



HHS Public Access

Author manuscript

Nat Hum Behav. Author manuscript; available in PMC 2017 October 05.

Published in final edited form as:

Nat Hum Behav. 2016 ; 1: . doi:10.1038/s41562-016-0006.

Fear reduction without fear through reinforcement of neural activity that bypasses conscious exposure

Ai Koizumi^{1,2,3,9}, **Kaoru Amano**^{3,9}, **Aurelio Cortese**^{1,3,4,5,9}, **Kazuhisa Shibata**^{1,6}, **Wako Yoshida**^{1,3,8}, **Ben Seymour**^{1,3,8,*}, **Mitsuo Kawato**^{1,4,*}, and **Hakwan Lau**^{2,5,7,*}

¹Dept. of Decoded Neurofeedback, ATR Cognitive Mechanisms Laboratories, Address: 2-2-2, Hikaridai, Seika-cho, Sorakugun, Kyoto, 619-0288, JAPAN

²Dept. of Psychology, Columbia University, Address: 1190 Amsterdam Ave. 370 Schermerhorn Ext. MC:5501, New York, 10027, USA

³Center for Information and Neural Networks (CiNet), NICT, Address: 1-4 Yamadaoka, Suita City, Osaka, 565-0871, JAPAN

⁴Graduate School of Information Science, Nara Institute of Science and Technology, Address: 8916-5 Takayama, Ikoma Nara, 630-0192, JAPAN

⁵Dept. of Psychology, UCLA, Address: BOX 951563, Los Angeles, CA 90095-1563, USA

⁶Dept. of Psychology, Graduate School of Environmental Studies, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, JAPAN 464-8601, JAPAN

⁷Brain Research Institute, UCLA, Address: Box 951761, Los Angeles, CA 90095-1761, USA

⁸Dept. of Engineering, University of Cambridge, Address: Trumpington St, Cambridge CB2 1PZ, UK

Keywords

Fear extinction; Multi-voxel decoding; fMRI decoded neurofeedback; Post-traumatic stress disorder

Fear conditioning is a fundamentally important and preserved process across species¹. In humans it is linked to fear-related disorders such as phobias and post-traumatic stress

Correspondence should be addressed to: Mitsuo Kawato (kawato@atr.jp), Ben Seymour (bjs49@cam.ac.uk), and Hakwan Lau (hakwan@gmail.com).

⁹These authors contributed equally to this work

Competing Interests

I have read the journal's policy and the authors of this manuscript have the following competing interests. KS and MK are the inventors of patents related to the DecNef method used in this study, while the original assignee of the patents is ATR, with which some of the authors are affiliated.

Data availability

There is no related open source for this study.

Author contributions

A.K., H.L., B.S., and M.K. designed the study while actively discussing with other co-authors, A.K., K.A., and A.C. implemented the experiment, A.K. conducted the experiment, A.K., K.S., A.C., H.L., and M.K. analyzed the results with support of K.A., and W.Y. Lastly, A.K., B.S., H.L., and M.K. wrote the manuscript.

disorder (PTSD)^{2,3}. Fear memories can be reduced by counter-conditioning, in which fear conditioned stimuli (CS+) are repeatedly reinforced with reward⁴ or with novel non-threatening stimuli⁵. However, this procedure involves explicit presentations of CS+, which is itself aversive before fear is successfully reduced. This aversiveness may be a problem when trying to translate such experimental paradigms into clinical settings⁶. It also raises the fundamental question as to whether explicit presentations of feared objects is necessary for fear reduction¹. While learning without explicit stimulus presentation has been previously demonstrated⁷⁻¹⁰, whether fear can be reduced while avoiding the explicit exposure to CS+ remains largely unknown. One recently developed approach employs an implicit method to induce learning by reinforcing stimulus-specific neural representations using real-time decoding of multivariate functional magnetic resonance imaging (fMRI) signals^{11,12,13}, in the absence of stimulus presentation; i.e. pairing rewards with the occurrences of multi-voxel brain activity patterns matching a specific stimulus (decoded fMRI neurofeedback: DecNef^{11,13}). It has been shown that participants exhibit perceptual learning for a specific visual stimulus feature through DecNef, without giving any strategy for induction of specific neural representations, and without awareness of the content of reinforced neural representations¹¹. Here, we examined whether a similar approach can be applied to counter-conditioning of fear (Figure 1a-e). We show that we can reduce fear towards CS+ by pairing rewards with the activation patterns in visual cortex representing a CS+, while participants remain unaware of the content and purpose of the procedure. This procedure may be an initial step towards novel treatments for fear-related disorders such as phobia and PTSD, via unconscious processing.

In the experiment (Figure 1a), participants first acquired fear response to two visual stimuli (Target CS+ and Control CS+) in the Acquisition session, and went through the three daily sessions of Neural Reinforcement by DecNef during which only the activation patterns for Target CS+ were reinforced to reduce its associated fear, without physical presentations of the Target CS+. On the subsequent day, participants were presented with two CS+s and their fear response was measured in the Test session. The details of these sessions are as follows:

In the *Acquisition session*, pavlovian aversive conditioning was performed in 17 healthy participants by pairing two visual cues (CS+) with uncomfortable but tolerable electrical shocks (Figure 1a). The CS+s were visual stimuli - vertical gratings of different colors (red and green), allowing them to be distinguished by multi-voxel pattern decoding in the visual cortex, V1/V2. Unbeknownst to the participant, one of the two CS+s was designated as Target CS+, meaning we intended to have its associated fear level subsequently reduced via DecNef. The other CS+ was designated as the Control CS+, as baseline comparison. The choice of stimuli was based on a previous study using similar procedures¹³. Towards the end of Acquisition (Figure 1a-d), both CS+s induced elevated skin conductance responses (SCR) in comparison to an unreinforced cue (CS-), indicating successful conditioning of fear (Figure 2a, Figure S1).

