

1 **SHORT TITLE: Test-retest reliability in fMRI of depression**

2

3 **Importance of test-retest reliability for promoting fMRI based screening and**
4 **interventions in major depressive disorder**

5 Laurie Compère, PhD ¹, Greg J. Siegle, PhD ¹, Kimberly Young, PhD ¹

6 ¹ Department of Psychiatry, University of Pittsburgh School of Medicine

7 Western Psychiatric Institute and Clinic

8 3811 O'Hara Street

9 Pittsburgh, PA 15213

10 ljc44@pitt.edu

11 gsiegle@pitt.edu

12 youngk@pitt.edu

13

14

15 **Corresponding author**

16 Laurie Compère, Department of Psychiatry, University of Pittsburgh School of
17 Medicine

18 Western Psychiatric Institute and Clinic

19 3811 O'Hara Street

20 Pittsburgh, PA 15213

21 Email: ljc44@pitt.edu

22

23 **Abstract**

24 Proponents of personalized medicine have promoted neuroimaging evaluation and
25 treatment of major depressive disorder in three areas of clinical application: clinical
26 prediction, outcome evaluation, and neurofeedback. Whereas psychometric
27 considerations such as test-retest reliability are basic precursors to clinical adoption
28 for most clinical instruments, they are often not considered for neuroimaging
29 assessments. As an example, we consider functional magnetic resonance imaging
30 (fMRI) of depression, a common and particularly well validated mechanistic technology
31 for understanding disorder and guiding treatment. In this article, we review work on
32 test-retest reliability for depression fMRI studies. We find that basic psychometrics
33 have not been regularly attended to in this domain. For instance, no fMRI
34 neurofeedback study has included measures of test-retest reliability despite the implicit

35 assumption that brain signals are stable enough to train. We consider several factors
36 that could be useful to aid clinical translation including 1) attending to how the BOLD
37 response is parameterized, 2) identifying and promoting regions or voxels with stronger
38 psychometric properties 3) accounting for within-individual changes (e.g., in
39 symptomatology) across time and 4) focusing on tasks and clinical populations that
40 are relevant for the intended clinical application. We apply these principles to published
41 prognostic and neurofeedback data sets. The broad implication of this work is that
42 attention to psychometrics is important for clinical adoption of mechanistic assessment,
43 is feasible, and may improve the underlying science.

44

45 **Keywords:** depression, fMRI, neurofeedback, psychometric, treatment prediction,
46 test-retest reliability

47

48

49

50

1. GENERAL INTRODUCTION

51 Proponents of personalized medicine have promoted mechanistic evaluation and
52 mechanistically targeted treatments for major depressive disorder (Hansen and Siegle,
53 2015). As an example, we consider functional magnetic resonance imaging (fMRI), a
54 common and particularly well validated mechanistic technology that represents a
55 promising proof-of-concept in this area. Longitudinal assessment of changes in
56 regional brain activity using functional magnetic resonance imaging (fMRI) has
57 increasingly been used in research on the treatment of psychiatric conditions including
58 major depressive disorder (MDD) (Fournier et al., 2014). As good psychometric
59 properties are essential for any measure to be considered for clinical adoption
60 (Pickford and Guilford, 2007), best-practice guidelines for increasing generalizability
61 and reproducibility of fMRI results are emerging (Nichols et al., 2017; Poldrack et al.,
62 2017). We focus here on test-retest reliability in task-based fMRI and neurofeedback
63 (fMRI-nf) designs, using MDD as a running case example. Ideally, our observations
64 can be applied to other technologies and across neuropsychiatric disorders.

65 To understand the current state of the field, we conducted literature reviews
66 quantifying how often test-retest reliability was reported in fMRI biomarker and real-
67 time fMRI neurofeedback (rtfMRI-nf) studies in MDD. As we will demonstrate below,
68 this was infrequent and the general literature has shown that when assessed, reliability
69 was generally low. We thus suggest a few analytic techniques for improving test-rest
70 reliability in fMRI and its clinical applicability. We focus on data analysis to make our
71 suggestions maximally applicable to already collected data. Finally, we test these
72 suggested principles on published MDD neuroimaging treatment outcome and
73 neurofeedback datasets as proofs of concept.

74 The idea that fMRI could have therapeutic utility is based on assumptions that
75 hemodynamic activity is reliable over time in the absence of intervention, and that
76 observed changes between one scan session and the next have significant and
77 interpretable values (Barch and Mathalon, 2011). The reliability of fMRI also affects its
78 criterion validity, as poor reliability limits the strength of association between the
79 instrument and other relevant measures (Vul et al., 2009).

80

81 1.1. On computing reliability of fMRI

82 Demonstrating ability to achieve similar results over time, or the reliability of
83 measures is considered critical to creating a clinically useful measure (Pickford and

84 Guilford, 2007). Reliability is a quantitative measure of stability of an individual's data
85 (Bennett and Miller, 2013). It refers to the ability of a measure to distinguish participants
86 from each other and to replicate the order of individuals' ranks during repeated
87 assessments, assuming they do not experience true signal change between
88 assessments (Barch and Mathalon, 2011).

89 Though stable regional hemodynamic activations at the group level can be
90 observed over time, there are significant changes in how each subject contributes
91 individually to the observed group activation (Caceres et al., 2009; Zandbelt et al.,
92 2008). Various approaches have been used to measure test-retest reliability for fMRI.
93 For example, a Pearson correlation between visits across time measure the degree
94 to which visits on two occasions are linearly related, where data from each visit are
95 independently scaled (e.g., Harrington et al., 2006). A more common approach, and
96 the measure we focus on in this manuscript, involves computing intra-class correlation
97 coefficients (ICC) that also reflect rank ordering of values across days (Bennett and
98 Miller, 2010) as a ratio of variance between values observed across subjects and sites
99 divided by the total visit variance (Bartko, 1966). Values range from 0 (no reliability) to
100 1 (perfect reliability). There are three different types of ICCs described by the princeps
101 article written by Shrout and Fleiss (1979). The ICC(1,1) index is similar to the Pearson
102 correlation but normalizes by the pooled mean and variance across visits. ICC(2,1) is
103 an agreement index that allows generalization of results across scanners while
104 ICC(3,1) works under the assumption that the variance is the same across scanners.
105 Therefore, the ICC(3,1), mostly used across studies, is a scanner consistency index
106 where the effect of scanner is considered a fixed effect (Shrout and Fleiss, 1979). In
107 order to match the literature in the field and because we considered the scanner as a
108 covariate of interest when investigating the impact of taking into account clinical and
109 design covariates when computing reliability indexes, we mainly used ICC(3,1) in our
110 analyses.

111 Interpretation of ICC values is subjective with no uniformly accepted standards; ICC
112 values of less than 0.4 are often considered poor, 0.4-0.59 fair, 0.60-0.74 good, and
113 above 0.75 excellent (Cicchetti, 1994; Plichta et al., 2012; Shrout and Fleiss, 1979),
114 though more stringent cutoffs have also been recommended (e.g., Portney and
115 Watkins, 2009). Negative ICC value are usually interpreted as no reliability (Bartko,
116 1976), since these values are outside the theoretical limits of ICC (although negative

117 values may appear when within-subject variance is greater than between-subject
118 variance) (Lahey et al., 1983).

119 Though the ICC has been recommended for use in fMRI (Caceres et al., 2009),
120 some fMRI analysis packages (SPM, FSL) do not inherently support computation of
121 this metric, potentially hinting at its perceived value in the field, though other packages
122 (e.g., AFNI, NIFTI) do provide for its computation, and add-on packages (e.g., reliability
123 toolbox for SPM or other packages on R) do allow such computations (see
124 Computation of voxelwise ICCs using different tools in Box 3 in supplementary
125 materials for more details). Indeed, reliability estimates have been rarely reported in
126 fMRI studies and usually reveal poor reliability when estimated (Elliott et al., 2020).
127 Non-clinical studies have generally found low to moderate test-retest reliability values
128 for regional fMRI activity, with ICCs ranging from 0.33-0.66 (reviewed in Bennett and
129 Miller, 2010).

130

131 **1.2. Biomarker/Prediction Studies Review**

132 Many studies suggest fMRI measurements can be used to predict treatment
133 outcome in MDD (for reviews, see Arnone, 2019; Fonseka et al., 2018; Phillips and
134 Swartz, 2014; Wessa and Lois, 2015). The underlying assumption is that biomarkers
135 in the brain are involved in the causal process of MDD. Therefore, it is expected that
136 the activity measured in these biomarkers is related to, and evolves over time with,
137 symptom changes in general and that for interventions targeting the biomarker the
138 more abnormal activity observed, the more effective the intervention will be. However,
139 clinical applications of these findings are limited by the possibility that these biomarkers
140 may have low test-retest reliability (Nord et al., 2017). If a biomarker is not reliable, it
141 is impractical to interpret its activation at the individual level (Fu et al., 2013; Guo et
142 al., 2012). Thus, despite strong predictive utility, researchers acknowledge that their
143 results might be limited by poor test-retest reliability (e.g., Fu et al., 2015). Of particular
144 interest, the amygdala, a commonly reported biomarker for MDD, shows poor to good
145 reliability when emotional stimuli are displayed, with great heterogeneity between
146 studies in healthy participants (Lois et al., 2018). Thus, we surveyed the predictive
147 fMRI literature in MDD to examine whether this first step was being taken.

148 **1.2.1. Method**

149 A PubMed search with the key words “fMRI AND biomarker OR prediction OR
150 predict AND depression OR MDD OR major depressive disorder NOT Rest NOT

151 Resting” produced 140,640 results. We combined this list with other articles
 152 discovered in our submitted fMRI meta-analysis of depression treatment outcome
 153 prediction studies (Strege et al., 2020) to complete the list of articles (Table 1).” After
 154 removing articles not including functional neuroimaging (i.e., studies focusing on
 155 volumetric measures or using PET) or human participants, we were left with 55 studies
 156 (Table 1).

157 **Table 1: Studies Examining neuroimaging biomarkers of pharmacotherapy and**
 158 **psychotherapy outcomes in Major Depressive Disorder and mention of test-**
 159 **retest reliability of the studies**

160

Reference	Treatment(s)	Biomarker	Findings	Mention of signal reliability	Possibility to test signal reliability
Sheline et al.(2001)	Sertraline	Amygdala	Decreased activation following treatment	No	Yes
Davidson, et al. (2003)	Venlafaxine	ACC	Greater activation at baseline associated with better treatment response	No	Yes
Fu et al. (2004)	Fluoxetine	ACC, ventral striatum, cerebellum	Reduction of dynamic range associated with symptoms improvement	No	Yes
Canli et al. (2005)	None	Amygdala	Amygdala activation at baseline predicts symptom improvement	No	No

Schaefer, et al. (2006)	Venlafaxine	Prefrontal, temporal and parietal cortices, insula, basal ganglia and hippocampus	Normalized activation after treatment	No	Yes
Siegle, Carter, & Thase (2006)	CBT	sgACC and amygdala	Low and high, respectively, activation is associated with greater symptom improvement after therapy	No	No
Anand, et al. (2007)	Sertaline	Amygdala and ACC	Decrease activation in limbic regions and increased connectivity with the ACC after treatment	No	Yes
Chen et al. (2007)	Fluoxetine	ACC	Greater activation at baseline predict faster rates of symptom improvement	No	No
Fales et al. (2007)	Escitalopram	DLPFC	Enhanced activation following treatment	No	Yes
Fitzgerald et al. (2007)	TMS	Middle frontal gyrus, left precuneus, left precentral gyrus, left medial frontal	Decreased activation after low frequency treatment in middle frontal gyrus and left	No	Yes

		gyrus, right inferior frontal gyrus	precuneus in respondents – Increased activation after high frequency treatment in left prefrontal gyrus, left medial frontal gyrus, right inferior frontal gyrus in respondents		
Fu et al. (2007)	Fluoxetine	Hippocampus and extrastriate cortex	Greater activation following treatment and associated with symptom improvement	No	Yes
Langenecker et al. (2007)	S-citalopram	Insula, right middle frontal gyrus, left inferior frontal gyrus, amygdala and cerebellar vermis	Greater activation at baseline associated with symptoms improvement	No	Yes
Robertson et al. (2007)	Bupropion	Amygdala	Reduced activation associated with symptom improvement	No	Yes
Walsh et al. (2007)	Fluoxetine	dACC, left middle frontal and lateral temporal cortices	Reduced activity at baseline associated with	Yes (discussion section) ^a	Yes

			symptom improvement		
Fu et al. (2008)	CBT	dACC	Reduced activation at baseline associated with symptom improvement	No	Yes
Benedetti et al. (2009)	Venlafaxine	Right medial frontal gyrus	Decreased activation following treatment was associated with symptom improvement	No	Yes
Costafreda, et al. (2009)	CBT	ACC, superior and middle frontal cortices, paracentral cortex, superior parietal cortex, precuneus and cerebellum	Activation contributed to prediction of remission	No	No
Dichter et al. (2010)	Behavioral Action Therapy	Paracingulate gyrus	Activation was prognostic for depressive symptom change after psychotherapy	No	Yes
Forbes et al., 2010	CBT and SSRI	Striatum and mPFC	Final levels of severity symptoms were related to pretreatment striatal reactivity	No	No

			and greater striatal and lower mPFC activity was prognostic for anxiety symptom reduction		
Keedwell et al. (2010)	Various antidepressants	Right visual cortex and right sgACC	Greater baseline activity associated with clinical improvement after treatment	No	Yes
Lemogne et al. (2010)	Various antidepressants	Left DLPFC	Reduced activation following treatment	No	Yes
López-Solà et al. (2010)	Duoxetine	pACC, right prefrontal cortex, pons	Clinical improvement associated with reduced activation	No	Yes
Roy et al. (2010)	Citalopram	Ventromedial prefrontal cortex and ACC	Greater activation at baseline associated with symptom improvement	No	Yes
Victor, et al. (2010)	Sertaline	Amygdala	Decreased activation after treatment	No	Yes
Wagner et al. (2010)	Citalopram, reboxetine	Amygdala, hippocampus	Decreased activation after citalopram treatment	No	Yes
Frodl et al. (2011)	Mirtazapine, venlafaxine	Left fusiform gyrus, right	Increased activation in the	No	Yes

		rolandic operculum	left fusiform gyrus at baseline was associated with a better response to venlafaxine and smaller activation in the right rolandic operculum was related to better response to mirtazapine		
Light et al. (2011)	Venlafaxine, fluoxetine	Ventrolateral prefrontal cortex	Reduced activity at baseline is associated with anhedonia reduction	No	Yes
Ritchey, et al. (2011)	CBT	Ventromedial prefrontal cortex	Increased activity at baseline associated with symptom improvement	No	Yes
Samson et al. (2011)	Mirtazapine, venlafaxine	dmPFC, posterior cingulate cortex, superior frontal gyrus, caudate nucleus and insula	Greater activation at treatment associated with better treatment response	No	Yes
Arnone et al. (2012)	Citalopram	Amygdala	Reduced activation	No	Yes

			following treatment		
Godlewska, et al. (2012)	Escitalopram	Amygdala	Reduced activity after treatment	No	No
Rosenblau et al. (2012)	Escitalopram	Amygdala, prefrontal cortex	Decreased activation following treatment	No	Yes
Ruhé, et al. (2012)	Paroxetine	Amygdala	Lower activation associated with better response to treatment after	No	Yes
Siegle et al. (2012)	CBT	sgACC	Reduced activation at baseline associated with greater symptom improvement	Yes	Yes
Stoy et al. (2012)	Escitalopram	Ventral striatum	Increased activation following treatment	No	Yes
Tao et al. (2012)	Fluoxetine	Amygdala, orbitofrontal cortex and sgACC	Decreased activation after treatment	Yes (discussion section) ^b	Yes
Wang et al. (2012)	Fluoxetine	Insula, left ACC and middle frontal gyrus	Decreased activation in insula and left ACC and greater in the middle frontal gyrus following treatment	No	Yes

