# Clinical Response to Neurofeedback in Major Depression Relates to Subtypes of Whole-Brain Activation Patterns During Training

Masaya Misaki Ph.D.[a,b*], Kymberly D. Young Ph.D.[c], Aki Tsuchiyagaito Ph.D.[a,b], Jonathan Savitz Ph.D.[a,b], and Salvador M. Guinjoan M.D. Ph.D.[a,b,d].

[a] Laureate Institute for Brain Research, Tulsa, OK, USA
[b] Oxley College of Health and Natural Sciences, The University of Tulsa, Tulsa, OK, USA
[c] Department of Psychiatry, University of Pittsburgh School of Medicine, Pittsburgh, PA, USA
[d] Department of Psychiatry, Oklahoma University Health Sciences Center at Tulsa, Tulsa, OK, USA

*Corresponding author:
Masaya Misaki

# Abstract

Major Depressive Disorder (MDD) poses a significant public health challenge due to its high prevalence and the substantial burden it places on individuals and healthcare systems. Real-time functional magnetic resonance imaging neurofeedback (rtfMRI-NF) shows promise as a treatment for this disorder, although its mechanisms of action remain unclear. This study investigated whole-brain response patterns during rtfMRI-NF training to explain interindividual variability in clinical efficacy in MDD. We analyzed data from 95 participants (67 active, 28 control) with MDD from previous rtfMRI-NF studies designed to increase left amygdala activation through positive autobiographical memory recall. Significant symptom reduction was observed in the active group ($t$=-4.404, $d$=-0.704, $p$<0.001) but not in the control group ($t$=-1.609, $d$=-0.430, $p$=0.111). However, left amygdala activation did not account for the variability in clinical efficacy. To elucidate the brain training process underlying the clinical effect, we examined whole-brain activation patterns during two critical phases of the neurofeedback procedure: activation during the self-regulation period, and transient responses to feedback signal presentations. Using a systematic process involving feature selection, manifold extraction, and clustering with cross-validation, we identified two subtypes of regulation activation and three subtypes of brain responses to feedback signals. These subtypes were significantly associated with the clinical effect (regulation subtype: $F$=8.735, $p$=0.005; feedback response subtype: $F$=5.326, $p$=0.008; subtypes' interaction: $F$=3.471, $p$=0.039). Subtypes associated with significant symptom reduction were characterized by selective increases in control regions, including lateral prefrontal areas, and decreases in regions associated with self-referential thinking, such as default mode areas. These findings suggest that large-scale brain activity during training is more critical for clinical efficacy than the level of activation in the neurofeedback target region itself. Tailoring neurofeedback training to incorporate these patterns could significantly enhance its therapeutic efficacy.

# Introduction

Major Depressive Disorder (MDD) presents a significant public health challenge, with approximately one-third of diagnosed patients not responding to first-line treatments such as antidepressants and psychotherapy. This results in substantial disability and economic losses due to treatment costs and lost productivity [1, 2]. Real-time functional magnetic resonance imaging neurofeedback (rtfMRI-NF) has emerged as a promising alternative, demonstrating large to medium effect sizes in treating depressive symptoms [3-5]. This noninvasive brain modulation technique involves the real-time analysis and visualization of brain activation signals, thereby enabling participants to self-regulate their brain activity. Its efficacy in training participants to modulate their brain activation is well-supported by many studies, including several meta-analyses [3-11]. However, the direct impact of rtfMRI-NF on symptom relief is not yet fully understood due to incomplete knowledge of the neural mechanisms by which this training alleviates symptoms through the regulation of specific brain activations.

Previous studies on the mechanisms of NF training [12-15], including investigations into brain responses to feedback signals [15-19], have identified a broad spectrum of brain activities associated with NF-mediated self-regulation training. This training process is considered to include aspects of reinforcement learning, and two types of brain activation - evaluation of feedback values and modulation of brain activation - are common components of the reinforcement learning process [20]. Thus, to elucidate the learning mechanisms of NF-mediated brain regulation, investigating these two epochs is crucial. Active regions during these epochs typically include the prefrontal cortex, salience network, and reward processing areas [12-15]. However, while many studies have focused on the success of regulating target brain activities, the relationship between these activities and subsequent symptom relief remains elusive. Furthermore, interindividual variability in the clinical efficacy of NF necessitates further investigation to identify brain response subtypes associated with therapeutic outcomes.

26    This study aimed to characterize whole-brain activation patterns during rtfMRI-NF training in

27    individuals with MDD, with the goal of identifying brain activation subtypes associated with

28    interindividual variability in therapeutic efficacy. To this end, we analyzed a large dataset from

29    rtfMRI-NF studies where participants with MDD were trained to regulate left amygdala activity

30    through neurofeedback [21-24]. These studies consistently observed significant reductions in

31    depressive symptoms on average post-training, albeit with variations in therapeutic outcomes

32    among participants.

33    We hypothesized that the observed variability in treatment efficacy could be explained by

34    variations in whole-brain activation patterns during self-regulation training, extending beyond the

35    NF target region (amygdala). The involvement of large-scale networks in NF training has been

36    demonstrated in a meta-analysis of NF studies [12], and burgeoning evidence suggests that the

37    effects of NF training may extend beyond the targeted brain region [25-27]. Specifically, we focused

38    on two types of brain activation critical for NF training: regulatory task activity throughout the

39    task block and instantaneous responses to neurofeedback signal presentations (Figure 1).

40    These two types of activation are thought to correspond to two critical components of the

41    reinforcement learning process [20] and should characterize the training process for each

42    individual.

## Methods

### Participants

45    The present study was a secondary analysis of data from our previously published studies [21,

46    22, 24] and a preliminary study utilizing the same rtfMRI-NF protocols in individuals diagnosed with

47    MDD. The University of Oklahoma Institutional Review Board (IRB) or the Western IRB

48    reviewed and approved the study protocols, ensuring adherence to the ethical principles of the

49    Declaration of Helsinki. Participants provided written informed consent prior to participation and

50    were financially compensated. Participants met DSM-IV-TR [28] criteria for MDD based on the

51     Structural Clinical Interview for DSM-IV disorders [29] or DSM-5 criteria for MDD based on the

52     Mini-International Neuropsychiatric Interview (MINI) [30]. Previous articles [21, 22, 24] detailed each

53     study's inclusion and exclusion criteria. Common inclusion criteria across studies are ages 18-

54     65, current diagnosis of MDD, and common exclusion criteria are current diagnosis of PTSD,

55     substance use disorder, bipolar disorder, active suicidal ideation or behavior within a year, a

56     history of psychosis, pregnancy, and MRI contraindicators.

57       The present analysis included data from 95 participants with MDD (68 females; mean age ±

58     SD = 33.6 ± 10.4 years), consisting of 67 in the active NF group (46 females) who received left

59     amygdala neurofeedback and 28 in the control group (22 females) who received neurofeedback

60     from a brain region not associated with emotional processing.

61 **Real-time fMRI Neurofeedback Paradigm**

62       The NF training was designed to enhance the activation signal in the left amygdala while

63     recalling happy autobiographical memories [31]. The task sequence utilized a blocked design

64     consisting of alternating periods of rest, self-regulation with NF, and number counting, each

65     lasting 40 seconds. This sequence was repeated four times in a single training run, with

66     participants undergoing three such runs per session. Because the number of sessions differed

67     among the studies, the present analyses were confined to data from the three runs of the initial

68     session. Depressive symptom severity was estimated with the Montgomery-Åsberg Depression

69     Rating Scale (MADRS) [32] immediately before the NF session, and one-week post-training.

70       MRI imaging in all experiments was conducted using the same 3T Discovery MR750 scanner

71     (GE Healthcare). Blood Oxygenation Level-Dependent (BOLD) fMRI data acquisition employed

72     a T2*-weighted gradient-echo planar imaging (EPI) sequence, with parameters set as follows:

73     TR/TE = 2000/30 ms, acquisition matrix = 96 × 96, field of view (FOV)/slice thickness = 240

74     mm/2.9 mm, flip angle = 90°, and 34 axial slices, using a SENSE acceleration factor of 2. EPI

75     images were then reconstructed to a 128x128 matrix, resulting in an fMRI voxel size of

76    1.875x1.875x2.9 mm³. Anatomical reference was obtained using a T1-weighted magnetization-

77    prepared rapid gradient-echo (MPRAGE) sequence.