The activation patterns discriminating the Target and non-target Control CS+s were determined by conventional multivariate decoding in another session (i.e., the fMRI session for MVPA, see Supplementary Method) before the Acquisition session (mean of decoding accuracy estimated with leave-one out cross validation; $72.1\% \pm 9.2$ s.d.), so that the

likelihood that Target CS+ is represented in V1/V2 activation patterns could be calculated in real-time during the subsequent Neural Reinforcement sessions.

In the three daily *Neural Reinforcement sessions* (Figure 1e), participants viewed achromatic visual gratings. They were asked to use any mental strategy they liked to try and increase the diameter of a disc on a monitor to earn monetary reward. They were unaware that the diameter reflected how the multi-voxel patterns reflected those induced by the Target CS+ in the visual cortex. The likelihood of occurrence of Target CS+ pattern was associated with magnitude of reward, leading to counter-conditioning for Target CS+⁴. We hypothesized this would reduce the ability of Target CS+ to elicit fear responses, relative to the Control CS+.

On Day 1 of the Neural Reinforcement session, the occurrence of Target CS+ pattern was around chance level ($50.2 \pm \text{s.e. } 4.1\%$; $t(16) = .05$, $P = .96$, two-tailed). Subsequently, Target CS+ likelihood exceeded chance level on Day 2 ($58.9 \pm 3.1\%$) and Day 3 ($57.2 \pm 3.3\%$) ($t(16) = 2.89$, $P = .01$, $t(16) = 2.20$, $P = .04$, respectively), providing evidence of successful DecNef of the Target CS+ pattern. The effect of day was significant (ANOVA; $F(2, 15) = 3.62$, $P = .038$), which was primarily due to the increase of Target CS+ likelihood from Day 1 to Day 2 ($P = .089$, Bonferroni-corrected). Importantly, although the overall likelihood of Target CS+ occurrence was modest, over the 3 days there was a sufficiently large variability in the trial-wise induction likelihood within each participant (across-participant mean of SD for the likelihood, 45.1%). Therefore participants were exposed to a full range of contingency between the induction likelihood and its corresponding reward, which is critical for the facilitation of reinforcement learning. More detailed progress of the occurrence of Target CS+ pattern is shown in Figure S2.

Post-experimental questionnaire confirmed that participants remained unaware of the association between the disc's diameter and occurrence of Target CS+ representation, and did not consciously use strategies related to color or grating (i.e. appropriate imagery-based strategies, Table S1). Neither were they able to guess the identity of the Target CS+ vs Control CS+ in a forced choice question afterwards (62.5% accuracy, Chi-square test; $\chi^2 = .780$, $P = .377$). This result is in agreement with the previous studies using a similar DecNef procedure^{11,13}. Furthermore, there were also no reports of subjective fear during the Neural Reinforcement session (Table S1), and SCR responses during these sessions were significantly less when compared to those associated with the CS+s at the end of Acquisition ($t(16) = 4.218$, $P = .001$), and did not correlate with induced Target CS+ pattern likelihood (Figure 2b&c).

Next, during the *Test session*, four unsignaled shocks (USs) were presented to reactivate the fear memory (i.e., reinstatement. Figure 1a) to assess potential fear reduction for the Target CS+¹⁴. Given that the Test session was performed several days after Acquisition, this reinstatement procedure helped to ensure observable fear responses for the baseline condition (i.e. Control CS+), to allow for a meaningful comparison. During the Test session, the two CS+s and the CS- were presented alone to evaluate the associated fear response (SCR). Critically, we found that the SCR to Target CS+ was significantly reduced when compared to the Control CS+ ($t(16) = 2.630$, $P = .018$, paired t-test, two-tailed) (Figure 2a, Figure S1), suggesting that fear towards the CS+ was reduced only when it went through the

DecNef procedure, where occurrence of the Target CS+ was paired with reward, effectively counter-conditioning the previously acquired fear. Interestingly, the magnitude of this effect was similar to what was observed following conventional extinction procedures^{15,16}, even though the underlying mechanism may be different, not least as participants were unaware of the occurrence of Target CS+ representations during Neural Reinforcement.

We also recorded fMRI responses during the Acquisition and Test sessions, focusing in particular on responses in the amygdala and the ventral medial prefrontal cortex (VMPFC), which have been implicated in acquisition and extinction of fear memory^{17,18}. The amygdala showed significant responses after conditioning for both CS+s, but a significant reduction in response to the Target CS+ compared to Control CS+ in the Test session (Figure 3a), mirroring the specific pattern of fear reduction seen in the SCRs. VMPFC responses were significantly negative for both CS+s during Acquisition as previously shown^{16,17,19}, but significantly less positive during Neural Reinforcement and Test sessions for the Target CS+ (Figure 3b). While responses in the amygdala and VMPFC were reduced for the Target CS+ following the Neural Reinforcement sessions, during these sessions the trial-wise response level in these regions did not correlate with the likelihood of Target CS+ induction, suggesting that fear memory was not strongly reactivated by the spontaneous occurrence of Target CS+ patterns in the visual cortex (Figure S3). The average responses in V1/V2 were similar between Target and Control CS+s both before and after the Neural Reinforcement sessions (Figure S4), undermining the possibility that the differential responses in the amygdala and VMPFC for the two CS+s during the Test session were merely due to the altered visual processing for Target CS+.

