

Quantifying Dysregulation of fMRI-Derived Control Circuits for Computational Psychiatry

Syed Sultan

State University of New York at Stony Brook

Steve Skiena

Stony Brook University

Lilianne Mujica-Parodi (✉ mujica@lcneuro.org)

State University of New York at Stony Brook

Article

Keywords: brain, circuit, dysregulation, fMRI, control system, trajectory, computational psychiatry, generative model

Posted Date: November 7th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1413254/v2>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Quantifying Dysregulation of fMRI-Derived Control Circuits for Computational Psychiatry

Syed Fahad Sultan¹, Steven Skiena², and Lilianne R. Mujica-Parodi^{3,4}✉

¹Computer Science Department, Furman University, Greenville, SC 29613

²Computer Science Department, Stony Brook University, Stony Brook, NY 11794

³Biomedical Engineering Department, Stony Brook University, Stony Brook, NY 11794

⁴Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital and Harvard Medical School, Charlestown, MA 02129

1 Psychiatric disorders are thought to result from *dysregulated*
2 *brain circuits*, yet human neuroimaging currently lacks stan-
3 *dardized methods for quantifying neural dysregulation*. Here,
4 *we present a scalable framework for extracting fMRI-derived*
5 *(generative) control circuits, then use circuit trajectories to es-*
6 *timate their control error*. Using synthetic circuits, we first
7 *demonstrate that our framework accurately identifies each cir-*
8 *cuit’s architecture and models its dynamics by estimation of*
9 *transfer functions*. As a use case, we then apply the frame-
10 *work to human task-based functional MRI data (UK Biobank,*
11 *N=19,831)*. In a purely data-driven manner, without priors,
12 *our framework identified thalamus-linked prefrontal-limbic and*
13 *ventral stream subcircuits, selectively engaged during sensori-*
14 *motor processing of affective and non-affective stimuli*. Finally,
15 *we demonstrate that circuit-wide dysregulation, defined by de-*
16 *gree of drift from healthy trajectories, tracks symptom sever-*
17 *ity for neuroticism (ventral subcircuit), depression (prefrontal-*
18 *limbic subcircuit), and bipolar disorder (full circuit)*.

19 brain | circuit | dysregulation | fMRI | control system | trajectory | computa-
20 tional psychiatry | generative model

21 Correspondence: mujica@lcneuro.org

22 Introduction

23 Psychiatric disorders are commonly understood to reflect
24 *dysregulation* of one or more brain circuits. Yet, clinical
25 neuroscience generally conflates the term *circuits* with co-
26 activated brain regions, the latter of which are more accu-
27 rately described as *networks*. Because neuroimaging-derived
28 networks are normally defined by linear regressions ($y =$
29 $b_0 + b_1x_1 + \dots$), they are capable of reliably modeling only
30 a very narrow range of topologies, in which one or more in-
31 puts leads to a single output (1). This limitation excludes the
32 capacity for positive and negative feedback loops, as required
33 for regulation.

34 To quantify brain circuit dysregulation we exploit the gen-
35 erative aspect of data-derived control circuits, which allows
36 us to predict how a circuit’s output time series will evolve
37 over time. In a classic engineering control application, such
38 as autopilot (Figure 1a), a vehicle corrects for deviations from
39 its desired trajectory through negative feedback (e.g., *as the*
40 *vehicle starts to drift to the right, the control circuit corrects*
41 *the drift by steering to the left*). As such, the difference be-
42 tween the autopilot’s actual versus desired trajectories pro-
43 vides a measure of its *control error*, or dysregulation (Fig-

ure 1b). Here, we use trajectory drift as a measure of control
error. We calculate circuit-wide dysregulation across fMRI-
derived control circuits, and demonstrate its clinical utility as
applied to three psychiatric use cases: *neuroticism*, *depres-*
sion, and *bipolar disorder* (See Methods for definitions).