Furey et al. (2013)	Scopolamine	Middle occipital cortex	Increased activation at baseline was prognostic for symptoms improvement	No	Yes
Heller et al. (2013)	Fluoxetine or venlafaxine	Nucleus accumbens	Greater activation following treatment associated with more self-reported positive affect	No	Yes
Miller et al. (2013)	Escitalopram	Midbrain, DLPFC, paracingulate, ACC, thalamus and caudate nuclei	Reduced activation at baseline correlated with greater improvement following treatment	No	No
Rizvi et al. (2013)	Fluoxetine and olanzapine	Premotor cortex	Increased activation at baseline in respondents was prognostic for symptom improvement	Yes but not reported (method section) ^c	Yes
Victor, et al. (2013)	Sertraline	pgACC	Increased correlation at baseline correlated with greater clinical improvement after treatment	No	Yes
Toki et al. (2014)	Various antidepressants	Left hippocampus	Increased activation	No	No

			associated with greater response treatment		
Yoshimura et al. (2014)	CBT	vACC	Improvements in depressive symptoms were negatively correlated with its activity	No	Yes
Fu et al. (2015)	Duloxetine	Posterior cingulate cortex	Increased activation following treatment	Yes (limitation section) ^d	Yes
Furey et al. (2015)	Scopolamine	sgACC and middle occipital cortex	Increased and decreased activation, respectively, associated with treatment response	No	Yes
Straub et al., 2015	CBT	sgACC	Activation before treatment related to therapeutic success	No	Yes
Williams et al. (2015)	Escitalopram, sertraline, venlafaxine	Amygdala	Decreased activation at baseline was associated with treatment response	No	Yes
Cullen et al. (2016)	Various antidepressants	Rostral and sgACC, insula, middle frontal cortex, right	Decreased activation in postral and sgACC and increased in ,	No	Yes

		hippocampus and left cerebellum	insula, middle frontal cortex, right hippocampus and left cerebellum associated with symptom improvement		
Delaveau et al. (2016)	Agomelatine	DLPFC and precuneus	Activation at baseline was related to treatment response	No	Yes
Doerig et al. (2016)	CBT	Amygdala	Activity in this region pre-intervention is negatively correlated with the outcome	No	No
Godlewska, et al. (2016)	Escitalopram	ACC, insula, amygdala and thalamus	Reduced activity after treatment associated with treatment response	No	Yes
Gyurak et al. (2016)	Escitalopram, sertraline and venlafaxine	DLPFC and inferior parietal cortex	Increased DLPFC activation at baseline associated with remission and increased inferior parietal activation associated with remission for SSRI and the	No	Yes

			opposite pattern for SNRI		
Opmeer et al. (2016)	-	Rostral ACC	Increased activation at baseline was prognostic for remission	No	Yes
Szczepanik et al. (2016)	Scopolamine	Amygdala	Increased activity at baseline was associated with symptoms improvement	Yes (limitation section) ^e	No
Fang et al., (2017)	Transcutaneous vagus nerve stimulation	Insula	Activation level at first stimulation session associated with clinical improvement	No	No
Sankar, et al. (2017)	Duloxetine	Left inferior frontal activity	Decreased activation following treatment	No	Yes
Spies et al. (2017)	Escitalopram	Precuneus and PCC	Deactivation before treatment was related to change in symptoms after 2 weeks of treatment	No	No
Godlewska et al. (2018)	Escitalopram	pgACC	Activity before treatment was able to predict response status (responder vs non-responder)	No	No

		at the level of individual participant			
			Increased activation following treatment associated with better treatment outcome	No	Yes
Rubin-Falcone et al. (2018)	CBT	sgACC, medial prefrontal cortex, lingual gyrus			

161

162 ACC: Anterior Cingulate Cortex; CBT: Cognitive Behavioral Therapy; dACC: dorsal Anterior Cingulate
 163 Cortex; DLPFC: Dorsolateral Prefrontal Cortex; dmPFC: dorsomedial Prefrontal Cortex; mPFC : medial
 164 Prefrontal Cortex ; MDD: Major Depressive Disorder; PCC : Posterior Cingulate Cortex; pgACC:
 165 pregenual Anterior Cingulate Cortex; sgACC: subgenual Anterior Cingulate Cortex ; SSRI : selective
 166 serotonin reuptake inhibitor

167

168 ^a "Test-retest effects were accounted for by the healthy control group, who underwent the same scans
 169 at the same time points"

170 ^b "repeat fMRI assessment of healthy comparison subjects, as well as repeat assessment of the
 171 depressed adolescents, thus providing assessment of expected test-retest reliability"

172 ^c "For analyses of change over time, a higher level fixed effects analysis was run for each subject,
 173 contrasting parameter estimates within subject for the response to slides at the two time points of
 174 interest."

175 ^d"perhaps in part reflecting the poor test-retest reliability of amygdala response to these emotional faces
 176 [54], while resting-state fMRI data show greater robustness and reproducibility [55]. Test-retest reliability
 177 of a neuroimaging measure becomes particularly important in the development of biomarkers for
 178 prognosis and diagnosis [44]."

179 ^e"some investigators have raised concerns regarding the reliability of the BOLD signal (Boubela et al.,
 180 2015). Nevertheless, studies have found that emotional stimuli evoke a consistent pattern of responsivity
 181 over repeated sessions (Johnstone et al., 2005)."

182

183 **1.2.2. Results**

184 Though most of the reviewed studies could have reported test-retest reliability (i.e.,
 185 participants performed two scans), most did not mention it. Seven mentioned reliability
 186 in the discussion and only one reported test-retest reliability at the subject level; Siegle
 187 et al. (2012) reported "sgACC z scores and reactivity had moderate test-retest
 188 reliability in controls undergoing testing approximately 16 weeks apart (N=27; r=0.39

189 [P=0.04]). All but 1 had a pretest z score less than 0.5, and all but 2 had a posttest z
190 score less than 0.5, suggesting stability within a restricted range.” Other studies that
191 mention reliability describe stability of group effects. For example, “Test-retest effects
192 were accounted for by the healthy control group, who underwent the same scans at
193 the same time points” (Walsh et al., 2007) is often reported in the discussion. This
194 technique, while valuable, does not yield estimates of test-retest reliability at the
195 individual subject level; the absence of a main effect of Time is evidence of the lack of
196 a mean shift, but not of the stability of participants ranks.

197

198 **1.3. rtfMRI-nf Studies Review**

199 Interventions that use biological measures as real-time targets, including rtfMRI-nf
200 also implicitly assume reliability. rtfMRI-nf trains patients to regulate the hemodynamic
201 activity in regions of interest (Decharms, 2008) with the hope that changing a causal
202 mechanism will result in symptom changes. rtfMRI-nf appears useful for several clinical
203 populations, including patients with MDD (Thibault et al., 2018). Most patients can
204 learn volitional control of hemodynamic activity in a targeted brain region (Fovet et al.,
205 2015) which has been associated with clinical improvements (Fovet et al., 2015;
206 Linden, 2014; Linden et al., 2012; Young et al., 2014) suggesting potential translational
207 applications (Decharms, 2008; Ruiz et al., 2014; Thibault et al., 2018). An implicit
208 assumption of rtfMRI-nf is that the signal measured on one day represents the same
209 quantity measured on subsequent days, and thus performance on that metric can be
210 trained over days. Consequently, test-retest reliability seems a strong prerequisite.
211 Thus, as for prediction studies, we considered whether test-retest reliability is being
212 reported in the fMRI neurofeedback literature.

213 **1.3.1. Method**

214 A PubMed search with the key words “(neurofeedback AND fMRI) OR rt-fMRI-nf)
215 AND (depression OR MDD OR major depressive disorder” provided 44 results. After
216 removing articles not including rtfMRI-nf or patients suffering from MDD, we were left
217 with 11 studies (Table 2).

218

219 **Table 2: rt-fMRI-nf studies in Major Depressive Disorder and mention and**
 220 **possibility of test-retest reliability**

Reference	Neurofeedback	ROI	Mention of the reliability of the signal	Possibility to test signal reliability	How could they look at reliability
Linden et al. (2012)	Upregulation	Functional localizer of brain areas involved in the generation of positive emotions (e.g., VLPLC, insula)	No	Yes	Same regions selected by the localizer on different sessions
Zotev, et al., (2014)	Upregulation	Left amygdala (anatomical)	No	No	-
Young et al. (2014) ^a	Upregulation	Left amygdala (anatomical)	No	Yes	Reliability of fMRI signal in ROI
Yuan et al. (2014) ^a	Upregulation	Left amygdala (anatomical)	No	No	-
Zotev et al. (2016) ^a	Upregulation	Left amygdala (anatomical)	No	No	-
Hamilton et al. (2016)	Downregulation	Functional localizer of the salience network	No	No	-
Young et al. (2017) ^b	Upregulation	Left amygdala (anatomical)	No	Yes	Reliability of fMRI signal in ROI

Young, Misaki, et al. (2017) ^b	Upregulation	Left amygdala (anatomical)	No	Yes	Reliability of fMRI signal in ROI
Young et al. (2018) ^b	Upregulation	Left amygdala (anatomical)	No	Yes	Reliability of fMRI signal in ROI
MacDuffie et al. (2018)	Upregulation and downregulation	Functional localizer of ACC	No	No	-
Mehler et al. (2018)	Upregulation	Functional localizer of brain areas involved in seeing positive versus neutral pictures (e.g., insula and striatum)	No	Yes	Same regions selected by the localizer on different sessions

221

222 ACC: Anterior Cingulate Cortex; VLPFC: Ventrolateral Prefrontal Cortex

223 *References associated with the same letter refer to the same data set*

224 **1.3.2. Results**

225 None of the examined fMRI-nf studies reported on the reliability of the signal being
 226 trained (Table 2 and specific discussion of functional localizers in Box 4 in
 227 supplements).

228

229 **1.4. Conclusions Thus Far**

230 MDD studies using fMRI for clinical prediction or treatment rarely mention reliability,
 231 mirroring the more general fMRI literature (for meta-analysis, see Elliott et al., 2020).
 232 This lack of reporting could be due to failure to consider psychometrics important, or
 233 systematic decisions not to report observed low reliabilities. Indeed, reliability in
 234 published fMRI research in non-clinical studies, across protocols, tasks, regions of
 235 interest, psychological functions, and retest intervals have been fairly low (ICC~0.50),

236 with most published studies reporting values between 0.33-0.66. These values are
237 mostly below “good” reliability thresholds for psychometrically sound clinical tests
238 (~0.6).

239

240 **2. POTENTIAL WAYS TO OPTIMIZE TEST-RETEST RELIABILITY IN fMRI/rtfMRI-** 241 **NF**

242 To facilitate reporting of reliability in clinical studies as part of every-day
243 neuroimaging-science, the remainder of this article is dedicated to introducing ways to
244 report, improve, and increase clinical applicability of test-retest reliability for fMRI in
245 clinical populations. We apply and evaluate these suggestions in two published data
246 sets (Siegle et al., 2012; Young et al., 2017b).

247 There is already a strong literature on optimizing preprocessing, which can increase
248 measurement of true signal, and thus reliability (Andersson et al., 2001; Miki et al.,
249 2000; Oakes et al., 2005; Zhilkin and Alexander, 2004). We therefore begin by
250 considering whether using alternate ways of indexing task-related reactivity in single-
251 subject data with optimized preprocessing lead to improved test-retest reliability.

252 As each combination of task, design, scanner, preprocessing and analysis strategy
253 has a unique value of reliability that cannot necessarily be generalized to other studies
254 (Braver et al., 2010), it may be useful to have standardized generally applicable
255 methods to find out which regions and analysis methods have sufficient psychometric
256 qualities to be used as biomarkers or in which the signal is stable enough to be able to
257 give relevant feedback of its activation.

258

259 **2.1. Optimize indices of task-related reactivity**

260 The first possibility we consider involves optimizing indices for task related
261 reactivity in fMRI. This Blood Oxygen Level Dependent (BOLD) response is generally
262 considered to be convolution of the time-course of neural activity with a physiological
263 hemodynamic response. Mis-specification of the shape of BOLD reactivity can
264 introduce inefficiency and noise into estimates, which decreases reliability in human
265 (Handwerker et al., 2012; Lindquist et al., 2009; Shan et al., 2014) and animal models
266 (Peng et al., 2019). If, for example, neural responses to task stimuli are sustained in
267 depression rather than increased in amplitude (e.g., Mandell, Siegle, Shutt, Feldmiller,
268 & Thase, 2014), standard indices such as the amplitude of the canonical BOLD
269 response may not capture relevant aspects of the pathology.

270 Thus, we propose evaluating indices such as the average amplitude, area under
271 the curve and timing/shape of the curve of the BOLD response in addition to its
272 canonical amplitude. Gamma variate models, in particular, yield parameters for onset,
273 rise and fall slopes, and magnitude of hemodynamic responses (e.g., Larson et al.,
274 2006), which can be evaluated for reliability. Similarly, including temporal and
275 dispersion derivatives can account for individual differences in peak response timing
276 and small differences in HRF length, providing larger test-retest reliability values
277 (Fournier et al., 2014).

278

279 **2.2. Examine Regions with Voxel-Wise High Test-Retest Reliability**

280 When considering task-related reactivity in a region of interest (ROI), it is useful
281 to reduce voxelwise reactivity to a single or few indices which capture reactivity across
282 the region as a whole. The same consideration applies for reliability. Caceres et. al.
283 (2009) suggest computing the ICC in each voxel within a region of interest (ROI) and
284 reporting the median ICC as an index of region's test-retest reliability. This approach
285 has been applied practically (Fournier et al., 2014; Lois et al., 2018). However, several
286 potential biomarkers and neurofeedback targets identified in the literature, including
287 the amygdala (Lebow and Chen, 2016; Young et al., 2014) and the sgACC (Siegle et
288 al., 2012), consist of subregions with anatomical and functional heterogeneity
289 (Hrybouski et al., 2016; Palomero-Gallagher et al., 2019). Their reliability may not be
290 the same across these sub-divisions (Brabec et al., 2010; Janak and Tye, 2015;
291 LeDoux, 2012). Therefore, it is possible that only some parts of ROIs may have
292 adequate reliability and that the median reliability will not capture the most reliable
293 parts of the signal. Just as questionnaires are traditionally constructed by eliminating
294 unreliable items from an initial theoretically plausible set (Sheatsley, 1983), an index
295 that inherits solely from the reliable voxels may increase psychometric properties of
296 the preserved portions of regions.

297 Ten years ago, Bennet and Miller (2010) suggested that voxelwise reliability
298 constitutes the most rigorous criteria of reliability since it implies that the level of activity
299 in all voxels should remain consistent between scans. Although few studies have used
300 this approach, we contend the available psychometric arguments weight in favor of
301 voxel-wise computation of ICCs, restricting "reliable" ROIs to those regions in which all
302 voxels have good or excellent reliability.