Although the likelihood of Target CS+ occurrence was estimated selectively from the activation patterns in V1/V2 during the Neural Reinforcement sessions, it remains unclear whether such information reflecting Target CS+ was confined to V1/V2 or other brain regions were also engaged to act in concert with V1/V2. To examine whether any brain regions outside the visual cortex were engaged, we conducted the whole-brain searchlight multi-voxel pattern analysis (MVPA)²⁰, which quantitatively measures the degree to which the activation patterns of other brain areas could predict the likelihood that the Target CS+ patterns were induced in V1/V2. Such predictability of the Target CS+ likelihood in V1/V2 would reflect the “information transmission” from V1/V2 to other areas^{11,13}. For more detailed advantages of this approach, see Supplementary Method. Specifically, the searchlight MVPA estimated the degree to which the trial-by-trial induction likelihood of Target CS+ in V1/V2 can be reconstructed from the multi-voxel patterns within a spherical region of interest (ROI) (radius=15 mm) centered at each voxel (see Supplementary Method). The analysis revealed that the Target CS+ likelihood (i.e., likelihood of red or green grating) was transmitted to many visual areas during the fMRI session for MVPA when the red and green gratings were physically presented (Figure 4a), but it was largely confined within V1/V2 during the Neural Reinforcement sessions, with a few exceptions (Figure 4b). In particular, the striatal area (caudate nucleus) had significant information transmission from V1/V2, in keeping with its possible role in reinforcement learning^{21,22}. Such engagement of the striatal area was further supported by a Psychophysiological Interaction (PPI) analysis²³ showing enhanced functional connectivity between V1/V2 and striatum, as a function of the increase of Target CS+ likelihood in V1/V2 (Supplementary

Method, Figure S6). These results suggest that, besides the visual cortex where Target CS+ was induced, the striatal area was engaged to some extent in the Neural Reinforcement sessions.

Previous studies with conventional extinction procedures have shown a role of VMPFC for successful extinction^{17,18}. Yet, the aforementioned searchlight analysis revealed no significant information transmission from V1/V2 to VMPFC, suggesting that VMPFC may not be actively engaged in DecNef on average within the participants group. Moreover, across participants the degree to which VMPFC patterns predicted V1/V2 patterns ('information transmission') was negatively correlated with the success of fear reduction (Spearman's $\rho = -.522$, $P = .034$; Figure 4c), suggesting that there was less VMPFC engagement for participants with more successful fear reduction. These results, as well as the whole brain analysis examining the correlation between information transmission and fear reduction (Figure S5), suggest that VMPFC disengagement may have led to larger reduction of fear. Thus, our results are consistent with the view that the fear reduction observed here depended on a possibly different mechanism in comparison to that of the conventional extinction procedures¹⁸ (see Discussion for details).

To summarize, we here provide behavioural and neurophysiological evidence that rewards can be directly paired with the patterns of neural activity in visual cortex to facilitate the reduction of fear. This DecNef procedure bypasses the requirement for physical visual presentation of a CS. As such, the procedure represents direct counter-conditioning of neural activity based on the fluctuations of the information content in the visual cortex activity. The results show that counter-conditioning of visual cortical-based fear is sufficient to bring about the reduction of behavioural and amygdalar fear responses, even in the absence of the participants' awareness of the content of neural induction and purpose of the procedure^{11,13,24}.

The pattern of activity in VMPFC contrasts with that observed in conventional studies of extinction, whereby it typically exhibits enhanced responses to extinguished fear cues^{16,17}. In conventional extinction, fear memory is thought to be inhibited by extinction learning of a context-specific 'safety' state, mediated by VMPFC, via the exertion of an inhibitory influence on amygdala-based fear memories^{16,17}. Similarly, in counter-conditioning, presenting rewards in place of previously paired aversive outcomes reduces conditional fear by a putatively similar inhibitory mechanism⁴. The fact that we see reduced VMPFC activity in association with reduced fear response for Target CS+ (Figure 3b) suggests that the mechanism of fear reduction demonstrated here may differ from conventional extinction procedures^{16,17,19}. This is further supported by the fact that greater interaction between the visual cortex and VMPFC during the Neural Reinforcement sessions as measured by "information transmission" is associated with less fear reduction (Figure 4c, Figure S5).

Given that successful fear reduction requires VMPFC involvement during conventional extinction training¹⁸, this result hints at the possibility that VMPFC disengagement may be a key to the success of our procedure¹⁹. Consistent with this view, some previous studies have also suggested that disengagement of VMPFC could in some cases counterintuitively lead to more robust extinction of fear. For example, a human lesion study has shown that VMPFC

damage prevents development of PTSD after traumatic experiences²⁵, implying robust fear extinction. Similarly, infant rats achieve more robust fear extinction than adult rats despite that infant rats do not rely on the medial prefrontal areas including an area homologous to the human VMPFC²⁶. Such VMPFC disengagement may have made the effect of DecNef robust, potentially making the fear reduction effect to be context-invariant (See Figure S7).

Although the exact neural mechanism underlying this neural counter-conditioning procedure still remains to be further elucidated, it is likely to critically involve the striatum, an area implicated in reinforcement learning^{21,22}. During the Neural Reinforcement sessions, occurrence of Target CS+ in V1/V2 could be predicted from the activation patterns of striatum, mainly caudate nucleus (Figure 4b). Similarly, the functional coupling between V1/V2 and striatum was modulated by the occurrence of Target CS+ in V1/V2 (Figure S6). Activation in striatal areas has been found in previous studies on neurofeedback training²⁷, which suggests that neurofeedback procedures in general may rely to some extent on striatal reinforcement.

While this is the first human study to show that fear can be reduced by directly pairing reward with the induced activation patterns for CS+ in visual cortex, an animal study has previously shown that pairing reward with the optogenetically reactivated hippocampal traces of a fear conditioned location (i.e., CS+) could reduce the associated fear response¹⁰. The study further showed that, after the optogenetic reactivation, the hippocampal memory trace lost its ability to activate the previously associated amygdala neurons. Similarly in the current study, reinforcing the Target CS+ pattern in visual cortex may have weakened its previously acquired fear association with amygdala, resulting in reduction of amygdala response to the physically presented Target CS+ (Figure 3a).

Overall, the fear reduction effect achieved with the DecNef procedures appears robust, as it tolerated the challenge provided by reinstatement to reactivate fear memory¹⁴. Tolerance of the fear reduction effect to such a challenge is ecologically meaningful as it may capture resistance to relapse of fear in real life fear memories¹⁴. Here we used reinstatement primarily to allow a sufficient magnitude of fear response after the multiple days of training period, but in principle it would be interesting to look at the effect on the fear response without reinstatement. From a theoretical perspective, showing fear reduction both with and without reinstatement would allow better clarification as to whether fear reduction was driven primarily by an inhibition of the original memory¹⁴. As fear reduction effects are often weakened after reinstatement²⁸⁻³¹, it is likely that the observed fear reduction effect would have been still present even if tested without reinstatement, but future studies could directly investigate this.