78    A detailed description of the rtfMRI-NF procedure can be found in our previous publications [21,

79    22, 24]. Briefly, a custom in-house rtfMRI system was utilized for the experiments [31]. For the active

80    group, the NF signal was extracted from the left amygdala, defined by 7-mm diameter spheres

81    centered at Talairach coordinates (x, y, z = -21, -5, -16 mm), and then mapped to each

82    participant's brain space. For the control group, the NF signal was sourced from the horizontal

83    segment of the intraparietal sulcus at Talairach coordinates (-42, -48, 48 mm), a region

84    suggested to be unrelated to emotion regulation [33]. During the happy memory recall block, the

85    NF signal was quantified as the percent signal change from the mean signal in the preceding

86    rest block. The signal was updated every 2 s and visually presented to participants as a red bar.

87    **MRI data processing**

88    The Analysis of Functional NeuroImages (AFNI; http://afni.nimh.nih.gov) software suite was

89    utilized for image processing. After discarding the first three volumes to achieve signal

90    equilibrium, preprocessing steps were conducted, including despike, RETROICOR [34] along with

91    respiratory volume per time (RVT) regression [35] for physiological noise correction, slice-timing

92    and motion corrections, nonlinear warping to the MNI template brain with resampling to 2 mm³

93    voxels using the Advanced Normalization Tools (ANTs; http://stnava.github.io/ANTs/), spatial

94    smoothing with a 6 mm full-width at half maximum (FWHM) Gaussian kernel, and scaling of

95    signals to percent change relative to the voxel-wise mean.

96    Activation during the self-regulation was assessed using a general linear model (GLM)

97    analysis. We extracted the two types of brain response in this process (Fig. 1). The GLM

98    regressors included response models for the regulation and counting blocks, each modeled with

99    a boxcar function convolved with the canonical hemodynamic response function (HRF). The

100    regressor for the regulation block was used for estimating the first type of brain activation,

101    regulation block activation (Fig. 1). The GLM regressors also included twelve motion parameters

102    (three rotations, three translations, and their temporal derivatives), three principal components

103    of ventricle signals, local white matter signals (ANATICOR) [36], and an event-related regressor

104    for the onset of any condition block (modeled as a delta function convolved with HRF) as

105    nuisance covariates. Volumes with frame-wise displacement greater than 0.3 and their

106    preceding volume were censored in the GLM.

107    Another type of brain response, the feedback event-related response (Figure 1), was

108    evaluated using additional event-related regressors for each feedback presentation. The

109    response was modeled as a delta function, modulated by the feedback amplitude normalized in

110    each run, and convolved with the hemodynamic response function (HRF). Since the feedback

111    signal was presented during the regulation block and this event regressor could be collinear with

112    the regulation block regressor, we orthogonalized the feedback-response-event regressor with

113    respect to the regulation block regressor. The beta coefficients from the general linear model

114    (GLM) analysis were used as response estimates.

115    These response estimates were assessed for each block independently using the least

116    squares - separate (LS-S) approach [37], in which separate regressors for the target block and all

117    other blocks were included to estimate the response in a block, and repeated for each block in

118    separate GLM analyses. The use of this block-wise response facilitates subsequent group

119    analysis using linear mixed-effect model and provides the estimates with better test-retest

120    reliability [38, 39].

**Clustering analysis**

121    

122    Clustering analysis was conducted on the whole-brain beta maps solely for the active group

123    participants to categorize their brain activation patterns. This approach was chosen because

124    including the control group could introduce brain activation patterns with various unspecific

125    effects. Such inclusion might increase the dimensionality of latent subtypes and complicate the

126    extraction of subtypes relevant to the active NF training process.

127    A significant challenge posed by this analysis is the high dimensionality of the whole-brain

128    beta maps; the problem often referred to as the "curse of dimensionality" [40]. This phenomenon

129    refers to the issues arising in high-dimensional spaces, where distances between points

130    become uniformly large, data points are sparsely distributed, and there is a high risk of

131    overfitting, leading to clustering solutions that are difficult to reproduce.

132    To address this challenge, we implemented a multi-step strategy to reduce dimensionality and

133    identify an informative space associated with treatment outcomes. Additionally, we employed

134    cross-validation to ensure the clustering solution's stability and reproducibility. Our approach

135    involved (1) aggregating voxel-wise responses into regional averages using a functional brain

136    atlas, (2) removing the regions irrelevant to treatment outcomes, (3) applying Uniform Manifold

137    Approximation and Projection (UMAP) [41] to extract a low-dimensional representation, (4)

138    applying k-means clustering in the UMAP space, and (5) employing repeated cross-validation to

139    assess the robustness of the clustering solution [42]. Figure 2 illustrates the flowchart of the

140    analysis performed to delineate subtypes of brain activation.

141    Initially, voxel-wise beta values were averaged within each functional region defined by the

142    Shen 268 atlas [43, 44]. This atlas was chosen for its demonstrated performance in various

143    predictive modeling studies [43, 45-49]. The data for each participant are averaged across the three

144    runs or for each run. We tested both approaches as a hyperparameter to test which one yielded

145    the most stable clustering solution. Subsequently, we identified regions correlating with changes

146    in the MADRS score (change relative to the baseline score expressed as a ratio) after adjusting

147    for the effects of age and sex. A threshold of $p < 0.1$ was used for this selection, prioritizing

148    dimensionality reduction rather than finding significantly related regions. After excluding the

149    regions irrelevant to the treatment outcome, we utilized UMAP for dimensionality reduction,

150    followed by k-means clustering. The UMAP parameters were set to optimize the clustered

151    distribution within the reduced-dimensional space, with 'n_neighbors' (the number of

152    neighboring points used in the local manifold approximation) set to 50, and 'min_dist' (the

8

153  minimum distance apart that points are allowed to be in the low-dimensional representation) set

154  to 0. The large 'n_neighbors' value and small 'min_dist' value encourage the UMAP to extract a

155  space with a clustered distribution. Subsequently, k-means clustering was executed within the

156  UMAP-defined space.

157      The stability of the clusters was rigorously evaluated through cross-validation, facilitated by

158  the 'reval' Python package [42]. This involved dividing the dataset into two halves, applying

159  clustering to one half, and then applying the derived clustering rule to the opposite half, and vice

160  versa, to assess the consistency of the cluster labels between the two rules (Fig. 2). This

161  procedure was repeated 10 times, each with a unique data split, to calculate the average

162  stability score. The stability measure is the mean normalized Hamming distance between the

163  two solutions, ranging from 0 to 1. A smaller value indicates a more stable and reproducible

164  solution. Since this measure is larger with a larger number of clusters, it was scaled by the

165  stability of random labeling of the same number of clusters [50].

166      In summary, these procedures were conducted across the hyperparameter space of response

167  summary (average across runs or a sequence of three runs), UMAP random seed (50 values),

168  UMAP dimension (ranging from 2 to 50), and the number of clusters (ranging from 2 to 5) to find

169  robust clustering with the smallest mean distance between the cross-validated solutions.

170  **Mapping the brain responses in each subtype**

171      After identifying the subtypes, we evaluated the voxel-wise response patterns for each

172  subtype using a linear mixed-effect (LME) model analysis performed voxel-wise with the 'lme4'

173  package [51] in R language and environment for statistical computing. This LME analysis was

174  applied to the beta values for each block, run, and participant, incorporating fixed effects for the

175  run, the identified subtype from the clustering analysis, age, and sex, as well as a random effect

176  for participants at the intercept level. The mean response for each subtype was calculated using

177  the 'emmeans' package [52] in R. The mean maps of the subtypes were thresholded at a voxel-

178  wise $p < 0.001$, corrected for cluster size with a $p < 0.05$ using AFNI 3dClustSim.