303

304 **2.3. Optimize Models to Account for Individual and Clinical Features**

305 Minimizing sources of non-interest that could vary between administrations
306 increases the reliability of acquired data (Lin and Monica Way, 2014). Some fMRI noise
307 sources such as differences in instrumentation, time of day, motion, etc. can be
308 controlled, to some extent, via design. Tasks can be selected which have few practice
309 effects and pre-baseline training can remove practice and strategy-development
310 effects (Barch and Mathalon, 2011; Palmer et al., 2018). Choosing as-simple-as-
311 possible tasks can minimize the impact of non-task cognitive processes. Standardizing
312 instructions and training procedures helps to ensure participants understand the task
313 before the first administration (Barch and Mathalon, 2011). Effects of other time-
314 varying noise sources, such as thermal and physiological noise, are routinely
315 minimized via preprocessing procedures (Krüger and Glover, 2001).

316 That said, if sources of variation across time, such as physiological or cognitive
317 features, cannot be fully managed within design or processing, statistical methods for
318 adjusting test-retest reliability estimates for them (Atri et al., 2011; Hsiao et al., 2011;
319 Laenen et al., 2006) may be important to consider. Indeed, individual differences in
320 state anxiety can account for amygdala activation (Calder et al., 2011) and habituation
321 (Sladky et al., 2012), and, variation in rumination in depression is continuously
322 associated with individual differences in amygdala, hippocampal, and prefrontal
323 reactivity to emotional stimuli (Mandell et al., 2014; Siegle et al., 2002). Thus, true
324 signal differences due to anxiety, mood or other symptoms between scans, especially
325 if test-retest reliability is being evaluated in the context of possible treatment-related
326 effects, might account for apparently unreliable neural responses, particularly to
327 emotional stimuli. Thus, it may be useful to account for individuals' differences that
328 could change across time statistically in estimating reliability, e.g., via the inclusion of
329 clinical covariates.

330

331 **2.4. Examine Reliability Within Relevant Tasks and Clinical Populations**

332 Estimating reliability in healthy participants or symptomatic individuals who do
333 not receive intervention may help separate effects of symptom change from practice
334 effects. Yet, these approaches can introduce other confounds (e.g., if a task is reliable
335 in patients but not controls or non-treatment seeking symptomatic individuals). The
336 majority of studies have examined reliability of fMRI data in homogenous samples of
337 healthy, often young, university students (Bennett and Miller, 2010; Lois et al., 2018).

338 Studies reviewed in Table 1 that discuss reliability in MDD generally restrict their
339 discussion to whether there was a main effect of Time in healthy controls. Generally,
340 BOLD response variability is greater in forms of between-subject responses than within
341 (Aguirre et al., 1998). A limitation of ICC is that simultaneous inclusion of within and
342 between subject variability causes estimators to be affected by sample composition.
343 As groups might differ in the degree to which regional signals are reliable between
344 measurements (Fournier et al., 2014), and because ICCs are proportional to between
345 subject variability, heterogeneous samples can produce different ICCs even with the
346 same degree of within-subject reliability of test-retest values. Using only healthy control
347 participants may underrepresent true variability or over represent measurement errors
348 in the population of interest, yielding inaccurate reliability estimates. Similarly, non-
349 treatment seeking patients differ from treatment seeking patients on many variables
350 that could affect test-retest reliability, such as symptomatology and comorbidity
351 (Galbaud Du Fort et al., 1999).

352 Thus, testing reliability in the population of interest may provide more accurate
353 estimates. We therefore recommend the use of representative samples to create a
354 voxel-wise, population- and task-specific map of test-retest reliability. For example, if
355 a task is to be used to distinguish symptomatic from healthy individuals, this method
356 should be applied to a mixed population of healthy and symptomatic participants prior
357 to the clinical application of the task. If the purpose is to distinguish respondents and
358 non-respondents to a treatment, we recommend assessing reliability among
359 treatment-seeking patients.

360

361 **3. EVALUATION OF SUGGESTED OPTIMIZATIONS IN A PROGNOSTIC** 362 **NEUROIMAGING TREATMENT OUTCOME DATASET**

363 We have described several approaches that could be useful when examining
364 and seeking to improve test-retest reliability in service of clinical translation including
365 R1) optimizing BOLD signal parameterization, R2) using regions or voxels with
366 stronger psychometric properties, R3) accounting for within-individual changes and
367 R4) studying relevant tasks and populations for the intended application. In this section
368 we demonstrate feasibility of these approaches and examine whether they are useful
369 when applied to a published clinical fMRI dataset (Siegle et al., 2012). Our code for
370 these analyses is freely available from https://github.com/PICANlab/Reliability_toolbox
371 in the folder named “activation_task_reliability”.

372 **3.1. Method**

373 The sample consisted of participants described in Siegle et al. (2012)
374 augmented by the addition of 8 patients who completed the same protocol after that
375 paper was submitted, yielding 57 patients with major depressive disorder (MDD), and
376 35 healthy control participants (see supplement for details of this dataset and its
377 relationship to Siegle et al 2012). Briefly, participants with MDD completed a slow
378 event-related task during 3T fMRI in which they labeled the valence of emotional words
379 (here, as in the published dataset, we analyzed only nominally negative words) before
380 and after 12-16 weeks of Cognitive Therapy.

381 We computed reliability estimates within 4 ROIs which the literature suggests
382 may function as biomarkers for treatment response including the amygdala (Arnone et
383 al., 2012; Godlewska et al., 2012; Sheline et al., 2001), dorsolateral prefrontal cortex
384 (DLPFC, Koenigs and Grafman, 2009), rostral anterior cingulate cortex (rACC, Hunter
385 et al., 2013) and subgenual cingulate cortex (sgACC, Siegle et al., 2012b; Straub et
386 al., 2015; Taylor et al., 2018) (our region-wise definitions are included in Box 1 in
387 Supplement).

388 **3.1.1. Optimize the BOLD Signal**

389 The BOLD response to negative words was modeled within participants using
390 4 different methods including 1) amplitude of a canonically shaped BOLD signal (using
391 AFNI's 3dDeconvolve with a narrow tent function ('BLOCK5(1,1)', Cox, 1996), 2) Area
392 under the curve (via multiple regression of a delta function across 8 TRs using
393 3dDeconvolve, i.e. computed with Finite Impulse Response/FIR basis, with sum of
394 betas as the parameter retained); 3) Peak amplitude from the same regressions as #2,
395 and 4) a gamma variate model with parameters for onset-delay, rise-decay rate, and
396 height. Voxelwise outliers outside the Tukey hinges were windsorized across
397 participants and ICCs (3,1) were computed (Shrout and Fleiss, 1979) within individuals
398 for each modeling method using custom Matlab code. While ICC(2,1) allows
399 generalizing results obtained from different scanners, we chose to use ICC(3,1) to be
400 able to compare with most of the literature, given that it is the most widely used ICC.
401 This approach also allowed us to examine the importance of including scanner as a
402 covariate in 3.1.3.

403 **3.1.2. Compute Voxelwise Reliability**

404 To measure the benefit of identifying reliable voxels, we calculated the mean,
405 median and standard deviation of the ICCs in each of the ROIs for each modeling
406 method and each group.

407 **3.1.3. Include Clinical and Design Related Measures**

408 We examined whether indices of reliability increased when clinical and design-
409 related measures were included. As the ICC does not easily allow inclusion of
410 covariates, we used semi partial correlations within the context of multiple regressions
411 with and without covariates to assess changes in reliability, where covariates were pre
412 and post clinical measures, as:

$$413 \widehat{Post} = \beta_0 + \beta(1 \rightarrow n)covariates + \beta(n + 1)Pre$$

414 This model accounts for the potential that participants who show little change in
415 symptoms may have better test-retest reliability. Modelling these clinical effects at the
416 group level should make it possible to identify variance unique to test-retest reliability.

417 We included indices of pre- and post-treatment depressive symptomatology
418 (Beck Depression Inventory; BDI, Beck et al., 1996), state and trait anxiety
419 (Spielberger, 1983), rumination (Nolen-Hoeksema et al., 1993), and sleepiness
420 (Johns, 1991) administered on the scan day, the scanner on which data were acquired,
421 and participant's group when patients and controls were considered in one sample,
422 coded as dummy variables, as covariates. Missing data were imputed via regression
423 from the other administered measures also used as covariates.

424 A primary question was whether any of the proposed techniques described
425 above, including different BOLD models, accounting for voxelwise variability, and the
426 use of covariates, would differentially affect reliability estimates (i.e., semi-partial
427 correlations). As such, after computing reliability estimates at each voxel, we rank
428 ordered them across all permutations of BOLD estimate parameters (6 parameters)
429 and the use or non-use of covariates (2 conditions) at each voxel per ROI, yielding 12
430 x #-voxels rankings per ROI. Following a Kolmogorov-Smirnov test justifying the need
431 to use non-parametric tests, we report a Kruskal-Wallis test to determine whether the
432 rankings differed across models in each ROI. If they did, as a simple effects test, we
433 generated confidence intervals around the mean of rankings for each of the 12
434 conditions via a one-way ANOVA (via Matlab's multcompare function). Non-
435 overlapping confidence intervals are interpretable as significant differences between
436 one condition and any other. To display them we generated figures showing the mean

437 of rankings for each condition, which will be numbers on the order of 1 to 12 x # voxels,
438 with higher means representing being at the top of the rankings across many voxels
439 within the ROI.

440 **3.1.4. Use Clinically Representative Samples**

441 All analyses were conducted on the whole sample (controls and patients) to
442 establish likely reliability of tests that could be used to discriminate groups, and on
443 patients only, to establish likely reliability of clinical prognostic and change indicators.
444 We considered multiple reliability effect size thresholds which might be used in other
445 studies (0.4 and 0.6 for fair and good reliability and 0.7, and 0.75 for traditional labels
446 of the data as “reliable” and clinically meaningful).

447 **3.1.5. Type 1 error control**

448 As 1) each of the hypotheses and regions examined for this manuscript was
449 considered a different family of tests and 2) we want our results to generalize to
450 reliability as it is reported in the confirmatory biomarker and neurofeedback literatures
451 where only one region is generally examined, consistent with the literature on test-
452 retest reliability in neuroimaging, type I error was not controlled across regions and
453 hypotheses for ROI-wise statistics. For simple-effects tests of differences in rankings
454 across conditions, we controlled for the number of conditions with a Bonferroni test.
455 For voxelwise statistics we subjected all voxelwise residual maps to empirical cluster
456 thresholding (AFNI's 3dFWHMx and 3dClustSim, acf model with small-volume
457 corrections for examined regions) using a p threshold (-pthr) based on each considered
458 effect size threshold (see in supplementary materials, table S3 for more details).

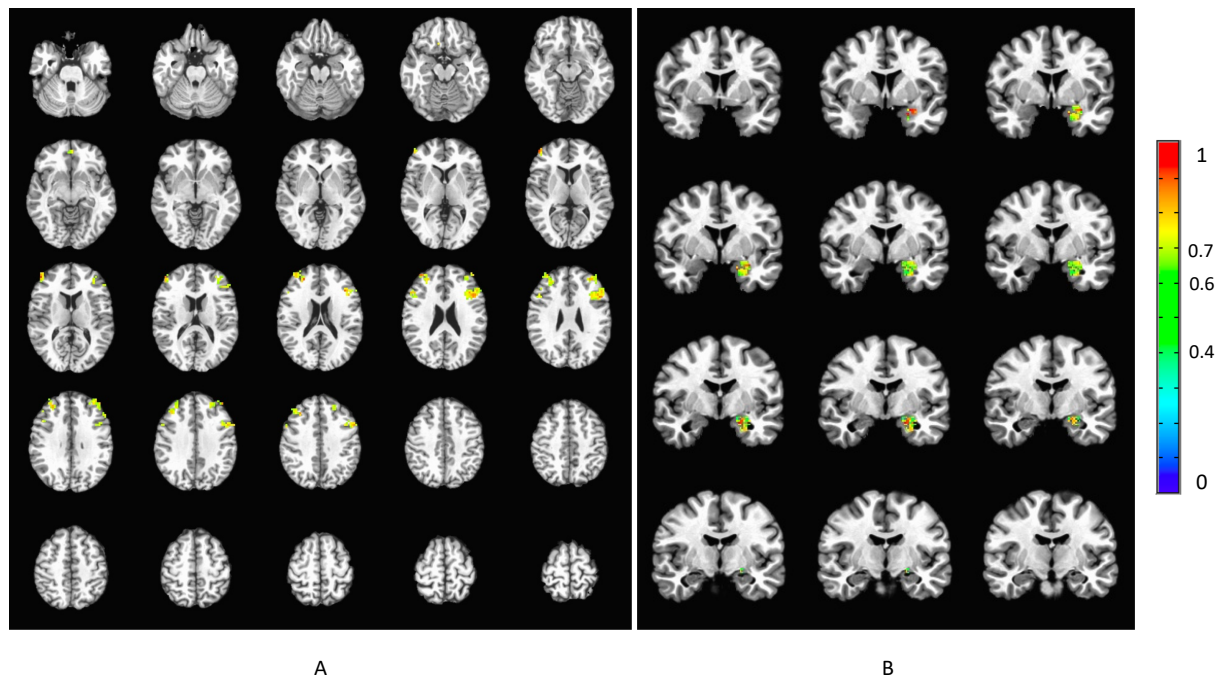
459

460 **3.2. Results and discussion**

461 **3.2.1. Optimizing the BOLD signal**

462 ICC's were uniformly low (<.3) for all BOLD parameterizations when entire ROIs
463 were considered (Table 4). Kruskal Wallis tests did suggest differential reliability across
464 our parameterizations (Table 5a). This held when the two outlying uniformly low
465 reliability parameterizations (rise decay with and without covariates) were removed
466 from consideration (Table 5b). Yet, there were non-overlapping confidence intervals
467 among counts of rank orderings of parameterizations for voxelwise tests, suggesting
468 that at least for some subsets of regions, some parameterizations were superior
469 (Supplement Figure S1, Table S1). For example, in the full sample, for the amygdala,
470 amplitude without covariates was superior to other parameters. Over all ROIs, the most

471 reliable parameters were amplitude, canonical amplitude, and height (Figure 1 shows
472 voxelwise variation within a Priori ROIs for the height parameter) for the whole sample
473 and amplitude, area under the curve, and height for only patients (Figure S1 and Table
474 S1). However, looking at ROIs and samples independently, the parameter offering the
475 highest levels of reliability varied.



476
477 **Figure 1: Test-retest reliability in ROIs estimated with voxel wise ICCs using**
478 **height parameter, a threshold of $ICC > 0.4$ and cluster correction applied for this**
479 **threshold in A. Siegle et al. (2012) dataset of patients and B. Young et al. (2017)**
480 **data set of the transfer run in the experimental group (signal with training)**
481 **preprocessed with the TBV style pipeline.**

482

483 3.2.2. Voxelwise reliability

484 In the whole sample, moderate reliability ($ICC > .4$) in clusters large enough to
485 infer significance was observed in the DLPFC using the canonical amplitude model
486 and in the amygdala using amplitude (Table 3). “Good” ($ICC > .6$) reliability was reached
487 in clusters large enough to infer significance when only the patients were considered,
488 using amplitude and height in the DLPFC. These levels of voxelwise test-retest
489 reliability were higher than using the median or mean value of ICCs within whole ROIs
490 (Table 4). Levels of generally accepted reliability for clinical measures ($ICC > .7$) were
491 not observed in clusters large enough to report.