In developing this framework, we started from several
desiderata: the ability to test homeostatic regulation in re-
sponse to driving inputs (perturbation), a fundamental re-
quirement of control theory(2–4); the ability to conduct
whole-brain circuit discovery, free of priors; and the ability to
scale, thereby leveraging the marked increase in both mega-
scale neuroimaging datasets made possible through open-
science initiatives, as well as high resolution, fine-granularity
parcellations of the brain.

To date, the only standardized method capable of esti-
mating fMRI-derived generative circuits is Dynamic Causal
Modeling (DCM) (5). DCM is normally used to estimate cir-
cuit architecture, in the form of a directed, weighted graph.
However, circuit architecture by itself is not sufficient to
provide quantitative estimation of control parameters such
as dysregulation. Moreover, DCM’s computationally ex-
pensive algorithms, even for faster (resting-state only) vari-
ants such as spectral DCM (6), result in convergence times
so lengthy that they remain impractical for extracting cir-
cuits from mega-scale datasets, using purely data-driven ap-
proaches (>100 brain regions), with short (e.g., 5 minute)
time series. More recent (driving input-compatible) variants
such as regression DCM (rDCM) (7) and sparse rDCM (8) re-
place the hemodynamic forward model with a fixed hemody-
namic response function (HRF). As a result, they fail to allow
for heterogeneity in the blood oxygenation level dependent
(BOLD) signal across brain regions and individuals (9, 10).
This biophysical constraint can lead to confounds, particu-
larly when applied to neurodevelopmental, aging, and patient
populations (11, 12).

Thus, we introduce a generalized framework based on state
space systems that bridges the gap between network the-
ory and control theory, with the scalability required to mine
mega-scale datasets such as UK Biobank (N=19,831) (13).
First, we confirm that using time-series to estimate systems
of differential equations without biophysical constraints still
permits recovery of complex causal relationships and feed-
back loops that characterize brain circuits. Using synthetic
data, we generate canonical circuit motifs to simulate circuits

