

# Common psychiatric treatments alter affective dynamics

Reviewed Preprint

v1 • September 5, 2025

Not revised

Quentin Dercon , Quentin JM Huys, Robb B Rutledge, Camilla L Nord

Applied Computational Psychiatry Lab, Mental Health Neuroscience Department, Division of Psychiatry and Max Planck UCL Centre for Computational Psychiatry and Ageing Research, Department of Imaging Neuroscience, Queen Square Institute of Neurology, UCL, London, United Kingdom • Department of Psychology, Yale University, New Haven, United States • Wu Tsai Institute, Yale University, New Haven, United States • Department of Psychiatry, Yale University, New Haven, United States • MRC Cognition and Brain Sciences Unit, University of Cambridge, Cambridge, United Kingdom • Department of Psychiatry, University of Cambridge, Cambridge, United Kingdom

 [https://en.wikipedia.org/wiki/Open\\_access](https://en.wikipedia.org/wiki/Open_access) Copyright information

## eLife Assessment

Dercon and colleagues combine common psychiatric treatments with a probabilistic reward learning task and trial-by-trial ratings of affect, confidence, and engagement. Using computational cognitive modeling, they show that, while both treatments serve to counter negative biases in affect and confidence, cognitive distancing and antidepressant medication have dissociable effects on subjective evaluations and reward-based choice behavior. This work provides **convincing** evidence regarding an **important** line of investigation into the dynamic integration of affect, cognition, and learning.

<https://doi.org/10.7554/eLife.107269.1.sa3>

## Abstract

Affective states are dynamic, fluctuating in response to recent events: an unexpected pleasure, a disappointing loss. Affective biases, which cause disruptions in these dynamics, are core components of mental ill-health, but the specific effects of treatments on these biases are poorly understood. Here, we investigate the impact of common psychiatric treatments on subjective assessments of happiness, confidence, and engagement during a reinforcement learning task ( $N = 935$ ; 130 taking antidepressant medications). Half ( $N = 459$ ) of the participants were randomised to practice a common psychotherapeutic technique—cognitive distancing—throughout the task. From a joint computational model of learning and affect, we find evidence for distinct and overlapping impacts of psychiatric treatments on affective dynamics. Cognitive distancing attenuates downward drift in happiness and engagement and increases recency bias in the affective impact of recent choices. Conversely, antidepressant use increases baseline happiness and confidence in individuals with similar levels of current symptoms, and decreases recency bias such that more past events influence affective states. Crucially, both cognitive distancing and antidepressant use converge to dampen negative biases in happiness and confidence specifically in participants experiencing higher levels of anxiety and depression symptoms. Together, our results indicate that common treatments for

mental ill-health may alter symptoms through their impact on affective dynamics, but via distinct mechanisms.

## 1 Background

How are you feeling right now? Research across economics, psychology, and health sciences suggests the answer to this question—your subjective well-being—is closely tied to objective quality of life<sup>1,2</sup> and health across the lifespan<sup>3</sup>. But feelings are far from static, momentarily fluctuating in response to recent events<sup>4–6</sup>, and even individual choices. Frequently asking participants to rate their feelings enables a read-out of moment-to-moment changes in subjective well-being, or their *affective dynamics*.

In influential work, Rutledge *et al.* (2014)<sup>5</sup> demonstrated momentary happiness ratings during a gambling task could be accurately predicted by a computational model incorporating the average reward for a gamble (*expected value*) and the outcome of the gamble minus this average (*prediction error*). Using a task where reward magnitude and probability were uncorrelated, Blain & Rutledge (2020)<sup>7</sup> subsequently showed that momentary happiness is particularly sensitive to changes in learning-related variables—specifically, prediction errors for reward probability—as compared to as compared to reward information that was relevant to behaviour but not learning. Links between happiness ratings and learning-related quantities may extend to subjective assessments of other affective states. Theoretical accounts posit that momentary confidence is the approximate probability a choice is correct<sup>8,9</sup> (though see<sup>10,11</sup>), while effort costs decrease the value of choices independently of reward probability<sup>12,13</sup>, in turn influencing momentary engagement. Together, these results suggest that affective dynamics are closely coupled with objective quantities that drive choices during learning.

Biases in subjective affect are a core feature of mental ill-health. Symptoms of depression have been consistently linked to lower<sup>7,14</sup> and less stable<sup>15,16</sup> momentary happiness, while transdiagnostic features of mental ill-health have been linked to biased confidence judgements at different timescales<sup>17–20</sup> and impairments in motivation and engagement<sup>13,21–23</sup>. Affective biases maintain symptoms of mental ill-health by inducing changes in emotion processing and perception<sup>24,25</sup>. For example, negatively biased perception—a common feature of depression<sup>26</sup>—may cause low mood by making outcomes appear less rewarding; low mood in turn further negatively biases perception, causing a positive feedback loop which spirals toward a depressive episode<sup>27</sup>. Successful psychiatric treatment may act to perturb these maladaptive cycles. Short-term selective serotonin reuptake inhibitor (SSRI) administration induces positive perceptual biases in healthy participants<sup>28</sup>, suggesting that affective biases may be an early target of antidepressant drugs, acting to shift choices away from those that maintain low mood<sup>29</sup>. Crucially, given that affective biases are precipitated and maintained by negative thinking patterns—a core target of cognitive psychotherapy<sup>30,31</sup>—they may represent a transdiagnostic treatment target of psychological and pharmacological interventions for symptoms of mental ill-health.

Here, we aimed to link choice behaviour to affective dynamics throughout a reinforcement learning (RL) task<sup>32–34</sup> (**Figure 1A**) and to relate this to mental ill-health symptoms and treatments. We asked online participants ( $N = 935$ ) to rate their feelings (from 0–100) on one of three different affect scales—happiness, confidence, and engagement—after receiving feedback on each choice they made. Half (49.1%) of the participants were randomised to an acute psychological intervention known as “cognitive distancing”, a common<sup>35</sup> and effective<sup>31</sup> component of psychological therapy which alters learning in this task<sup>34</sup>. We also collected information on current antidepressant medication use in a demographic questionnaire (reported by 13.9% of participants), and derived transdiagnostic mental health symptom factor scores from a psychiatric questionnaire battery<sup>36,37</sup>. We then assessed how participants’ affect ratings across each of

the three scales covaried with learning-related outcomes throughout the task, accounting for underlying affective biases, using computational modelling. By quantifying the associations between model-derived measures of affective dynamics and transdiagnostic features of mental ill-health, the cognitive distancing intervention, and self-reported antidepressant medication use, we asked whether affective dynamics may be a common target of treatments for symptoms of mental ill-health.

## 2 Methods

### 2.1 Online experiment and sample

A total of 995 participants were recruited via Prolific<sup>38</sup> over three weeks in April–May 2021. Participants were recruited in batches with fixed pre-screeners for age range, gender, and history of any mental health diagnosis, which resulted in a sample broadly representative of the UK population in terms of these characteristics (see<sup>34</sup> for details). After completing a demographic questionnaire, which included questions on current medication usage and mental health diagnoses, participants completed the RL task described below. They then took a short test of working memory (visual digit span), and answered questions from several psychiatric questionnaires, answers from which were used to derive three validated transdiagnostic features of mental health (anxiety/depression, compulsive behaviour, and social withdrawal)<sup>36</sup>, using methods for computational factor modelling<sup>37</sup> described in detail elsewhere<sup>34,39</sup>. Sixty participants were excluded for meeting pre-registered criteria<sup>34</sup>. The study was approved by the University of Cambridge Human Biology Research Ethics Committee (HBREC) (HBREC.2020.40) and jointly sponsored by the University of Cambridge and Cambridge University Hospitals National Health Service (NHS) Foundation Trust (IRAS ID 289980). All participants provided written informed consent through an online form, in line with University of Cambridge HBREC procedures for online studies.

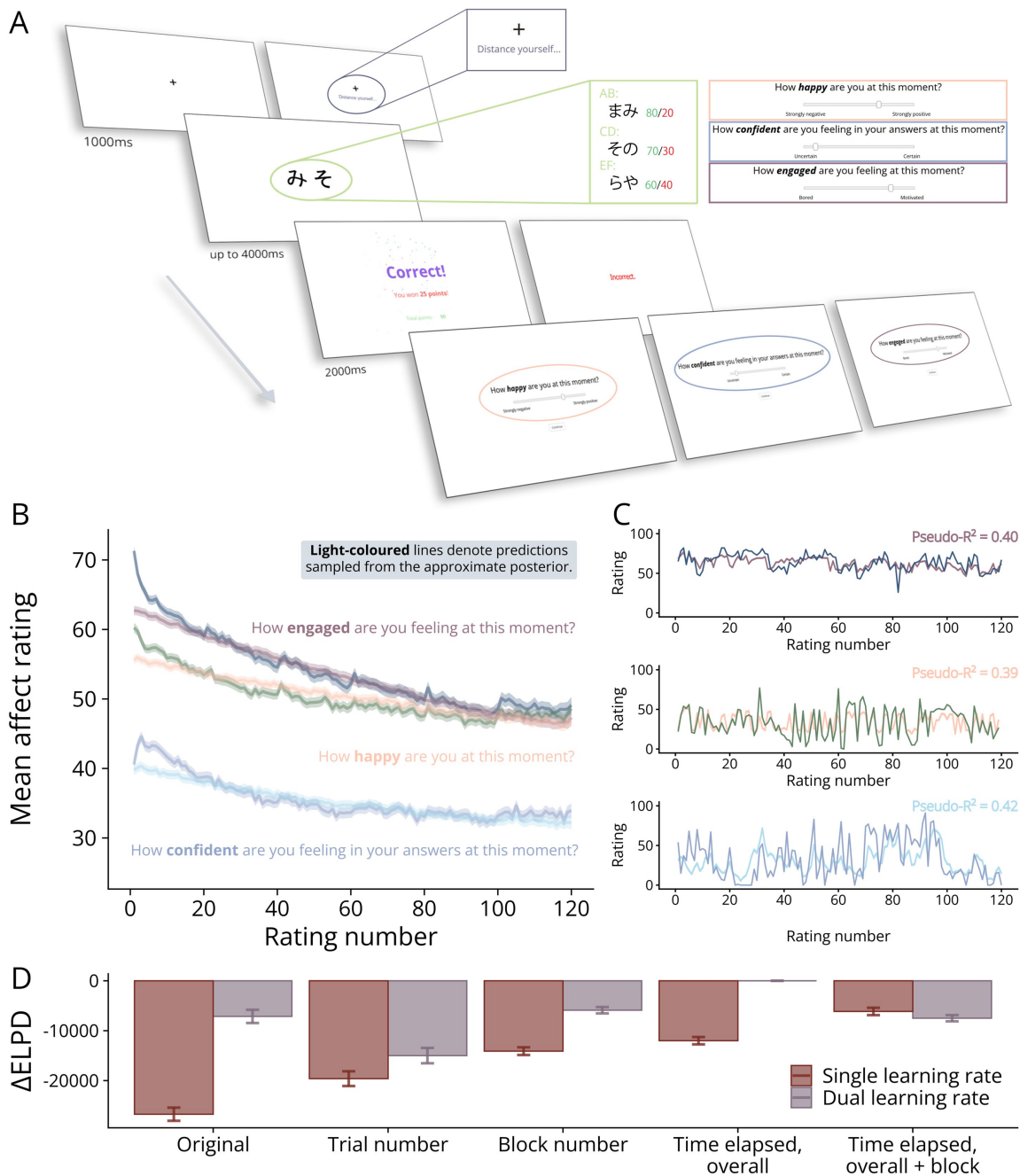
### 2.2 Reinforcement learning task

The reinforcement learning task in the present study—the probabilistic selection task<sup>32,33</sup>—involved learning which symbol in each of three pairs was more likely correct. Consistent choice of the “better” symbol in each pair enabled a participant to accumulate more points (and maximise their chances of winning a monetary bonus). One symbol in each pair was always more likely correct, but the contingencies varied across the pairs, from 0.8/0.2 (‘AB’) to 0.7/0.3 (‘CD’), to 0.6/0.4 (‘EF’). All participants saw the same six symbols, but the pairs were randomised across individuals and counterbalanced across trials. After making a choice, participants received feedback (“Correct!” or “Incorrect.”), and then rated their subjective happiness, confidence, or engagement (**Figure 1A**).

One of the three questions was asked after each trial outcome, with each question asked twenty times per block of sixty trials, and never more than twice in a row. Participants were also asked to rate (again from 0–100) how fatigued they felt compared to the beginning of the block after the end of each of the sixty-trial training blocks. After six training blocks, participants were tested on all fifteen unique pairs without feedback. We previously reported the effects of cognitive distancing on task performance and learning, including results of the test phase, in the same sample<sup>34</sup>.

### 2.3 Acute psychological intervention and antidepressant use

Half of the participants ( $n=459$ ; 49.1%) were randomised to be taught, and then practice throughout the task, a psychotherapeutic technique termed “cognitive distancing”, which encouraged them to “take a step back” from their emotional reactions to feedback throughout the task (see here<sup>34</sup> for further details). Apart from an additional instructional video before the task



**Figure 1**

**Task design, affect model posterior predictions and model comparison.**

**A.** The task design. **B.** Mean affect ratings for each rating type (engaged, happy, or confident), by distancing group, compared to model predictions (light-coloured lines). **C.** Example comparison of predictions from the best-fitting model (light-coloured lines) to raw affect ratings from three different individuals with the median pseudo-R<sup>2</sup> for each rating type. **D.** Model fit compared to the best-fitting model (time elapsed, overall with dual learning rate) in terms of their ELPD (i.e., higher ELPD [or less negative ELPD compared to the best-fitting model] is better), estimated via Bayesian approximate LOO cross-validation<sup>47</sup>. Ribbons in B-C and error bars in D denote standard errors.

started and a small prompt to “Distance yourself...” which appeared with each fixation cross (Figure 1A), the task was identical for distanced and non-distanced participants. To explore similarities between the effects of this psychological intervention, and a pharmacological treatment for mental ill-health, antidepressant medication, we also asked participants to report their current medication use: 130 participants (13.9%) reported current antidepressant use, with the majority ( $n=94$ ; 72.3%) taking an SSRI.