492 **Table 3: Table of number of voxels reaching different reliability thresholds for**
 493 **each sample, first level parameter, and ROI with cluster correction applied.**

ROI		Amygdala (242 voxels)		DLPFC (2675 voxels)		rACC (865 voxels)		sgACC liberally thresholded (33 voxels)		sgACC conservatively thresholded (18 voxels)	
Population	Reactivity model	ICC thresholds		ICC thresholds		ICC thresholds		ICC thresholds		ICC thresholds	
		0.4	0.6	0.4	0.6	0.4	0.6	0.4	0.6	0.4	0.6
Controls & patients	Canonical amplitude	0	0	465	0	0	0	0	0	0	0
	Amplitude	66	0	5	0	0	0	0	0	0	0
	Area under the curve	10	0	0	0	0	0	0	0	0	0
	Onset delay	0	0	0	0	0	0	0	0	0	0
	Rise decay	0	0	0	0	0	0	0	0	0	0
	Height	0	0	290	2	0	0	0	0	0	0
	Patients	Canonical amplitude	0	0	299	6	6	0	0	0	0
Amplitude		24	0	0	0	0	0	0	0	0	0
Area under the curve		0	0	0	0	0	0	0	0	0	0
Onset delay		0	0	0	0	0	0	0	0	0	0
Rise decay		0	0	0	0	0	0	0	0	0	0
Height		0	0	374	5	5	0	2	0	1	0

494

495 **Table 4: Table of mean, standard deviation and median values of ICCs for each**
 496 **sample, reactivity model, and ROI.**

Population	Reactivity model	Amygdala	DLPFC	rACC	sgACC liberally thresholded	sgACC conservatively thresholded
Controls & patients	Canonical amplitude	0.11 (± 0.09); 0.11	0.24 (± 0.16);	0.09 (± 0.10);	0.15 (± 0.08); 0.13	0.17 (± 0.09); 0.18

		0.26	0.09		
Amplitude	0.23(±0.14);	0.12	0.11	-0.01 (±0.13);	-0.04 (±0.14); -
	0.22	(±0.11);	(±0.10);	-0.04	0.08
Area under the curve	0.13 (±0.14);	0.08	0.03	-0.03 (±0.09);	-0.06 (±0.10); -
	0.12	(±0.10);	(±0.11);	-0.04	0.07
Onset delay	0 (±0.09); -	0.01	0 (±0.10);	0 (±0.08);	-0.01 (±0.10); 0
	0.01	(±0.09); 0	0	0.01	
Rise decay	0 (±0); 0	0 (±0); 0	0 (±0); 0	0 (±0); 0	0 (±0); 0
Height	0.08 (±0.10);	0.21	0.13	0.16 (±0.12);	0.18 (±0.12);
	0.09	(±0.15);	(±0.12);	0.17	0.23
Canonical amplitude	0.09(±0.11);	0.22	0.08	0.10 (±0.12);	0.14 (±0.15);
	0.11	(±0.16);	(±0.14);	0.07	0.12
Amplitude	0.22 (±0.15);	0.11	0.10	-0.06 (±0.15);	-0.08 (±0.14); -
	0.22	(±0.13);	(±0.13);	-0.07	0.08
Area under the curve	0.13(±0.14);	0.6	0.03	-0.08 (±0.13);	-0.10 (±0.13); -
	0.12	(±0.12);	(±0.13);	-0.08	0.09
Onset delay	-0.01 (±0.12);	0.01	-0.01	0.02 (±0.11);	0.01 (±0.12);
	-0.01	(±0.12); 0	(±0.13); -	0.02	0.05
Rise decay	0 (±0); 0	0 (±0); 0	0 (±0); 0	0 (±0); 0	0 (±0); 0
Height	0.09 (±0.12);	0.22	0.12	0.16 (±0.17);	0.17 (±0.17);
	0.08	(±0.16);	(±0.15);	0.18	0.21
		0.23	0.13		

497 Mean (±standard deviation); median

498

499 **Table 5a: Table of Kruskal Wallis tests' output for each sample, reactivity model**
 500 **with and without covariates, and ROI with Bonferroni correction applied.**

Population	Amygdala	DLPFC	rACC	sgACC	sgACC
				liberally thresholded	conservatively thresholded
Controls & patients	H(11)=1414.67, p<0.001	H(11)=12717.07, p<0.001	H(11)=4794.14, p<0.001	H(11)=206.47, p<0.001	H(11)=118.32, p<0.001

Patients	H(11)=1233.13, p<0.001	H(11)=10371.75, p<0.001	H(11)=4477.55, p<0.001	H(11)=240.89, p<0.001	H(11)=136.93, p<0.001
-----------------	---------------------------	----------------------------	---------------------------	--------------------------	--------------------------

501 **Note: Applying Bonferroni correction for 6 reactivity models with and without**
 502 **covariates ($p < 0.05/12 = 0.004$).**

503

504 **Table 5b: Table of Kruskal Wallis tests' output for each sample, reactivity model**
 505 **with and without covariates, and ROI with Bonferroni correction applied, without**
 506 **rise decay.**

Population	Amygdala	DLPFC	rACC	sgACC liberally thresholded	sgACC conservatively thresholded
Controls & patients	H(9)=285.88, p<0.001	H(9)=4876.99, p<0.001	H(9)=644.19, p<0.001	H(9)=58.90, p<0.001	H(9)=40.55, p<0.001
Patients	H(9)=25.68, p=0.002	H(9)=1588.15, p<0.001	H(9)=190.42, p<0.001	H(9)=108.21, p<0.001	H(9)=67.20, p<0.001

507 **Note: Applying Bonferroni correction for 6 reactivity models with and without**
 508 **covariates ($p < 0.05/10 = 0.005$).**

509

510 **3.2.3. Clinical and Design Related Measures**

511 The addition of covariates never resulted in significantly higher average ranks
 512 for semi partial correlations in any ROI, in the whole sample or just the patients (Figure
 513 S1). In other words, adding covariates did not improve the reliability, and in some
 514 instances made it worse.

515

516 **4. EVALUATION OF SUGGESTED OPTIMIZATIONS IN AN EMPIRICAL** 517 **NEUROFEEDBACK DATASET**

518 To further support the feasibility of applying these recommendations and to
 519 evaluate the consistency of their performance in a second dataset, we consider a
 520 published fMRI neurofeedback dataset (Young, Siegle, et al., 2017, code available
 521 from https://github.com/PICANlab/Reliability_toolbox in the folder named “rtfMRI-
 522 nf_reliability”).

523

524 **4.1. Method**

525 This dataset constituted 18 patients in the experimental group who received
 526 amygdala neurofeedback and 16 patients in the control group who received parietal

527 neurofeedback. Briefly, participants completed two training scans on different days
528 within 2 weeks, each including a “baseline” and “transfer” runs during which no
529 feedback was presented. The analyzed task was a 40-second per block design during
530 which participants alternately rested, worked to upregulate a target region during recall
531 of positive memories, and did a distraction (counting) task (see supplement Box 5 for
532 details of this dataset). Here, we focus on a) the baseline data on the two training days
533 in control-feedback participants during recall of positive autobiographical memories
534 prior to neurofeedback training. As their amygdala signal did not change over the
535 course of the study at the group level (Young et al., 2017b), this allows us to examine
536 test-retest reliability of the left amygdala signal without the influence of neurofeedback.
537 b) the left amygdala signal during the two transfer runs in the experimental group, as
538 this represents the effect of neurofeedback training. Activity during the two post-
539 training transfer runs did not differ at the group level, allowing us to examine the test-
540 retest reliability of the amygdala signal after neurofeedback training. Because this
541 dataset only included patients with MDD, only the first 3 principles (i.e., optimization of
542 the BOLD signal, computation of voxelwise reliability, and inclusion of clinical and
543 design related measures) are evaluated in this dataset.

544 *Feedback signal*

545 To analyze the feedback signal averaged over the left amygdala we used the
546 output of the script used in Young, Siegle, et al. (2017) that allowed computation of the
547 feedback signal in real-time before considering the voxel-wise signal.

548 *Voxel-wise*

549 As rtfMRI-nf involves real-time preprocessing of the data, we sought to examine
550 whether this kind of preprocessing could affect the test-retest reliability of the signal.
551 We therefore performed data preprocessing emulating the real-time data processing
552 performed by the commercially available neurofeedback software Turbo BrainVoyager
553 (BrainVoyager, The Netherlands; henceforth “TBV style”) and a more classic
554 contemporary post-hoc preprocessing stream (here referred to as “standard
555 preprocessing”). Both streams were implemented using AFNI.

556 - TBV style preprocessing

557 Turbo BrainVoyager performs the following functions in real-time: 3D motion
558 correction, spatial smoothing, and drift removal via the design matrix. We used AFNI
559 to approximate these steps. After spatially transforming the anatomical then functionals
560 to the International Consortium for Brain Mapping 152 template, we then rescaled them

561 to conform to the Talairach atlas dimensions and then performed motion correction to
562 the first image, spatial smoothing 4mm FWHM smoothing kernel and fourth order
563 detrend for drift removal.

564 - Standard preprocessing

565 MRI pre-processing included despiking, volume registration and slice timing correction
566 for all EPI volumes in a given exam. After applying an intensity uniformity correction
567 on the anatomical, the anatomical was spatially transformed to the International
568 Consortium for Brain Mapping 152 template and rescaled to conform to the Talairach
569 atlas dimensions. Then, the fMRI data for each run were warped nonlinearly and the
570 same spatial transformations were applied. The fMRI run was spatially smoothed
571 within the grey matter mask using a Gaussian kernel with full width at half maximum
572 (FWHM) of 4 mm. A first standard GLM analysis was then applied separately for each
573 of the fMRI runs. The following regressors were included in the GLM model: six motion
574 parameters and their derivatives as nuisance covariates to take into account possible
575 artifacts caused by head motion, white matter and cerebrospinal fluid signals, and five
576 polynomial terms for modeling drift.

577 **4.1.1. Optimize the BOLD Signal**

578 **4.1.1.1. Amygdala signal**

579 From each participant's real-time left amygdala signal we calculated an
580 "amygdala signal" for each positive recall block minus the mean of the preceding rest
581 block from the output of previously used scripts for real-time preprocessing (Young et
582 al., 2017b), and recreated the feedback signal by taking the amount of activation at
583 every TR during the experimental condition minus the mean activation in the previous
584 rest condition, on the baseline run of control participants at visits 1 and 2 (signal without
585 training) and on the transfer run of experimental participants at visits 1 and 2 (signal
586 with training), independently. We then averaged the time course of the feedback signal
587 over all happy blocks. We summarized the activation for each participant for each visit
588 by either a mean of the amygdala signal or by fitting the time course with a gamma
589 variate model with parameters for onset-delay, rise-decay-rate, and height (see
590 Methodological choice to fit gamma variates in Supplement Box 2 for more information
591 of this methodological choice).

592 ICC(3,1) estimates were computed (Shrout and Fleiss, 1979) independently on the
593 estimates of the feedback signal with and without training.

594 **4.1.1.2. Voxelwise signal**

595 The same reactivity models as in the treatment outcome dataset were applied (see
596 part 3.2.1.1) to data preprocessed with both types of preprocessing but adapted to this
597 design (AFNI tent parameters to accommodate 40 s blocks as BLOCK(40,1), and area
598 under the curve across entire blocks).

599 **4.1.2. Compute Voxelwise Reliability**

600 As in the treatment outcome data set, to measure the benefit of identifying
601 reliable voxels, we calculated the mean, median and standard deviation of the ICCs in
602 the left amygdala for each model, group, and additionally for both preprocessing
603 pipelines.

604 **4.1.3. Include Clinical and Design Related Measures**

605 As in the treatment outcome data set, semi partial correlations were computed
606 with and without covariates. We included indices of depressive symptomatology (Beck
607 Depression Inventory; BDI, Beck et al., 1996), state and trait anxiety (Spielberger,
608 1983), sleepiness and drowsiness administered on the scan day, and the scanner on
609 which data were acquired coded as dummy variables, as covariates. There was no
610 missing data. We then compared the semi-partial correlations across all models of
611 individual responses with and without covariates for each group and preprocessing
612 pipeline as in section 3.1.3, to understand which models offered adequate test-retest
613 reliability and whether there were differences between them.

614 **4.1.4. Type 1 error control**

615 As discussed in section 3.1.5, cluster correction was applied on voxelwise
616 statistics (further details in supplement table S4).

617 **4.2. Results and discussion**

618 **4.2.1. Optimizing the BOLD signal**

619 **4.2.1.1. Amygdala signal**

620 The mean amygdala signals with and without training showed poor reliability
621 (ICCs<0.1). When the amygdala signal within the left amygdala was fit using a gamma
622 variate function, the onset-delay and height parameters showed fair reliability for the
623 signal without training (ICC=0.54 and ICC=0.47, respectively), with all other models,
624 including those with training, showing minimal reliability (ICC<.1). Therefore, it appears
625 that the shape of the signal without training is consistent across sessions and that the
626 signal in the left amygdala is more reliable when unchanged by training, which is
627 consistent with the assumption that training is changing the signal over time.

628 **4.2.1.2. Voxel-wise signal**

629 Kruskal Wallis tests suggested there were differences between the parameters
630 in reliability (Tables 6a and 6b). In particular, reliability for the height parameter (as well
631 as amplitude for the signal without training) was higher than for other parameters
632 (Figure S1). The height parameter also yielded a large enough cluster to infer
633 significance for “excellent” ($ICC > .7$) reliability in both samples (Table 7, Figure 1 for
634 illustration).

635 The use of the standard preprocessing stream had non-significantly-different
636 reliabilities from the stream emulating the real-time preprocessing run by Turbo
637 BrainVoyager over all parameters with or without covariates, with the exception of the
638 height parameter without covariates, which showed higher reliability with TBV style
639 preprocessing than with standard preprocessing in the signal without training (see
640 Figure S1).