9

**Statistical testing**

LME analysis was performed to investigate changes in MADRS scores one week after a NF session for both active and control groups. The LME model included fixed effects for time (pre-/post-NF), group (active/control), age, and sex, as well as a random effect for participants at the intercept. We also examined whether treatment efficacy was correlated with NF training performance measures, including the mean NF amplitude, mean left amygdala response during the regulation block, and the change in left amygdala response between the first and the last training runs using linear model analysis.

We then examined whether the identified subtypes of brain activation during NF training were associated with demographic and training performance variables, as well as the depressive symptom score (MADRS) at the baseline and its change post-NF. Each variable was examined using linear model analysis, with subtypes serving as independent variables. The post-hoc analysis for the mean response for each subtype and the difference between the subtypes was performed using the 'emmeans' package [52], and *p*-values for the post-hoc comparisons were corrected by multivariate *t* adjustment.

In addition, the subtyping rules established in the active group were applied to the control group to assess if similar patterns in MADRS score changes could be observed across the subtypes. The presence of consistent subtype associations with symptom reduction in both the active and control groups would suggest the operation of a common treatment mechanism. This could indicate that the efficacy of NF training is not specific to the neurofeedback signal from an emotion-related region.

## Results

### Data selection

Post-training MADRS scores were not available for two participants in the active group. fMRI data from five participants were excluded from the analysis due to problems with physiological

10

204    signal acquisition (two from the active group and one from the control group) or excessive head

205    motion (two from the active group) with more than 30% of time points censored (Framewise

206    Displacement [FD] > 0.3) in all training runs. Additionally, four participants in the active group

207    who exhibited excessive head motion in any one of the training runs were excluded only when

208    analyzing data with a series of three training runs. Analyses were performed using the largest

209    sample size available for each test measure, regardless of missing data on other variables (see

210    Supplementary Information [SI], Section 1). No significant differences in age, sex, and baseline

211    MADRS scores were observed between the active and control groups in any of these selected

212    datasets.

**Reduction of depressive symptoms following NF training**

214    Figure 3 shows the MADRS scores before and after NF training for each group, with average

215    scores represented by the height of bars and individual participant scores depicted as points

216    connected by lines. LME analysis identified a significant main effect of time. Post-hoc analysis

217    revealed a significant decrease in post-session scores ($t$ = -4.404, $d$ = -0.704, $p$ < 0.001). SI

218    Table S2a provides the ANOVA tables for the LME analysis. While the time by group interaction

219    was not statistically significant ($p$ = 0.090), the active group showed a significant symptom

220    reduction ($t$ = -5.576, $d$ = -0.978, $p$ < 0.001) but not the control group ($t$ = -1.609, $d$ = -0.430, $p$ =

221    0.111).

222    However, left amygdala activation had no significant main effect of the group ($F$ = 2.870, $p$ =

223    0.094) and the interaction between the group and run ($F$ = 0.597, $p$ = 0.551). Furthermore,

224    within the active group, no significant associations were found between changes in MADRS

225    scores and measures of NF training performance, including mean NF amplitude, average left

226    amygdala response during regulation blocks, or changes in left amygdala response from the

227    first to the last training runs (see SI Section 2 for comprehensive details).

**Clustering results**

In clustering of activation maps related to the regulation block, the optimal stability score of 0.217 was obtained using the following hyperparameters: averaging responses across runs for each participant, extracting a 40-dimensional UMAP space, and clustering into two subtypes, which included 32 and 31 participants respectively. Similarly, in the clustering of activation maps associated with the feedback event, an optimal stability score of 0.229 was reached by analyzing a series of responses across three runs, extracting a 14-dimensional UMAP space, and forming three subtypes consisting of 23, 15, and 22 participants each. However, no significant association was found between the subtypes identified in the regulation block and those in the feedback event ($\chi^2 = 3.143$, $p = 0.208$).

Figure 4 shows *t*-maps of the mean beta values for the subtypes of regulation block activations (REG subtypes), with detailed peak coordinates available in SI, Section 3, Table S3. The first subtype, REG-A, exhibited increased BOLD signals in regions commonly identified in NF studies [12], including the lateral prefrontal cortex, superior and inferior parietal lobules, supplementary motor area (SMA), anterior insula, thalamus, and cerebellum. In contrast, decreased BOLD signals were observed in Heschl's gyrus, posterior insula, middle cingulate cortex, and cuneus. The second subtype, REG-B, demonstrated increased BOLD signals in the lateral prefrontal cortex, SMA, and anterior insula, with decreased activity in default mode network (DMN) areas such as the precuneus, posterior cingulate cortex, mediodorsal prefrontal cortex, Rolandic operculum, and fusiform gyrus.

Figure 5 presents *t*-maps of the mean beta values for subtypes of the brain responses to the feedback event (FBR subtypes), with peak coordinates detailed in SI, Section 4, Table S4. The first subtype, labeled FBR0, demonstrated a limited response, particularly in the Rolandic operculum and right supramarginal gyrus. The second subtype, labeled FBR-, was characterized by a negative association of BOLD signal changes in response to NF signals, affecting regions including the SMA, middle cingulate cortex, Heschl's gyrus, Rolandic

254    operculum, and visual cortex. In contrast, the third subtype, labeled FBR+, displayed a positive

255    association of BOLD signal changes in areas such as the precuneus, middle cingulate cortex,

256    lateral prefrontal cortex, nucleus accumbens, cerebellum, and inferior occipital regions.

257    **Subtype associations with demographics, NF-training characteristics, and the**
258    **depressive symptom in the active group**

259    Linear model analyses examining associations of the REG and FBR subtypes with age, mean

260    NF signal amplitude, mean left amygdala activation across runs, and baseline MADRS scores

261    revealed no significant differences between the subtypes. Similarly, $\chi^2$ tests showed no

262    significant associations of these subtypes with sex or study participation (refer to SI section 5 for

263    further detailed statistics).

264    However, changes in MADRS scores, indicative of clinical efficacy, were significantly

265    associated with these subtypes. The main effects of both REG and FBR subtypes, as well as

266    their interaction effect on changes in MADRS scores, were statistically significant (REG, $F =$

267    8.735, $p = 0.005$; FBR, $F = 5.326$, $p = 0.008$; interaction $F = 3.471$, $p = 0.039$). Post-hoc

268    analysis revealed a significant decrease in MADRS scores within the REG-B subtype ($t = -6.077$,

269    $d = -1.702$, $p < 0.001$), in contrast to REG-A ($t = -2.170$, $d = -0.608$, $p = 0.068$). For the FBR

270    subtypes, significant decreases in MADRS scores were observed in both FBR- ($t = -4.995$, $d = -$

271    $1.399$, $p < 0.001$) and FBR+ ($t = -3.818$, $d = -1.069$, $p = 0.001$), whereas the FBR0 subtype

272    exhibited no significant change ($t = -0.839$, $d = -0.235$, $p = 0.786$).

273    Moreover, changes in MADRS scores varied by subtype combination, as illustrated in Figure

274    6. A significant reduction in scores was noted when REG-A was paired with FBR- ($t = -3.327$, $d$

275    $= -0.931$, $p = 0.003$). Within the REG-B group, no significant reduction was observed with the

276    FBR0 subtype ($t = -1.221$, $d = -0.342$, $p = 0.401$). However, significant reductions were seen

277    when REG-B was paired with FBR- ($t = -3.789$, $d = -1.061$, $p < 0.001$) and FBR+ ($t = -6.018$, $d =$

278    $-1.685$, $p < 0.001$). In summary, participants in the active group classified as REG-A with FBR-,

279    or REG-B with FBR- or FBR+ experienced significant reductions in MADRS scores one week

280    after NF training.