## 2.4 Computational modelling: joint RL-affect models

For consistency with previous literature, we used the well-characterised model of momentary happiness first described by Rutledge *et al.* (2014) as a baseline model. This model assumes fluctuations around a baseline (i.e., longer-term mean) can be captured by a weighted sum of recent expected values and prediction error and, importantly, does not condition on previous ratings.

The Rutledge *et al.* (2014) model has been primarily validated in (e.g., gambling) tasks where expected values and prediction errors are explicitly available to participants. As such, we extended it to account for learning in this task. Specifically, our model—which we term a joint RL-affect model—comprised two components: (1) a  $Q$ -learning model to infer expected values and prediction errors from participants’ choices (which has been shown to accurately capture choice behaviour in this task); and (2) a model for momentary affect which assumes fluctuations around a baseline can be captured by a recency-weighted sum of  $Q$ -learning model-derived expected values and prediction errors. Hierarchical models were simultaneously fitted to task choices plus happiness, confidence, and engagement self-reports, assuming different parameter weights across all scales and participants.

### 2.4.1 RL models

A  $Q$ -learning model infers expected values (termed  $Q$ -values) from participant choices—the action ( $a_t$ ) of choosing one symbol over the other in each of the three pairs (denoted as states,  $s_t$ )—by assuming they update at each trial  $t$  based on prediction errors  $\delta_t$ , with the update magnitude controlled by a learning rate  $\alpha \in [0,1]$ :

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha \delta_t \quad \text{where} \quad \delta_t = r_t - Q_t(s_t, a_t). \quad (1)$$

We additionally considered a dual learning rate  $Q$ -learning model for choices, in which the learning rate is assumed to differ depending on whether the outcome was rewarding ( $\alpha_{\text{reward}}$ ) or not ( $\alpha_{\text{loss}}$ ):

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha_{\text{reward}}[\delta_t]_+ + \alpha_{\text{loss}}[\delta_t]_-. \quad (2)$$

In both cases, the difference in  $Q$ -values between the chosen ( $a_t$ ) and avoided action ( $\bar{a}_t$ ) is transformed to a choice probability using a softmax function, weighted by an inverse temperature  $\beta$ :

$$P_t(s_t, a_t) = \left( 1 + \exp \left\{ -\beta [Q_t(s_t, a_t) - Q_t(s_t, \bar{a}_t)] \right\} \right)^{-1}. \quad (3)$$

### 2.4.2 Affect models

## Baseline model

Affect ratings scaled to [0, 1] for each participant and rating type (i.e., happiness, confidence or engagement) were assumed to be drawn from independent Beta distributions with a mean-variance reparameterization which models the shape parameters as functions of a (conditional) mean and precision<sup>42</sup>. Extreme ratings (0 or 1) were allowed in the task, so we transformed ratings to the (0,1) interval using a simple transformation<sup>43</sup> ( $y' = \frac{y(N-1)+0.5}{N}$ ; here  $N$  is the total number of participants but could be any large number). A logit link function was used, so the Beta regression weights (excluding the intercept  $w_0$ ) can be interpreted as the log odds ratio per unit change in the covariate for an increase in affect rating, with other terms held constant (note that our approach differs slightly from that of Rutledge *et al.* (2014)<sup>5</sup>, who assume momentary happiness ratings follow a Gaussian distribution; see Forbes & Bennett (2024)<sup>41</sup> for a validation of this Beta regression approach). Following<sup>5,41</sup> we write the  $i$ th participant's affect rating  $Y^{(i)}$  for rating type  $p$  at trial  $t$  as,

$$Y_{t,p}^{(i)} \sim \text{Beta}\left(\mu_t^{(i)} \phi_p, (1 - \mu_t^{(i)}) \phi_p\right) \quad (4)$$

$$\log\left(\frac{\mu_t^{(i)}}{1 - \mu_t^{(i)}}\right) = w_0^p + w_2^p \sum_{t'=1}^t \gamma_p^{t-t'} Q_{t'}(s_{t'}, a_{t'}) + w_3^p \sum_{t'=1}^t \gamma_p^{t-t'} [r_{t'} - Q_{t'}(s_{t'}, a_{t'})], \quad (5)$$

where  $t$  is the overall trial number at rating number  $i$  for rating type  $p$ ,  $\gamma$  is the discount or 'forgetting' factor which imposes a strict weighting on recent trials, and  $w_k^p$  the weights on each of the  $k$  quantities of interest for the  $p$ th rating type; and  $Q_t$  and  $r_t$  are the  $Q$ -learning model-derived expected values for the chosen symbol  $a_t$  in the state  $s_t$  (i.e., the pair of symbols presented on trial  $t$ ) and the feedback valence ( $\pm 1$  for correct/incorrect to allow for negative  $Q$ -values) respectively, both at trial  $t$ . Note that both sums are over all previous trials, not just those of rating type  $p$ .

## Accounting for drift over time

We also fit models including an extra weight  $w_1^p$  to account for potential "drift over time" in affect<sup>44</sup>, by modifying (5) as follows,

$$\log\left(\frac{\mu_t^{(i)}}{1 - \mu_t^{(i)}}\right) = w_0^p + w_1^p \tau_t + w_2^p \sum_{t'=1}^t \gamma_p^{t-t'} Q_{t'}(s_{t'}, a_{t'}) + w_3^p \sum_{t'=1}^t \gamma_p^{t-t'} [r_{t'} - Q_{t'}(s_{t'}, a_{t'})], \quad (6)$$

where  $\tau_t$  is some measure of time elapsed up to trial  $t$ : either trial number, block number, overall time elapsed, or time elapsed since the start of that block (see Model comparison).

## 2.4.3 Fits to data

Models were fitted in a hierarchical Bayesian manner, with approximate posteriors derived via automatic differentiation variational inference (ADVI)<sup>45</sup> implemented in CmdStan<sup>46</sup>. All models were fit to choices and ratings on all three affect scales simultaneously across both distancing and non-distancing participants, with separate weights and decay factors assumed for each person and question, and separate group-level (hyper)priors on each parameter. In other words, participant-level parameter distributions are assumed to be conditionally independent given the group-level distribution over that parameter. Individual-level predictive accuracy was

assessed by comparing responses predicted from each participant's approximate posterior to their observed affect ratings via pseudo- $R^2$ , defined following previous work<sup>42</sup> as the squared correlation between observed and mean posterior predictions.

#### 2.4.4 Model comparison

In the affect model, we tested for “drift over time”<sup>44</sup>; and in the Q-learning model we tested for separate learning rates for rewarding and non-rewarding outcomes. We assumed the Rutledge *et al.* (2014)<sup>5</sup> model (equation 5) to be the baseline model, and so the parameters in this model were included in all models.

Drift over time in affect may be particularly relevant to our task, as participants were able to take as long as they wished to rate their subjective feelings, and time between blocks was additionally unconstrained. As such, we compared four models with different measures of time elapsed (i.e., variants of equation 6) to the baseline model (equation 5), and either single or dual learning rates. The extra parameter added linear weights on either trial number, block number, or total time elapsed. We also tested a final model with two extra parameters: weights on both total time elapsed and time elapsed since the beginning of that block. We then compared all ten models in terms of their approximate leave one out (LOO) expected log pointwise predictive density (ELPD), a metric of estimated out-of-sample predictive accuracy<sup>47</sup>, corrected for the use of variational approximations to the true posterior<sup>48</sup>.

### 2.5 Statistical analysis

For consistency with computational modelling, we adopt a fully Bayesian approach for statistical analyses where possible. As such, results are given as estimates with a highest density interval (HDI), which (unlike a confidence interval (CI)<sup>49</sup>) can be interpreted as the probability that the true value falls within a given range. We report 95% HDIs to align with convention but interpret results in terms of strength of evidence throughout; an overlap with the null value should not be seen as evidence for lack of an effect, but rather weakened evidence for it.

#### 2.5.1 Associations between model parameters and mental health symptoms & treatments

We tested the effect of differences in transdiagnostic mental health symptoms, current self-reported antidepressant use, or cognitive distancing on model parameters using generalised linear models (GLMs), adjusted for age, gender, and working memory capacity (measured with visual digit span), separately for each rating type. GLMs relating model parameters to antidepressant use were run with and without adjustment for concurrent anxiety/depression symptoms (i.e., factor score), as medication use was not randomised. We also considered whether the effect of cognitive distancing and current antidepressant use on affect model parameters may differ in relation to their association with transdiagnostic mental health, by including factor score as an interaction term in GLMs.

Posterior samples for GLM coefficients were obtained via Markov chain Monte Carlo (MCMC) implemented in CmdStan<sup>46</sup>, with 2,000 warm-up and 10,000 sampling iterations for each of four chains, using models and priors from the rstanarm R package<sup>50</sup>. Response distributions were assumed Gaussian for all parameters except for learning and decay rates, which were modelled via Gamma family GLMs (with log link functions). See Interpretation of model-derived parameters in the Supplementary Methods for intuition as to how these regression coefficients are interpreted.

## 2.5.2 Differences in associations with baseline affect between transdiagnostic factors

To account for potential collinearities in the three transdiagnostic mental health symptom factor scores obtained via computational factor modelling, we used partial least squares (PLS) regression to test which of the transdiagnostic factor scores was most strongly associated with baseline affect ( $w_0$ ). PLS regression is a data-driven method which identifies latent components linking multivariate responses to predictors based on shared covariance, and so is well-suited to the problem of multicollinearity<sup>51</sup>.

In line with best-practices<sup>39</sup>, we first identified the number of components that best described our data in a training set (80% of participants) in terms of mean squared error (MSE) and  $R^2$  via ten-fold cross-validation. We then validated the predictive accuracy of this number of components in held-out test data (20% of participants), and formally tested this using a permutation test, where the PLS regression model was re-trained on 10,000 training datasets with shuffled outcome labels (providing a null distribution), and the fraction of these datasets where the MSE between the test data and the predictions from the permuted datasets was lower than the true train-test MSE taken as the  $p$ -value. The PLS regression with the chosen number of components was then refitted in the whole dataset, to obtain component loadings on each of the responses and predictors. Lastly, we computed bias-corrected and accelerated (BCa) CIs for each of the loadings, plus the differences in loadings between transdiagnostic factors, from 10,000 bootstrap samples.

## 3 Results

### 3.1 A computational model of subjective happiness accounting for learning and affective drift also captures momentary confidence and engagement

Following model comparison, we found that the best-fitting RL-affect model included separate learning rates for rewarding and non-rewarding outcomes and a linear effect of time elapsed since the beginning of the task (equation 6; **Figure 1D**). This model, when fitted to all affect ratings and participants simultaneously, explained participants' variance in happiness, confidence, and engagement assessments with similar accuracy (mean [standard deviation (SD)] pseudo- $R^2 = 0.40$ – $0.42$  [0.23–0.26]; see **Figure 1C** for example individuals).

### 3.2 Baseline affect is negatively associated with transdiagnostic mental health

We first assessed whether individuals' estimated model parameters from the best-fitting joint RL-affect model were associated with transdiagnostic features of mental health.

In line with previous work<sup>7,14</sup>, we found strong evidence for a negative association between baseline happiness ( $w_0^{happ}$ ) and anxiety/depression symptoms from a linear model adjusting for age, gender, digit span, and distancing (mean 4.44-point lower baseline happiness rating per SD increase in factor score; 95% HDI = [−5.46, −3.41]). Higher anxiety/depression factor scores were additionally associated with lower baseline confidence (mean 3.75-point lower  $w_0^{conf}$  per SD increase in score; 95% HDIs = [−5.01, −2.46]) and lower baseline engagement (mean 2.90-point lower  $w_0^{eng}$  per SD increase in score; 95% HDIs = [−4.12, −1.66]; **Figure 2A**). We also found that higher anxiety/depression factor scores were associated with lower odds of increases over time in happiness (mean 11.2% lower  $w_1^{happ}$  per SD increase in anxiety/depression score per hour; 95% HDI for multiplier = [0.808, 0.971]) and engagement (17.0% lower  $w_1^{eng}$  for a unit increase in

anxiety/depression score per hour; 95% HDI for multiplier = [0.726, 0.943]; **Figure 2A** [ii](#)), suggesting a higher rate of decline (i.e., more drift) in happiness and engagement in participants with higher anxiety/depression scores.

Higher compulsive behaviour and social withdrawal factor scores were also each associated with lower baseline happiness (e.g., mean 2.20-point lower  $\mu_{\text{happiness}}$  [95% HDI = (-3.31, -1.06)] per SD increase in compulsive behaviour score; **Figure 2B** [i](#)) and confidence (e.g., mean 3.59-point lower  $\mu_{\text{confidence}}$  [95% HDI = (-3.98, -1.38)] per SD increase in social withdrawal score; **Figure 2C** [i](#)), but were not associated with greater drift in affect over time (**Figure 2B** [ii](#) & **Figure 2C** [ii](#)). Associations between transdiagnostic symptom scores and choice-related affect measures (i.e.,  $\mu_{\text{affect}}$  and  $\mu_{\text{confidence}}$  were also observed (**Figure S3** [i](#)), as were strong associations between post-block fatigue ratings and both baseline and drift in affect over time (**Figure S5** [i](#); see Supplementary Results for more details).