641

642 **Table 6a: Table of Kruskal Wallis tests’ output for each sample, reactivity model**
643 **with and without covariates in the left amygdala with Bonferroni correction**
644 **applied.**

Population	Amygdala
Without training - Control - Baseline	H(23)=2964.56, p<0.001
With training - Experimental - Transfer	H(23)=3142.17, p<0.001

645 **Note: Applying Bonferroni correction for 6 reactivity models with and without**
646 **covariates ($p < 0.05/12 = 0.004$).**

647 **Table 6b: Table of Kruskal Wallis tests’ output for each sample, reactivity model**
648 **with and without covariates in the left amygdala with Bonferroni correction**
649 **applied, without rise decay.**

Population	Amygdala
Without training - Control - Baseline	H(19)=1397.84, p<0.001
With training - Experimental - Transfer	H(19)=1702.57, p<0.001

650 **Note: Applying Bonferroni correction for 6 reactivity models with and without**
651 **covariates ($p < 0.05/10 = 0.005$).**

652

653 **Table 7: Table of number of voxels reaching different reliability thresholds for**
 654 **each sample, preprocessing, and first level parameter with cluster correction**
 655 **applied.**

ROI		Amygdala (214 voxels)							
Preprocessing		BV style				Standard			
Population	First level model	ICC thresholds				ICC thresholds			
		0.4	0.6	0.7	0.75	0.4	0.6	0.7	0.75
Without training - Control - Baseline	Canonical amplitude	0	0	0	0	0	0	0	0
	Amplitude	52	16	6	2	35	0	0	0
	Area under the curve	0	0	0	0	40	0	0	0
	Onset-delay	0	0	0	0	0	0	0	0
	Rise-decay	0	0	0	0	0	0	0	0
	Height	78	26	13	13	53	24	9	5
With training - Experimental - Transfer	Canonical amplitude	0	0	0	0	0	0	0	0
	Amplitude	66	4	2	2	42	11	3	2
	Area under the curve	0	0	0	0	0	0	0	0
	Onset-delay	0	4	4	4	0	5	5	5
	Rise-decay	0	0	0	0	0	0	0	0
	Height	159	81	25	16	73	47	24	21

656

657 **4.2.2. Voxelwise reliability**

658 Some voxelwise ICC values obtained were higher than those computed on the
 659 real-time signal covering the entire left amygdala or mean or median ICC values
 660 computed over the entire left amygdala (Table 5 vs statistics reported in 4.2.1.1 and
 661 Table 6), with some clusters achieving an excellent level of reliability (ICC>.7, see
 662 Table 5) for standard and TBV-like preprocessing both for the trained and untrained
 663 signals, which did not occur for the region as a whole.

664

665 **Table 8: Table of mean, standard deviation and median values of ICCs for each**
 666 **sample, preprocessing, and first level parameter with cluster correction applied.**

Preprocessing	TBV style	Standard
---------------	-----------	----------

Without training - Control - Baseline	Canonical amplitude	-0.07 (± 0.21); -0.09	0.01 (± 0.24); 0
	Amplitude	0.29 (± 0.2); 0.3	0.26 (± 0.22); 0.27
	Area under the curve	0.02(± 0.21); 0.01	0.21 (± 0.23); 0.18
	Onset-delay	-0.03 (± 0.23); -0.05	-0.11 (± 0.20); -0.14
	Rise-decay	NA ($\pm NA$); NA	NA ($\pm NA$); NA
	Height	0.36 (± 0.23); 0.33	0.17 (± 0.38); 0.24
With training - Experimental - Transfer	Canonical amplitude	-0.11 (± 0.21); -0.12	0.08 (± 0.21); 0.09
	Amplitude	0.3 (± 0.18); 0.31	0.26 (± 0.21); 0.25
	Area under the curve	0.06 (± 0.20); 0.07	0.13 (± 0.18); 0.13
	Onset-delay	0.02 (± 0.24); -0.02	-0.05 (± 0.24); -0.13
	Rise-decay	NA ($\pm NA$); NA	NA ($\pm NA$); NA
	Height	0.52 (± 0.19); 0.56	0.35 (± 0.28); 0.34

667 Mean (\pm standard deviation); median

668

669 **4.2.3. Clinical and Design Related Measures**

670 **4.2.3.1. Amygdala signal**

671 Adding covariates when computing semi-partial correlations over the mean
 672 amygdala signal improved descriptively reliability estimates for the signal without
 673 training (from $sr=0.11$ to $sr=0.14$) as well as parameters tested to fit the signal with
 674 training (mean: from $sr=0.06$ to $sr=0.12$, onset-delay: from $sr=0.14$ to $sr=0.21$, rise-
 675 decay: from $sr=0.03$ to $sr=0.14$, height: from $sr=0.16$ to $sr=0.29$) although in no case
 676 we did achieve a fair level of reliability ($sr < 0.4$).

677 **4.2.3.2. Voxelwise signal**

678 The addition of covariates in never resulted in higher average ranks of
 679 semipartial correlation distributions on the untrained or trained signal preprocessed
 680 with the TBV-like or standard pipeline (Figure S1).

681

682

5. GENERAL DISCUSSION

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

As stated in a recent meta-analysis (Elliott et al., 2020), task fMRI reliability is not systematically evaluated and when it is, task-related fMRI measures show poor reliability. Our literature review shows that both prognostic and interventional fMRI studies in MDD, which might otherwise be poised for clinical translation, also do not attend to reliability. We demonstrate that by attending to some fairly simple principles, we can achieve fair to good reliability in a clinical prediction outcome dataset and excellent reliability in a neurofeedback fMRI study dataset (Figure 1). These principles include careful modeling of the BOLD signal, identification of reliable voxels within regions of interest, and calculation of reliability in the population for which translational applications are being considered. Across both datasets, the height parameter from a gamma variate function was the most reliable way to model the BOLD signal, especially among patients with MDD, in some regions of interest, and was, in some combinations of region and population or training condition, more reliable than canonical amplitude, though in other cases the reverse was true (Table 3 and 5 and Figure S1). Consequently, we recommend that researchers explore multiple ways of modeling the BOLD signal, particularly including gamma variate modeling in MDD, before concluding their experiment has low reliability. It may also be helpful for software for real-time analysis of fMRI data to implement alternative, potentially more reliable ways of characterizing BOLD responses in real-time.

Increasingly, functional differentiation of sub-regions of subcortical structures such as the amygdala has been acknowledged as important for fMRI (Balderston, Schultz, Hopkins, & Helmstetter, 2014; Ball et al., 2007; Michely, Rigoli, Rutledge, Hauser, & Dolan, 2020; Roy et al., 2009). The comparison of test-retest reliability estimates obtained on the feedback signal averaged over the whole amygdala versus these same estimates computed voxelwise in the neurofeedback dataset suggest non-uniformity across the amygdala in signal reliability as well; the extent to which these differences explain previous results localizing function to subregions is unclear. Thus, we suggest it may be useful to use a voxel-wise or subregion approach to estimating test-retest reliability. Indeed, this method reveals significantly large clusters of voxels with excellent test-retest reliability in the left amygdala which could be used as masks for neurofeedback targets; our method is easily feasible for new studies. Such excellent reliability, which is a prerequisite for clinical translation, was not attained in

715 our dataset, using the more common computation of median ICCs for each ROI (e.g.,
716 as recommended by Caceres et al., 2009) (see Tables 4 and 6).

717 Contrary to our hypotheses, we did not find that adding covariates to the model,
718 including the scanner on which participants were run and severity, which did change
719 as a function of intervention, improved test-retest reliability in these datasets (Figure
720 S1) in ROI-based or whole-brain analyses (see Figure S2). That said, covariates may
721 still be useful to include in other datasets – we recommend exploring this option further
722 before dismissing their utility.

723 Reliability did vary by whether the entire sample or only patient's data were
724 included and by whether or not participants were trained on the task, supporting the
725 potential utility of quantifying reliability on tasks and populations that are relevant for
726 the clinical application intended (Tables 3 and 5 and Figure S1).

727 There are several limitations of this review and analyses. As we have focused
728 only on MDD, it is unclear whether our conclusions apply transdiagnostically.
729 Improving reliability may require different strategies in other diseases, such as
730 Parkinsons, due to age-related atrophy, increased movement, and differences in
731 neurovascular coupling (Lecrux et al., 2019; Paek et al., 2019). There are many fMRI-
732 based metrics we could have examined, including functional connectivity, volumetric
733 measures, and resting state designs, which all provoke unique considerations for
734 optimizing test-retest reliability, some of which have been explored elsewhere (e.g.,
735 Noble et al., 2019). Here, we focused on regional BOLD activity as it is a common
736 feature of prediction and neurofeedback studies. Our published data sets had relatively
737 small number of subjects. This is typical for most clinical fMRI studies, but does raise
738 the concern that the sample is too small and underpowered. Therefore, we strongly
739 encourage the replication of these results and that is also why we have applied these
740 suggestions to two different data sets.

741

742

6. CONCLUSIONS

743 To summarize, demonstrating that mechanistic indices are reliable is important
744 before their clinical adoption in prediction or treatment-development. The literature in
745 these areas has implicitly accepted this assumption without testing it. Other non-clinical
746 fMRI studies have shown many of the regions targeted in clinical fMRI studies have
747 fairly low test-retest reliability, which was largely replicated using the most common
748 analytic techniques in our datasets. Yet, we have suggested a few principles that

749 appear to improve the test-retest reliability of the obtained mechanistic signals, have
750 shown their feasibility in two previously published fMRI data sets, and have made code
751 publicly available so that researchers with minimal mathematical and programming
752 knowledge can implement them. Wider adoption of these methods could help to realize
753 the potential of clinical fMRI and could extend to improving psychometrics for other
754 time-varying mechanistic indices.

755

756

Acknowledgements

757 This work was supported by the National Institutes of Health/ National Institute of
758 Mental Health [grant numbers MH115927, MH106591, MH074807, MH58356,
759 MH69618] and through the Pittsburgh Foundation Emmerling Fund [grant number
760 M2007-0114]. The content is solely the responsibility of the authors and does not
761 necessarily represent the official views of the funding agencies.

762

763

Summary declaration of interest

764

Declarations of interest: none.

765

766

767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800

REFERENCES

- Aguirre, G.K., Zarahn, E., D'esposito, M., 1998. The variability of human, BOLD hemodynamic responses. *Neuroimage*. <https://doi.org/10.1006/nimg.1998.0369>
- Anand, A., Li, Y., Wang, Y., Gardner, K., Lowe, M.J., 2007. Reciprocal Effects of Antidepressant Treatment on Activity and Connectivity of the Mood Regulating Circuit: An fMRI Study. *J. Neuropsychiatry Clin. Neurosci.* <https://doi.org/10.1176/jnp.2007.19.3.274>
- Andersson, J.L.R., Hutton, C., Ashburner, J., Turner, R., Friston, K., 2001. Modeling geometric deformations in EPI time series. *Neuroimage*. <https://doi.org/10.1006/nimg.2001.0746>
- Arnone, D., 2019. Functional MRI findings, pharmacological treatment in major depression and clinical response. *Prog. Neuro-Psychopharmacology Biol. Psychiatry*. <https://doi.org/10.1016/j.pnpbp.2018.08.004>
- Arnone, D., McKie, S., Elliott, R., Thomas, E.J., Downey, D., Juhasz, G., Williams, S.R., Deakin, J.F.W., Anderson, I.M., 2012. Increased amygdala responses to sad but not fearful faces in major depression: Relation to mood state and pharmacological treatment. *Am. J. Psychiatry*. <https://doi.org/10.1176/appi.ajp.2012.11121774>
- Atri, A., O'Brien, J.L., Sreenivasan, A., Rastegar, S., Salisbury, S., DeLuca, A.N., O'Keefe, K.M., LaViolette, P.S., Rentz, D.M., Locascio, J.J., Sperling, R.A., 2011. Test-retest reliability of memory task functional magnetic resonance imaging in alzheimer disease clinical trials. *Arch. Neurol.* <https://doi.org/10.1001/archneurol.2011.94>
- Balderston, N.L., Schultz, D.H., Hopkins, L., Helmstetter, F.J., 2014. Functionally distinct amygdala subregions identified using DTI and high-resolution fMRI. *Soc. Cogn. Affect. Neurosci.* 10, 1615–1622. <https://doi.org/10.1093/scan/nsv055>
- Ball, T., Rahm, B., Eickhoff, S.B., Schulze-Bonhage, A., Speck, O., Mutschler, I., 2007. Response Properties of Human Amygdala Subregions: Evidence Based on Functional MRI Combined with Probabilistic Anatomical Maps. *PLoS One* 2, e307. <https://doi.org/10.1371/journal.pone.0000307>
- Barch, D.M., Mathalon, D.H., 2011. Using brain imaging measures in studies of procognitive pharmacologic agents in schizophrenia: Psychometric and quality assurance considerations. *Biol. Psychiatry*.

- 801 <https://doi.org/10.1016/j.biopsych.2011.01.004>
- 802 Bartko, J.J., 1976. On various intraclass correlation reliability coefficients. *Psychol.*
803 *Bull.* <https://doi.org/10.1037/0033-2909.83.5.762>
- 804 Bartko, J.J., 1966. The Intraclass Correlation Coefficient as a Measure of Reliability.
805 *Psychol. Rep.* <https://doi.org/10.2466/pr0.1966.19.1.3>
- 806 Beck, A., Steer, R., Brown, G., 1996. Beck Depression Inventory II manual (2nd ed.
807 Ed.), ... for Beck Depression Inventory-II.
- 808 Benedetti, F., Radaelli, D., Bernasconi, A., Dallaspezia, S., Colombo, C., Smeraldi,
809 E., 2009. Changes in medial prefrontal cortex neural responses parallel
810 successful antidepressant combination of venlafaxine and light therapy. *Arch.*
811 *Ital. Biol.*
- 812 Bennett, C.M., Miller, M.B., 2013. fMRI reliability: Influences of task and
813 experimental design. *Cogn. Affect. Behav. Neurosci.*
814 <https://doi.org/10.3758/s13415-013-0195-1>
- 815 Bennett, C.M., Miller, M.B., 2010. How reliable are the results from functional
816 magnetic resonance imaging? *Ann. N. Y. Acad. Sci.*
817 <https://doi.org/10.1111/j.1749-6632.2010.05446.x>
- 818 Brabec, J., Rulseh, A., Hoyt, B., Vizek, M., Horinek, D., Hort, J., Petrovicky, P., 2010.
819 Volumetry of the human amygdala - An anatomical study. *Psychiatry Res. -*
820 *Neuroimaging.* <https://doi.org/10.1016/j.psychres.2009.11.005>
- 821 Braver, T.S., Cole, M.W., Yarkoni, T., 2010. Vive les differences! Individual variation
822 in neural mechanisms of executive control. *Curr. Opin. Neurobiol.*
823 <https://doi.org/10.1016/j.conb.2010.03.002>
- 824 Caceres, A., Hall, D.L., Zelaya, F.O., Williams, S.C.R., Mehta, M.A., 2009. Measuring
825 fMRI reliability with the intra-class correlation coefficient. *Neuroimage.*
826 <https://doi.org/10.1016/j.neuroimage.2008.12.035>
- 827 Calder, A.J., Ewbank, M., Passamonti, L., 2011. Personality influences the neural
828 responses to viewing facial expressions of emotion. *Philos. Trans. R. Soc. B*
829 *Biol. Sci.* <https://doi.org/10.1098/rstb.2010.0362>
- 830 Canli, T., Cooney, R.E., Goldin, P., Shah, M., Sivers, H., Thomason, M.E., Whitfield-
831 Gabrieli, S., Gabriels, J.D.E., Gotlib, I.H., 2005. Amygdala reactivity to emotional
832 faces predicts improvement in major depression. *Neuroreport.*
833 <https://doi.org/10.1097/01.wnr.0000174407.09515.cc>
- 834 Chen, C.H., Ridler, K., Suckling, J., Williams, S., Fu, C.H.Y., Merlo-Pich, E.,