Our current findings may eventually benefit clinical treatments of fear related disorders. From a translational perspective, the traditional application of fear extinction to anxiety-related disorders faces several challenges. One difficulty in applying traditional associative learning concepts to exposure therapy is that some participants may not comply with explicit encounters with feared objects in the first place⁶, because of their intrinsic aversiveness. Here, the induction of brain patterns are not accompanied by conscious awareness of the relevant content, and so this may alleviate the problem of patient attrition. However, to

achieve this ultimate goal, one needs to pass several technical hurdles. For instance, for ecological validity, one would need to develop procedures to decode images with rich real-life content, and the learning of the relevant multi-voxel patterns for individual patients would also need to be done without conscious presentation to the patients. These may be overcome by building decoders with subliminal presentation, or adapting them from other individuals' brain activity³². Despite these challenges, the present results hopefully represent an initial step towards a potential new avenue for treatment.

Methods

The entire experiment consisted of the five main sessions (Figure 1a), Acquisition, Neural Reinforcement x 3, and Test, which were conducted after the two preparatory sessions, Retinotopy, and fMRI session for MVPA. All sessions were conducted on different days separated by at least 24 hrs. All of the experiments were conducted with fMRI measurement.

Participants

Twenty-four participants (15 males, 23.2±2.5 years old) were initially enrolled. Seven participants were eliminated prior to participating in Neural Reinforcement session because six of them failed to show measurable fear response (see Acquisition session) and one participant did not complete Acquisition session due to excessive anxiety. The remaining 17 participants (11 males, 23.5±2.8 mean years old) completed all experimental sessions. We predetermined the number of participants to complete all the sessions based on our pilot study. Participants gave a written consent prior to participating in each session. The study was approved by the Institutional Review Board of ATR, Japan.

Retinotopy session

We first conducted a standard retinotopic mapping experiment to localize V1/V2 in each participant³³ (see Supplementary Information for complete method).

fMRI session for MVPA

The aim of the fMRI session for MVPA was to obtain fMRI data for constructing a decoder to classify the activation patterns in V1/V2 (see Retinotopy session) evoked by isoluminant red versus green vertical gratings (Figure 1b), which were to serve as the CS+s in the subsequent Acquisition session. The decoder was used in the following Neural Reinforcement sessions to evaluate the trial-by-trial likelihood that participants could induce brain activation patterns for Target CS+ (red or green grating, counterbalanced across participants). During this session each trial consisted of a fixation disc (6 sec), followed by a grating which flickered at 0.5 Hz (6 sec total).

The preprocessed fMRI signals from the localized V1/V2 subregions were then used to construct a decoder to classify the activation patterns for red versus green grating (see Supplementary Information). We used sparse logistic regression (SLR)³⁴ to automatically select the voxels that were relevant for classification. We trained the decoder using 192 data points obtained from 192 trials (across all 12 fMRI runs).

Acquisition session

The aim of the Acquisition session was to establish fear memory for red and green gratings (conditioned stimulus, CS+) by pairing them with an uncomfortable but tolerable electric shock (unconditioned stimulus, US). These two gratings were identical to the fMRI session for MVPA. A grating with novel color (blue or yellow) was introduced as CS-, which was never paired with US. The choice of color for CS- was counterbalanced across participant in orthogonal to the choice of color for Target CS+ (red or green) (e.g., approximately half of the participants with green Target CS+ were assigned blue CS-, while the other half participants were assigned yellow CS-). With such counterbalancing, we avoided the situation where the activation patterns for CS- would be always more similar to one of CS+s (e.g., Target) than to the other CS+ (e.g., Control). The experimenter was not blind to the color assignments in the Acquisition session and the subsequent sessions, as the handling with blindness was difficult due to the complexity of our procedures. Two CS+s were presented either with or without US (5 and 8 times, respectively), and CS- was always presented without US (8 times)¹⁹. Trial order was randomized. Each trial started with a presentation of a CS (4 sec) followed by a fixation disc (12 sec). On trials with US, a CS+ co-terminated with a burst of electric shocks (36 impulses across 200ms total). Skin conductance response (SCR) was recorded using BrainAmp Ag/AgCl sintered MR electrodes (Brain Products) attached to the distal phalanges of index and middle fingers of right hand. Among twenty-four participants who completed Acquisition session, six participants were excluded because no SCR was detected for the CS+s. Another participant did not complete Acquisition session due to excessive anxiety. The remaining 17 participants proceeded to the subsequent sessions. To estimate fear response in late Acquisition, we calculated the mean SCR during the last 2 trials for each CS (Figure 2a, left panel).

Neural Reinforcement session

The Neural Reinforcement sessions were conducted for three consecutive days. The aim of the session was to repetitively induce V1/V2 activation patterns for one of the CS+s (red, N=9; green, N=8) without participants' awareness of the induced Target CS+. We reinforced participants with monetary reward for inducing the patterns for one of the CS+s, given the capacity of reward to reinforce behaviour³⁵ as well as neural activity³⁶. Participants were not attached to an electrode for electric shock.

Each trial had a sequence of an induction period (6 sec), a fixation period (7 sec), a feedback period (1 sec), and an inter-trial interval (6 sec) (Figure 1e). During the induction period, participants were instructed to somewhat regulate their brain activity so as to maximize the size of white disc which served as feedback. Feedback was presented after 6 sec of the fixation period following the induction period. In the induction period, a gray vertical grating was presented. The gray grating flickered at 0.5 Hz (6 sec total; three repetitions of a grating (1.5 sec) and a fixation (0.5 sec)). Participants were not informed as to what the feedback disc size represented (i.e., Target CS+ likelihood in V1/V2).

V1/V2 activation pattern during the induction period was analyzed online to estimate the likelihood that the currently achieved brain activation patterns represented the patterns for

Target CS+ (red or green) that were previously decoded from the fMRI session for MVPA. Hemodynamic delay of 6 s was taken into account.

