between the SRI

2 dose and the imbalanced settings of $\nu$ ($\nu^+ - \nu^-$).

3

4 **Imaging data acquisition, preprocessing, and statistical analysis**

5 rs-fMRI data were collected using three different MRI scanners: 49 $OCD_{NoSRI}$ participants and

6 53 HCs at Kajiicho Medical Imaging Center (dataset A), 10 $OCD_{NoSRI}$ participants and 18 HCs at

7 Kyoto Prefectural University of Medicine (dataset B) for replication of the findings of dataset A,

8 and 20 HCs with the delayed feedback task (10 HCs for the imbalanced cluster and the

9 remaining 10 HCs for the balanced cluster) in ATR (dataset C). There was no overlap of

10 participants among datasets. All demographic distributions were matched between groups

11 (dataset A, Supplementary Table 4; dataset B, Supplementary Table 5; dataset C, Supplementary

12 Table 6). Some of the participants in dataset A and B were included in our previous studies

13 conducted for a different purpose using a different method[14, 15, 28, 29]. All fMRI imaging protocols

14 using gradient EPI sequences are summarized in Supplementary Table 3. High-resolution T1-

15 weighted structural images were also acquired.

16     Preprocessing of rs-fMRI data was conducted using fmriprep_ciftify 1.3.0.post2-2.3.0[53,

17 54]. Briefly, typical preprocessing steps such as slice timing correction, motion correction, and

18 spatial normalization into Montreal Neurological Institute space were conducted using fmriprep.

19 We then converted the data from the volumetric NIfTI format to the surface-based CIFTI format

20 (https://www.nitrc.org/projects/cifti/) using ciftify[54]. To remove artifacts and increase the

21 signal/noise ratio, the time course of the rs-fMRI data was detrended, bandpass-filtered (0.01–

22 0.08 Hz), and linearly regressed out of nuisance variables (the temporal fluctuations of the entire

23 brain and six head motion parameters, their derivatives, and six principal components of

1    anatomical CompCor[55]). With respect to motion artifacts, framewise displacement (FD) was not

2    significantly different between groups in all datasets [median (interquartile range); dataset A:

3    OCD patients, 0.082 (0.070–0.097) mm; HCs, 0.087 (0.067–0.10) mm; Brunner-Munzel test, p >

4    0.05; dataset B: OCD patients, 0.095 (0.067–0.12) mm; HCs, 0.079 (0.062–0.094) mm; Brunner-

5    Munzel test, p > 0.05; dataset C: balanced cluster, 0.13 (0.11–0.18) mm; imbalanced cluster, 0.11

6    (0.094–0.14) mm; Brunner-Munzel test, p > 0.05]. The first six functional scans were discarded

7    to allow magnetization to reach equilibrium. For each participant, mean time series were

8    extracted from 360 cortical and 358 subcortical parcels using CAB-NP[31], which is the

9    comprehensive whole-brain solution for large-scale functional networks based on the cortical

10   parcellation developed by Glasser et al.[56] (Human Connectome Project Multi-Modal

11   Parcellation). Pearson correlation coefficients were calculated between each pair of parcels and

12   transformed to Fisher's Z scores to obtain the FC matrix.

13          To evaluate the neural substrate of nonmedicated patients with OCD hypothetically

14   related to $v^+ > v^-$, we compared the FC matrices between 49 $OCD_{NoSRI}$ patients and 53 HCs

15   (dataset A) using NBS[32]. We chose NBS because it facilitates the detection of a subnetwork

16   related to the condition of interest while controlling the family-wise error rate. Briefly, NBS was

17   performed in the following two steps. First, the between-group comparison for every possible FC

18   was conducted and thresholded at t = 3.87 (corresponding to p = 0.0001) by considering the FD

19   as a nuisance variable. We detected the thresholded subnetworks (networks of nodes

20   interconnected by significant FCs), and their network size was calculated. Second, the

21   significance of subnetworks was tested using 10,000 random permutations of groups, which

22   determined the null distribution of the largest subnetwork size. Only subnetworks whose network

23   size exceeded the estimated family-wise error-corrected p-value 0.05 were identified as the

1    network that was significantly different between $OCD_{NoSRI}$ patients and HCs. The detected

2    subnetwork was visualized using BrainNet Viewer[57]. To confirm the robustness of the OCD

3    network, we compared the mean FC of the detected OCD network between 10 $OCD_{NoSRI}$ patients

4    and 18 HCs in the entirely independent dataset (dataset B) using the Brunner-Munzel test. To

5    further explore whether the imbalanced $(v^+ > v^-)$ HC cluster showed OCD-like characteristics

6    also regarding a functional network, we compared every FC of the OCD network between 10

7    HCs each in the imbalanced and balanced clusters (dataset C) detected in our delayed feedback

8    task using the Brunner-Munzel test.