### 3.3 Baseline affect is most strongly associated with anxiety/depression symptoms

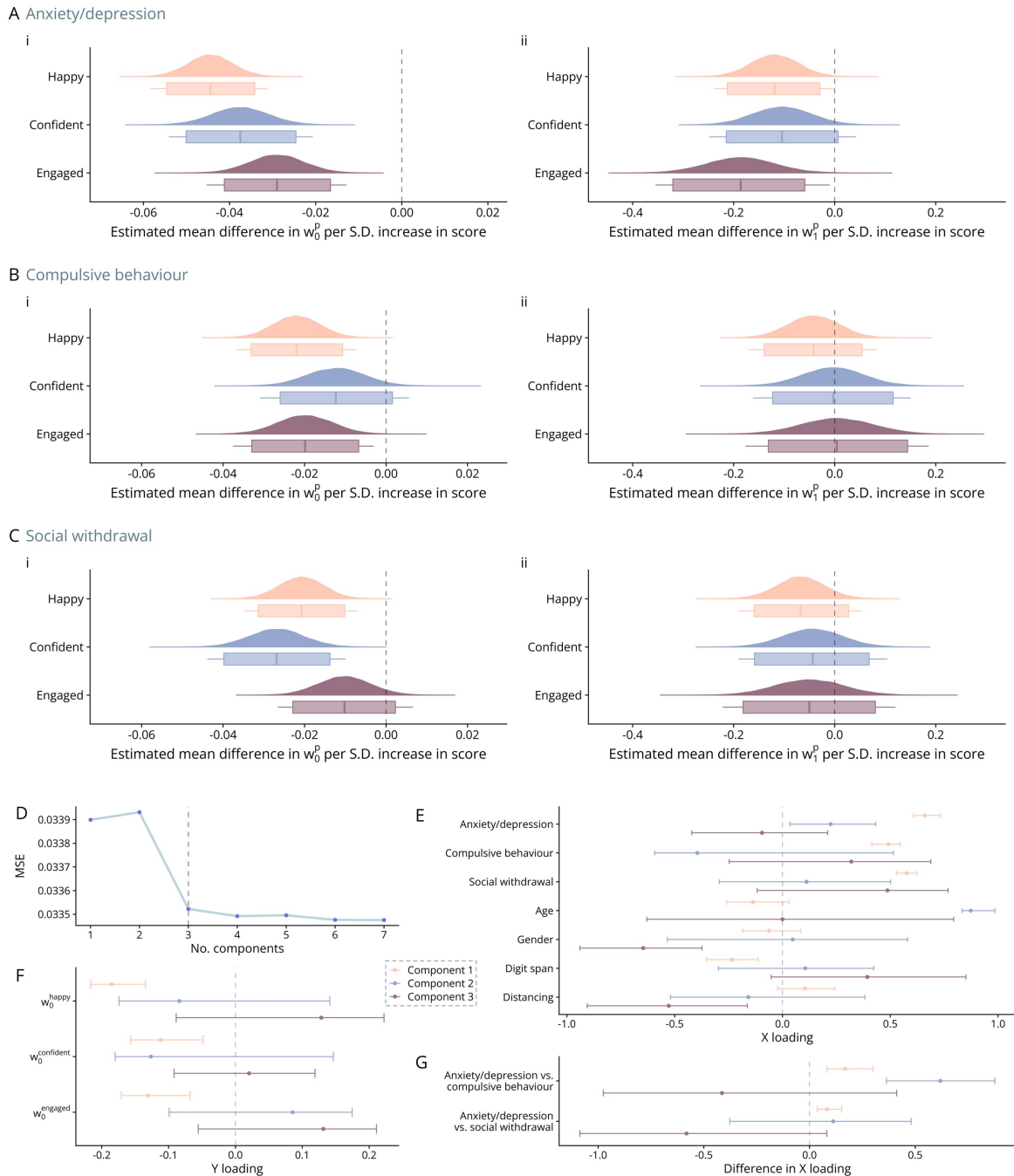
The three transdiagnostic mental health symptom factor scores obtained via computational factor modelling were highly correlated ( $r = 0.47$  [95% CI: (0.42, 0.52)] between anxiety/ depression and compulsive behaviour;  $r = 0.61$  [95% CI: (0.57, 0.65)] between anxiety/ depression and social withdrawal;  $r = 0.42$  [95% CI: (0.37, 0.47)] between compulsive behaviour and social withdrawal). As such, to compare the strength of associations between baseline affect and transdiagnostic symptom factors, we used a partial least squares regression model to relate baseline affect to the three scores plus age, gender, digit span, and distancing.

We found that a three-component model represented the best compromise between predictive accuracy (in training data) and parsimony (**Figure 2D** [i](#)), and there was strong statistical evidence that this model could accurately predict responses in held-out test data (permutation test  $p < 0.0001$ ). The first component of the model negatively loaded on baseline happiness (loading = -0.185, BCa bootstrapped 95% CI = [-0.217, -0.135]), confidence (loading = -0.112, BCa bootstrapped 95% CI = [-0.157, -0.049]), and engagement (loading = -0.131, BCa bootstrapped 95% CI = [-0.171, -0.068]) (**Figure 2F** [i](#)). It also positively loaded on each of the transdiagnostic symptom factors: anxiety/depression (loading = 0.660, BCa bootstrapped 95% CI = [0.608, 0.733]), compulsive behaviour (loading = 0.491, BCa bootstrapped 95% CI = [0.415, 0.545]), and social withdrawal (loading = 0.577, BCa bootstrapped 95% CI = [0.529, 0.623]). The other two components did not show this pattern (**Figure 2E** [i](#)). There was strong evidence that the first component's loading was higher for anxiety/depression than both compulsive behaviour and social withdrawal (component 1 loading difference [BCa bootstrapped 95% CI] = 0.169 [0.083, 0.299] and 0.083 [0.037, 0.152] respectively; **Figure 2G** [i](#)).

### 3.4 Cognitive distancing slows affective drift and antidepressant use positively modulates baseline affect

We then assessed the evidence for effects of cognitive distancing and antidepressant use on choice-independent affective dynamics: baseline affect and its drift.

Previously, in the same sample, we reported evidence from linear mixed models that participants practising cognitive distancing declined slightly less in happiness and engagement, but not confidence, over the course of the task<sup>34</sup> [i](#). Evidence from our RL-affect model was consistent with a decrease in affect drift over time: distancing individuals on average drifted less across the task in happiness (estimated mean 21.3% higher odds of increase in happiness over per hour; 95% HDI for multiplier = [1.01, 1.46]; **Figure 3A** [ii](#)), in spite of lower baseline engagement (estimated



**Figure 2**

**Associations between affect parameters and transdiagnostic mental health dimensions, and results of PLS regression.**

**A-C.** Estimated differences in baseline affect ( $w_0^p$ ; *i*) or drift in affect over time ( $w_1^p$ ; *ii*) for a unit increase in anxiety/depression (A), compulsive behaviour (B), or social withdrawal (C) transdiagnostic symptom factor score. **D-G.** Results of partial least squares regression: elbow plot of ten-fold cross-validated mean squared error for models with increasing numbers of components (D); loadings of the three-factor model on independent (E) and response (F) variables; and loading differences between anxiety/depression and other transdiagnostic factors (G). Boxplot boxes in A-C denote 95% HDIs and lines denote 99% HDIs; error bars in E-G denote BCa bootstrapped 95% CIs.

mean = -2.89 points; 95% HDI = [-5.42, -0.464]; **Figure 3A** [↗ i](#)). There was also some weak evidence of less drift in engagement in participants randomised to the intervention (17.4% higher  $\gamma_{\text{engaged}}$  in distancing individuals; 95% HDI for multiplier = [0.91, 1.52]; **Figure 3A** [↗ ii](#)).

There was limited evidence for any difference in affective drift in participants taking antidepressants (**Figure 3B** [↗ ii](#)). However, there was evidence that participants self-reporting current antidepressant use had 3.50-point higher baseline happiness (95% HDI = [0.311, 6.60]) and 2.69-point higher baseline confidence (with much weaker evidence: 95% HDI = [-1.16, 6.54]), after adjusting for anxiety/depression symptom scores (**Figure 3B** [↗ i](#)).

### 3.5 Cognitive distancing and antidepressant use have opposite effects on the weighting of choices in subjective happiness

Next, we quantified the associations between treatments—cognitive distancing and antidepressant use—and choice-dependent parameters (i.e.,  $\omega_{\text{e}}$  and  $\omega_{\text{t}}$ ) which control the extent of trial-to-trial fluctuations in affect.

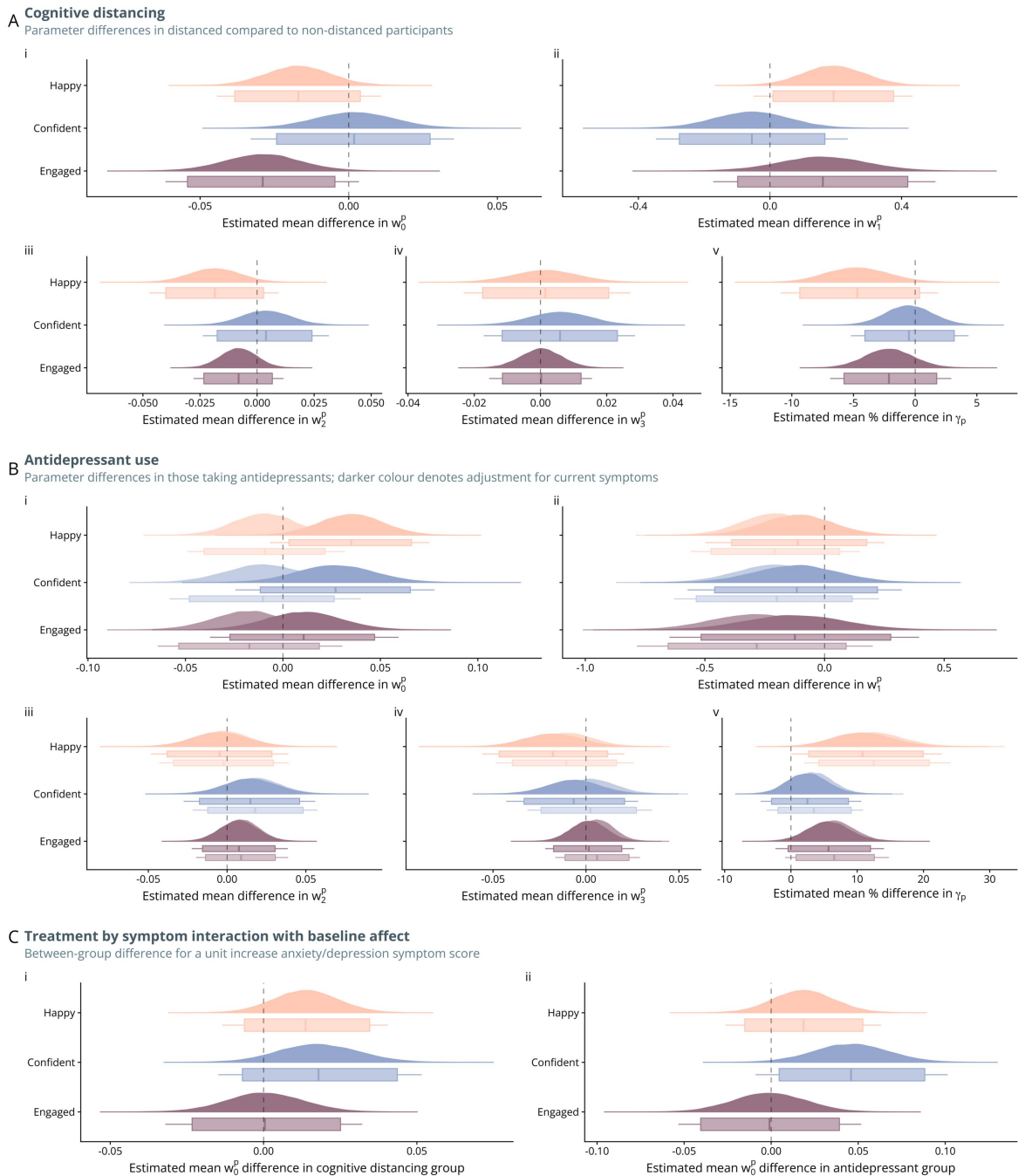
There was limited evidence for an association between either treatment and weights on recent prediction errors ( $\omega_{\text{e}}$ ; **Figure 3A** [↗ iv](#) & **Figure 3B** [↗ iv](#)). However, there was some evidence that cognitive distancing lowered weighting of recent expected (choice) values in happiness ratings (**Figure 3A** [↗ iii](#)), with 1.83% (95% HDI for multiplier=[0.961, 1.003]) lower odds of increased happiness ratings for the same weighted sum of recent expected values in distanced participants. Meanwhile, exploratory analyses with an extended between-rating model showed a specific effect of antidepressant use on the weighting of recent expected values in engagement and confidence ratings (see Supplementary Results & **Figure S4B** [↗ iii-iv](#)).