Test session

A day after the last day of the Neural Reinforcement session, we conducted the Test session to measure fear responses to Target CS+, Control CS+, and CS-. Based on our preliminary studies, we presented four unsigned USs before the Test session to activate the fear memory (i.e., Reinstatement) in a similar manner as a previous study¹⁹. Following reinstatement, each CS was presented for 11 times in a semi-randomized order: CS- was always presented on the first trial to capture irrelevant SCR due to orienting effect¹⁹. Data of this first CS- was discarded from the subsequent analyses. A trial sequence was identical to Acquisition session, except that there was no trial with US. SCR was recorded in the same manner as in the Acquisition session.

MRI parameters

Participants were scanned in a 3T MRI scanner (Trio, Siemens) with a head coil at the ATR Brain Activation Imaging Center. See supplementary method for more detailed parameters.

Definition of ROIs

Along with SCR, we measured response in the amygdala and VMPFC to track the fear related activity in these areas with fMRI. To determine the amygdala ROI, we first defined anatomical boundary of the amygdala with freesurfer segmentation, and selected voxels within this anatomical boundary that showed greater response for all US trials and the last 2 trials of each CS+ (i.e., fear relevant trials) relative to fixation during the Acquisition session. To define the VMPFC ROI, we first created an anatomical mask of a sphere with 15 mm radius centered around previously reported MNI coordinates [0, 40, -12]³⁷, which was estimated based on the representative literature^{16,17,38}. We then selected voxels within the sphere ROI that showed smaller response for all US trials and the last 2 trials of each CS+ relative to fixation during the Acquisition session, which was the expected direction of activity based on previous literature¹⁹. Caudate and ventral striatum were defined using FSL Structural Striatum Atlas⁴⁰.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

The study was conducted in ImPACT Program of Council for Science, Technology and Innovation (Cabinet Office, Government of Japan). This work was partially supported by "Brain machine Interface Development" under the Strategic Research Program for Brain Sciences supported by AMED of Japan, the ATR entrust research contract from National Institute of Information and Communications Technology, and the US National Institute of Neurological Disorders and Stroke of the National Institutes of Health (Grant No. R01NS088628 to H.L.). BS is funded by the Wellcome Trust, UK. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. We thank Kaori Nakamura for her help in scheduling and conducting the experiment, Nobuo Hiroe for the assistance with equipment, and Yasuhiro Shimada and Akikakazu Nishikido for operating the fMRI scanner, Hiroshi Ban for technical advices, Michelle Craske, Michael Treanor, Michael Sun, Alicia Izquierdo, and Frank Krasne for their comments on the manuscript.

References

1. LeDoux, J. *Anxious*. Oneworld Publications; 2015.
2. Lissek S, et al. Classical fear conditioning in the anxiety disorders: a meta-analysis. *Behav Res Ther*. 2005; 43:1391–1424. [PubMed: 15885654]
3. Yehuda R, LeDoux J. Response variation following trauma: a translational neuroscience approach to understanding PTSD. *Neuron*. 2007; 56:19–32. [PubMed: 17920012]
4. Dickinson, A., Dearing, M. Mechanism of learning and motivation. Dickinson, RA., editor. 1979.
5. Dunsmoor JE, Campese VD, Ceceli AO, LeDoux JE, Phelps EA. Novelty-facilitated extinction: providing a novel outcome in place of an expected threat diminishes recovery of defensive responses. *Biol Psychiatry*. 2015; 78:203–209. [PubMed: 25636175]
6. Schnurr PP, et al. Cognitive behavioral therapy for posttraumatic stress disorder in women: a randomized controlled trial. *JAMA*. 2007; 297:820–830. [PubMed: 17327524]
7. Esteves F, Parra C, Dimberg U, Ohman A. Nonconscious associative learning: Pavlovian conditioning of skin conductance responses to masked fear-relevant facial stimuli. *Psychophysiology*. 1994; 31:375–385. [PubMed: 10690918]
8. Knight DC, Nguyen HT, Bandettini PA. Expression of conditional fear with and without awareness. *Proc Natl Acad Sci U S A*. 2003; 100:15280–15283. [PubMed: 14657356]
9. Raio CM, Carmel D, Carrasco M, Phelps EA. Nonconscious fear is quickly acquired but swiftly forgotten. *Curr Biol*. 2012; 22:R477–9. [PubMed: 22720676]
10. Redondo RL, et al. Bidirectional switch of the valence associated with a hippocampal contextual memory engram. *Nature*. 2014; 513:426–430. [PubMed: 25162525]
11. Shibata K, Watanabe T, Sasaki Y, Kawato M. Perceptual Learning Incepted by Decoded fMRI Neurofeedback Without Stimulus Presentation. *Science*. 2011; 334:1413–1415. [PubMed: 22158821]
12. deBettencourt MT, Cohen JD, Lee RF, Norman KA, Turk-Browne NB. Closed-loop training of attention with real-time brain imaging. *Nat Neurosci*. 2015; 18:470–475. [PubMed: 25664913]
13. Amano K, Shibata K, Kawato M, Sasaki Y, Watanabe T. Learning to associate orientation with color in early visual areas by associative decoded fMRI neurofeedback. *26*, 1861–1866. *Curr Biol*. 2016
14. Bouton ME. Context and Behavioral Processes in Extinction. *Learn Mem*. 2004; 11:485–494. [PubMed: 15466298]
15. Milad MR, et al. Neurobiological basis of failure to recall extinction memory in posttraumatic stress disorder. *Biol Psychiatry*. 2009; 66:1075–1082. [PubMed: 19748076]
16. Milad MR, et al. Recall of fear extinction in humans activates the ventromedial prefrontal cortex and hippocampus in concert. *Biol Psychiatry*. 2007; 62:446–454. [PubMed: 17217927]
17. Phelps EA, Delgado MR, Nearing KI, LeDoux JE. Extinction learning in humans: role of the amygdala and vmPFC. *Neuron*. 2004; 43:897–905. [PubMed: 15363399]
18. Do-Monte FH, Manzano-Nieves G, Quiñones-Laracuente K, Ramos-Medina L, Quirk GJ. Revisiting the role of infralimbic cortex in fear extinction with optogenetics. *J Neurosci*. 2015; 35:3607–3615. [PubMed: 25716859]
19. Schiller D, Kanen JW, LeDoux JE, Monfils MH, Phelps EA. Extinction during reconsolidation of threat memory diminishes prefrontal cortex involvement. *Proc Natl Acad Sci U S A*. 2013; 110:20040–20045. [PubMed: 24277809]
20. Kriegeskorte N, Goebel R, Bandettini P. Information-based functional brain mapping. *Proceedings of the National Academy of Sciences*. 2006; 103:3863–3868.
21. Haruno M, Kawato M. Different neural correlates of reward expectation and reward expectation error in the putamen and caudate nucleus during stimulus-action-reward association learning. *J Neurophysiol*. 2006; 95:948–959. [PubMed: 16192338]
22. O’Doherty J, et al. Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science*. 2004; 304:452–454. [PubMed: 15087550]
23. Friston KJ, et al. Psychophysiological and modulatory interactions in neuroimaging. *Neuroimage*. 1997; 6:218–229. [PubMed: 9344826]