**Application of clustering rules to the control group**

282    When the clustering rules derived from the active group were applied to the control group,

283    participants were classified as follows: 13 into REG-A and 14 into REG-B for the REG subtypes;

284    and 11 into FBR0, 5 into FBR-, and 11 into FBR+ for the FBR subtypes. No significant

285    association was observed between the REG and FBR subtypes within the control group ($\chi^2$ =

286    0.345, $p$ = 0.842). Additionally, the distribution of participants across subtypes did not

287    significantly differ from that observed in the active group for both REG ($\chi^2$ = 0, $p$ = 1.000) and

288    FBR ($\chi^2$ = 0.526, $p$ = 0.768) subtypes.

289    In the control group, no significant associations were found between the subtypes and

290    variables such as age, sex, study participation, mean NF signal amplitude, changes in left

291    amygdala activation, or baseline MADRS scores, as detailed in SI Section 6. Contrary to the

292    active group, there were no significant differences in MADRS score changes across subtypes

293    within the control group.


## Discussion

295    The primary objective of this study was to identify subtypes of brain activation patterns during

296    NF training that could explain interindividual differences in clinical response. Our analysis

297    revealed that training characteristics within the target brain region, such as mean NF signal

298    amplitude, mean left amygdala activation, and signal changes in the left amygdala, were not

299    associated with changes in depressive symptoms. However, we found significant associations

300    between subtypes of whole-brain activation patterns during NF training and changes in MADRS

301    scores. In contrast, the control group showed no significant associations between the subtypes

302    and changes in MADRS scores, underscoring that the effects of NF treatment are distinct from

303    placebo effects.

14

304     Significant symptom reduction was observed across multiple subtype combinations, indicating

305     the existence of multiple brain functional pathways to successful treatment. Notably, individuals

306     classified within the FBR0 subtype, characterized by their non-responsiveness to feedback

307     signals, did not experience significant symptom reduction, highlighting the critical role of brain

308     response to feedback signals in effective NF training. Significant symptom reduction was

309     evident in individuals exhibiting brain activation correlated with feedback signals, either in a

310     positive (FBR+) or negative (FBR-) manner. However, the efficacy of this response pattern was

311     dependent on their regulation subtype. Individuals exhibiting increased activation across broad

312     brain areas during self-regulation (REG-A), which overlap with areas commonly reported in NF

313     studies [12], required a negative response to NF signals (FBR-) to achieve significant symptom

314     reduction. Conversely, individuals in the REG-B subtype, who showed increased activation in

315     executive control areas such as the lateral prefrontal cortex and SMA and salience network

316     regions such as the anterior insula (areas that overlap with those involved in general skill

317     learning and emotion regulation [53, 54]), alongside suppressive activation in DMN regions during

318     self-regulation, achieved symptom reduction regardless of the polarity of their brain response to

319     feedback signals. These findings suggest that nonspecific increases in brain activation do not

320     contribute to treatment efficacy; rather, selective modulation of brain responses is crucial for

321     successful NF therapy.

322     Increased activation of the DMN, often associated with self-referential thoughts [55, 56], has

323     been frequently observed in individuals with MDD compared to healthy controls, and is linked to

324     repetitive negative thinking [57-65]. Particularly relevant to the pathophysiology of MDD and its

325     treatment mechanisms could be the significant role of deactivating the posterior DMN,

326     especially the posterior cingulate cortex (PCC). Current evidence suggests that the PCC acts as

327     a major cortical hub for various self-referential phenomena, ranging from spatial body

328     localization to self-talk and autobiographical information processing [66-69]. In this context, the

329     characteristic deactivation of the PCC observed in brain activation subtypes associated with

330    clinical response in our study supports the hypothesis that this region may underlie the

331    persistent and intensified negative self-referential mentation that characterizes depression.

332    The absence of suppressive activation during the regulation task can be compensated for by

333    negative responses to feedback signals in the SMA, middle cingulate cortex (MCC), auditory

334    cortex (including Heschl's gyrus and Rolandic operculum), and visual cortex. Although the

335    functional mechanism of these negative responses to NF signals remains unclear, they may

336    counteract maladaptive brain activation patterns associated with MDD. Notably, alterations in

337    SMA activity, frequently reported in MDD [70-73], and changes in middle cingulate activation during

338    reward anticipation have also been highlighted in a meta-analysis [74]. The association between

339    less MCC activation during NF training and greater symptom reduction was also reported [26].

340    Furthermore, increased connectivity in the auditory cortex was associated with repetitive

341    negative thinking [75], and heightened activity or connectivity in the visual cortex have also been

342    reported in individuals with MDD [48, 76].

343    Selective increases and decreases in brain activation associated with successful NF training

344    have also been documented in a meta-analysis of amygdala NF studies [4]. This study reported

345    that successful amygdala modulation was linked to deactivation in the posterior insula and DMN

346    regions, including the ventromedial prefrontal cortex and PCC, during the regulation period.

347    Additionally, negative activations in successful modulators were noted in the left

348    parahippocampal gyrus and cuneus, aligning with the present results observed in the FBR-

349    subtype. Notably, this meta-analysis [4] defined training success as the ability to down-regulate

350    amygdala activation, in contrast to our study, which focused on up-regulating amygdala

351    activation and evaluated training success by changes in depressive symptoms. This suggests a

352    common underlying mechanism in emotion regulation training that leads to depressive symptom

353    reduction, despite varying approaches to NF signal definition and operational definitions of

354    success.

355     In light of our findings, the therapeutic effects of NF targeting the left amygdala cannot be

356    solely attributed to changes in its activity alone. This observation aligns with previous research

357    showing that activations in brain regions beyond the intended target during NF training can

358    mediate symptom reduction [26]. Such evidence supports the notion that brain self-regulation

359    through NF involves a large-scale network, extending well beyond a single target region.

360    Although amygdala activity is a reliable marker of emotional states and can indicate successful

361    regulatory efforts, as demonstrated in many NF protocols [11], our findings underscore the

362    importance of the regulation process itself in achieving treatment success and explaining the

363    observed variability in treatment outcomes.

364    To our knowledge, this study is the first to explore interindividual differences in brain activation

365    during NF training and its association with therapeutic effects on depressive symptoms. This

366    investigation marks a significant step toward understanding the mechanisms of NF as a

367    treatment for psychiatric disorders. Our findings, which highlight the importance of the whole-

368    brain regulation process over the mere amplitude of amygdala activity, could pave the way for

369    improvements in NF protocols. Thus, incorporating feedback that targets both regulation-related

370    activities and responses to feedback within large-scale brain networks could potentially enhance

371    treatment efficacy. As NF methodologies continue to evolve, focusing on the process of

372    regulating affective states rather than merely activating specific brain regions might yield more

373    effective treatments. In this context, pattern-based NF approaches such as DecNef [77],

374    connectome-based neurofeedback [78], and semantic neurofeedback [79, 80], which emphasize

375    extracting feedback signals from multivariate brain activation patterns, offer a promising

376    direction for refining training protocols.

377    The limitations of this study warrant careful consideration. Although the clustering analysis

378    was designed to classify individuals into distinct subtypes, it is possible that these subtypes do

379    not represent distinct groups per se. Rather, the subtypes may be part of a continuous

380    distribution with no clear demarcation between groups. This is suggested by the fact that the

17

381  obtained optimal stability scores, which indicate the discrepancy between cross-validated

382  solutions, were greater than 0.2, indicating potential ambiguity in cluster definition. In scenarios

383  where the samples have a continuous distribution, splitting into an equal number of participants

384  could emerge as a more stable approach. This phenomenon may explain why our analysis

385  resulted in nearly equal numbers of participants being categorized into each subtype,

386  highlighting that these identified subtypes may not delineate distinct groups but rather capture

387  different facets of the continuous response spectrum. It is also noteworthy that subtype

388  classification was not related to age, sex, or baseline severity of depressive symptoms.