- 835 Bullmore, E., 2007. Brain Imaging Correlates of Depressive Symptom Severity
836 and Predictors of Symptom Improvement After Antidepressant Treatment. *Biol.*
837 *Psychiatry*. <https://doi.org/10.1016/j.biopsych.2006.09.018>
- 838 Cicchetti, D. V., 1994. Guidelines, Criteria, and Rules of Thumb for Evaluating
839 Normed and Standardized Assessment Instruments in Psychology. *Psychol.*
840 *Assess*. <https://doi.org/10.1037/1040-3590.6.4.284>
- 841 Costafreda, S.G., Khanna, A., Mourao-Miranda, J., Fu, C.H.Y., 2009. Neural
842 correlates of sad faces predict clinical remission to cognitive behavioural therapy
843 in depression. *Neuroreport*. <https://doi.org/10.1097/WNR.0b013e3283294159>
- 844 Cox, R.W., 1996. AFNI: Software for analysis and visualization of functional magnetic
845 resonance neuroimages. *Comput. Biomed. Res*.
846 <https://doi.org/10.1006/cbmr.1996.0014>
- 847 Cullen, K.R., Klimes-Dougan, B., Vu, D.P., Westlund Schreiner, M., Mueller, B.A.,
848 Eberly, L.E., Camchong, J., Westervelt, A., Lim, K.O., 2016. Neural Correlates of
849 Antidepressant Treatment Response in Adolescents with Major Depressive
850 Disorder. *J. Child Adolesc. Psychopharmacol*.
851 <https://doi.org/10.1089/cap.2015.0232>
- 852 Davidson, R.J., Irwin, W., Anderle, M.J., Kalin, N.H., 2003. The neural substrates of
853 affective processing in depressed patients treated with venlafaxine. *Am. J.*
854 *Psychiatry*. <https://doi.org/10.1176/appi.ajp.160.1.64>
- 855 Decharms, R.C., 2008. Applications of real-time fMRI. *Nat. Rev. Neurosci*.
856 <https://doi.org/10.1038/nrn2414>
- 857 Delaveau, P., Jabourian, M., Lemogne, C., Allaïli, N., Choucha, W., Girault, N.,
858 Lehericy, S., Laredo, J., Fossati, P., 2016. Antidepressant short-term and long-
859 term brain effects during self-referential processing in major depression.
860 *Psychiatry Res. - Neuroimaging*.
861 <https://doi.org/10.1016/j.psychres.2015.11.007>
- 862 Dichter, G.S., Felder, J.N., Smoski, M.J., 2010. The effects of Brief Behavioral
863 Activation Therapy for Depression on cognitive control in affective contexts: An
864 fMRI investigation. *J. Affect. Disord*. <https://doi.org/10.1016/j.jad.2010.03.022>
- 865 Doerig, N., Krieger, T., Altenstein, D., Schlumpf, Y., Spinelli, S., Späti, J., Brakowski,
866 J., Quednow, B.B., Seifritz, E., Holtforth, M.G., 2016. Amygdala response to self-
867 critical stimuli and symptom improvement in psychotherapy for depression. *Br. J.*
868 *Psychiatry*. <https://doi.org/10.1192/bjp.bp.114.149971>

- 869 Elliott, M., Knodt, A., Ireland, D., Morris, M., Poulton, R., Ramrakha, S., Sison, M.,
870 Moffitt, T., Caspi, A., Hariri, A., 2020. What is the Test-Retest Reliability of
871 Common Task-fMRI Measures? New Empirical Evidence and a Meta-Analysis.
872 *Biol. Psychiatry*. <https://doi.org/10.1016/j.biopsych.2020.02.356>
- 873 Fales, C.L., Barch, D.M., Rundle, M.M., Mintun, M.A., Mathews, J., Snyder, A.Z.,
874 Sheline, Y.I., 2009. Antidepressant treatment normalizes hypoactivity in
875 dorsolateral prefrontal cortex during emotional interference processing in major
876 depression. *J. Affect. Disord.* <https://doi.org/10.1016/j.jad.2008.04.027>
- 877 Fang, J., Egorova, N., Rong, P., Liu, J., Hong, Y., Fan, Y., Wang, X., Wang, H., Yu,
878 Y., Ma, Y., Xu, C., Li, S., Zhao, J., Luo, M., Zhu, B., Kong, J., 2017. Early cortical
879 biomarkers of longitudinal transcutaneous vagus nerve stimulation treatment
880 success in depression. *NeuroImage Clin.*
881 <https://doi.org/10.1016/j.nicl.2016.12.016>
- 882 Fitzgerald, P.B., Sritharan, A., Daskalakis, Z.J., De Castella, A.R., Kulkarni, J., Egan,
883 G., 2007. A functional magnetic resonance imaging study of the effects of low
884 frequency right prefrontal transcranial magnetic stimulation in depression. *J.*
885 *Clin. Psychopharmacol.* <https://doi.org/10.1097/jcp.0b013e318151521c>
- 886 Fonseka, T.M., MacQueen, G.M., Kennedy, S.H., 2018. Neuroimaging biomarkers as
887 predictors of treatment outcome in Major Depressive Disorder. *J. Affect. Disord.*
888 <https://doi.org/10.1016/j.jad.2017.10.049>
- 889 Forbes, E.E., Olino, T.M., Ryan, N.D., Birmaher, B., Axelson, D., Moyles, D.L., Dahl,
890 R.E., 2010. Reward-related brain function as a predictor of treatment response
891 in adolescents with major depressive disorder. *Cogn. Affect. Behav. Neurosci.*
892 <https://doi.org/10.3758/CABN.10.1.107>
- 893 Fournier, J.C., Chase, H.W., Almeida, J., Phillips, M.L., 2014. Model specification
894 and the reliability of fMRI results: Implications for longitudinal neuroimaging
895 studies in psychiatry. *PLoS One*. <https://doi.org/10.1371/journal.pone.0105169>
- 896 Fovet, T., Jardri, R., Linden, D., 2015. Current Issues in the Use of fMRI-Based
897 Neurofeedback to Relieve Psychiatric Symptoms. *Curr. Pharm. Des.*
898 <https://doi.org/10.2174/1381612821666150619092540>
- 899 Frodl, T., Scheuerecker, J., Schoepf, V., Linn, J., Koutsouleris, N., Bokde, A.L.W.,
900 Hampel, H., Möller, H.J., Brückmann, H., Wiesmann, M., Meisenzahl, E., 2011.
901 Different effects of mirtazapine and venlafaxine on brain activation: An open
902 randomized controlled fMRI study. *J. Clin. Psychiatry*.

- 903 <https://doi.org/10.4088/JCP.09m05393blu>
- 904 Fu, C.H.Y., Costafreda, S.G., Sankar, A., Adams, T.M., Rasenick, M.M., Liu, P.,
905 Donati, R., Maglanoc, L.A., Horton, P., Marangell, L.B., 2015. Multimodal
906 functional and structural neuroimaging investigation of major depressive disorder
907 following treatment with duloxetine. *BMC Psychiatry*.
908 <https://doi.org/10.1186/s12888-015-0457-2>
- 909 Fu, C.H.Y., Steiner, H., Costafreda, S.G., 2013. Predictive neural biomarkers of
910 clinical response in depression: A meta-analysis of functional and structural
911 neuroimaging studies of pharmacological and psychological therapies.
912 *Neurobiol. Dis.* <https://doi.org/10.1016/j.nbd.2012.05.008>
- 913 Fu, C.H.Y., Williams, S.C.R., Brammer, M.J., Suckling, J., Kim, J., Cleare, A.J.,
914 Walsh, N.D., Mitterschiffthaler, M.T., Andrew, C.M., Pich, E.M., Bullmore, E.T.,
915 2007. Neural responses to happy facial expressions in major depression
916 following antidepressant treatment. *Am. J. Psychiatry*.
917 <https://doi.org/10.1176/ajp.2007.164.4.599>
- 918 Fu, C.H.Y., Williams, S.C.R., Cleare, A.J., Brammer, M.J., Walsh, N.D., Kim, J.,
919 Andrew, C.M., Pich, E.M., Williams, P.M., Reed, L.J., Mitterschiffthaler, M.T.,
920 Suckling, J., Bullmore, E.T., 2004. Attenuation of the neural response to sad
921 faces in major depression by antidepressant treatment: A prospective, event-
922 related functional magnetic resonance imaging study. *Arch. Gen. Psychiatry*.
923 <https://doi.org/10.1001/archpsyc.61.9.877>
- 924 Fu, C.H.Y., Williams, S.C.R., Cleare, A.J., Scott, J., Mitterschiffthaler, M.T., Walsh,
925 N.D., Donaldson, C., Suckling, J., Andrew, C., Steiner, H., Murray, R.M., 2008.
926 Neural Responses to Sad Facial Expressions in Major Depression Following
927 Cognitive Behavioral Therapy. *Biol. Psychiatry*.
928 <https://doi.org/10.1016/j.biopsych.2008.04.033>
- 929 Furey, M.L., Drevets, W.C., Hoffman, E.M., Frankel, E., Speer, A.M., Zarate, C.A.,
930 2013. Potential of pretreatment neural activity in the visual cortex during
931 emotional processing to predict treatment response to scopolamine in major
932 depressive disorder. *JAMA Psychiatry*.
933 <https://doi.org/10.1001/2013.jamapsychiatry.60>
- 934 Furey, M.L., Drevets, W.C., Szczepanik, J., Khanna, A., Nugent, A., Zarate, C.A.,
935 2015. Pretreatment differences in BOLD response to emotional faces correlate
936 with Antidepressant response to scopolamine. *Int. J. Neuropsychopharmacol.*

- 937 <https://doi.org/10.1093/ijnp/pyv028>
- 938 Galbaud Du Fort, G., Newman, S.C., Boothroyd, L.J., Bland, R.C., 1999. Treatment
939 seeking for depression: Role of depressive symptoms and comorbid psychiatric
940 diagnoses. *J. Affect. Disord.* [https://doi.org/10.1016/S0165-0327\(98\)00052-4](https://doi.org/10.1016/S0165-0327(98)00052-4)
- 941 Godlewska, B.R., Browning, M., Norbury, R., Cowen, P.J., Harmer, C.J., 2016. Early
942 changes in emotional processing as a marker of clinical response to SSRI
943 treatment in depression. *Transl. Psychiatry.* <https://doi.org/10.1038/tp.2016.130>
- 944 Godlewska, B.R., Browning, M., Norbury, R., Igoumenou, A., Cowen, P.J., Harmer,
945 C.J., 2018. Predicting Treatment Response in Depression: The Role of Anterior
946 Cingulate Cortex. *Int. J. Neuropsychopharmacol.*
947 <https://doi.org/10.1093/ijnp/pyy069>
- 948 Godlewska, B.R., Norbury, R., Selvaraj, S., Cowen, P.J., Harmer, C.J., 2012. Short-
949 term SSRI treatment normalises amygdala hyperactivity in depressed patients.
950 *Psychol. Med.* <https://doi.org/10.1017/S0033291712000591>
- 951 Guo, C.C., Kurth, F., Zhou, J., Mayer, E.A., Eickhoff, S.B., Kramer, J.H., Seeley,
952 W.W., 2012. One-year test-retest reliability of intrinsic connectivity network fMRI
953 in older adults. *Neuroimage.* <https://doi.org/10.1016/j.neuroimage.2012.03.027>
- 954 Gyurak, A., Patenaude, B., Korgaonkar, M.S., Grieve, S.M., Williams, L.M., Etkin, A.,
955 2016. Frontoparietal activation during response inhibition predicts remission to
956 antidepressants in patients with major depression. *Biol. Psychiatry.*
957 <https://doi.org/10.1016/j.biopsych.2015.02.037>
- 958 Hamilton, J.P., Glover, G.H., Bagarinao, E., Chang, C., Mackey, S., Sacchet, M.D.,
959 Gotlib, I.H., 2016. Effects of salience-network-node neurofeedback training on
960 affective biases in major depressive disorder. *Psychiatry Res. - Neuroimaging.*
961 <https://doi.org/10.1016/j.psychres.2016.01.016>
- 962 Handwerker, D.A., Gonzalez-Castillo, J., D'Esposito, M., Bandettini, P.A., 2012. The
963 continuing challenge of understanding and modeling hemodynamic variation in
964 fMRI. *Neuroimage.* <https://doi.org/10.1016/j.neuroimage.2012.02.015>
- 965 Hansen, N.S., Siegle, G., 2015. Paving the road to the neurocognitive clinic of
966 tomorrow: Appealing to standards, in: *From Symptom to Synapse: A
967 Neurocognitive Perspective on Clinical Psychology.*
968 <https://doi.org/10.4324/9780203507131>
- 969 Harrington, G.S., Tomaszewski Farias, S., Buonocore, M.H., Yonelinas, A.P., 2006.
970 The intersubject and intrasubject reproducibility of fMRI activation during three

- 971 encoding tasks: Implications for clinical applications. *Neuroradiology*.
972 <https://doi.org/10.1007/s00234-006-0083-2>
- 973 Heller, A.S., Johnstone, T., Light, S.N., Peterson, M.J., Kolden, G.G., Kalin, N.H.,
974 Davidson, R.J., 2013. Relationships between changes in sustained fronto-striatal
975 connectivity and positive affect in major depression resulting from antidepressant
976 treatment. *Am. J. Psychiatry*. <https://doi.org/10.1176/appi.ajp.2012.12010014>
- 977 Hrybouski, S., Aghamohammadi-Sereshki, A., Madan, C.R., Shafer, A.T., Baron,
978 C.A., Seres, P., Beaulieu, C., Olsen, F., Malykhin, N. V., 2016. Amygdala
979 subnuclei response and connectivity during emotional processing. *Neuroimage*.
980 <https://doi.org/10.1016/j.neuroimage.2016.02.056>
- 981 Hsiao, C.K., Chen, P.C., Kao, W.H., 2011. Bayesian random effects for interrater and
982 test-retest reliability with nested clinical observations. *J. Clin. Epidemiol*.
983 <https://doi.org/10.1016/j.jclinepi.2010.10.015>
- 984 Hunter, A.M., Korb, A.S., Cook, I.A., Leuchter, A.F., 2013. Rostral anterior cingulate
985 activity in major depressive disorder: State or trait marker of responsiveness to
986 medication? *J. Neuropsychiatry Clin. Neurosci*.
987 <https://doi.org/10.1176/appi.neuropsych.11110330>
- 988 Janak, P.H., Tye, K.M., 2015. From circuits to behaviour in the amygdala. *Nature*.
989 <https://doi.org/10.1038/nature14188>
- 990 Johns, M.W., 1991. A new method for measuring daytime sleepiness: The Epworth
991 sleepiness scale. *Sleep*. <https://doi.org/10.1093/sleep/14.6.540>
- 992 Keedwell, P.A., Drapier, D., Surguladze, S., Giampietro, V., Brammer, M., Phillips,
993 M., 2010. Subgenual cingulate and visual cortex responses to sad faces predict
994 clinical outcome during antidepressant treatment for depression. *J. Affect.*
995 *Disord*. <https://doi.org/10.1016/j.jad.2009.04.031>
- 996 Koenigs, M., Grafman, J., 2009. The functional neuroanatomy of depression: Distinct
997 roles for ventromedial and dorsolateral prefrontal cortex. *Behav. Brain Res*.
998 <https://doi.org/10.1016/j.bbr.2009.03.004>
- 999 Krüger, G., Glover, G.H., 2001. Physiological noise in oxygenation-sensitive
1000 magnetic resonance imaging. *Magn. Reson. Med*.
1001 <https://doi.org/10.1002/mrm.1240>
- 1002 Laenen, A., Vangeneugden, T., Geys, H., Molenberghs, G., 2006. Generalized
1003 reliability estimation using repeated measurements. *Br. J. Math. Stat. Psychol*.
1004 <https://doi.org/10.1348/000711005X66068>