9

**Data Availability**

The patients' data supporting the conclusion of this paper are not publicly available due to them containing information that could compromise research participant privacy or consent. The theoretical derivation of our computational model is in the Supplementary Note along with the MATLAB source code.

**Acknowledgments**

**Author Contributions**

Yuki S., Yutaka S., J.N., and S.C.T designed the study. Yutaka S. developed the theory and performed computational modeling. Yuki S., Y.A., J.N., and S.C.T. collected the data. Yuki S. conducted the computational modeling and statistical analysis of the data. Yuki S. and Yutaka S.

35

1    wrote the manuscript, which was edited by all authors. S.C.T. acquired funding to support theory

2    development and data analysis.

3

4    **Competing Interests Statement**

5    The authors declare no competing financial interests.

6

## References

1. Karni, A. & Sagi, D. Where practice makes perfect in texture discrimination: evidence for primary visual cortex plasticity. *Proc Natl Acad Sci U S A* **88**, 4966-4970 (1991).

2. Seitz, A.R. & Watanabe, T. Psychophysics: Is subliminal learning really passive? *Nature* **422**, 36 (2003).

3. Wichers, M.*, et al.* From affective experience to motivated action: tracking reward-seeking and punishment-avoidant behaviour in real-life. *PLoS One* **10**, e0129722 (2015).

4. Sutton, R.S. & Barto, A.G. Reinforcement Learning: An Introduction. (MIT Press, Cambridge, MA. 1998).

5. Schultz, W., Dayan, P. & Montague, P.R. A neural substrate of prediction and reward. *Science* **275**, 1593-1599 (1997).

6. Graybiel, A.M. The basal ganglia. *Curr Biol* **10**, R509-511 (2000).

7. Kravitz, A.V., Tye, L.D. & Kreitzer, A.C. Distinct roles for direct and indirect pathway striatal neurons in reinforcement. *Nat Neurosci* **15**, 816-818 (2012).

8. Nonomura, S.*, et al.* Monitoring and updating of action selection for goal-directed behavior through the striatal direct and indirect pathways. *Neuron* **99**, 1302-1314. e1305 (2018).

9. Shen, W., Flajolet, M., Greengard, P. & Surmeier, D.J. Dichotomous dopaminergic control of striatal synaptic plasticity. *Science* **321**, 848-851 (2008).

10. Yagishita, S.*, et al.* A critical time window for dopamine actions on the structural plasticity of dendritic spines. *Science* **345**, 1616-1620 (2014).

11. Izhikevich, E.M. Solving the distal reward problem through linkage of STDP and dopamine signaling. *Cereb Cortex* **17**, 2443-2452 (2007).

12. Gerstner, W., Lehmann, M., Liakoni, V., Corneil, D. & Brea, J. Eligibility traces and plasticity on behavioral time scales: experimental support of neoHebbian three-factor learning rules. *Frontiers in neural circuits* **12**, 53 (2018).

13. Pauls, D.L., Abramovitch, A., Rauch, S.L. & Geller, D.A. Obsessive-compulsive disorder: an integrative genetic and neurobiological perspective. *Nat Rev Neurosci* **15**, 410-424 (2014).

14. Sakai, Y.*, et al.* Corticostriatal functional connectivity in non-medicated patients with obsessive-compulsive disorder. *Eur Psychiatry* **26**, 463-469 (2011).

15. Abe, Y.*, et al.* Hyper-influence of the orbitofrontal cortex over the ventral striatum in obsessive-compulsive disorder. *Eur Neuropsychopharmacol* **25**, 1898-1905 (2015).

16. Salkovskis, P.M. Understanding and treating obsessive-compulsive disorder. *Behav Res Ther* **37 Suppl 1**, S29-52 (1999).

17. Cavedini, P., Gorini, A. & Bellodi, L. Understanding obsessive-compulsive disorder: focus on decision making. *Neuropsychol Rev* **16**, 3-15 (2006).