389  However, the relationship between the subtypes and other demographic factors or

390  neurobiological indicators before and after treatment remains unexplored. Further investigation

391  may reveal subtypes of treatment response based on specific pretreatment characteristics.

392   In conclusion, the present study demonstrates that the therapeutic outcomes of NF training

393  are significantly influenced by whole-brain activation patterns both during the process of affect

394  regulation and during the response to feedback signals. Our findings reveal multiple patterns of

395  brain activity associated with significant therapeutic effects, suggesting a variety of potential

396  pathways to recovery through NF training. In future studies, the optimization of NF training may

397  involve real-time monitoring of brain activity during regulation efforts and in response to

398  feedback signals, as well as tailoring feedback signals to the individual's current state of training

399  progress. Such an approach could enhance treatment precision, adjusting it to match each

400  participant's unique neural response patterns, thus potentially increasing the effectiveness of NF

401  training.

## Acknowledgement

406    the study design, data collection and analysis, decision to publish, or preparation of the

407    manuscript. The authors are solely responsible for the content. The authors would like to

408    acknowledge Jerzy Bodurka, Ph.D. (1964–2021), for his intellectual and scientific contributions

409    to the designing the left amygdala neurofeedback protocols and development of the real-time

410    fMRI neurofeedback systems, which provided the foundation for the present work. During the

411    preparation of this work the authors used ChatGPT (https://chat.openai.com/) and DeepL

412    (https://www.deepl.com/write) in order to improve language and readability. After using these

413    tools, the authors reviewed and edited the content as needed and take full responsibility for the

414    content of the publication.

## Conflict of Interest

416    The authors report no financial relationships with commercial interests in relation to the work

417    described.

# References

1.  Institute for Health Metrics and Evaluation (IHME). Global Burden of Disease Study 2019. IHME, https://www.healthdata.org/research-analysis/gbd2020.

2.  Greenberg PE, Fournier A-A, Sisitsky T, Simes M, Berman R, Koenigsberg SH *et al.* The Economic Burden of Adults with Major Depressive Disorder in the United States (2010 and 2018). *PharmacoEconomics* 2021; **39**(6)**:** 653-665.

3.  Pindi P, Houenou J, Piguet C, Favre P. Real-time fMRI neurofeedback as a new treatment for psychiatric disorders: A meta-analysis. *Progress in neuro-psychopharmacology & biological psychiatry* 2022; **119:** 110605.

4.  Goldway N, Jalon I, Keynan JN, Hellrung L, Horstmann A, Paret C *et al.* Feasibility and utility of amygdala neurofeedback. *Neurosci Biobehav Rev* 2022; **138:** 104694.

5.  Fernandez-Alvarez J, Grassi M, Colombo D, Botella C, Cipresso P, Perna G *et al.* Efficacy of bio- and neurofeedback for depression: a meta-analysis. *Psychological medicine* 2022; **52**(2)**:** 201-216.

6.  Taylor SF, Martz ME. Real-time fMRI neurofeedback: the promising potential of brain-training technology to advance clinical neuroscience. *Neuropsychopharmacology : official publication of the American College of Neuropsychopharmacology* 2023; **48**(1)**:** 238-239.

7.  Girges C, Vijiaratnam N, Zrinzo L, Ekanayake J, Foltynie T. Volitional Control of Brain Motor Activity and Its Therapeutic Potential. *Neuromodulation* 2022; **25**(8)**:** 1187-1196.

8.  Dudek E, Dodell-Feder D. The efficacy of real-time functional magnetic resonance imaging neurofeedback for psychiatric illness: A meta-analysis of brain and behavioral outcomes. *Neurosci Biobehav Rev* 2021; **121:** 291-306.

9.  Trambaiolli LR, Kohl SH, Linden DEJ, Mehler DMA. Neurofeedback training in major depressive disorder: A systematic review of clinical efficacy, study quality and reporting practices. *Neurosci Biobehav Rev* 2021; **125:** 33-56.

10. Martz ME, Hart T, Heitzeg MM, Peltier SJ. Neuromodulation of brain activation associated with addiction: A review of real-time fMRI neurofeedback studies. *NeuroImage Clinical* 2020; **27:** 102350.

11. Barreiros AR, Almeida I, Baia BC, Castelo-Branco M. Amygdala Modulation During Emotion Regulation Training With fMRI-Based Neurofeedback. *Frontiers in human neuroscience* 2019; **13:** 89.

12. Emmert K, Kopel R, Sulzer J, Bruhl AB, Berman BD, Linden DEJ *et al.* Meta-analysis of real-time fMRI neurofeedback studies using individual participant data: How is brain regulation mediated? *Neuroimage* 2016; **124**(Pt A)**:** 806-812.

13.    Zotev V, Phillips R, Misaki M, Wong CK, Wurfel BE, Krueger F *et al.* Real-time fMRI neurofeedback training of the amygdala activity with simultaneous EEG in veterans with combat-related PTSD. *NeuroImage Clinical* 2018; **19:** 106-121.

14.    Herwig U, Lutz J, Scherpiet S, Scheerer H, Kohlberg J, Opialla S *et al.* Training emotion regulation through real-time fMRI neurofeedback of amygdala activity. *Neuroimage* 2019; **184:** 687-696.

15.    Paret C, Zahringer J, Ruf M, Gerchen MF, Mall S, Hendler T *et al.* Monitoring and control of amygdala neurofeedback involves distributed information processing in the human brain. *Hum Brain Mapp* 2018; **39**(7)**:** 3018-3031.

16.    Lubianiker N, Paret C, Dayan P, Hendler T. Neurofeedback through the lens of reinforcement learning. *Trends Neurosci* 2022; **45**(8)**:** 579-593.

17.    Sitaram R, Ros T, Stoeckel L, Haller S, Scharnowski F, Lewis-Peacock J *et al.* Closed-loop brain training: the science of neurofeedback. *Nat Rev Neurosci* 2017; **18**(2)**:** 86-100.

18.    Lawrence EJ, Su L, Barker GJ, Medford N, Dalton J, Williams SC *et al.* Self-regulation of the anterior insula: Reinforcement learning using real-time fMRI neurofeedback. *Neuroimage* 2014; **88:** 113-124.

19.    Skottnik L, Sorger B, Kamp T, Linden D, Goebel R. Success and failure of controlling the real-time functional magnetic resonance imaging neurofeedback signal are reflected in the striatum. *Brain Behav* 2019; **9**(3)**:** e01240.

20.    Sutton RS, Barto AG. *Reinforcement learning: An introduction*. 2nd edn. MIT Press Cambridge, MA2018.

21.    Young KD, Zotev V, Phillips R, Misaki M, Yuan H, Drevets WC *et al.* Real-time FMRI neurofeedback training of amygdala activity in patients with major depressive disorder. *PLoS One* 2014; **9**(2)**:** e88785.

22.    Young KD, Siegle GJ, Bodurka J, Drevets WC. Amygdala Activity During Autobiographical Memory Recall in Depressed and Vulnerable Individuals: Association With Symptom Severity and Autobiographical Overgenerality. *Am J Psychiatry* 2016; **173**(1)**:** 78-89.

23.    Yuan H, Young KD, Phillips R, Zotev V, Misaki M, Bodurka J. Resting-state functional connectivity modulation and sustained changes after real-time functional magnetic resonance imaging neurofeedback training in depression. *Brain connectivity* 2014; **4**(9)**:** 690-701.

24.    Tsuchiyagaito A, Smith JL, El-Sabbagh N, Zotev V, Misaki M, Al Zoubi O *et al.* Real-time fMRI neurofeedback amygdala training may influence kynurenine pathway metabolism in major depressive disorder. *NeuroImage Clinical* 2021; **29:** 102559.