24. Sarrazin JC, Cleeremans A, Haggard P. How do we know what we are doing? Time, intention and awareness of action. *Conscious Cogn.* 2008; 17:602–615. [PubMed: 17468011]
25. Koenigs M, et al. Focal brain damage protects against post-traumatic stress disorder in combat veterans. *Nat Neurosci.* 2008; 11:232–237. [PubMed: 18157125]
26. Kim JH, Hamlin AS, Richardson R. Fear extinction across development: the involvement of the medial prefrontal cortex as assessed by temporary inactivation and immunohistochemistry. *J Neurosci.* 2009; 29:10802–10808. [PubMed: 19726637]
27. Emmert K, et al. Meta-analysis of real-time fMRI neurofeedback studies using individual participant data: How is brain regulation mediated? *Neuroimage.* 2016; 124:806–812. [PubMed: 26419389]
28. Brooks DC, Beth H, Nelson JB, Bouton ME. Reinstatement after counterconditioning. *Anim Learn Behav.* 1995; 23:383–390.
29. LaBar KS, Phelps EA. Reinstatement of conditioned fear in humans is context dependent and impaired in amnesia. *Behav Neurosci.* 2005; 119:677–686. [PubMed: 15998188]
30. Norrholm SD, et al. Conditioned fear extinction and reinstatement in a human fear-potentiated startle paradigm. *Learn Mem.* 2006; 13:681–685. [PubMed: 17142300]
31. Hermans D, et al. Reinstatement of fear responses in human aversive conditioning. *Behav Res Ther.* 2005; 43:533–551. [PubMed: 15701362]
32. Haxby JV, et al. A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron.* 2011; 72:404–416. [PubMed: 22017997]
33. Engel SA, Glover GH, Wandell BA. Retinotopic organization in human visual cortex and the spatial precision of functional MRI. *Cereb Cortex.* 1997; 7:181–192. [PubMed: 9087826]
34. Yamashita O, Sato MA, Yoshioka T, Tong F, Kamitani Y. Sparse estimation automatically selects voxels relevant for the decoding of fMRI activity patterns. *Neuroimage.* 2008; 42:1414–1429. [PubMed: 18598768]
35. Kobayashi S, et al. Influences of rewarding and aversive outcomes on activity in macaque lateral prefrontal cortex. *Neuron.* 2006; 51:861–870. [PubMed: 16982429]
36. Bray S, Shimojo S, O’Doherty JP. Direct instrumental conditioning of neural activity using functional magnetic resonance imaging-derived reward feedback. *J Neurosci.* 2007; 27:7498–7507. [PubMed: 17626211]
37. Lonsdorf TB, Haaker J, Kalisch R. Long-term expression of human contextual fear and extinction memories involves amygdala, hippocampus and ventromedial prefrontal cortex: a reinstatement study in two independent samples. *Soc Cogn Affect Neurosci.* 2014; 9:1973–1983. [PubMed: 24493848]
38. Kalisch R, et al. Context-dependent human extinction memory is mediated by a ventromedial prefrontal and hippocampal network. *J Neurosci.* 2006; 26:9503–9511. [PubMed: 16971534]
39. Schiller D, et al. Preventing the return of fear in humans using reconsolidation update mechanisms. *Nature.* 2010; 463:49–53. [PubMed: 20010606]
40. Tziortzi AC, et al. Imaging dopamine receptors in humans with [11C]-()-PHNO: Dissection of D3 signal and anatomy. *Neuroimage.* 2011; 54:264–277. [PubMed: 20600980]
41. Nichols TE, Holmes AP. Nonparametric permutation tests for functional neuroimaging: a primer with examples. *Hum Brain Mapp.* 2002; 15:1–25. [PubMed: 11747097]