18. Koran, L.M.*, et al.* Practice guideline for the treatment of patients with obsessive-compulsive disorder. *Am J Psychiatry* **164**, 5-53 (2007).

19. Bloch, M.H., McGuire, J., Landeros-Weisenberger, A., Leckman, J.F. & Pittenger, C. Meta-analysis of the dose-response relationship of SSRI in obsessive-compulsive disorder. *Mol Psychiatry* **15**, 850-855 (2010).

20. Inada, T. & Inagaki, A. Psychotropic dose equivalence in Japan. *Psychiatry Clin Neurosci* **69**, 440-447 (2015).

21. Nakajima, T.*, et al.* Reliability and validity of the Japanese version of the Yale-Brown Obsessive-Compulsive Scale. *Psychiatry Clin Neurosci* **49**, 121-126 (1995).

1    22.     Hamilton, M. Development of a rating scale for primary depressive illness. *Br J Soc Clin*

2    *Psychol* **6**, 278-296 (1967).

3    23.     Tanaka, S.C.*, et al.* Serotonin affects association of aversive outcomes to past actions. *J*

4    *Neurosci* **29**, 15669-15674 (2009).

5    24.     Daw, N.D. Trial-by-trial data analysis using computational models. in *Decision Making,*

6    *Affect, and Learning: Attention and performance XXIII* (eds. Delgado, M.R., Phelps, E.A. &

7    Robbins, T.W.) 3-38 (Oxford Scholarship Online, Oxford, 2011).

8    25.     Campello, R.J., Moulavi, D., Zimek, A. & Sander, J. Hierarchical density estimates for

9    data clustering, visualization, and outlier detection. *ACM T Knowl Discov D* **10**, 5 (2015).

10   26.     Sugiura, Y. & Tanno, Y. Self-report inventory of obsessive-compulsive symptoms:

11   Reliability and validity of the Japanese version of the Padua Inventory. *Arch Pschyiatr Diagn*

12   *Clin Eval* **11**, 175-189 (2000).

13   27.     Wakabayashi, A., Tojo, Y., Baron-Cohen, S. & Wheelwright, S. [The Autism-Spectrum

14   Quotient (AQ) Japanese version: evidence from high-functioning clinical group and normal

15   adults]. *Shinrigaku Kenkyu* **75**, 78-84 (2004).

16   28.     Takagi, Y.*, et al.* A neural marker of obsessive-compulsive disorder from whole-brain

17   functional connectivity. *Sci Rep* **7**, 7538 (2017).

18   29.     Takagi, Y.*, et al.* A common brain network among state, trait, and pathological anxiety

19   from whole-brain functional connectivity. *Neuroimage* **172**, 506-516 (2018).

20   30.     Smith, S.M.*, et al.* A positive-negative mode of population covariation links brain

21   connectivity, demographics and behavior. *Nat Neurosci* **18**, 1565-1567 (2015).

22   31.     Ji, J.L.*, et al.* Mapping the human brain's cortical-subcortical functional network

23   organization. *Neuroimage* **185**, 35-57 (2019).

32. Zalesky, A., Fornito, A. & Bullmore, E.T. Network-based statistic: identifying differences in brain networks. *Neuroimage* **53**, 1197-1207 (2010).

33. He, K.*, et al.* Distinct eligibility traces for LTP and LTD in cortical synapses. *Neuron* **88**, 528-538 (2015).

34. Brzosko, Z., Schultz, W. & Paulsen, O. Retroactive modulation of spike timing-dependent plasticity by dopamine. *Elife* **4**, e09685 (2015).

35. Brzosko, Z., Zannone, S., Schultz, W., Clopath, C. & Paulsen, O. Sequential neuromodulation of Hebbian plasticity offers mechanism for effective reward-based navigation. *Elife* **6**, e27756 (2017).

36. Hauser, T.U.*, et al.* Increased decision thresholds enhance information gathering performance in juvenile obsessive-compulsive disorder (OCD). *PLoS Comput Biol* **13**, e1005440 (2017).

37. Vaghi, M.M.*, et al.* Compulsivity reveals a novel dissociation between action and confidence. *Neuron* **96**, 348-354.e4 (2017).

38. Gillan, C.M. & Robbins, T.W. Goal-directed learning and obsessive-compulsive disorder. *Philos Trans R Soc Lond B Biol Sci* **369** (2014).

39. Gillan, C.M. & Sahakian, B.J. Which is the driver, the obsessions or the compulsions, in OCD? *Neuropsychopharmacology* **40**, 247-248 (2015).