25.    Misaki M, Mulyana B, Zotev V, Wurfel BE, Krueger F, Feldner M *et al.* Hippocampal volume recovery with real-time functional MRI amygdala neurofeedback emotional training for posttraumatic stress disorder. *Journal of affective disorders* 2021; **283:** 229-235.

26. Misaki M, Phillips R, Zotev V, Wong CK, Wurfel BE, Krueger F *et al.* Brain activity mediators of PTSD symptom reduction during real-time fMRI amygdala neurofeedback emotional training. *NeuroImage Clinical* 2019; **24:** 102047.

27. Misaki M, Phillips R, Zotev V, Wong CK, Wurfel BE, Krueger F *et al.* Real-time fMRI amygdala neurofeedback positive emotional training normalized resting-state functional connectivity in combat veterans with and without PTSD: a connectome-wide investigation. *NeuroImage Clinical* 2018; **20:** 543-555.

28. American Psychiatric Association. *Diagnostic and Statistical Manual of Mental Disorders DSM-IV-TR Fourth Edition (Text Revision).* American Psychiatric Publishing: Washington, DC, 2000.

29. First MB. Structured clinical interview for DSM‑IV‑TR axis I disorders, research version, patient edition (SCID‑I/P). *Biometrics research*: New York, NY, 2002.

30. Sheehan DV, Lecrubier Y, Sheehan KH, Amorim P, Janavs J, Weiller E *et al.* The Mini-International Neuropsychiatric Interview (M.I.N.I.): the development and validation of a structured diagnostic psychiatric interview for DSM-IV and ICD-10. *The Journal of clinical psychiatry* 1998; **59 Suppl 20:** 22-33;quiz 34-57.

31. Zotev V, Krueger F, Phillips R, Alvarez RP, Simmons WK, Bellgowan P *et al.* Self-regulation of amygdala activation using real-time FMRI neurofeedback. *PLoS One* 2011; **6**(9)**:** e24522.

32. Montgomery SA, Åsberg M. A new depression scale designed to be sensitive to change. *The British journal of psychiatry : the journal of mental science* 1979; **134:** 382-389.

33. Fias W, Lammertyn J, Caessens B, Orban GA. Processing of Abstract Ordinal Knowledge in the Horizontal Segment of the Intraparietal Sulcus. *The Journal of Neuroscience* 2007; **27**(33)**:** 8952-8956.

34. Glover GH, Li TQ, Ress D. Image-based method for retrospective correction of physiological motion effects in fMRI: RETROICOR. *Magn Reson Med* 2000; **44**(1)**:** 162-167.

35. Birn RM, Smith MA, Jones TB, Bandettini PA. The respiration response function: the temporal dynamics of fMRI signal fluctuations related to changes in respiration. *Neuroimage* 2008; **40**(2)**:** 644-654.

36. Jo HJ, Saad ZS, Simmons WK, Milbury LA, Cox RW. Mapping sources of correlation in resting state FMRI, with artifact detection and removal. *Neuroimage* 2010; **52**(2)**:** 571-582.

37. Mumford JA, Turner BO, Ashby FG, Poldrack RA. Deconvolving BOLD activation in event-related designs for multivoxel pattern classification analyses. *Neuroimage* 2012; **59**(3)**:** 2636-2643.

38. Chen G, Taylor PA, Stoddard J, Cox RW, Bandettini PA, Pessoa L. Sources of Information Waste in Neuroimaging: Mishandling Structures, Thinking Dichotomously, and Over-Reducing Data. *Aperture Neuro* 2022; **2:** 1-22.

39. Chen G, Padmala S, Chen Y, Taylor PA, Cox RW, Pessoa L. To pool or not to pool: Can we ignore cross-trial variability in FMRI? *Neuroimage* 2021; **225:** 117496.

40. Altman N, Krzywinski M. The curse(s) of dimensionality. *Nature methods* 2018; **15**(6)**:** 399-400.

41. McInnes L, Healy J, Melville J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:180203426* 2018.

42. Landi I, Mandelli V, Lombardo MV. reval: A Python package to determine best clustering solutions with stability-based relative clustering validation. *Patterns (N Y)* 2021; **2**(4)**:** 100228.

43. Shen X, Finn ES, Scheinost D, Rosenberg MD, Chun MM, Papademetris X *et al.* Using connectome-based predictive modeling to predict individual behavior from brain connectivity. *Nature protocols* 2017; **12**(3)**:** 506-518.

44. Shen X, Tokoglu F, Papademetris X, Constable RT. Groupwise whole-brain parcellation from resting-state fMRI data for network node identification. *Neuroimage* 2013; **82:** 403-415.

45. Greene AS, Gao S, Scheinost D, Constable RT. Task-induced brain state manipulation improves prediction of individual traits. *Nature communications* 2018; **9**(1)**:** 2807.

46. Greene AS, Gao S, Noble S, Scheinost D, Constable RT. How Tasks Change Whole-Brain Functional Organization to Reveal Brain-Phenotype Relationships. *Cell Rep* 2020; **32**(8)**:** 108066.

47. Ju S, Horien C, Shen X, Abuwarda H, Trainer A, Constable RT *et al.* Connectome-based predictive modeling shows sex differences in brain-based predictors of memory performance. *Frontiers in Dementia* 2023; **2**.

48. Misaki M, Tsuchiyagaito A, Guinjoan SM, Rohan ML, Paulus MP. Trait repetitive negative thinking in depression is associated with functional connectivity in negative thinking state rather than resting state. *Journal of affective disorders* 2023; **340:** 843-854.

49. Yoo K, Rosenberg MD, Hsu WT, Zhang S, Li CR, Scheinost D *et al.* Connectome-based predictive modeling of attention: Comparing different functional connectivity features and prediction methods across datasets. *Neuroimage* 2018; **167:** 11-22.

50. Lange T, Roth V, Braun ML, Buhmann JM. Stability-Based Validation of Clustering Solutions. *Neural Computation* 2004; **16**(6)**:** 1299-1323.

51. Bates D, Mächler M, Bolker B, Walker S. Fitting Linear Mixed-Effects Models Usinglme4. *Journal of Statistical Software* 2015; **67**(1)**:** 1–48.
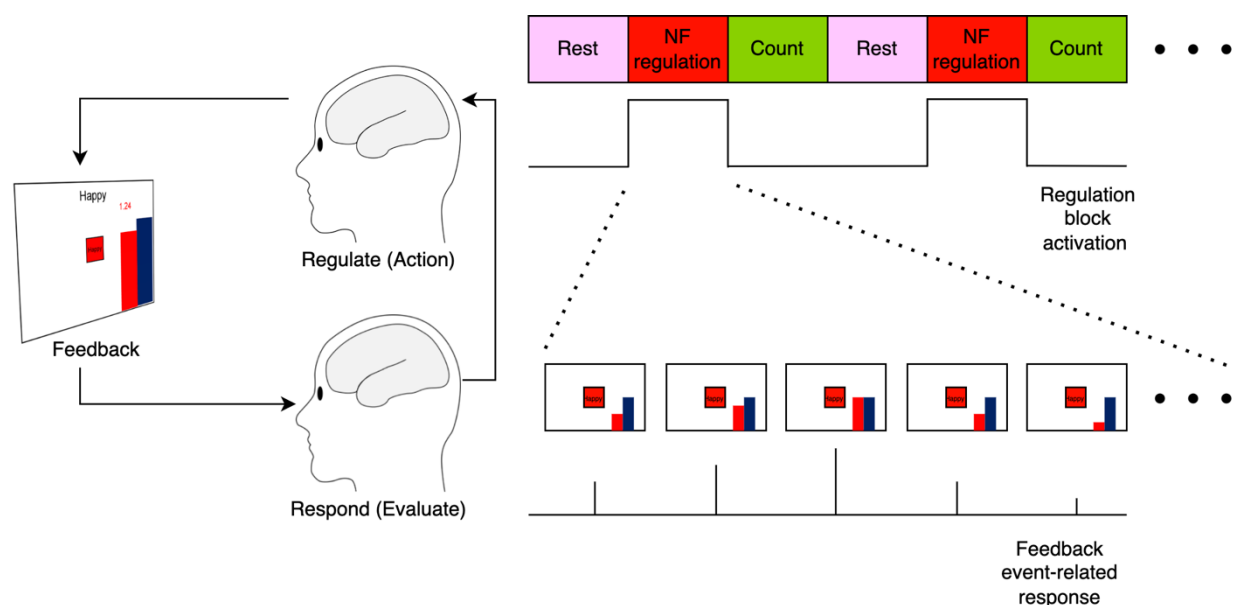
52.  emmeans: Estimated Marginal Means, aka Least-Squares Means. https://CRAN.R-project.org/package=emmeans, 2022, Accessed Date Accessed 2022 Accessed.