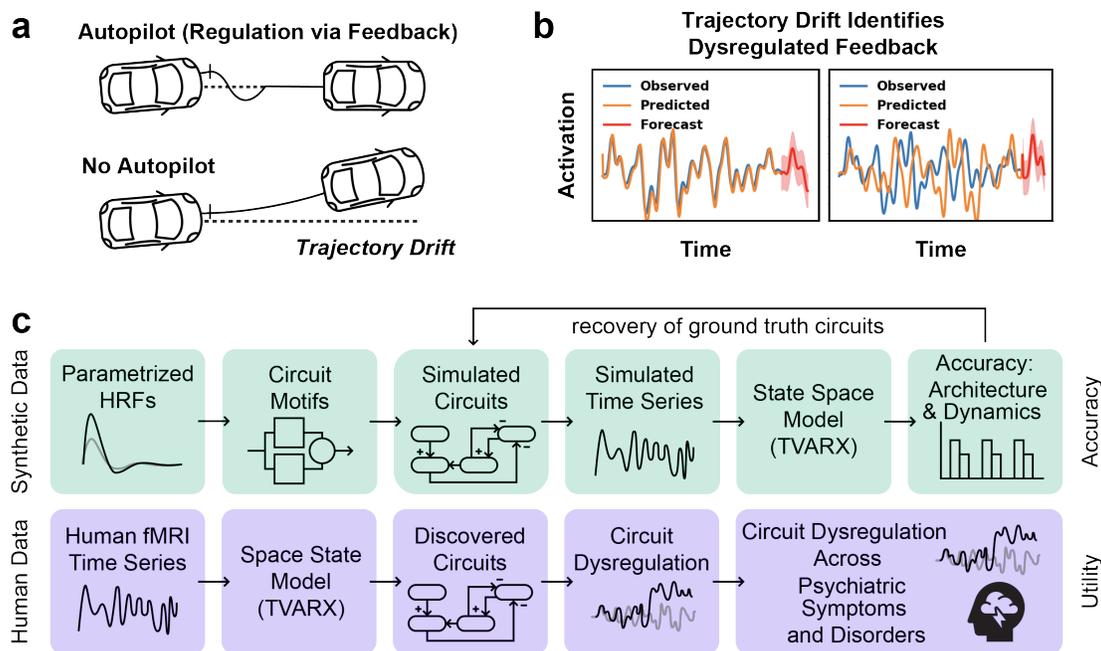


Fig. 1. Trajectory Drift as a Measure of Feedback Control Error, and Thus Circuit Dysregulation. (a) In a classic engineering control application, such as autopilot, a vehicle corrects for deviations from its desired trajectory through negative feedback. As such, the difference between the autopilot's actual versus desired trajectories provides a measure of its *control error*, or dysregulation. (b) As per the autopilot example, we use trajectory drift as control error to calculate dysregulation across fMRI-derived control circuits, and demonstrate its application for three psychiatric use cases: *neuroticism*, *depression*, and *bipolar disorder*. (c) Schematic of the pipeline for the discovery of circuit architecture and dynamics from human fMRI and simulated time series, using Time-Varying Autoregressive Model with Exogenous Inputs (TVARX) and other state space models. We use trajectory drift between predicted and actual trajectories to quantify circuit dysregulation across subjects with varying degrees of severity for psychiatric symptoms.

with varying architectures and dynamics to test our framework's ability to recover both (Figure 1c top row). Second, having validated the framework on synthetic data, we then apply the framework to UK Biobank fMRI data. Using tasks designed to dissociate processing of affective versus non-affective stimuli (14, 15), we extract the control circuit selectively engaged by each. Third, from each individual's circuits we calculate the circuit's trajectory control error, which quantifies its degree of dysregulation. From these control errors, we statistically test the relationship between circuit-wide dysregulation and psychiatric symptoms (Figure 1c bottom row).

Results

Recovering circuit motifs from dynamic outputs

We first evaluate our framework using circuit motifs. Five thousand synthetic circuits are constructed by connecting nodes, each with its own transfer function, according to three basic motifs: *series*, *parallel* and *feedback*. These motifs are then combined in a modular fashion, to create larger circuits of varying levels of complexity (See Methods, Figure 4a).

The transfer functions used in our experiments were designed to resemble the hemodynamic response function (HRF) (16) extensively used to model blood oxygen level dependent (BOLD) (17) signals measured using functional magnetic resonance imaging (fMRI) (Figure 4b). The HRF function is parameterized by response height, time-to-peak

and full-width at half-max. In our simulations, each node had different parameter values for the HRF, as previously shown for human data (9, 12).

The transfer function for each motif is an algebraic combination of node transfer functions (Figure 4c). Each motif also had an inverse variant. Serial and parallel connections each had both excitatory and inhibitory variants, while feedback loops had both positive and negative variants. The inverse variants are obtained by inverting the sign of their corresponding algebraic expressions (See Methods, Figure 4c). Note that although the node transfer functions are parameterized HRFs, successive connections and their corresponding algebraic expressions applied to HRFs can result in complex transfer functions.

We evaluate the ability to recover the canonical circuits both in terms of architecture and trajectory dynamics using a range of models with varying levels of complexity (Figure 2a). All of these models, with the exception of DCM, can be generalized by our state space system equations (See Methods, equation 3, Table 2). A detailed discussion regarding comparison with DCM is presented in a subsequent subsection.

The architecture is estimated through classifying relationship between each pair of nodes. For each node, parameters of its state space model are learned against time series of all other nodes (details in Methods). Based on the learned feed-forward matrix, a time-varying causality graph is established