- 1005 Lahey, M.A., Downey, R.G., Saal, F.E., 1983. Intraclass correlations: There's more
1006 there than meets the eye. *Psychol. Bull.* [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-2909.93.3.586)
1007 [2909.93.3.586](https://doi.org/10.1037/0033-2909.93.3.586)
- 1008 Langenecker, S.A., Kennedy, S.E., Guidotti, L.M., Briceno, E.M., Own, L.S., Hooven,
1009 T., Young, E.A., Akil, H., Noll, D.C., Zubieta, J.K., 2007. Frontal and Limbic
1010 Activation During Inhibitory Control Predicts Treatment Response in Major
1011 Depressive Disorder. *Biol. Psychiatry.*
1012 <https://doi.org/10.1016/j.biopsych.2007.02.019>
- 1013 Larson, C.L., Schaefer, H.S., Siegle, G.J., Jackson, C.A.B., Anderle, M.J., Davidson,
1014 R.J., 2006. Fear Is Fast in Phobic Individuals: Amygdala Activation in Response
1015 to Fear-Relevant Stimuli. *Biol. Psychiatry.*
1016 <https://doi.org/10.1016/j.biopsych.2006.03.079>
- 1017 Lebow, M.A., Chen, A., 2016. Overshadowed by the amygdala: The bed nucleus of
1018 the stria terminalis emerges as key to psychiatric disorders. *Mol. Psychiatry.*
1019 <https://doi.org/10.1038/mp.2016.1>
- 1020 Lecrux, C., Bourourou, M., Hamel, E., 2019. How reliable is cerebral blood flow to
1021 map changes in neuronal activity? *Auton. Neurosci. Basic Clin.*
1022 <https://doi.org/10.1016/j.autneu.2019.01.005>
- 1023 LeDoux, J., 2012. Rethinking the Emotional Brain. *Neuron.*
1024 <https://doi.org/10.1016/j.neuron.2012.02.004>
- 1025 Lemogne, C., Mayberg, H., Bergouignan, L., Volle, E., Delaveau, P., Lehericy, S.,
1026 Allilaire, J.F., Fossati, P., 2010. Self-referential processing and the prefrontal
1027 cortex over the course of depression: A pilot study. *J. Affect. Disord.*
1028 <https://doi.org/10.1016/j.jad.2009.11.003>
- 1029 Light, S.N., Heller, A.S., Johnstone, T., Kolden, G.G., Peterson, M.J., Kalin, N.H.,
1030 Davidson, R.J., 2011. Reduced right ventrolateral prefrontal cortex activity while
1031 inhibiting positive affect is associated with improvement in hedonic capacity after
1032 8 weeks of antidepressant treatment in major depressive disorder. *Biol.*
1033 *Psychiatry.* <https://doi.org/10.1016/j.biopsych.2011.06.031>
- 1034 Lin, A.L., Monica Way, H.Y., 2014. Functional Magnetic Resonance Imaging, in:
1035 *Pathobiology of Human Disease: A Dynamic Encyclopedia of Disease*
1036 *Mechanisms.* <https://doi.org/10.1016/B978-0-12-386456-7.07610-3>
- 1037 Linden, D.E.J., 2014. Neurofeedback and networks of depression. *Dialogues Clin.*
1038 *Neurosci.*

- 1039 Linden, D.E.J., Habes, I., Johnston, S.J., Linden, S., Tatineni, R., Subramanian, L.,
1040 Sorger, B., Healy, D., Goebel, R., 2012. Real-time self-regulation of emotion
1041 networks in patients with depression. PLoS One.
1042 <https://doi.org/10.1371/journal.pone.0038115>
- 1043 Lindquist, M.A., Meng Loh, J., Atlas, L.Y., Wager, T.D., 2009. Modeling the
1044 hemodynamic response function in fMRI: efficiency, bias and mis-modeling.
1045 Neuroimage. <https://doi.org/10.1016/j.neuroimage.2008.10.065>
- 1046 Lois, G., Kirsch, P., Sandner, M., Plichta, M.M., Wessa, M., 2018. Experimental and
1047 methodological factors affecting test-retest reliability of amygdala BOLD
1048 responses. Psychophysiology. <https://doi.org/10.1111/psyp.13220>
- 1049 López-Solà, M., Pujol, J., Hernández-Ribas, R., Harrison, B.J., Contreras-Rodríguez,
1050 O., Soriano-Mas, C., Deus, J., Ortiz, H., Menchón, J.M., Vallejo, J., Cardoner,
1051 N., 2010. Effects of duloxetine treatment on brain response to painful stimulation
1052 in major depressive disorder. Neuropsychopharmacology.
1053 <https://doi.org/10.1038/npp.2010.108>
- 1054 MacDuffie, K.E., MacInnes, J., Dickerson, K.C., Eddington, K.M., Strauman, T.J.,
1055 Adcock, R.A., 2018. Single session real-time fMRI neurofeedback has a lasting
1056 impact on cognitive behavioral therapy strategies. NeuroImage Clin.
1057 <https://doi.org/10.1016/j.nicl.2018.06.009>
- 1058 Mandell, D., Siegle, G.J., Shutt, L., Feldmiller, J., Thase, M.E., 2014. Neural
1059 substrates of trait ruminations in depression. J. Abnorm. Psychol. 123, 35–48.
1060 <https://doi.org/10.1037/a0035834>
- 1061 Mehler, D.M.A., Sokunbi, M.O., Habes, I., Barawi, K., Subramanian, L., Range, M.,
1062 Evans, J., Hood, K., Lührs, M., Keedwell, P., Goebel, R., Linden, D.E.J., 2018.
1063 Targeting the affective brain—a randomized controlled trial of real-time fMRI
1064 neurofeedback in patients with depression. Neuropsychopharmacology.
1065 <https://doi.org/10.1038/s41386-018-0126-5>
- 1066 Michely, J., Rigoli, F., Rutledge, R.B., Hauser, T.U., Dolan, R.J., 2020. Distinct
1067 Processing of Aversive Experience in Amygdala Subregions. Biol. Psychiatry
1068 Cogn. Neurosci. Neuroimaging 5, 291–300.
1069 <https://doi.org/10.1016/j.bpsc.2019.07.008>
- 1070 Miki, A., Raz, J., Van Erp, T.G.M., Liu, C.S.J., Haselgrove, J.C., Liu, G.T., 2000.
1071 Reproducibility of visual activation in functional MR imaging and effects of
1072 postprocessing. Am. J. Neuroradiol.

- 1073 Miller, J.M., Schneck, N., Siegle, G.J., Chen, Y., Ogden, R.T., Kikuchi, T., Oquendo,
1074 M.A., Mann, J.J., Parsey, R. V., 2013. fMRI response to negative words and
1075 SSRI treatment outcome in major depressive disorder: A preliminary study.
1076 *Psychiatry Res. - Neuroimaging*.
1077 <https://doi.org/10.1016/j.pscychresns.2013.08.001>
- 1078 Nichols, T.E., Das, S., Eickhoff, S.B., Evans, A.C., Glatard, T., Hanke, M.,
1079 Kriegeskorte, N., Milham, M.P., Poldrack, R.A., Poline, J.B., Proal, E., Thirion,
1080 B., Van Essen, D.C., White, T., Yeo, B.T.T., 2017. Best practices in data
1081 analysis and sharing in neuroimaging using MRI. *Nat. Neurosci*.
1082 <https://doi.org/10.1038/nn.4500>
- 1083 Noble, S., Scheinost, D., Constable, R.T., 2019. A decade of test-retest reliability of
1084 functional connectivity: A systematic review and meta-analysis. *Neuroimage*.
1085 <https://doi.org/10.1016/j.neuroimage.2019.116157>
- 1086 Nolen-Hoeksema, S., Morrow, J., Fredrickson, B.L., 1993. Response styles and the
1087 duration of episodes of depressed mood. *J. Abnorm. Psychol.* 102, 20–28.
1088 <https://doi.org/10.1037/0021-843X.102.1.20>
- 1089 Nord, C.L., Gray, A., Charpentier, C.J., Robinson, O.J., Roiser, J.P., 2017.
1090 Unreliability of putative fMRI biomarkers during emotional face processing.
1091 *Neuroimage*. <https://doi.org/10.1016/j.neuroimage.2017.05.024>
- 1092 Oakes, T.R., Johnstone, T., Ores Walsh, K.S., Greischar, L.L., Alexander, A.L., Fox,
1093 A.S., Davidson, R.J., 2005. Comparison of fMRI motion correction software
1094 tools. *Neuroimage*. <https://doi.org/10.1016/j.neuroimage.2005.05.058>
- 1095 Opmeer, E.M., Kortekaas, R., Van Tol, M.J., Renken, R.J., Demenescu, L.R.,
1096 Woudstra, S., Ter Horst, G.J., Van Buchem, M.A., Van Der Wee, N.J.A.,
1097 Veltman, D.J., Aleman, A., 2016. Changes in regional brain activation related to
1098 depressive state: A 2-year longitudinal functional MRI study. *Depress. Anxiety*.
1099 <https://doi.org/10.1002/da.22425>
- 1100 Paek, E.J., Murray, L.L., Newman, S.D., Kim, D.J., 2019. Test-retest reliability in an
1101 fMRI study of naming in dementia. *Brain Lang*.
1102 <https://doi.org/10.1016/j.bandl.2019.02.002>
- 1103 Palmer, C.E., Langbehn, D., Tabrizi, S.J., Papoutsis, M., 2018. Test-retest reliability of
1104 measures commonly used to measure striatal dysfunction across multiple testing
1105 sessions: A longitudinal study. *Front. Psychol*.
1106 <https://doi.org/10.3389/fpsyg.2017.02363>

- 1107 Palomero-Gallagher, N., Hoffstaedter, F., Mohlberg, H., Eickhoff, S.B., Amunts, K.,
1108 Zilles, K., 2019. Human Pregenual Anterior Cingulate Cortex: Structural,
1109 Functional, and Connectional Heterogeneity. *Cereb. Cortex*.
1110 <https://doi.org/10.1093/cercor/bhy124>
- 1111 Peng, S.-L., Chen, C.-M., Huang, C.-Y., Shih, C.-T., Huang, C.-W., Chiu, S.-C.,
1112 Shen, W.-C., 2019. Effects of Hemodynamic Response Function Selection on
1113 Rat fMRI Statistical Analyses. *Front. Neurosci*.
1114 <https://doi.org/10.3389/fnins.2019.00400>
- 1115 Phillips, M.L., Swartz, H.A., 2014. A Critical Appraisal of Neuroimaging Studies of
1116 Bipolar Disorder: Toward a New Conceptualization of Underlying Neural Circuitry
1117 and a Road Map for Future Research. *Am. J. Psychiatry* 171, 829–843.
1118 <https://doi.org/10.1176/appi.ajp.2014.13081008>
- 1119 Pickford, R.W., Guilford, J.P., 2007. Psychometric Methods. *Br. J. Sociol*.
1120 <https://doi.org/10.2307/586971>
- 1121 Plichta, M.M., Schwarz, A.J., Grimm, O., Morgen, K., Mier, D., Haddad, L., Gerdes,
1122 A.B.M., Sauer, C., Tost, H., Esslinger, C., Colman, P., Wilson, F., Kirsch, P.,
1123 Meyer-Lindenberg, A., 2012. Test-retest reliability of evoked BOLD signals from
1124 a cognitive-emotive fMRI test battery. *Neuroimage*.
1125 <https://doi.org/10.1016/j.neuroimage.2012.01.129>
- 1126 Poldrack, R.A., Baker, C.I., Durnez, J., Gorgolewski, K.J., Matthews, P.M., Munafò,
1127 M.R., Nichols, T.E., Poline, J.B., Vul, E., Yarkoni, T., 2017. Scanning the
1128 horizon: Towards transparent and reproducible neuroimaging research. *Nat.*
1129 *Rev. Neurosci.* <https://doi.org/10.1038/nrn.2016.167>
- 1130 Portney, L., Watkins, M., 2009. *Foundations of Clinical Research: Applications to*
1131 *Practice*, Upper Saddle River, NJ: Pearson/Prentice Hall.
- 1132 Ritchey, M., Dolcos, F., Eddington, K.M., Strauman, T.J., Cabeza, R., 2011. Neural
1133 correlates of emotional processing in depression: Changes with cognitive
1134 behavioral therapy and predictors of treatment response. *J. Psychiatr. Res.*
1135 <https://doi.org/10.1016/j.jpsychires.2010.09.007>
- 1136 Rizvi, S.J., Salomons, T. V., Konarski, J.Z., Downar, J., Jacobbe, P., McIntyre, R.S.,
1137 Kennedy, S.H., 2013. Neural response to emotional stimuli associated with
1138 successful antidepressant treatment and behavioral activation. *J. Affect. Disord.*
1139 <https://doi.org/10.1016/j.jad.2013.06.050>
- 1140 Robertson, B., Wang, L., Diaz, M.T., Aiello, M., Gersing, K., Beyer, J., Mukundan, S.,

- 1141 McCarthy, G., Doraiswamy, P.M., 2007. Effect of bupropion extended release on
1142 negative emotion processing in major depressive disorder: A pilot functional
1143 magnetic resonance imaging study. *J. Clin. Psychiatry*.
1144 <https://doi.org/10.4088/JCP.v68n0212>
- 1145 Rosenblau, G., Sterzer, P., Stoy, M., Park, S., Friedel, E., Heinz, A., Pilhatsch, M.,
1146 Bauer, M., Ströhle, A., 2012. Functional neuroanatomy of emotion processing in
1147 major depressive disorder is altered after successful antidepressant therapy. *J.*
1148 *Psychopharmacol.* <https://doi.org/10.1177/0269881112450779>
- 1149 Roy, A.K., Shehzad, Z., Margulies, D.S., Kelly, A.M.C., Uddin, L.Q., Gotimer, K.,
1150 Biswal, B.B., Castellanos, F.X., Milham, M.P., 2009. Functional connectivity of
1151 the human amygdala using resting state fMRI. *Neuroimage*.
1152 <https://doi.org/10.1016/j.neuroimage.2008.11.030>
- 1153 Roy, M., Harvey, P.O., Berlim, M.T., Mamdani, F., Beaulieu, M.M., Turecki, G.,
1154 Lepage, M., 2010. Medial prefrontal cortex activity during memory encoding of
1155 pictures and its relation to symptomatic improvement after citalopram treatment
1156 in patients with major depression. *J. Psychiatry Neurosci.*
1157 <https://doi.org/10.1503/jpn.090010>
- 1158 Rubin-Falcone, H., Weber, J., Kishon, R., Ochsner, K., Delaparte, L., Doré, B.,
1159 Zanderigo, F., Oquendo, M.A., Mann, J.J., Miller, J.M., 2018. Longitudinal effects
1160 of cognitive behavioral therapy for depression on the neural correlates of
1161 emotion regulation. *Psychiatry Res. - Neuroimaging*.
1162 <https://doi.org/10.1016/j.psychres.2017.11.002>
- 1163 Ruhé, H.G., Booij, J., Veltman, D.J., Michel, M.C., Schene, A.H., 2012. Successful
1164 pharmacologic treatment of major depressive disorder attenuates amygdala
1165 activation to negative facial expressions: A functional magnetic resonance
1166 imaging study. *J. Clin. Psychiatry*. <https://doi.org/10.4088/JCP.10m06584>
- 1167 Ruiz, S., Buyukturkoglu, K., Rana, M., Birbaumer, N., Sitaram, R., 2014. Real-time
1168 fMRI brain computer interfaces: Self-regulation of single brain regions to
1169 networks. *Biol. Psychol.* <https://doi.org/10.1016/j.biopsycho.2013.04.010>
- 1170 Samson, A.C., Meisenzahl, E., Scheuerecker, J., Rose, E., Schoepf, V., Wiesmann,
1171 M., Frodl, T., 2011. Brain activation predicts treatment improvement in patients
1172 with major depressive disorder. *J. Psychiatr. Res.*
1173 <https://doi.org/10.1016/j.jpsychires.2011.03.009>
- 1174 Sankar, A., Adams, T.M., Costafreda, S.G., Marangell, L.B., Fu, C.H.Y., 2017.