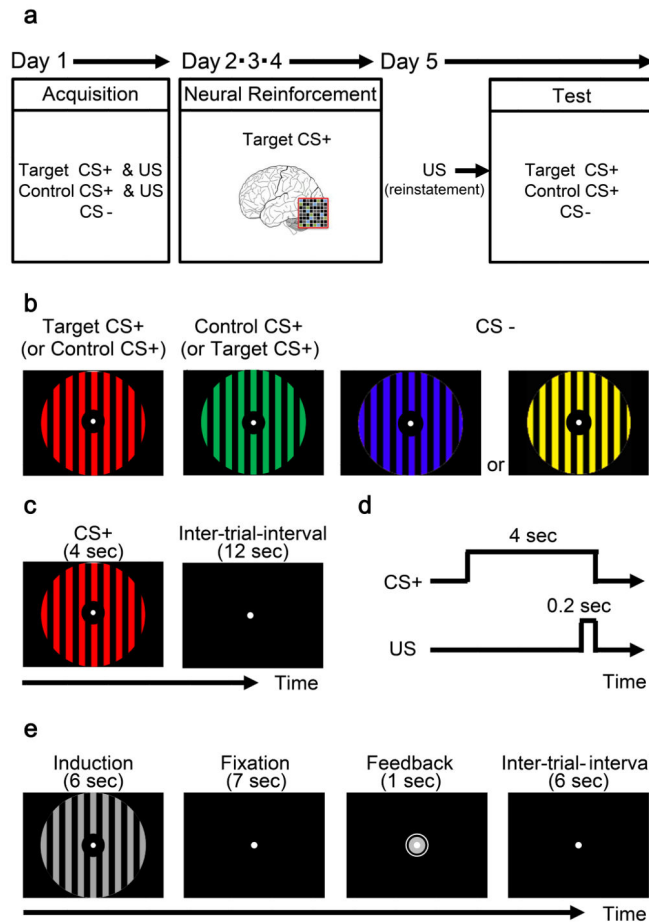
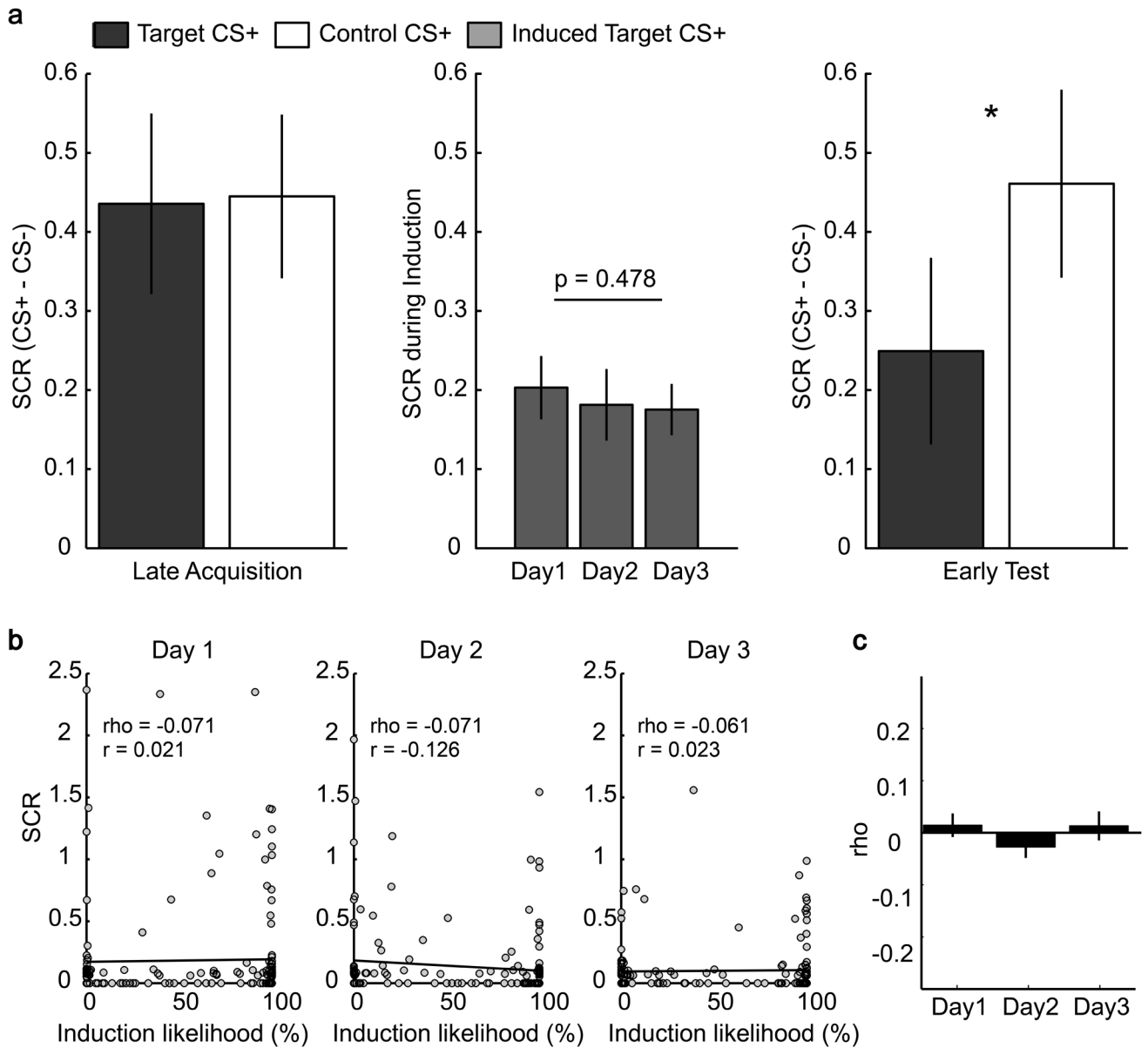


Figure 1.

Overall experimental design. **a**) After retinotopy and the fMRI session for MVPA (see Supplementary Method), participants went through 5 days of main experimental sessions in the MRI scanner. **b**) The stimuli used as CSs were colored vertical grating patterns, with choice of colors (red and green) for Target CS+ and Control CS+ counterbalanced across participants. A choice of color (blue or yellow) for CS- was also counterbalanced across participants. **c**) Timeline for a single trial in the Acquisition or Test session. In an Acquisition trial, both Target and Control CS+s (red/green) were paired with US (electric shock). **d**) CS+s were paired with the co-terminating US during Acquisition at the contingency rate of 38%. **e**) During a Neural Reinforcement trial, participants were required to somehow regulate their brain activity, upon seeing a gray vertical grating (Induction cue). The size of the disc during the feedback period indicated the online-calculated likelihood of the Target CS+ patterns in V1/V2. The disc size was proportional to the amount of monetary reward earned.