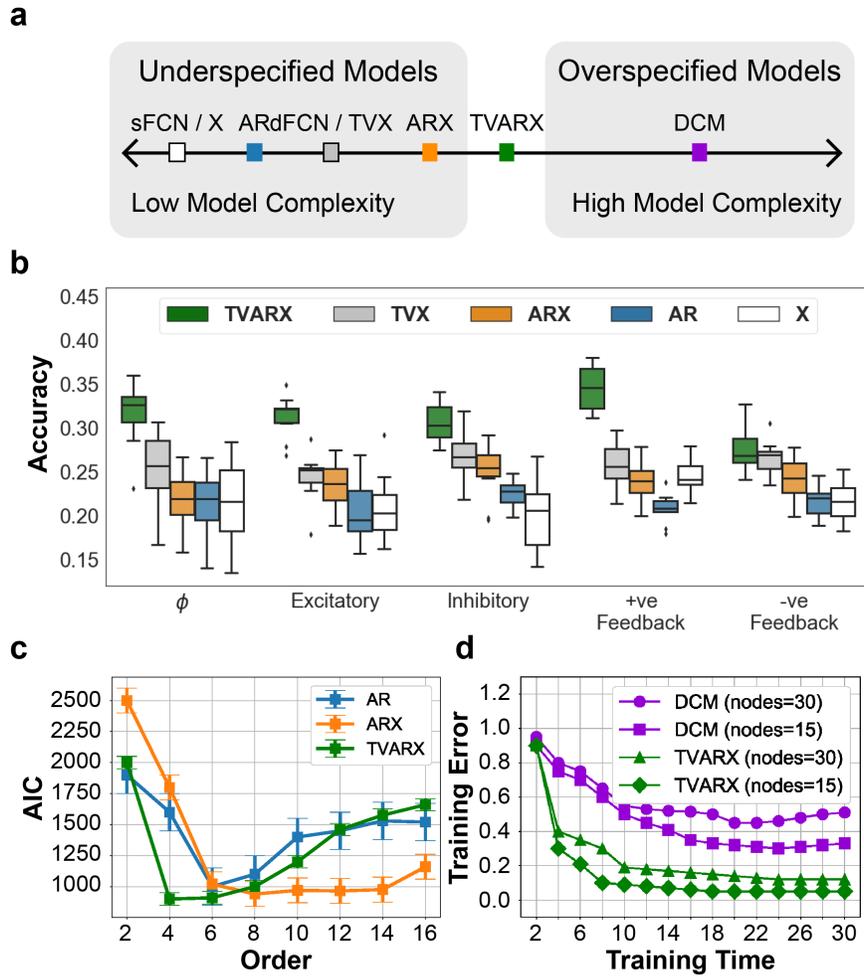


Fig. 2. Using Synthetic Data, we Compare the Performance (Accuracy and Speed) of System Identification Algorithms in Recovering Control Circuit Architecture and Dynamics (a) In this work, we fill the gap in literature between under-specified models that model networks, rather than circuits, and over-specified models that cannot scale computationally and therefore fail to extract circuits for large number of regions using shorter (~ 5 minute) time series. (b) Accuracy scores for classification of 5000 simulated circuits and their components, including correct classification of no connection: ϕ using different state space models discussed in Table 2. TVARX performed best at recovering the original circuit topology (c) AIC scores across models with respect to order of autoregressive component, which accounts for increasing complexity. TVARX performs best even when penalized for having the larger number of parameters. (d) Comparison of TVARX with (stochastic) Dynamic Causal Modeling (DCM) on human task-based fMRI. DCM fails to converge for shorter (~ 5 minute) fMRI time series as well as for larger number of nodes. *sFCN*=static Functional Connectivity Networks, *X*=eXogenous inputs, *AR*=AutoRegressive, *dFCN*=dynamic Functional Connectivity Networks, *TVX*=Time Varying with eXogenous inputs, *ARX*=AutoRegressive with eXogenous inputs, *TVARX*=Time Varying AutoRegressive with eXogenous inputs, *DCM*=Dynamic Causal Modeling.

140 (Figure 4e) and Eulerian/elementary circuits are identified.
 141 The connections of these circuits are further classified into
 142 one of the four connection types.

143 Figure 2b provides the accuracy scores for each model with
 144 respect to classification of each connection type. The absence
 145 of DCM in Figure 2b reflects the fact that it failed to converge
 146 and thus did not yield meaningful results for our synthetic
 147 circuits. The Time-Varying autoRegressive with eXogenous
 148 inputs (TVARX) model outperforms other, simpler, models
 149 in identifying each connection type and thus recovering the
 150 overall circuit architecture.