- 1175 Effects of antidepressant therapy on neural components of verbal working
1176 memory in depression. *J. Psychopharmacol.*
1177 <https://doi.org/10.1177/0269881117724594>
- 1178 Schaefer, H.S., Putnam, K.M., Benca, R.M., Davidson, R.J., 2006. Event-Related
1179 Functional Magnetic Resonance Imaging Measures of Neural Activity to Positive
1180 Social Stimuli in Pre- and Post-Treatment Depression. *Biol. Psychiatry.*
1181 <https://doi.org/10.1016/j.biopsych.2006.03.024>
- 1182 Shan, Z.Y., Wright, M.J., Thompson, P.M., McMahon, K.L., Blokland, G.G.A.M., De
1183 Zubicaray, G.I., Martin, N.G., Vinkhuyzen, A.A.E., Reutens, D.C., 2014.
1184 Modeling of the hemodynamic responses in block design fMRI studies. *J. Cereb.*
1185 *Blood Flow Metab.* <https://doi.org/10.1038/jcbfm.2013.200>
- 1186 Sheatsley, P.B., 1983. Questionnaire Construction and Item Writing, in: *Handbook of*
1187 *Survey Research.* <https://doi.org/10.1016/b978-0-12-598226-9.50012-4>
- 1188 Sheline, Y.I., Barch, D.M., Donnelly, J.M., Ollinger, J.M., Snyder, A.Z., Mintun, M.A.,
1189 2001. Increased amygdala response to masked emotional faces in depressed
1190 subjects resolves with antidepressant treatment: An fMRI study. *Biol. Psychiatry.*
1191 [https://doi.org/10.1016/S0006-3223\(01\)01263-X](https://doi.org/10.1016/S0006-3223(01)01263-X)
- 1192 Shrout, P.E., Fleiss, J.L., 1979. Intraclass correlations: Uses in assessing rater
1193 reliability. *Psychol. Bull.* <https://doi.org/10.1037/0033-2909.86.2.420>
- 1194 Siegle, G.J., Carter, C.S., Thase, M.E., 2006. Use of fMRI to predict recovery from
1195 unipolar depression with cognitive behavior therapy. *Am. J. Psychiatry.*
1196 <https://doi.org/10.1176/ajp.2006.163.4.735>
- 1197 Siegle, G.J., Steinhauer, S.R., Thase, M.E., Stenger, V.A., Carter, C.S., 2002. Can't
1198 shake that feeling: Event-related fMRI assessment of sustained amygdala
1199 activity in response to emotional information in depressed individuals. *Biol.*
1200 *Psychiatry.* <https://doi.org/10.4324/9780203706978>
- 1201 Siegle, G.J., Thompson, W.K., Collier, A., Berman, S.R., Feldmiller, J., Thase, M.E.,
1202 Friedman, E.S., 2012. Toward Clinically Useful Neuroimaging in Depression
1203 Treatment. *Arch. Gen. Psychiatry.*
1204 <https://doi.org/10.1001/archgenpsychiatry.2012.65>
- 1205 Sladky, R., Höflich, A., Atanelov, J., Kraus, C., Baldinger, P., Moser, E.,
1206 Lanzenberger, R., Windischberger, C., 2012. Increased Neural Habituation in
1207 the Amygdala and Orbitofrontal Cortex in Social Anxiety Disorder Revealed by
1208 fMRI. *PLoS One.* <https://doi.org/10.1371/journal.pone.0050050>

- 1209 Spielberger, C., 1983. Manual for the State-Trait Anxiety Inventory (STAI). Consult.
1210 Psychol. Press 4–26. https://doi.org/10.1007/978-0-387-78665-0_6709
- 1211 Spies, M., Kraus, C., Geissberger, N., Auer, B., Klöbl, M., Tik, M., Störkat, I.L., Hahn,
1212 A., Woletz, M., Pfabigan, D.M., Kasper, S., Lamm, C., Windischberger, C.,
1213 Lanzenberger, R., 2017. Default mode network deactivation during emotion
1214 processing predicts early antidepressant response. *Transl. Psychiatry*.
1215 <https://doi.org/10.1038/tp.2016.265>
- 1216 Stoy, M., Schlagenhaut, F., Sterzer, P., BERPohl, F., Hägele, C., Suchotzki, K.,
1217 Schmack, K., Wrase, J., Ricken, R., Knutson, B., Adli, M., Bauer, M., Heinz, A.,
1218 Ströhle, A., 2012. Hyporeactivity of ventral striatum towards incentive stimuli in
1219 unmedicated depressed patients normalizes after treatment with escitalopram. *J.*
1220 *Psychopharmacol.* <https://doi.org/10.1177/02698811111416686>
- 1221 Straub, J., Plener, P.L., Sproeber, N., Sprenger, L., Koelch, M.G., Groen, G., Abler,
1222 B., 2015. Neural correlates of successful psychotherapy of depression in
1223 adolescents. *J. Affect. Disord.* <https://doi.org/10.1016/j.jad.2015.05.020>
- 1224 Strege, M. V, Siegle, G.J., Young, K., 2020. Cingulate prediction of response to
1225 antidepressant and cognitive behavioral therapies for depression: Theory, meta-
1226 analysis, and empirical application. *bioRxiv* 2020.12.02.407841.
1227 <https://doi.org/10.1101/2020.12.02.407841>
- 1228 Szczepanik, J., Nugent, A.C., Drevets, W.C., Khanna, A., Zarate, C.A., Furey, M.L.,
1229 2016. Amygdala response to explicit sad face stimuli at baseline predicts
1230 antidepressant treatment response to scopolamine in major depressive disorder.
1231 *Psychiatry Res. - Neuroimaging*.
1232 <https://doi.org/10.1016/j.pscychresns.2016.06.005>
- 1233 Tao, R., Calley, C.S., Hart, J., Mayes, T.L., Nakonezny, P.A., Lu, H., Kennard, B.D.,
1234 Tamminga, C.A., Emslie, G.J., 2012. Brain activity in adolescent major
1235 depressive disorder before and after fluoxetine treatment. *Am. J. Psychiatry*.
1236 <https://doi.org/10.1176/appi.ajp.2011.11040615>
- 1237 Taylor, S.F., Ho, S.S., Abagis, T., Angstadt, M., Maixner, D.F., Welsh, R.C.,
1238 Hernandez-Garcia, L., 2018. Changes in brain connectivity during a sham-
1239 controlled, transcranial magnetic stimulation trial for depression. *J. Affect.*
1240 *Disord.* <https://doi.org/10.1016/j.jad.2018.02.019>
- 1241 Thibault, R.T., MacPherson, A., Lifshitz, M., Roth, R.R., Raz, A., 2018.
1242 Neurofeedback with fMRI: A critical systematic review. *Neuroimage*.

- 1243 <https://doi.org/10.1016/j.neuroimage.2017.12.071>
- 1244 Toki, S., Okamoto, Y., Onoda, K., Matsumoto, T., Yoshimura, S., Kunisato, Y.,
1245 Okada, G., Shishida, K., Kobayakawa, M., Fukumoto, T., Machino, A., Inagaki,
1246 M., Yamawaki, S., 2014. Hippocampal activation during associative encoding of
1247 word pairs and its relation to symptomatic improvement in depression: A
1248 functional and volumetric MRI study. *J. Affect. Disord.*
1249 <https://doi.org/10.1016/j.jad.2013.07.021>
- 1250 Victor, T.A., Furey, M.L., Fromm, S.J., Öhman, A., Drevets, W.C., 2013. Changes in
1251 the neural correlates of implicit emotional face processing during antidepressant
1252 treatment in major depressive disorder, in: *International Journal of*
1253 *Neuropsychopharmacology*. <https://doi.org/10.1017/S146114571300062X>
- 1254 Victor, T.A., Furey, M.L., Fromm, S.J., Öhman, A., Drevets, W.C., 2010. Relationship
1255 between amygdala responses to masked faces and mood state and treatment in
1256 major depressive disorder. *Arch. Gen. Psychiatry*.
1257 <https://doi.org/10.1001/archgenpsychiatry.2010.144>
- 1258 Vul, E., Harris, C., Winkielman, P., Pashler, H., 2009. Puzzlingly High Correlations in
1259 fMRI Studies of Emotion, Personality, and Social Cognition (a.k.a. Voodoo
1260 Correlations in Social Neuroscience). *Perspect. Psychol. Sci.*
- 1261 Wagner, G., Koch, K., Schachtzabel, C., Sobanski, T., Reichenbach, J.R., Sauer, H.,
1262 Schlösser, R.G.M., 2010. Differential effects of serotonergic and noradrenergic
1263 antidepressants on brain activity during a cognitive control task and
1264 neurofunctional prediction of treatment outcome in patients with depression. *J.*
1265 *Psychiatry Neurosci*. <https://doi.org/10.1503/jpn.090081>
- 1266 Walsh, N.D., Williams, S.C.R., Brammer, M.J., Bullmore, E.T., Kim, J., Suckling, J.,
1267 Mitterschiffthaler, M.T., Cleare, A.J., Pich, E.M., Mehta, M.A., Fu, C.H.Y., 2007.
1268 A Longitudinal Functional Magnetic Resonance Imaging Study of Verbal
1269 Working Memory in Depression After Antidepressant Therapy. *Biol. Psychiatry*.
1270 <https://doi.org/10.1016/j.biopsych.2006.12.022>
- 1271 Wang, Y., Xu, C., Cao, X., Gao, Q., Li, J., Liu, Z., Sun, N., Ren, Y., Zhang, K., 2012.
1272 Effects of an antidepressant on neural correlates of emotional processing in
1273 patients with major depression. *Neurosci. Lett.*
1274 <https://doi.org/10.1016/j.neulet.2012.08.034>
- 1275 Wessa, M., Lois, G., 2015. Brain Functional Effects of Psychopharmacological
1276 Treatment in Major Depression: a Focus on Neural Circuitry of Affective

- 1277 Processing. *Curr. Neuropharmacol.*
1278 <https://doi.org/10.2174/1570159x13666150416224801>
- 1279 Williams, L.M., Korgaonkar, M.S., Song, Y.C., Paton, R., Eagles, S., Goldstein-
1280 Piekarski, A., Grieve, S.M., Harris, A.W.F., Usherwood, T., Etkin, A., 2015.
1281 Amygdala Reactivity to Emotional Faces in the Prediction of General and
1282 Medication-Specific Responses to Antidepressant Treatment in the Randomized
1283 iSPOT-D Trial. *Neuropsychopharmacology*. <https://doi.org/10.1038/npp.2015.89>
- 1284 Yoshimura, S., Okamoto, Y., Onoda, K., Matsunaga, M., Okada, G., Kunisato, Y.,
1285 Yoshino, A., Ueda, K., Suzuki, S. ichi, Yamawaki, S., 2014. Cognitive behavioral
1286 therapy for depression changes medial prefrontal and ventral anterior cingulate
1287 cortex activity associated with self-referential processing. *Soc. Cogn. Affect.*
1288 *Neurosci.* <https://doi.org/10.1093/scan/nst009>
- 1289 Young, K.D., Misaki, M., Harmer, C.J., Victor, T., Zotev, V., Phillips, R., Siegle, G.J.,
1290 Drevets, W.C., Bodurka, J., 2017a. Real-Time fMRI Amygdala Neurofeedback
1291 Changes Positive Information Processing in Major Depressive Disorder. *Biol.*
1292 *Psychiatry*. <https://doi.org/10.1016/j.biopsych.2017.03.013>
- 1293 Young, K.D., Siegle, G.J., Misaki, M., Zotev, V., Phillips, R., Drevets, W.C., Bodurka,
1294 J., 2018. Altered task-based and resting-state amygdala functional connectivity
1295 following real-time fMRI amygdala neurofeedback training in major depressive
1296 disorder. *NeuroImage Clin.* <https://doi.org/10.1016/j.nicl.2017.12.004>
- 1297 Young, K.D., Siegle, G.J., Zotev, V., Phillips, R., Misaki, M., Yuan, H., Drevets, W.C.,
1298 Bodurka, J., 2017b. Randomized clinical trial of real-time fMRI amygdala
1299 neurofeedback for major depressive disorder: Effect on symptoms and
1300 autobiographical memory recall, in: *American Journal of Psychiatry*.
1301 <https://doi.org/10.1176/appi.ajp.2017.16060637>
- 1302 Young, K.D., Zotev, V., Phillips, R., Misaki, M., Yuan, H., Drevets, W.C., Bodurka, J.,
1303 2014. Real-time fMRI neurofeedback training of amygdala activity in patients
1304 with major depressive disorder. *PLoS One*.
1305 <https://doi.org/10.1371/journal.pone.0088785>
- 1306 Yuan, H., Young, K.D., Phillips, R., Zotev, V., Misaki, M., Bodurka, J., 2014. Resting-
1307 State Functional Connectivity Modulation and Sustained Changes After Real-
1308 Time Functional Magnetic Resonance Imaging Neurofeedback Training in
1309 Depression. *Brain Connect.* <https://doi.org/10.1089/brain.2014.0262>
- 1310 Zandbelt, B.B., Gladwin, T.E., Raemaekers, M., van Buuren, M., Neggers, S.F.,

- 1311 Kahn, R.S., Ramsey, N.F., Vink, M., 2008. Within-subject variation in BOLD-
1312 fMRI signal changes across repeated measurements: Quantification and
1313 implications for sample size. *Neuroimage*.
1314 <https://doi.org/10.1016/j.neuroimage.2008.04.183>
- 1315 Zhilkin, P., Alexander, M.E., 2004. Affine registration: A comparison of several
1316 programs. *Magn. Reson. Imaging*. <https://doi.org/10.1016/j.mri.2003.05.004>
- 1317 Zotev, V., Phillips, R., Yuan, H., Misaki, M., Bodurka, J., 2014. Self-regulation of
1318 human brain activity using simultaneous real-time fMRI and EEG
1319 neurofeedback. *Neuroimage*. <https://doi.org/10.1016/j.neuroimage.2013.04.126>
- 1320 Zotev, V., Yuan, H., Misaki, M., Phillips, R., Young, K.D., Feldner, M.T., Bodurka, J.,
1321 2016. Correlation between amygdala BOLD activity and frontal EEG asymmetry
1322 during real-time fMRI neurofeedback training in patients with depression.
1323 *NeuroImage Clin*. <https://doi.org/10.1016/j.nicl.2016.02.003>
1324