53.  Linhartova P, Latalova A, Kosa B, Kasparek T, Schmahl C, Paret C. fMRI neurofeedback in emotion regulation: A literature review. *Neuroimage* 2019; **193:** 75-92.

54.  Etkin A, Buchel C, Gross JJ. The neural bases of emotion regulation. *Nat Rev Neurosci* 2015; **16**(11)**:** 693-700.

55.  Axelrod V, Rees G, Bar M. The default network and the combination of cognitive processes that mediate self-generated thought. *Nature Human Behaviour* 2017; **1**(12)**:** 896-910.

56.  Andrews-Hanna JR, Smallwood J, Spreng RN. The default network and self-generated thought: component processes, dynamic control, and clinical relevance. *Annals of the New York Academy of Sciences* 2014; **1316**(1)**:** 29-52.

57.  Hamilton JP, Farmer M, Fogelman P, Gotlib IH. Depressive Rumination, the Default-Mode Network, and the Dark Matter of Clinical Neuroscience. *Biol Psychiatry* 2015; **78**(4)**:** 224-230.

58.  Wise T, Marwood L, Perkins AM, Herane-Vives A, Joules R, Lythgoe DJ *et al.* Instability of default mode network connectivity in major depression: a two-sample confirmation study. *Translational psychiatry* 2017; **7**(4)**:** e1105.

59.  Zhu X, Zhu Q, Shen H, Liao W, Yuan F. Rumination and Default Mode Network Subsystems Connectivity in First-episode, Drug-Naive Young Patients with Major Depressive Disorder. *Sci Rep* 2017; **7:** 43105.

60.  Bessette KL, Jenkins LM, Skerrett KA, Gowins JR, DelDonno SR, Zubieta JK *et al.* Reliability, Convergent Validity and Time Invariance of Default Mode Network Deviations in Early Adult Major Depressive Disorder. *Front Psychiatry* 2018; **9:** 244.

61.  Jacob Y, Morris LS, Huang KH, Schneider M, Rutter S, Verma G *et al.* Neural correlates of rumination in major depressive disorder: A brain network analysis. *NeuroImage Clinical* 2020; **25:** 102142.

62.  Makovac E, Fagioli S, Rae CL, Critchley HD, Ottaviani C. Can't get it off my brain: Meta-analysis of neuroimaging studies on perseverative cognition. *Psychiatry Res Neuroimaging* 2020; **295:** 111020.

63.  Misaki M, Tsuchiyagaito A, Al Zoubi O, Paulus M, Bodurka J, Tulsa I. Connectome-wide search for functional connectivity locus associated with pathological rumination as a target for real-time fMRI neurofeedback intervention. *NeuroImage Clinical* 2020; **26:** 102244.

64.  Stern ER, Eng GK, De Nadai AS, Iosifescu DV, Tobe RH, Collins KA. Imbalance between default mode and sensorimotor connectivity is associated with perseverative thinking in obsessive-compulsive disorder. *Translational psychiatry* 2022; **12**(1)**:** 19.
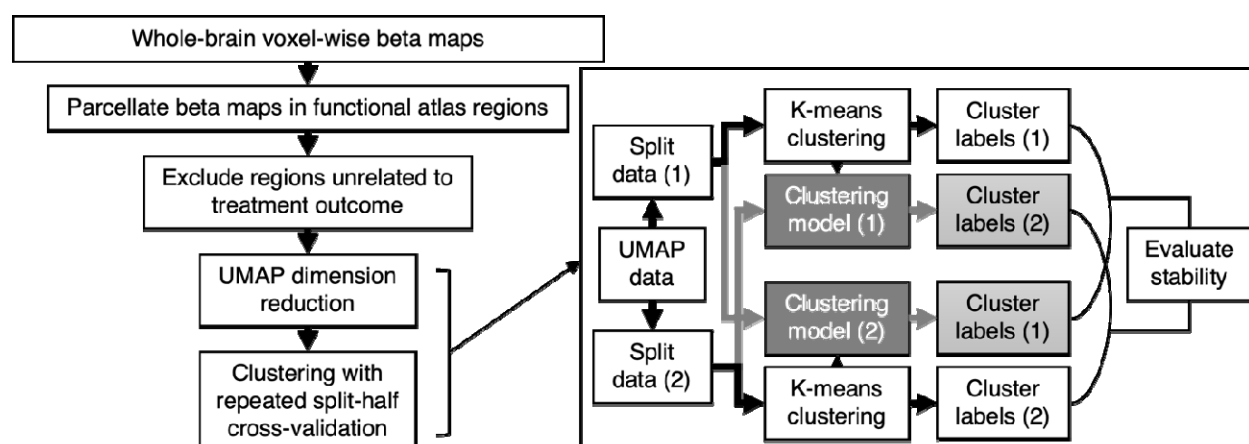
65.   Yang MH, Guo ZP, Lv XY, Zhang ZQ, Wang WD, Wang J *et al.* BMRMI Reduces Depressive Rumination Possibly through Improving Abnormal FC of Dorsal ACC. *Neural Plast* 2022; **2022:** 8068988.

66.   Agathos J, Steward T, Davey CG, Felmingham KL, Ince S, Moffat BA *et al.* Differential engagement of the posterior cingulate cortex during cognitive restructuring of negative self- and social beliefs. *Social cognitive and affective neuroscience* 2023; **18**(1).

67.   Natu VS, Lin JJ, Burks A, Arora A, Rugg MD, Lega B. Stimulation of the Posterior Cingulate Cortex Impairs Episodic Memory Encoding. *J Neurosci* 2019; **39**(36)**:** 7173-7182.

68.   Guterstam A, Bjornsdotter M, Gentile G, Ehrsson HH. Posterior cingulate cortex integrates the senses of self-location and body ownership. *Curr Biol* 2015; **25**(11)**:** 1416-1425.

69.   Foster BL, Koslov SR, Aponik-Gremillion L, Monko ME, Hayden BY, Heilbronner SR. A tripartite view of the posterior cingulate cortex. *Nat Rev Neurosci* 2023; **24**(3)**:** 173-189.

70.   Andreescu C, Butters M, Lenze EJ, Venkatraman VK, Nable M, Reynolds III CF *et al.* fMRI activation in late-life anxious depression: a potential biomarker. *International Journal of Geriatric Psychiatry* 2009; **24**(8)**:** 820-828.

71.   Ionescu DF, Niciu MJ, Mathews DC, Richards EM, Zarate Jr CA. NEUROBIOLOGY OF ANXIOUS DEPRESSION: A REVIEW. *Depression and Anxiety* 2013; **30**(4)**:** 374-385.

72.   Sarkheil P, Odysseos P, Bee I, Zvyagintsev M, Neuner I, Mathiak K. Functional connectivity of supplementary motor area during finger-tapping in major depression. *Comprehensive psychiatry* 2020; **99:** 152166.

73.   Northoff G, Hirjak D, Wolf RC, Magioncalda P, Martino M. All roads lead to the motor cortex: psychomotor mechanisms and their biochemical modulation in psychiatric disorders. *Molecular psychiatry* 2021; **26**(1)**:** 92-102.