Current antidepressant use was associated with less forgetting of choices and outcomes in happiness ratings (12.5% higher  $\gamma_{\text{happy}}$ ; 95% HDI for multiplier = [1.04, 1.21]) and engagement (6.52% higher  $\gamma_{\text{engaged}}$ ; 95% HDI for multiplier = [1.01, 1.13]); **Figure 3B** [↗ iii](#)), suggesting higher weighting of earlier trials' expected values and prediction errors in subsequent affect ratings. Evidence for this association remained, albeit slightly weakened, after additionally adjusting for current anxiety/depression symptoms (**Figure 3B** [↗ iii](#)), which were themselves positively associated with  $\gamma_{\text{happy}}$  and (to a lesser extent)  $\gamma_{\text{engaged}}$  (**Figure S3A** [↗ iii](#)). Notably, the converse effect was seen in distancing participants, with the psychological intervention associated with lower happiness forgetting factors, albeit with weak evidence (4.69% lower  $\gamma_{\text{happy}}$ ; 95% HDI for multiplier = [0.906, 1.004]; **Figure 3A** [↗ iii](#)).

### 3.6 Cognitive distancing and antidepressant use dampen negative associations between baseline affect and anxiety/depression symptoms

Lastly, we explored whether the negative associations between choice-independent affective dynamics and transdiagnostic anxiety/depression symptoms were altered by cognitive distancing or current antidepressant use, by including treatment by symptom interactions in outcome GLMs.

We found that both the distancing intervention and antidepressant use weakened the negative associations between baseline happiness and confidence, but not engagement. Specifically, distancing individuals with higher anxiety/depression scores had higher baseline happiness and confidence (mean 1.37-point higher  $\omega_{\text{happy}}$  and 1.79-point higher  $\omega_{\text{confidence}}$  per SD increase in anxiety/depression score respectively) relative to non-distancing participants with the same symptom scores, though with weak evidence (95% HDIs for this distancing by symptom interaction = [-0.63, 3.46] for baseline happiness and [-1.11, 6.12] for confidence (**Figure 3C** [↗ i](#))). These effects were mirrored in participants taking antidepressants. The evidence for an antidepressant by anxiety/depression symptom interaction with respect to baseline happiness was



**Figure 3**

**Associations between treatments and affect model parameters.**

**A-B.** Estimated mean differences in individuals' baseline affect ( $w_0^d$ ), drift in affect over time ( $w_1^d$ ) and forgetting rate of previous trials' expected values and prediction errors ( $\gamma_p$ ), in participants practicing cognitive distancing (A) and taking antidepressants (B; darker colour denotes adjustment for current anxiety/depression symptoms). **C.** Interactions between cognitive distancing (i) and antidepressant use (ii) and higher anxiety/depression symptom scores, with respect to baseline affect. **D.** Associations between antidepressant use and expected value weights in affect rating computation: main effect (i), interaction with trial lag (ii), and posterior mean parameters for weighting of previous choices' expected values in engagement ratings ( $w_{t-1}^d$ ) by antidepressant group (iii). In all plots, boxplot boxes denote 95% HDIs, and lines denote 99% HDIs.

fairly weak (1.86-point higher  $w_0^{\text{happy}}$  per SD increase in anxiety/depression score; 95% HDI = [-1.52, 5.27]), but the evidence for the corresponding interaction effect on baseline confidence was stronger (mean 4.59-point higher  $w_0^{\text{confidence}}$  per SD increase in anxiety/depression score; 95% HDI = [0.559, 12.53]; **Figure 3C** [it](#)).

## 4 Discussion

Here, we applied a computational model of momentary happiness which assumes fluctuations in affect ratings depend simply on baseline affect, its drift over time, and recency-decayed expected and received outcomes. By extending this model to also capture learning, we were able to link objective behaviour to subjective feelings across distinct affective states—happiness, confidence, and engagement ratings—and show that a core component of psychological therapy, cognitive distancing, and antidepressant medication use have different effects on affective dynamics, but converge to alter affective biases associated with symptoms of mental ill-health.

There were distinct effects of both treatments on affective dynamics. Randomisation to a psychotherapeutic intervention, cognitive distancing, attenuated declines in happiness and engagement over time, adding to our previously reported findings that this psychotherapeutic technique alters aspects of reward learning<sup>34</sup>. Self-reported antidepressant use, meanwhile, was associated with higher baseline happiness and confidence after adjustment for current anxiety/depression symptoms (as antidepressant use was not randomised), which is consistent with evidence that antidepressants exert positive affective biases<sup>28</sup>. Subsequent exploration of changes in affect revealed further mechanistic divergence: current antidepressant use was associated with lower recency biases across all scales (i.e., forgetting factors closer to one), and cognitive distancing reduced the weighting of expectations and higher recency bias in happiness ratings. Together, these results suggest that psychiatric treatments act to alter the contribution of objective learning-related quantities to subjective value judgements.

Consistent with extensive evidence<sup>14,18,19,26</sup>, we found negative associations between model-derived baseline affect (across all scales) and transdiagnostic psychiatric symptom measures derived from computational factor modelling<sup>36,37</sup>. This effect, which was strongest for anxiety/depression symptom scores, indicates a consistent, time-invariant negative affective bias which scales with mental ill-health symptom load. Critically, we found evidence for a convergent treatment by symptom interaction with baseline affect across both cognitive distancing and antidepressant use: negative associations between anxiety/depression symptoms and baseline happiness and confidence were attenuated in participants with higher anxiety/depression symptom scores. These results support the cognitive neuropsychological model of antidepressant action, whereby antidepressants are proposed to act acutely to revert negative or maladaptive affective biases<sup>29,52</sup>, and suggest that cognitive distancing<sup>31</sup> and other components of psychotherapy may also act clinically to alter affective biases contributing to symptoms of mental ill-health. We propose that changes in affective dynamics should be investigated further in longitudinal studies as a computational predictor of subsequent symptom change.

Our methodological approach also extends previous work in two ways. First, we applied a theory-driven computational model which allows for fluctuations in momentary happiness as a function of the history of expected values and prediction errors resulting from those expectations, which has been primarily validated in tasks where learning is not required<sup>5</sup>. We not only show that this model can capture happiness ratings in a task where expected values are never explicitly available and have to be learned from experience, but can also accurately capture variation in subjective ratings of confidence and engagement. Second, we found evidence for drift in happiness over time, replicating recent work which characterised ‘mood drift over time’<sup>44</sup>, and extended this to both confidence and engagement. We also found that this drift was strongly

associated with self-reported fatigue (**Figure S5**; see Supplementary Results for more details). We note that we did not find evidence of a previously reported effect—reduced mood drift over time with increased depressive symptoms<sup>44</sup>—instead finding evidence to the contrary (**Figure 1A ii**). However, this work primarily reported evidence from short gambling tasks<sup>44</sup>; our results indicate that associations between affective drift and mental ill-health symptoms are not task-invariant.

We note several limitations. Firstly, there were limitations in our outcome measures. Transdiagnostic measures of mental health psychopathology were estimated using questionnaire subsets<sup>34,39</sup>, precluding investigation of associations between parameters and individual diagnostic scales. Antidepressant use, meanwhile, was self-reported, non-randomised, and we did not collect information concerning length of treatment. Secondly, our use of a single computational model for all three affect scales is a powerful approach, but limited in its ability to truly contrast trial-to-trial fluctuations in each individual rating scale, as we are only comparing the contribution of a small number of computational components (i.e., affective weights) to a fraction of their variation; the residual (scale-specific) variation is likely also important in explaining how these ratings overlap and differ. Third, as model complexity meant only approximate (mean-field variational, rather than sampling-based) inference was viable, we were unable to account for uncertainty in estimates of individual-level posterior mean parameters in associations with quantities of interest (e.g., by using precision-weighted GLMs), as the posterior covariance matrix cannot accurately capture local interdependencies, meaning that parameter precisions are not reliable enough for uncertainty-weighted outcome models<sup>45</sup>.

To conclude, we integrated objective choice behaviour in a learning task with trial-to-trial affect ratings across three distinct states—happiness, confidence, and engagement—within a unified computational model. This enabled us to uncover associations between model parameters and treatments for mental health conditions, offering new insights into their underlying mechanisms-of-action. Our results demonstrate the critical importance of affective biases in the maintenance and updating of affective states in mental ill-health, and indicate that existing, effective treatments can be understood at least in part as acting to shift these biases towards the healthy range.

## Supplementary Material

### Supplementary Methods

#### Interpretation of model-derived parameters

In the Results, we report intercept (i.e., baseline affect) parameters ( $w_0^a$ ) following an inverse logit transformation to allow interpretation of the GLM coefficient on the original (point) scale (from 0 to 1) as the difference in baseline affect rating between individuals differing only in the covariate of interest (by one unit). Other GLM-estimated weight parameter differences are interpreted as the difference in the (log) odds of an increase in affect rating between individuals differing only in the covariate that the parameter weights.

For intuition, consider an estimated GLM coefficient of 0.1 for the distancing group in relation to baseline happiness ( $w_0^{\text{happ}}$ ) inverse logit transformed to a 0 to 1 scale. This result indicates an estimated 10-point (on the 0–100 scale) higher baseline happiness rating in distanced participants after covariate adjustment. Meanwhile, a GLM coefficient for the distancing group of 0.1 in relation to  $w_1^{\text{happ}}$  for the time elapsed (overall) model (i.e., where  $w_1^{\text{happ}}$  is the adjusted log odds ratio for an increase in happiness for a one hour increase in elapsed time) suggests that the distanced group were 10.5% more likely than non-distanced participants to have an increase in happiness after an hour (i.e., the estimated multiplier is  $\exp(0.1) = 1.105$ ) (note the exact interpretation (i.e.,

greater increase vs. slower decline in odds) depends on the GLM intercept term, which in this case would be the estimated  $w_1^{\text{happy}}$  in non-distanced participants after covariate adjustment). Similarly, a GLM coefficient for the distancing group of 0.1 in relation to  $w_2^{\text{happy}}$  suggests that a unit increase in  $Q$ -value is 10.5% more likely to produce an increase in happiness rating, suggesting a higher contribution of (recent) expected values in the affect rating computation.

## Between-rating RL-affect model

### Model definition and fit

In an exploratory analysis, we modified (6) to partition the weights on expected values and prediction errors ( $w_2$  and  $w_3$ ) into individual weights on the outcomes of each previous trial since the previous rating (in our task, this could be up to four intervening outcomes plus the current trial). As such, the lagged trial weight parameters  $w_{k,t-1}$  and  $w_{k,t-2}$  can be seen as capturing the between-rating change in affective dynamics. We hence modify (6) as follows,

$$\log\left(\frac{\mu_t^{(i)}}{1 - \mu_t^{(i)}}\right) = w_0^p + w_1^p \tau_t + \sum_{t'=0}^I w_{2(-t')}^p Q_{t-t'}(s_{t-t'}, a_{t-t'}) + \sum_{t'=0}^I w_{3(-t')}^p [r_{t-t'} - Q_{t-t'}(s_{t-t'}, a_{t-t'})], \quad (7)$$

where  $I$  denotes the number of ratings between the current rating and the last rating of the same type; in this task,  $I \leq 4$ .  $w_{k,t-1}$  are the weights on the outcomes  $t'$  trials back, so  $w_{k,t}$  is the weight on the outcome of the current trial,  $w_{k,t-1}$  is the weight on the outcome of the previous trial, and so on. This model was fit to choices and ratings on all three affect scales simultaneously in a hierarchical Bayesian manner via ADVI<sup>45</sup> as explained in Methods. The only difference was that lagged trial weights were assumed drawn from  $w_2$  and  $w_3$  priors specific to each individual which in turn were drawn from overall group-level hyperpriors (i.e., a three-level hierarchy).

### Associations between between-rating RL-affect model parameters and treatments

The associations between the weights on expected values and prediction errors for individual trials and cognitive distancing and antidepressant use were quantified using multilevel Bayesian linear regression models implemented in the brms R package<sup>53</sup> and CmdStan<sup>46</sup>. All models controlled for the same covariates (i.e., age, gender, digit span), but also included a participant-level random intercept and slopes (on trial lag) to account for the fact there were five parameters per person for each affect rating, as well as the main effect of trial lag and its interaction with each treatment. Separate models were fit for each affect rating and parameter (i.e.,  $w_{2(-t)}$  and  $w_{3(-t)}$ ).

### Parameter recovery

#### Joint RL-affect model with mood drift over time

To test whether we could recover known parameter values from the best-fitting model (i.e., dual learning rate, time elapsed over time), we simulated one hundred datasets (including choices and affect ratings) with parameters drawn from the following distributions:

We then fit the model to these simulated data with approximate inference, and compared the posterior mean parameter estimated to those known to have generated the data. We found that all parameters could be recovered with high accuracy ( $r > 0.87$ ; **Figure S1A**).