40. Brzosko, Z., Mierau, S.B. & Paulsen, O. Neuromodulation of spike-timing-dependent plasticity: past, present, and future. *Neuron* **103**, 563-581 (2019).

41. Cavaccini, A.*, et al.* Serotonergic signaling controls input-specific synaptic plasticity at striatal circuits. *Neuron* **98**, 801-816 e807 (2018).

1  42.    Gursel, D.A., Avram, M., Sorg, C., Brandl, F. & Koch, K. Frontoparietal areas link

2  impairments of large-scale intrinsic brain networks with aberrant fronto-striatal interactions in

3  OCD: a meta-analysis of resting-state functional connectivity. *Neurosci Biobehav Rev* **87**, 151-

4  160 (2018).

5  43.    Ambrose, R.E., Pfeiffer, B.E. & Foster, D.J. Reverse replay of hippocampal place cells is

6  uniquely modulated by changing reward. *Neuron* **91**, 1124-1136 (2016).

7  44.    Gersch, T.M., Foley, N.C., Eisenberg, I. & Gottlieb, J. Neural correlates of temporal

8  credit assignment in the parietal lobe. *PLoS One* **9**, e88725 (2014).

9  45.    Asaad, W.F., Lauro, P.M., Perge, J.A. & Eskandar, E.N. Prefrontal neurons encode a

10  solution to the credit-assignment problem. *J Neurosci* **37**, 6995-7007 (2017).

11  46.    Vann, S.D., Aggleton, J.P. & Maguire, E.A. What does the retrosplenial cortex do? *Nat

12  Rev Neurosci* **10**, 792-802 (2009).

13  47.    Chamberlain, S.R. & Menzies, L. Endophenotypes of obsessive-compulsive disorder:

14  rationale, evidence and future potential. *Expert Rev Neurother* **9**, 1133-1146 (2009).

15  48.    Gillan, C.M.*, et al.* Comparison of the association between goal-directed planning and

16  self-reported compulsivity vs obsessive-compulsive disorder diagnosis. *JAMA Psychiatry*, 1-10

17  (2019).

18

19  **References for Online Methods and Supplementary Information**

20  49.    First, M.B., Spizer, R.L., Gibbon, M. & Williams, J.B.W. *Structures Clinical Interview

21  for Axis I DSM-IV Disorders - Patient Edition (CID-I/P)* (Biometrics Research Department, New

22  York State Psychiatric Institute, New York, 1994).

1  50.     Bergstra, J.S., Bardenet, R., Bengio, Y. & Kégl, B. Algorithms for hyper-parameter

2  optimization. in *Advances in Neural Information Processing Systems* (eds. Shawe-Taylor, J.,

3  Zemel, R.S., Bartlett, P.L., Pereira, F. & Weinerger, K.Q.) 2546-2554 (Neural Information

4  Processing Systems, Inc., San Diego, 2011).

5  51.     Yazdani, S., Vahabie, A.-H., Araabi, B.N. & Ahmadabadi, M.N. Better than maximum

6  likelihood estimation of model-based and model-free learning style. Preprint at

7  https://www.biorxiv.org/content/10.1101/296335v1 (2018).

8  52.     Anderson, M.J., Walsh, D.C.I., Robert Clarke, K., Gorley, R.N. & Guerra-Castro, E.

9  Some solutions to the multivariate Behrens-Fisher problem for dissimilarity-based analyses. *Aust*

10  *NZ J Stat* **59**, 57-79 (2017).

11  53.     Esteban, O.*, et al.* fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat*

12  *Methods* **16**, 111-116 (2019).

13  54.     Dickie, E.W.*, et al.* Ciftify: A framework for surface-based analysis of legacy MR

14  acquisitions. *Neuroimage* **197**, 818-826 (2019).

15  55.     Behzadi, Y., Restom, K., Liau, J. & Liu, T.T. A component based noise correction

16  method (CompCor) for BOLD and perfusion based fMRI. *Neuroimage* **37**, 90-101 (2007).

17  56.     Glasser, M.F.*, et al.* A multi-modal parcellation of human cerebral cortex. *Nature* **536**,

18  171-178 (2016).

19  57.     Xia, M., Wang, J. & He, Y. BrainNet Viewer: a network visualization tool for human

20  brain connectomics. *PLoS One* **8**, e68910 (2013).

21