74.   Jauhar S, Fortea L, Solanes A, Albajes-Eizagirre A, McKenna PJ, Radua J. Brain activations associated with anticipation and delivery of monetary reward: A systematic review and meta-analysis of fMRI studies. *PLOS ONE* 2021; **16**(8)**:** e0255292.

75.   Tsuchiyagaito A, Sanchez SM, Misaki M, Kuplicki R, Park H, Paulus MP *et al.* Intensity of repetitive negative thinking in depression is associated with greater functional connectivity between semantic processing and emotion regulation areas. *Psychological medicine* 2023; **53**(12)**:** 5488-5499.

76.   Wu F, Lu Q, Kong Y, Zhang Z. A Comprehensive Overview of the Role of Visual Cortex Malfunction in Depressive Disorders: Opportunities and Challenges. *Neuroscience Bulletin* 2023; **39**(9)**:** 1426-1438.

77.   Shibata K, Lisi G, Cortese A, Watanabe T, Sasaki Y, Kawato M. Toward a comprehensive understanding of the neural mechanisms of decoded neurofeedback. *Neuroimage* 2019; **188:** 539-556.
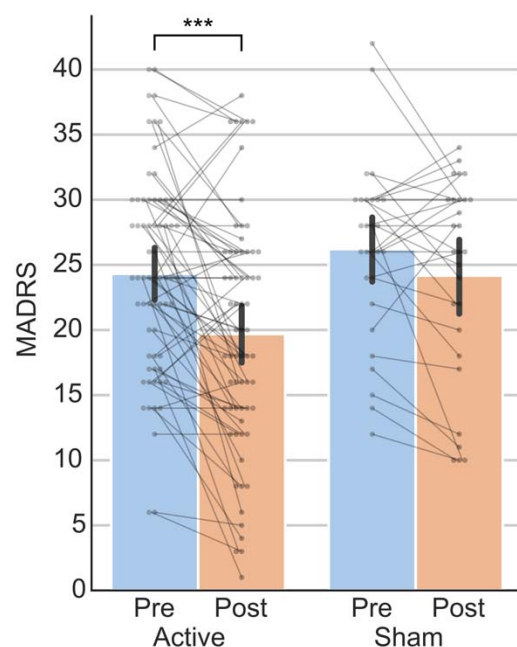
78.  Scheinost D, Hsu TW, Avery EW, Hampson M, Constable RT, Chun MM *et al.* Connectome-based neurofeedback: A pilot study to improve sustained attention. *Neuroimage* 2020; **212:** 116684.

79.  Ciarlo A, Russo AG, Ponticorvo S, di Salle F, Lührs M, Goebel R *et al.* Semantic fMRI neurofeedback: a multi-subject study at 3 tesla. *Journal of Neural Engineering* 2022; **19**(3)**:** 036020.

80.  Russo AG, Luhrs M, Di Salle F, Esposito F, Goebel R. Towards semantic fMRI neurofeedback: navigating among mental states using real-time representational similarity analysis. *J Neural Eng* 2021; **18**(4)**:** 046015.

**Figure 1**. Schematic diagram of the neurofeedback-based brain regulation training loop (left panel) and the models used to estimate brain activation at each epoch (right panel). The reinforcement learning process is typically characterized by the evaluation of a feedback signal followed by the adjustment of action (mental regulation) to increase the reward (feedback) signal. Brain activations corresponding to these two components were analyzed using a block-wise response model during the neurofeedback regulation blocks and an event-related response model for each neurofeedback presentation at every TR, modulated by feedback amplitude.

**Figure 2**. Procedures of subtyping whole-brain beta maps.

**Figure 3:** MADRS scores before and after NF training in the active and control groups. The bars show the mean MADRS scores for each group, with the error bars representing their standard error. Individual participant scores are shown as dots, with lines connecting pre- and post-training scores to highlight changes for each individual. A statistically significant decrease in post-training MADRS scores was observed in the active group (***, $p < 0.001$).

**REG-A**



**REG-B**



**Figure 4**. *t* maps of mean activation in the regulation block for the REG subtypes. The maps are thresholded at a voxel-wise p < 0.001, with cluster-size correction at p < 0.05.
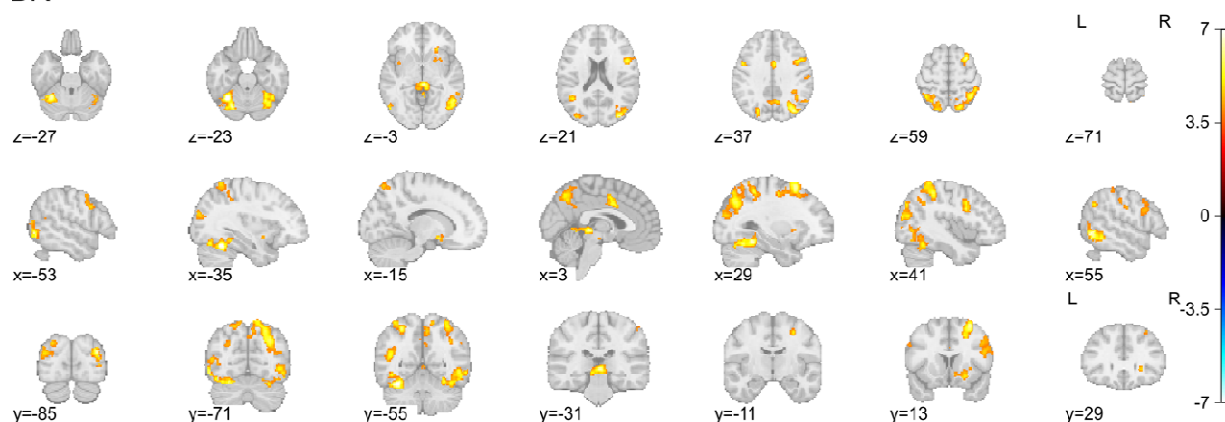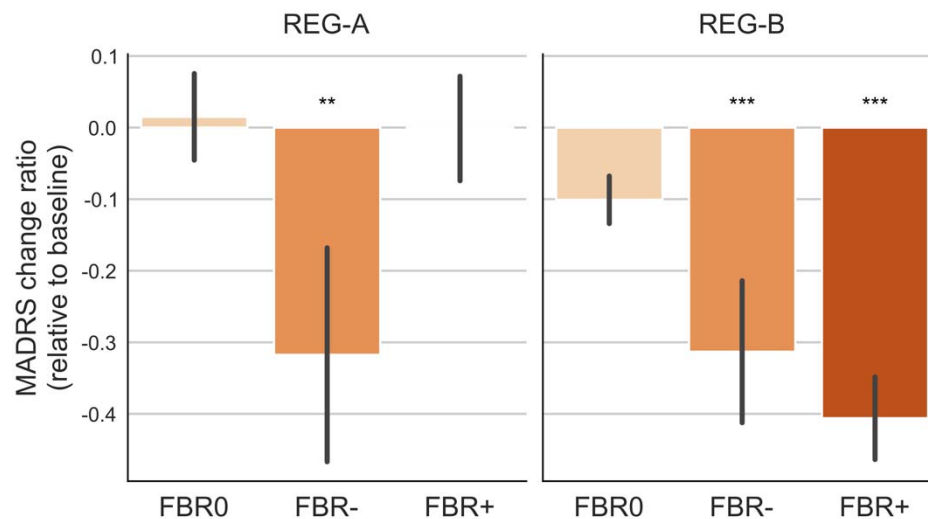
**Figure 5**. *t* maps of mean activation in response to NF signals for the FBR subtypes. The maps are thresholded at a voxel-wise *p* < 0.001, with cluster-size correction at *p* < 0.05.

31

**Figure 6**. Ratio of change in MADRS score from baseline for each REG and FBR subtype. **, $p$ < 0.005; ***, $p$ < 0.001.