Table S1

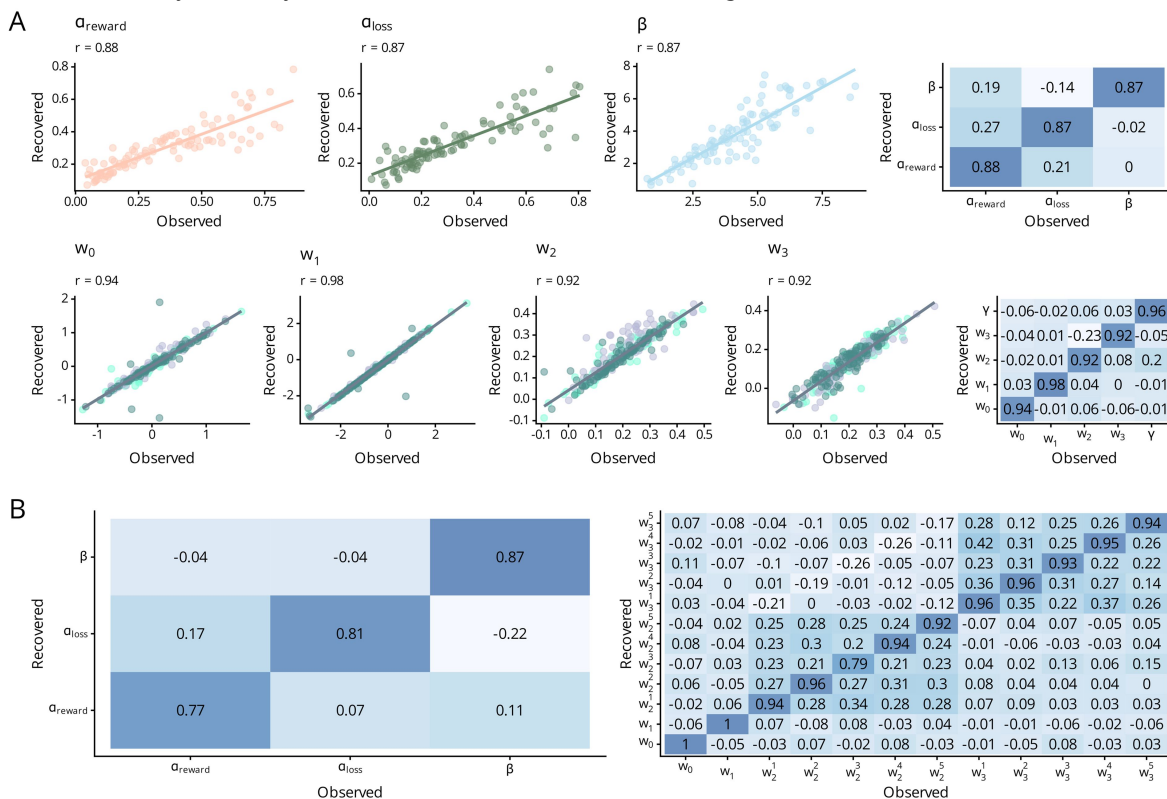
Parameter distributions used to simulate data and test parameter recovery.

<sup>†</sup>i.e., so  $\beta \in [0,10]$

Parameter	Distribution
$\alpha_{reward}, \alpha_{loss}$	Beta(1.5, 3)
$0.1\beta^{\dagger}$	Beta(3, 4)
$w_0^p$	Normal(0, 0.5)
$w_1^p$	Normal(-0.5, 1)
$w_2^p, w_3^p$	Normal(0.2, 0.1)
$\gamma_p$	Beta(2, 2)

Figure S1

Parameter recovery for A) the joint RL-affect model, and B) the between-rating RL-affect model.



### Between-rating RL-affect model

Parameter recovery for the between-rating RL-affect model was tested similarly, with one hundred simulated datasets. The parameter settings were identical to the above except for the time-dependent parameters (and the absence of  $\gamma_p$  in the model). Specifically,  $w_{t-t'}^a$  and  $w_{t-t'}^r$  were sampled from Beta(1,  $j$ ) distributions, where  $t'$  is the trial lag, so that weights for earlier trials were weighted (on average) lower, as we see in the real dataset. Again, even with the high complexity of the model, we found that all parameters could be recovered with reasonably high accuracy ( $r > 0.77$ ; [Figure S1B](#) [↗](#)).

### Comparison between results from models fit to choices alone (with sampling-based inference) vs. in the joint RL-affect model (with variational inference)

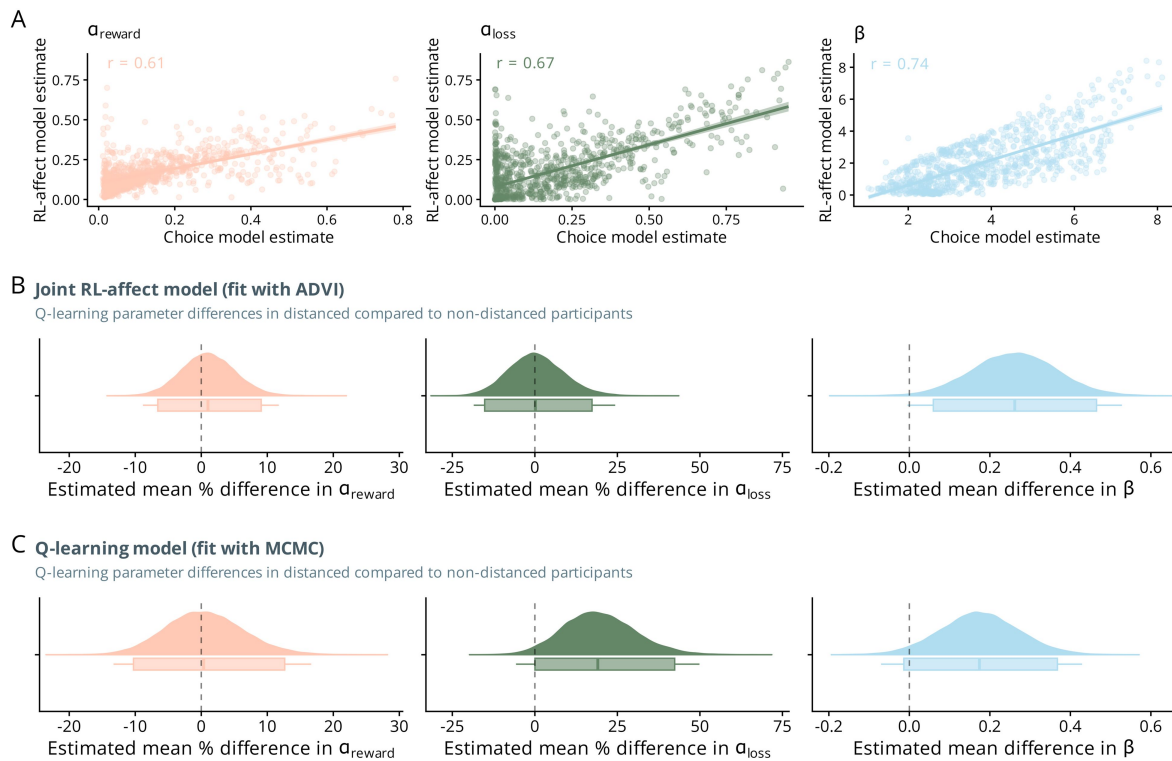
We previously reported results in this dataset where we fit Q-learning models to choices alone and compared parameters in distanced and non-distanced participants [34](#) [↗](#).

Besides the obvious difference—that the models were fit to choices alone as opposed to both choices and affect ratings—the models in our previous work [34](#) [↗](#) were fit to data via sampling-based inference (MCMC), as opposed to variational inference (ADVI [47](#) [↗](#)). However, despite these differences, we find that individuals' posterior mean Q-learning parameters from (i) dual learning rate models fit to choices alone with MCMC, and (ii) the best-fitting joint RL-affect model fit choices and affect ratings simultaneously with ADVI are highly correlated with those we previously reported from Q-learning models fit to choices alone [34](#) [↗](#) in the same sample ( $\alpha_{reward}$ :  $r=0.61$  [95% CI = (0.57, 0.65)];  $\alpha_{loss}$ :  $r=0.67$  [95% CI: 0.63, 0.70];  $\beta$ :  $r = 0.74$  [95% CI = (0.70, 0.76)]; [Figure S2A](#) [↗](#)). We also replicate a key result from our earlier work: higher inverse temperatures ( $\beta$ ) in the distancing group ([Figure S2B-C](#) [↗](#)).

## Supplementary Results

### Group-level parameter estimates for the best-fitting joint RL-affect model

At the group-level, model-predicted baseline affect ( $w_0^a$ ) was highest for engagement (group-level mean = 65.9 points; 95% HDI = [65.9, 66.0]), followed by happiness (group-level mean = 56.5, 95% HDI = [56.4, 56.6]), and lowest for confidence (group-level mean = 32.2, 95% HDI = [32.1, 32.2]). Affect drift over time ( $w_0^r$ ) was steepest for engagement: an estimated 81.2% lower odds of increased engagement per hour (all other terms being equal; 95% HDI for multiplier = [0.187, 0.189]), compared to 54.1% lower confidence (95% HDI for multiplier = [0.457, 0.462]) and 37.7% lower happiness (95% HDI for multiplier = [0.620, 0.625]) per hour. Group-level means for weights on recent expected values ( $w_1^a$ ) and prediction errors ( $w_1^r$ ) were positive for all three ratings, showing higher affect with increased recent rewards, and were lowest for engagement (posterior mean [95% HDI] for multiplier:  $w_{engagement}^a = 1.200$  [1.198, 1.201],  $w_{engagement}^r = 1.216$  [1.216, 1.218],  $w_{engagement}^e = 1.116$  [1.114, 1.116];  $w_{happiness}^a = 1.136$  [1.134, 1.136],  $w_{happiness}^r = 1.208$  [1.207, 1.210],  $w_{happiness}^e = 1.109$  [1.108, 1.109]). Lastly, the average decay factor was highest for confidence (group-level mean = 0.563; 95% HDI = [0.561, 0.565]), and lowest for happiness (group-level mean=0.407; 95% HDI = [0.405, 0.409]), suggesting higher weighting of outcomes of earlier trials (i.e., prior to the most recent trial) in confidence judgements compared to happiness ratings.



**Figure S2**

Comparison of Q-learning parameters and effects of distancing between dual learning rate models fit to choices alone and the joint RL-affect model additionally fit to affect ratings.

## Compulsive behaviour and social withdrawal factor scores are also associated with altered affective dynamics

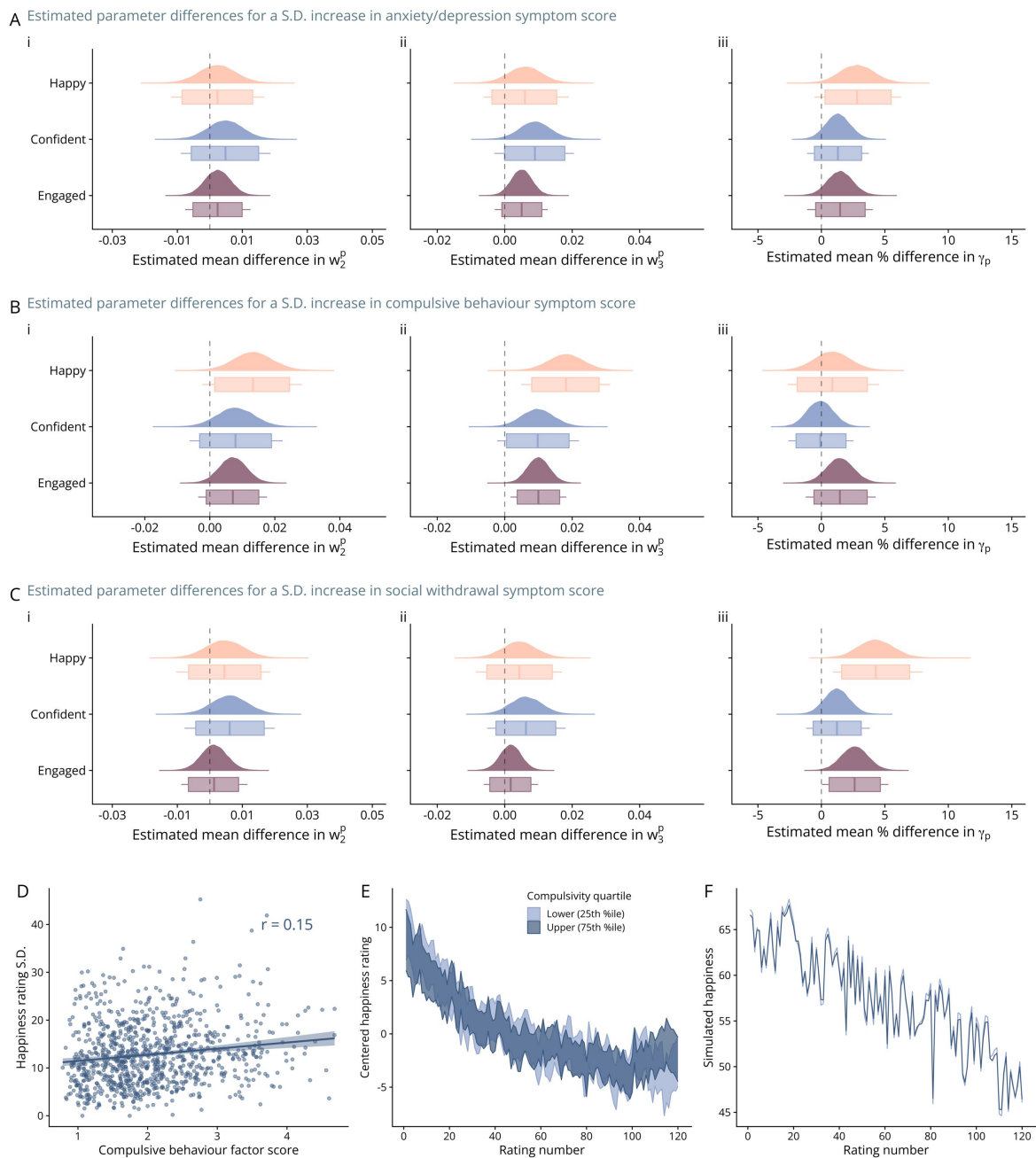
In addition to differences in baseline affect and its drift over time detailed in the main text, there was evidence that participants with higher compulsive behaviour scores placed more weight on recent expected values (higher  $w_{t-1}^{\text{happy}}$ ; 95% HDI excluding zero for  $w_{t-1}^{\text{happy}}$  only) and prediction errors (higher  $w_{t-1}^{\text{err}}$ ) in their subjective affect judgements (**Figure S3** [i–iii](#)). The consequences of this were evident in the observed data. For example, there was evidence of a weak positive correlation between happiness rating variability and compulsive behaviour factor scores (correlation coefficient [95% CI]  $r = 0.15$  [0.086, 0.212],  $p < 0.0001$ ; **Figure S3D** [i](#)), which can also be qualitatively observed by comparing mean-centred happiness ratings for participants in the bottom quartile versus the upper quartile of compulsive behaviour factor scores (**Figure S3E** [i](#)). That said, the actual effect of the estimated differences in  $w_{t-1}^{\text{happy}}$  and  $w_{t-1}^{\text{err}}$  on ratings is small, as shown by simulating happiness ratings for individuals who differ only in having the estimated  $w_{t-1}^{\text{happy}}$  and  $w_{t-1}^{\text{err}}$  for those with the 25<sup>th</sup> percentile versus the 75<sup>th</sup> percentile compulsive behaviour factor score (**Figure S3E** [ii](#)).