**Figure 2.**

Reduction of fear response as measured by SCR. **a** In late Acquisition (last 2 trials), participants developed positive responses for both Target and Control CS+. During the Neural Reinforcement session (middle panel), such responses were lower than those associated with the CS+s at the end of Acquisition ($t(16)=4.218$, $P=.001$). In early Test (first 2 trials; see Figure 1a), response to Target CS+ was reduced compared to Control CS+ ($t(16)=2.630$, $P=.018$, paired t-test, two-tailed). **b** During Neural Reinforcement, trial-wise correlation between SCR and induction likelihood (i.e. degree to which activity in V1/V2 resembled the multi-voxel pattern for Target CS+) was negligible, indicating that induced Target CS+ did not lead to fear response. Shown are scatter plots of a representative participant, where each dot represents a single induction trial. **c** Plotted is the Fisher-

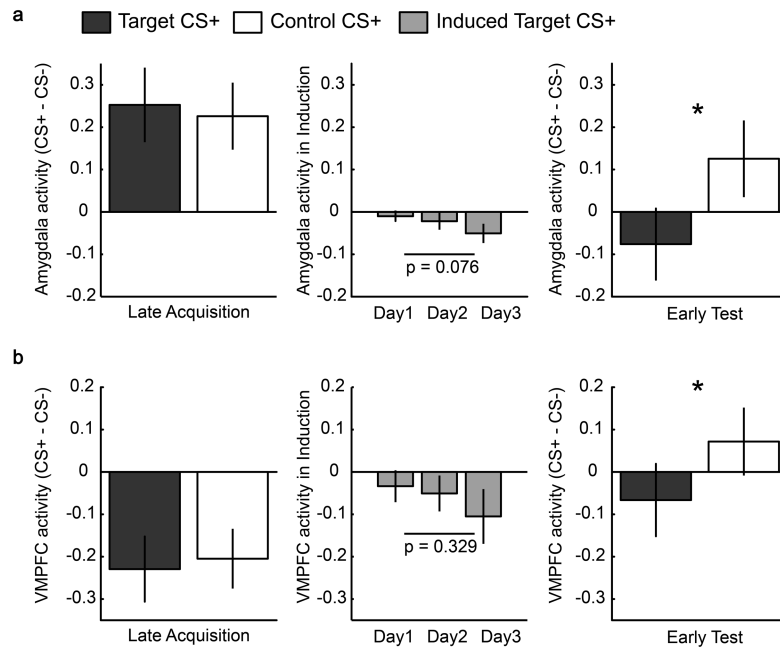
transformed correlation coefficients for each day of the Neural Reinforcement session averaged across participants. Error bars represent standard errors.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Figure 3.**

Brain activity in the amygdala and VMPFC. Same labeling as in Figure 2a except the dependent measure here is the average level of activity (arbitrary unit) in the amygdala (**a**) and VMPFC (**b**). **a**) Amygdala activity was reduced for Target CS+ compared with Control CS+ ($t(16)=2.21$, $P=.042$; two-tailed) in early Test. **b**) VMPFC activity was reduced for Target CS+ relative to Control CS+ ($t(16)=2.13$, $P=.049$; two-tailed) in early Test. While activity in the amygdala and VMPFC numerically decreased across the three days of the Neural Reinforcement sessions (middle panels of **a**&**b**, respectively), these decreases were not significant ($P=.076$, $P=.329$, respectively). Error bars represent standard errors. Please also see Figure S3.

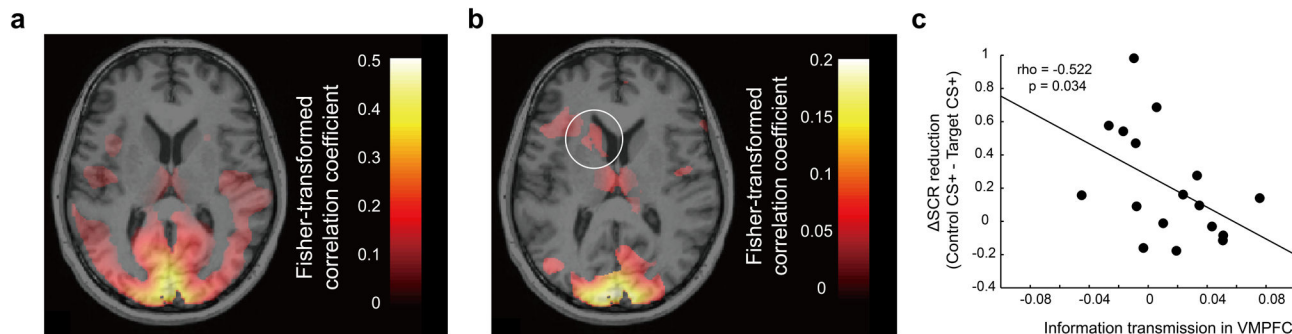


Figure 4.

Engagement of striatum (**b**) and disengagement of VMPFC (**c**) during Neural Reinforcement. **a**) Whole-brain searchlight MVPA quantitatively evaluated information transmission from V1/V2 to the whole brain, during the fMRI session for MVPA. Information transmission was estimated as the correlation between the Target CS+ likelihood in V1/V2 and its reconstructed value obtained from the multi-voxel patterns within each searchlight sphere ROI (radius=15 mm). Information was transmitted mostly within, although not confined to, many visual areas. This result ensures the sensitivity and power of the whole-brain searchlight MVPA to detect information transmission outside V1/V2, if there is any. **b**) Whole-brain MVPA during the Neural Reinforcement sessions. Overall, information transmission from V1/V2 was mostly confined to the early visual cortex. However, there was a notable transmission to striatum, mostly caudate, consistent with its role in reinforcement learning^{21,22}. A white circle highlights the significant information transmission in striatum. In **a** & **b**, Fisher-transformed correlation coefficient for the significant voxels are shown ($P < 0.05$; multiple corrections with permutation procedure⁴¹). **c**) Disengagement of VMPFC and successful fear reduction. Less information transmission between V1/V2 and VMPFC during Neural Reinforcement was related with larger reduction of SCR (Control CS+ - Target CS+) in early Test (Spearman's $\rho = -.522$, $P = .034$). Each data point corresponds to each participant. Solid line represents a least-square regression line. See also Figure S5.