We additionally found some weak evidence that increases in anxiety/depression factor were associated with slightly higher weighting of previous trials' expected values and prediction errors for happiness ( $\gamma_{\text{happy}}$ ; e.g., trial-before-last weighted an estimated 4.04% higher; 95% HDI for multiplier = [1.003, 1.079]; **Figure S3A** [i](#)). There was also some stronger evidence for a positive association between social withdrawal factor score and decay factors for happiness and engagement (**Figure S3C** [i](#)), suggesting marginally higher weighting of previous trials' expected values and prediction errors in the computation of affect ratings in those with higher levels social withdrawal symptoms.

## Antidepressant use is associated with increased weighting of previous choices' expected values in affect ratings

To further unpick the effects of treatments on the weighting of previous outcomes, we fit a more flexible between-rating RL-affect model (*equation 7* [i](#)). This model, which allowed for different weights on previous expected values ( $w_{t-1}^{\text{happy}}$ ) and prediction errors ( $w_{t-1}^{\text{err}}$ ) since the previous rating (up to five trials back), was able to capture the ratings well, albeit with marginally worse accuracy than the winning drift over time model (mean [SD] pseudo- $R^2$  for between-rating RL-affect model = 0.39 [0.22–0.25] across all three ratings). We then related  $w_{t-1}^{\text{happy}}$  and  $w_{t-1}^{\text{err}}$  from each rating type separately to both treatments via multilevel GLMs with participant-level random intercepts and slopes (on trial lag), adjusting for age, gender, and digit span as before.

Parameters from this between-rating model suggested limited evidence for a difference between distancing and non-distancing participants in the weighting of the most recent or intervening outcomes in their affective judgements (**Figure S4A** [i](#)). There was also no evidence of an effect of either treatment on weightings of prediction errors from previous trials (**Figure S4Aiii–iv** & **Figure S4Biii–iv** [i](#)). There was, however, some evidence of a small effect of antidepressant use on between-rating changes in affect: higher weighting of the most recent expected value in subjective affect ratings (**Figure S4B** [i](#)). The evidence for this was strongest for engagement, with a unit increase in the most recent  $Q$ -value associated with 4.33% higher odds of an increase in engagement rating (95% HDI for multiplier = [1.001, 1.089]). Furthermore, there was limited evidence of an accompanying interaction effect (**Figure S4B** [ii](#)), suggesting the contribution of less recent expected values to engagement ratings was also marginally higher in participants taking antidepressants, which may in turn explain the higher forgetting factor  $\gamma_{\text{engaged}}$  (**Figure 3B** [v](#)).

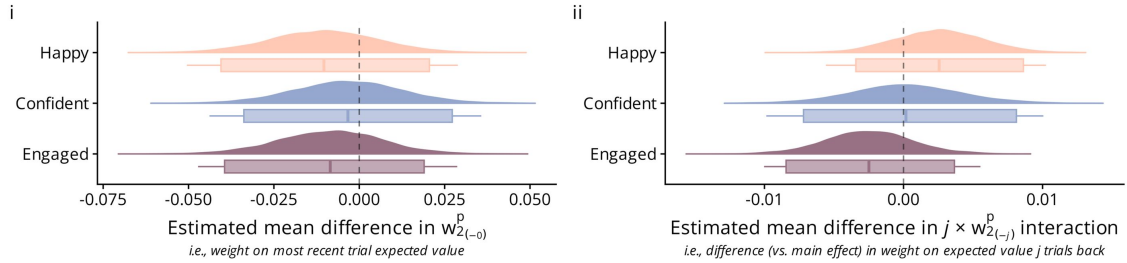


**Figure S3**

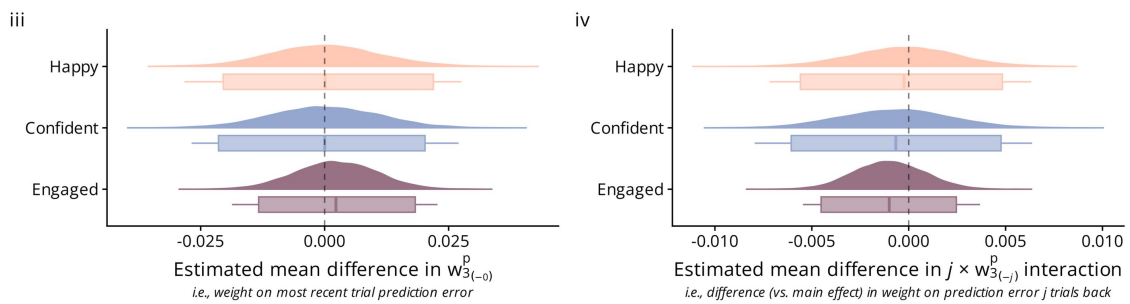
Associations between higher transdiagnostic psychiatric symptom factor scores and additional affect parameters (**A-C**), correlation between variance in happiness rating and compulsive behaviour score (**D**), and the simulated effect on happiness ratings of higher  $w_2^{happy}$  and  $w_3^{happy}$  (**E-F**).

## A Cognitive distancing

Differences in weights on expected values: most recent trial and change with increasing lag

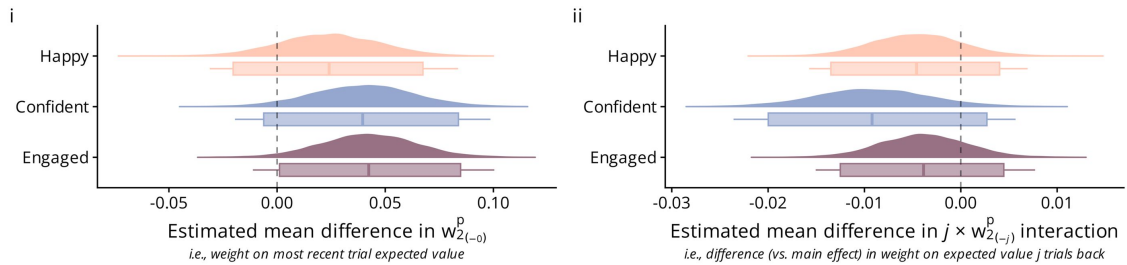


Differences in weights on prediction errors: most recent trial and change with increasing lag

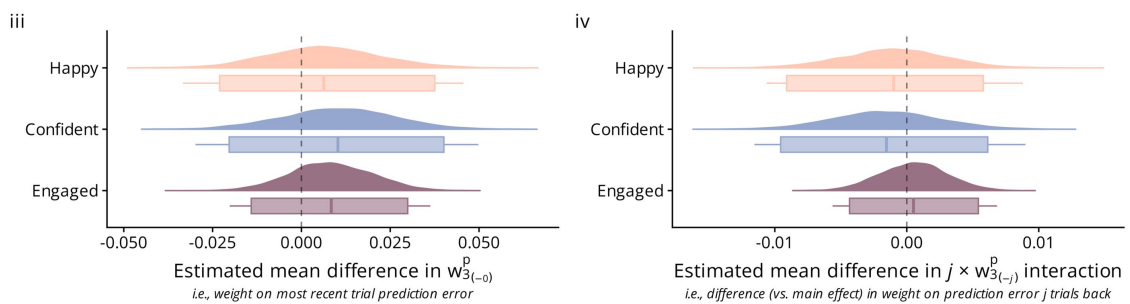


## B Antidepressant use

Differences in weights on expected values: most recent trial and change with increasing lag



Differences in weights on prediction errors: most recent trial and change with increasing lag



**Figure S4**

Effects of cognitive distancing (A) and antidepressant use (B) on expected value and prediction error parameters, derived from the between-rating RL-affect model.

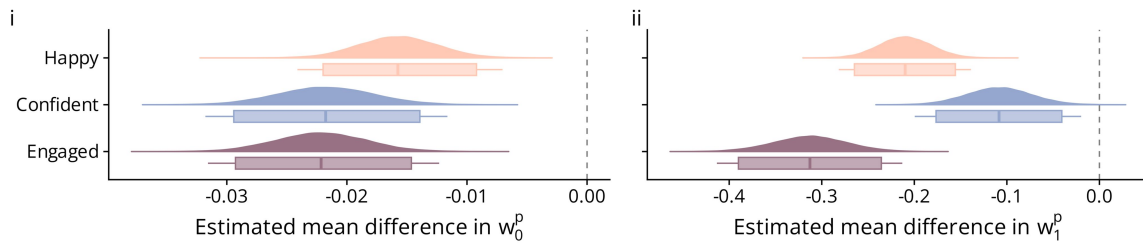
## Affective drift over time is associated with self-reported fatigue

Previous work on ‘mood drift over time’ has suggested it is mostly distinct from boredom and mind wandering<sup>44</sup>. Here, we were able to test an additional aspect of this phenomenon, namely its relation to fatigue, as we asked participants the following question after the end of each of the six blocks: “How fatigued do you feel compared to the beginning of the block?”. We hence examined the association between  $\omega_{\text{fatigue}}$  and both participants’ overall mean post-block fatigue and their change in post-block fatigue ratings over the course of the six training blocks. To quantify change in post-block fatigue, the six ratings were regressed on block number for each participant, with the regression coefficient ( $\beta_{\Delta\text{fatigue}}$ ) on block number taken as the quantity of interest (i.e., higher values suggest increases in post-block fatigue ratings over time).

We found strong evidence, after adjusting for age, gender, digit span, and distancing group, that higher mean post-block fatigue ratings were associated with lower baseline affect (lower  $\omega_{\text{affect}}$ ; **Figure S5A** [i](#)) and decreased odds of higher affect ratings across the task (lower  $\omega_{\text{affect}}$ ; **Figure S5A** [ii](#)), across all three affect rating types (e.g., estimated mean 18.9% lower odds of increased happiness across the task with a ten-point increase in mean post-block fatigue; 95% HDI for multiplier = [0.767, 0.856]). In addition, regression coefficients capturing change in post-block fatigue were strongly negatively associated with  $\omega_{\text{affect}}$  across all rating types after adjusting for mean post-block fatigue, most strongly for engagement (estimated mean 5.74% lower odds of increase in engagement for a one-point per block in fatigue rating rate-of-change, 95% HDI for multiplier = [0.944, 0.968]; **Figure S5B** [i](#)). Notably, associations in the opposite (positive) direction were observed between  $\beta_{\Delta\text{fatigue}}$  and baseline affect, again most strongly for engagement (estimated 0.714-point increase in engagement rating for a one-point per block increase in fatigue rating rate-of-change; 95% HDI = [0.452, 0.975]; **Figure S5B** [i](#)). Speculatively, this may represent an effect of motivation, where participants who were more engaged towards the beginning of the task also exerted more effort, resulting in larger overall increases in fatigue and decreases in engagement.

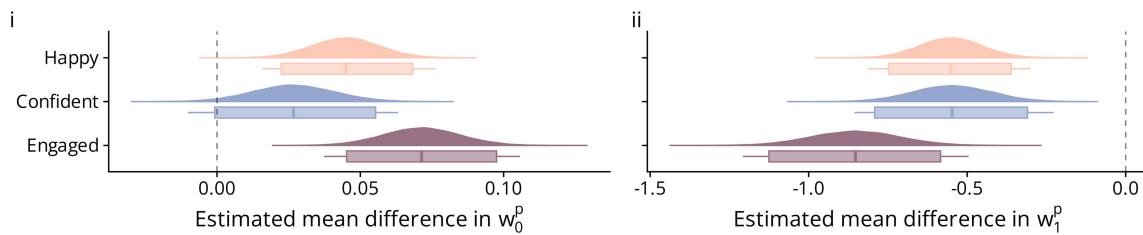
### A Mean post-block fatigue rating

Differences correspond to a 10-point increase in fatigue



### B Change in post-block fatigue rating

Differences correspond to a unit increase in  $\beta_{\Delta \text{ fatigue}}$



**Figure S5**

Associations between baseline affect and affective drift, and self-reported fatigue.

## Data and code availability

All code to replicate the analyses here can be found in accompanying Jupyter notebooks, alongside cleaned, anonymised data. See the GitHub repository for more details.

## Acknowledgements

This study was funded by an AXA Research Fund Fellowship awarded to C.L.N. (G102329) and the Medical Research Council (MC\_UU\_00030/12). C.L.N. is funded by a Wellcome Career Development Award (226490/Z/22/Z) and acknowledges support by the NIHR Cambridge NIHR Biomedical Research Centre (BRC-1215-20014). Q.D. is funded by a Wellcome Trust PhD studentship. Q.D. and Q.J.M.H acknowledge support by the NIHR UCLH BRC. Q.J.M.H. acknowledges grant funding from the NIHR, Wellcome Trust, Carigest S.A. and Koa Health. R.B.R. is supported by the National Institute of Mental Health (R01MH124110). R.B.R. holds equity in Maia.

## Additional information

### Rights retention

For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from this submission.

### Acronyms

- ADVI: Automatic Differentiation Variational Inference
- BCa: Bias-Corrected and accelerated
- CI: Confidence Interval
- ELPD: Expected Log Pointwise Predictive Density
- GLM: Generalised Linear Model
- HDI: Highest Density Interval
- HBREC: Human Biology Research Ethics Committee
- LOO: Leave One Out
- MSE: Mean Squared Error
- MCMC: Markov Chain Monte Carlo
- NHS: National Health Service
- PLS: Partial Least Squares
- RL: Reinforcement Learning
- SD: Standard Deviation
- SSRI: Selective Serotonin Reuptake Inhibitor

## References

1. Oswald A. J., Wu S. (2010) **Objective Confirmation of Subjective Measures of Human Well-Being: Evidence from the U.S.A** *Science* **327**:576–79 <https://doi.org/10.1126/science.1180606> | Google Scholar
2. Diener E., Oishi S., Tay L. (2018) **Advances in Subjective Well-Being Research** *Nature Human Behaviour* **2**:253–60 <https://doi.org/10.1038/s41562-018-0307-6> | Google Scholar
3. Steptoe A., Deaton A., Stone A. A. (2015) **Subjective Wellbeing, Health, and Ageing** *The Lancet* **385**:640–48 [https://doi.org/10.1016/S0140-6736\(13\)61489-0](https://doi.org/10.1016/S0140-6736(13)61489-0) | Google Scholar
4. Suh E., Diener E., Fujita F. (1996) **Events and Subjective Well-Being: Only Recent Events Matter** *Journal of Personality and Social Psychology* **70**:1091–1102 <https://doi.org/10.1037/0022-3514.70.5.1091> | Google Scholar
5. Rutledge R. B., Skandali N., Dayan P., Dolan R. J. (2014) **A Computational and Neural Model of Momentary Subjective Well-Being** *Proceedings of the National Academy of Sciences of the United States of America* **111**:12252–57 <https://doi.org/10.1073/pnas.1407535111> | Google Scholar
6. Taquet M., Quoidbach J., Montjoye Y.-A., Desseilles M., Gross J. J. (2016) **Hedonism and the Choice of Everyday Activities** *Proceedings of the National Academy of Sciences of the United States of America* **113**:9769–73 <https://doi.org/10.1073/pnas.1519998113> | Google Scholar
7. Blain B., Rutledge R. B. (2020) **Momentary Subjective Well-Being Depends on Learning and Not Reward** *eLife* **9**:e57977 <https://doi.org/10.7554/eLife.57977> | Google Scholar
8. Pouget A., Drugowitsch J., Kepecs A. (2016) **Confidence and Certainty: Distinct Probabilistic Quantities for Different Goals** *Nature Neuroscience* **19**:366–74 <https://doi.org/10.1038/nn.4240> | Google Scholar
9. Adler W. T., Ma W. J. (2018) **Comparing Bayesian and Non-Bayesian Accounts of Human Confidence Reports** *PLOS Computational Biology* **14**:1006572 <https://doi.org/10.1371/journal.pcbi.1006572> | Google Scholar
10. Navajas J., et al. (2017) **The Idiosyncratic Nature of Confidence** *Nature Human Behaviour* **1**:810–18 <https://doi.org/10.1038/s41562-017-0215-1> | Google Scholar
11. Li H.-H., Ma W. J. (2004) **Confidence Reports in Decision-Making with Multiple Alternatives Violate the Bayesian Confidence Hypothesis** *Nature Communications* **11** <https://doi.org/10.1038/s41467-020-15581-6> | Google Scholar
12. Walton M. E., Kennerley S. W., Bannerman D. M., Phillips P., Rushworth M. F. S. (2006) **Weighing up the Benefits of Work: Behavioral and Neural Analyses of Effort-Related Decision Making** *Neural Networks* **19**:1302–14 <https://doi.org/10.1016/j.neunet.2006.03.005> | Google Scholar
13. Ang Y.-S., Gelda S. E., Pizzagalli D. A. (2022) **Cognitive Effort-Based Decision-Making in Major Depressive Disorder** *Psychological Medicine* :1–8 <https://doi.org/10.1017/S0033291722000964>

Google Scholar

14. Rutledge R. B., et al. (2017) **Association of Neural and Emotional Impacts of Reward Prediction Errors With Major Depression** *JAMA Psychiatry* **74**:790–97 <https://doi.org/10.1001/jamapsychiatry.2017.1713> | Google Scholar
15. Barge-Schaapveld D. Q., Nicolson N. A., Berkhof J., deVries M. (1999) **Quality of Life in Depression: Daily Life Determinants and Variability** *Psychiatry Research* **88**:173–89 [https://doi.org/10.1016/s0165-1781\(99\)00081-5](https://doi.org/10.1016/s0165-1781(99)00081-5) | Google Scholar
16. Taquet M., Quoidbach J., Gross J. J., Saunders K. E., Goodwin G. M. (2020) **Mood Homeostasis, Low Mood, and History of Depression in 2 Large Population Samples** *JAMA Psychiatry* **77**:944–51 <https://doi.org/10.1001/jamapsychiatry.2020.0588> | Google Scholar
17. Rouault M., Seow T., Gillan C. M., Fleming S. M. (2018) **Psychiatric Symptom Dimensions Are Associated With Dissociable Shifts in Metacognition but Not Task Performance** *Biological Psychiatry* **84**:443–51 <https://doi.org/10.1016/j.biopsych.2017.12.017> | Google Scholar
18. Hoven M., et al. (2019) **Abnormalities of Confidence in Psychiatry: An Overview and Future Perspectives** *Translational Psychiatry* **9**:268 <https://doi.org/10.1038/s41398-019-0602-7> | Google Scholar
19. Hoven M., Denys D., Rouault M., Luigjes J., Holst R. (2022) **How do confidence and self-beliefs relate in psychopathology: a transdiagnostic approach** *Nature Mental Health* **1**:337–345 <https://doi.org/10.1038/s44220-023-00062-8> | Google Scholar
20. Katyal S., Huys Q. J., Dolan R. J., Fleming S. M. (2025) **Distorted learning from local metacognition supports transdiagnostic underconfidence** *Nature Communications* **16**:1854 <https://doi.org/10.1038/s41467-025-57040-0> | Google Scholar
21. Clery-Melin M.-L., et al. (2011) **Why Don't You Try Harder? An Investigation of Effort Production in Major Depression** *PLoS One* **6**:23178 <https://doi.org/10.1371/journal.pone.0023178> | Google Scholar
22. Pessiglione M., Vinckier F., Bouret S., Daunizeau J., Bouc R. L. (2018) **Why Not Try Harder? Computational Approach to Motivation Deficits in Neuro-Psychiatric Diseases** *Brain* **141**:629–50 <https://doi.org/10.1093/brain/awx278> | Google Scholar
23. Husain M., Roiser J. P. (2018) **Neuroscience of Apathy and Anhedonia: A Transdiagnostic Approach** *Nature Reviews Neuroscience* **19**:470–84 <https://doi.org/10.1038/s41583-018-0029-9> | Google Scholar
24. Disner S. G., Beevers C. G., Haigh E. A. P., Beck A. T. (2011) **Neural mechanisms of the cognitive model of depression** *Nature Reviews Neuroscience* **12**:467–477 <https://doi.org/10.1038/nrn3027> | Google Scholar
25. Lewis G., et al. (2017) **Variation in the recall of socially rewarding information and depressive symptom severity: a prospective cohort study** *Acta Psychiatrica Scandinavica* **135**:489–498 <https://doi.org/10.1111/acps.12729> | Google Scholar

26. Bylsma L. M., Taylor-Clift A., Rottenberg J. (2011) **Emotional Reactivity to Daily Events in Major and Minor Depression** *Journal of Abnormal Psychology* **120**:155–67 <https://doi.org/10.1037/a0021662> | [Google Scholar](#)
27. Eldar E., Rutledge R. B., Dolan R. J., Niv Y. (2016) **Mood as Representation of Momentum** *Trends in Cognitive Sciences* **20**:15–24 <https://doi.org/10.1016/j.tics.2015.07.010> | [Google Scholar](#)
28. Harmer C. J. (2008) **Serotonin and Emotional Processing: Does It Help Explain Antidepressant Drug Action?** *Neuropharmacology* **55**:1023–28 <https://doi.org/10.1016/j.neuropharm.2008.06.036> | [Google Scholar](#)
29. Harmer C. J., Duman R. S., Cowen P. J. (2017) **How Do Antidepressants Work? New Perspectives for Refining Future Treatment Approaches** *The Lancet Psychiatry* **4**:409–18 [https://doi.org/10.1016/S2215-0366\(17\)30015-9](https://doi.org/10.1016/S2215-0366(17)30015-9) | [Google Scholar](#)
30. Beck A. T. (1964) **Thinking and depression: II. Theory and therapy** *Archives of General Psychiatry* **10**:561–571 <https://doi.org/10.1001/archpsyc.1964.01720240015003> | [Google Scholar](#)
31. Kross E., Gard D., Deldin P., Clifton J., Ayduk O. (2012) **“Asking Why” from a Distance: Its Cognitive and Emotional Consequences for People with Major Depressive Disorder** *Journal of Abnormal Psychology* **121**:559–69 <https://doi.org/10.1037/a0028808> | [Google Scholar](#)
32. Frank M. J., Seeberger L. C., O’Reilly R. C. (2004) **By Carrot or by Stick: Cognitive Reinforcement Learning in Parkinsonism** *Science* **306**:1940–43 <https://doi.org/10.1126/SCIENCE.1102941> | [Google Scholar](#)
33. Frank M. J., Moustafa A. A., Haughey H. M., Curran T., Hutchison K. E. (2007) **Genetic Triple Dissociation Reveals Multiple Roles for Dopamine in Reinforcement Learning** *Proceedings of the National Academy of Sciences of the United States of America* **104**:16311–16 <https://doi.org/10.1073/pnas.0706111104> | [Google Scholar](#)
34. Dercon Q., Mehrhof S. Z., et al. (2023) **A Core Component of Psychological Therapy Causes Adaptive Changes in Computational Learning Mechanisms** *Psychological Medicine* :1–11 <https://doi.org/10.1017/S0033291723001587> | [Google Scholar](#)
35. Powers J. P., LaBar K. S. (2019) **Regulating emotion through distancing: A taxonomy, neurocognitive model, and supporting meta-analysis** *Neuroscience & Biobehavioral Reviews* **96**:155–173 <https://doi.org/10.1016/j.neubiorev.2018.04.023> | [Google Scholar](#)
36. Gillan C. M., Kosinski M., Whelan R., Phelps E. A., Daw N. D. (2016) **Characterizing a Psychiatric Symptom Dimension Related to Deficits in Goal-Directed Control** *eLife* **5**:e11305 <https://doi.org/10.7554/eLife.11305> | [Google Scholar](#)
37. Wise T., Robinson O., Gillan C. (2022) **Identifying Transdiagnostic Mechanisms in Mental Health Using Computational Factor Modeling** *Biological Psychiatry* **93**:690–703 <https://doi.org/10.1016/j.biopsych.2022.09.034> | [Google Scholar](#)
38. Palan S., Schitter C. (2018) **Prolific.Ac—A Subject Pool for Online Experiments** *Journal of Behavioral and Experimental Finance* **17**:22–27 <https://doi.org/10.1016/J.JBEF.2017.12.004> | [Google Scholar](#)

39. Wise T., Dolan R. J. (2020) **Associations between Aversive Learning Processes and Transdiagnostic Psychiatric Symptoms in a General Population Sample** *Nature Communications* **11**:4179 <https://doi.org/10.1038/s41467-020-17977-w> | Google Scholar
40. Kao C.-H., Feng G. W., Hur J. K., Jarvis H., Rutledge R. B. (2023) **Computational Models of Subjective Feelings in Psychiatry** *Neuroscience & Biobehavioral Reviews* **145**:105008 <https://doi.org/10.1016/j.neubiorev.2022.105008> | Google Scholar
41. Forbes L., Bennett D. (2024) **The effect of reward prediction errors on subjective affect depends on outcome valence and decision context** *Emotion* **24**:894–911 <https://doi.org/10.1037/emo0001310> | Google Scholar
42. Ferrari S., Cribari-Neto F. (2004) **Beta Regression for Modelling Rates and Proportions** *Journal of Applied Statistics* **31**:799–815 <https://doi.org/10.1080/0266476042000214501> | Google Scholar
43. Smithson M., Verkuilen J. (2006) **A Better Lemon Squeezer? Maximum-Likelihood Regression with Beta-Distributed Dependent Variables** *Psychological Methods* **11**:54–71 <https://doi.org/10.1037/1082-989X.11.1.54> | Google Scholar
44. Jangraw D. C., et al. (2023) **A Highly Replicable Decline in Mood during Rest and Simple Tasks** *Nature Human Behaviour* **7**:596–610 <https://doi.org/10.1038/s41562-023-01519-7> | Google Scholar
45. Kucukelbir A., Tran D., Ranganath R., Gelman A., Blei D. M. (2016) **Automatic Differentiation Variational Inference** *arXiv* <https://doi.org/10.48550/arXiv.1603.00788> | Google Scholar
46. Stan Development Team (2022) **Stan Modelling Language Users Guide and Reference Manual** [https://mc-stan.org/docs/2\\_31/cmdstan-guide-2\\_31.pdf](https://mc-stan.org/docs/2_31/cmdstan-guide-2_31.pdf)
47. Vehtari A., Gelman A., Gabry J. (2016) **Practical Bayesian Model Evaluation Using Leave-One-out Cross-Validation and WAIC** *Statistics and Computing* **27**:1413–32 <https://doi.org/10.1007/S11222-016-9696-4> | Google Scholar
48. Magnusson M., Andersen M. R., Jonasson J., Vehtari A. (2020) **Leave-One-Out Cross-Validation for Bayesian Model Comparison in Large Data** *arXiv* <https://doi.org/10.48550/arXiv.2001.00980> | Google Scholar
49. Greenland S., et al. (2016) **Statistical Tests, P Values, Confidence Intervals, and Power: A Guide to Misinterpretations** *European Journal of Epidemiology* **31**:337–50 <https://doi.org/10.1007/s10654-016-0149-3> | Google Scholar
50. Goodrich B., Gabry J., Ali I., Brilleman S. (2020) **rstanarm: Bayesian applied regression modeling via Stan R package** <https://mc-stan.org/rstanarm/>
51. Wold S., Ruhe A., Wold H., Dunn W. J. (1984) **The Collinearity Problem in Linear Regression. The Partial Least Squares (PLS) Approach to Generalized Inverses** *SIAM Journal on Scientific and Statistical Computing* **5**:735–43 <https://doi.org/10.1137/0905052> | Google Scholar
52. Harmer C. J., Goodwin G. M., Cowen P. J. (2009) **Why do antidepressants take so long to work? A cognitive neuropsychological model of antidepressant drug action** *British Journal of Psychiatry* **195**:102–108 <https://doi.org/10.1192/bjp.bp.108.051193> | Google Scholar

53. Burkner P.-C. (2017) **brms: An R Package for Bayesian Multilevel Models Using Stan** *Journal of Statistical Software* **80**:1–28 <https://doi.org/10.18637/jss.v080.i01> | [Google Scholar](#)

## Author information

### Quentin Dercon

Applied Computational Psychiatry Lab, Mental Health Neuroscience Department, Division of Psychiatry and Max Planck UCL Centre for Computational Psychiatry and Ageing Research, Department of Imaging Neuroscience, Queen Square Institute of Neurology, UCL, London, United Kingdom

**For correspondence:** [quentin.dercon.22@ucl.ac.uk](mailto:quentin.dercon.22@ucl.ac.uk)

### Quentin JM Huys

Applied Computational Psychiatry Lab, Mental Health Neuroscience Department, Division of Psychiatry and Max Planck UCL Centre for Computational Psychiatry and Ageing Research, Department of Imaging Neuroscience, Queen Square Institute of Neurology, UCL, London, United Kingdom

### Robb B Rutledge

Department of Psychology, Yale University, New Haven, United States, Wu Tsai Institute, Yale University, New Haven, United States, Department of Psychiatry, Yale University, New Haven, United States

### Camilla L Nord

MRC Cognition and Brain Sciences Unit, University of Cambridge, Cambridge, United Kingdom, Department of Psychiatry, University of Cambridge, Cambridge, United Kingdom

## Editors

Reviewing Editor

### Mimi Liljeholm

University of California, Irvine, Irvine, United States of America

Senior Editor

### Christian Büchel

University Medical Center Hamburg-Eppendorf, Hamburg, Germany

### Reviewer #1 (Public review):

Summary:

This study examines how two common psychiatric treatments, antidepressant medication and cognitive distancing, influence baseline levels and moment-to-moment changes in happiness, confidence, and engagement during a reinforcement learning task. Combining a probabilistic selection task, trial-by-trial affect ratings, psychiatric questionnaires, and computational modeling, the authors demonstrate that each treatment has distinct effects on affective dynamics. Notably, the results highlight the key role of affective biases in how people with mental health conditions experience and update their feelings over time, and

suggest that interventions like cognitive distancing and antidepressant medication may work, at least in part, by shifting these biases.

#### Strengths:

- (1) Addresses an important question: how common psychiatric treatments impact affective biases, with potential translational relevance for understanding and improving mental health interventions.
- (2) The introduction is strong, clear, and accessible, making the study approachable for readers less familiar with the underlying literature.
- (3) Utilizes a large sample that is broadly representative of the UK population in terms of age and psychiatric symptom history, enhancing generalizability.
- (4) Employs a theory-driven computational modeling framework that links learning processes with subjective emotional experiences.
- (5) Uses cross-validation to support the robustness and generalizability of model comparisons and findings.

#### Weaknesses:

The authors acknowledge the limitations in the discussion section.

#### Additional questions:

- (1) Group Balance & Screening for Medication Use: How many participants in the cognitive distancing and control groups were taking antidepressant medication? Why wasn't medication use included as part of the screening to ensure both groups had a similar number of participants taking medication?
- (2) Assessment of the Practice of Cognitive Distancing: Is there a direct or more objective method to evaluate whether participants actively engaged in cognitive distancing during the task, and to what extent? Currently, the study infers engagement indirectly through the outcomes, but does not include explicit measures of participants' use of the technique. Would including self-report check-ins throughout the task, asking participants whether they were actively engaging in cognitive distancing, have been useful? However, including frequent self-report check-ins would increase procedural differences between groups, making perhaps the tasks less comparable beyond the intended treatment manipulation. Maybe incorporating a question at the end of the task, asking how much they engaged in cognitive distancing, could offer a useful measure of subjective engagement without overly disrupting the task flow.

#### Conclusion:

This study advances our understanding of the mechanisms underlying mental health interventions. The combination of computational modeling with behavioral and affective data offers a powerful framework for understanding how treatments influence affective biases and dynamics. These findings are of broad interest across clinical and mental health sciences, cognitive and affective research, and applied translational fields focused on improving psychological well-being.

<https://doi.org/10.7554/eLife.107269.1.sa2>

**Reviewer #2 (Public review):**

In this paper, Dercon and colleagues report on affective changes related to components of reinforcement learning and on the effects of brief training in psychological distancing and participants' self-reported antidepressant use. About 1,000 participants were assessed online, with half randomized to a brief training in psychological distancing with reminders to distance during the subsequent reinforcement learning (RL) task. Participants completed a battery of psychiatric questionnaires and answered questions about medication use, with about 14% of participants reporting current antidepressant use. All participants completed the RL task and rated their happiness, confidence, engagement, and (at the end of each block of trials) fatigue throughout the task. Computational models were used to estimate trial-by-trial values of expected value and prediction error and to assess the effects of these values on self-reported affect. Participants' affect ratings decreased over time, and participants with higher psychiatric symptoms (particularly anxiety/depressive symptoms) showed lower baseline affect and greater decreases in affect. Participants randomized to the distancing intervention and who reported antidepressant use differed in their affective ratings: distancing reduced the reductions in happiness over time, while antidepressant use was related to higher baseline happiness. Distancing also reduced the effects of trial-level expected value on happiness, while antidepressant use was related to a more enduring effect of trial-level values on happiness.

Overall, this is an interesting paper with strong methods and an interesting approach. That psychiatric symptoms and cognitive distancing are related to affective ratings is not terribly novel; the relationship with antidepressant use is a bit more novel. The extension of the mood model to an RL task is a new contribution, as is the relationship of these effects with psychologically related manipulations.

One major concern is the inference that can be drawn from the two "treatments": one is a brief instruction in a component of psychotherapy, and one is ongoing use of medication. The former is not a treatment in and of itself, but a (presumably) active ingredient of one. How to interpret antidepressant use as measured is unclear, e.g., are the residual symptoms in these participants an early indicator of treatment resistance? Are these participants with better access to health care? Are they receiving antidepressants for a mental health issue?

There are some clarifications needed in the affect model as well.

<https://doi.org/10.7554/eLife.107269.1.sa1>

**Reviewer #3 (Public review):**

Summary:

The present manuscript investigates and proposes different mechanisms for the effects of two therapeutic approaches - cognitive distancing technique and use of antidepressants - on subjective ratings of happiness, confidence, and task engagement, and on the influence of such subjective experiences on choice behavior. Both approaches were found to link to changes in affective state dynamics in a choice task, specifically reduced drift (cognitive distancing) and increased baseline (antidepressant use). Results also suggest that cognitive distancing may reduce the weighing of recent expected values in the happiness model, while antidepressant use may reduce forgetting of choices and outcomes.

Strengths:

This is a timely topic and a significant contribution to ongoing efforts to improve our mechanistic understanding of psychopathology and devise effective novel interventions. The

relevance of the manuscript's central question is clear, and the links to previous literature and the broader field of computational psychiatry are well established. The modelling approaches are thoughtful and rigorously tested, with appropriate model checks and persuasive evidence that modelling complements the theoretical argument and empirical findings.

Weaknesses:

Some vagueness and lack of clarity in theoretical mechanisms and interpretation of results leave outstanding questions regarding (a) the specific links drawn between affective biases, therapies aimed at mitigating them, and mental health function, and (b) the structure and assumptions of the modelling, and how they support the manuscript's central claims. Broadly, I do not fully understand the distinction between how choice behavior vs. affect are impacted separately or together by cognitive distancing. Clarification on this point is needed, possibly through a more explicit proposal of a mechanism (or several alternative mechanisms?) in the introduction and more explicit interpretation of the modelling results in the context of the cyclical choice-affect mechanism.

#### (1) Theoretical framework and proposed mechanisms

The link between affective biases and negative thinking patterns is a bit unclear. The authors seem to make a causal claim that "affective biases are precipitated and maintained by negative thinking patterns", but it is unclear what precisely these negative patterns are; earlier in the same paragraph, they state that affective biases "cause low mood" and possibly shift choices toward those that maintain low mood. So the directionality of the mechanism here is unclear - possibly explaining a bit more of the cyclic nature of this mechanism, and maybe clarifying what "negative thinking patterns" refer to will be helpful.

More generally, this link between affect and choices, especially given the modelling results later on, should be clarified further. What is the mechanism by which these two impact each other? How do the models of choice and affect ratings in the RL task test this mechanism? I'm not quite sure the paper answers these questions clearly right now.

The authors also seem to implicitly make the claim that symptoms of mental ill-health are at least in part related to choice behavior. I find this a persuasive claim generally; however, it is understated and undersupported in the introduction, to the point where a reader may need to rely on significant prior knowledge to understand why mitigating the impact of affective biases on choice behavior would make sense as the target of therapeutic interventions. This is a core tenet of the paper, and it would be beneficial to clarify this earlier on.

It would be helpful to interpret a bit more clearly the findings from 3.4. on decreased drift in all three subjective assessments in the cognitive distancing group. What is the proposed mechanism for this? The discussion mentions that "attenuated declines [...] over time, [add] to our previously reported findings that this psychotherapeutic technique alters aspects of reward learning" - but this is vague and I do not understand, if an explanation for how this happens is offered, what that explanation is. Given the strong correlation of the drift with fatigue, is the explanation that cognitive distancing mitigates affect drift under fatigue? Or is this merely reporting the result without an interpretation around potential mechanisms?

(Relatedly, aside from possibly explaining the drift parameter, do the fatigue ratings link with choice behavior in any way? Is it possible that the cognitive distancing was helping participants improve choices under fatigue?)

#### (2) Task Structure and Modelling

It is unclear what counted as a "rewarding" vs. "unrewarding" trial in the model. From my understanding of the task description, participants obtained positive or no reward (no losses),

and verbal feedback, Correct/Incorrect. But given the probabilistic nature of the task, it follows that even some correct choices likely had unrewarding results. Was the verbal feedback still "Correct" in those cases, but with no points shown? I did not see any discussion on whether it is the #points earned or the verbal feedback that is considered a reward in the model. I am assuming the former, but based on previous literature, likely both play a role; so it would be interesting - and possibly necessary to strengthen the paper's argument - to see a model that assigns value to positive/negative feedback and earned points separately.

From a theory perspective, it's interesting that the authors chose to assume separate learning rates for rewarding and non-rewarding trials. Why not, for example, separate reward sensitivity parameters? E.g., rather than a scaling parameter on the PE, a parameter modifying the  $r$  term inside the PE equation to, perhaps, assign different values to positive and zero points? (While I think overall the math works out similarly at the fitting time, this type of model should be less flexible on scaling the expected value and more flexible on scaling the actual #points / the subjective experience of the obtained verbal feedback, which seems more in line with the theoretical argument made in the introduction). The introduction explicitly states that negative biases "may cause low mood by making outcomes appear less rewarding" - which in modelling equations seems more likely to translate to different reward-perception biases, and not different learning rates. Alternatively, one might incorporate a perseveration parameter (e.g., similar to Collins et al. 2014) that would also accomplish a negative bias. Either of these two mechanisms seems perhaps worth testing out in a model - especially in a model that defines more clearly what rewarding vs. unrewarding may mean to the participant.

If I understand correctly, the affect ratings models assume that the Q-value and the PE independently impact rating (so they have different weights,  $w_2$  and  $w_3$ ), but there is no parameter allowing for different impact for perceived rewarding and unrewarding outcomes? (I may be misreading equations 4-5, but if not, Q-value and PE impact the model via static rather than dynamic parameters.) Given the joint RL-affect fit, this seems to carry the assumption that any perceptual processing differences leading to different subjective perceptions of reward associated with each outcome only impact choice behavior, but not affect? (whereas affect is more broadly impacted, if I'm understanding this correctly, just by the magnitude of the values and PEs?) This is an interesting assumption, and the authors seem to have tested it a bit more in the Supplementary material, as shown in Figure S4. I'm wondering why this was excluded from the main text - it seems like the more flexible model found some potentially interesting differences which may be worth including, especially as they might shed additional insight into the influence of cognitive distancing on the cyclical choice-affect mechanisms proposed.

Minor comments:

If fatigue ratings were strongly associated with drift in the best-fitting model (as per page 13), I wonder if it would make sense to use those fatigue ratings as a proxy rather than allow the parameter to vary freely? (This does not in any way detract from the winning model's explanatory power, but if a parameter seems to be strongly explained by a variable we have empirical data for, it's not clear what extra benefit is earned by having that parameter in the model).

<https://doi.org/10.7554/eLife.107269.1.sa0>