151 To account for varying model complexities, in our evalua-
 152 tions for predicting trajectories we compute the Akaike Infor-
 153 mation Criterion (AIC) for different models as we increase

154 the autoregressive order (Figure 2c). Only models that in-
 155 clude past states were included in these comparisons, as ma-
 156 jority of parameters are part of the autoregressive component
 157 and only these models are capable of generating future pre-
 158 dicted non-linear trajectories. Here again, the TVARX model
 159 outperforms other models we evaluated in our experiments,
 160 even after accounting for the larger number of parameters.

Data-Driven Circuit Discovery Using Human fMRI

161 Participants from UK-Biobank (N=19,831) (13) were
 162 scanned while engaged in a task designed to elicit affective
 163 and non-affective sensorimotor processing. These were used
 164 to identify circuits selectively activated for each type of pro-
 165 cessing.
 166

167 Task Design

168 During fMRI scans, participants were administered the Hariri
169 faces/shapes "emotion" task (14, 15), as also implemented in
170 the Human Connectome Project (HCP) (18) but with shorter
171 overall duration and hence fewer total stimulus block repeats.
172 Participants were presented with alternating blocks of trials
173 with visual stimuli consisting of human faces and geomet-
174 ric shapes (circles, ellipses), with brief periods of rest in be-
175 tween. For facial expressions, 12 different images were used:
176 six of each gender and affect (angry or fearful), all from a
177 standardized set of pictures of facial affect (19). In each trial,
178 either three faces or three shapes were presented in a triang-
179 ular configuration: one centered above the other two. The par-
180 ticipants were asked to indicate, by pressing a button, which
181 stimulus on the bottom row matched the stimulus on the top
182 row. The response triggered either the next trial or an eight
183 second period of rest. The total length of the scan for each
184 subject was 4 minutes; we obtained data for N=19,831 par-
185 ticipants. The data were acquired on harmonized Siemens
186 3T Skyra scanners. The scans are 2.4mm isotropic with TR
187 of 0.735s and 332 frames per run. Each subject had one run.
188 The parcellation used in our experiments was provided by
189 UK Biobank and included 139 regions of interest (ROIs).
190 These ROIs are defined in MNI152 space, combining par-
191 cellations from several atlases: the Harvard-Oxford cortical
192 and subcortical atlases (20, 21) and the Diedrichsen cerebel-
193 lar atlas (22). Further information for the dataset is provided
194 in Methods.

195 Comparison of TVARX with Dynamic Causal Modeling

196 In Figure 2d, we present a comparison between TVARX and
197 DCM with respect to training time and corresponding train-
198 ing error. Even for small circuits (nodes ≤ 30), DCM fails
199 to converge for time series of 4 minutes (332 timepoints) as
200 available in UK-Biobank. As shown, not only does TVARX
201 converge considerably faster, but training error for DCM does
202 not decrease monotonically, indicating failure to converge.

203 This is not surprising, since DCM is designed to be
204 hypothesis-driven, testing competing models from a pre-
205 specified set of nodes(6). Several competing hypotheses that
206 constitute a model space are specified in the form of sub-
207 graphs, which are then compared using Bayesian model se-
208 lection. Increasing the number of nodes is challenging be-
209 cause the number of extrinsic (between-node) connections or
210 edges increases with the square of the number of nodes. This
211 can lead to models with an enormous number of free param-
212 eters and profound conditional dependencies among the pa-
213 rameters. Furthermore, the computational time required to
214 invert these models grows exponentially with the number of
215 free parameters. More recent variants have been developed
216 to successfully address this issue, but were not compared
217 to TVARX because of other limitations: spectral DCM is
218 not appropriate for measuring homeostatic regulation in re-
219 sponse to driving inputs, and regression/sparse DCM (7, 8)
220 constrains the HRF in ways that can introduce confounds in
221 clinical populations (9–12).

222 Prefrontal-Limbic and Ventral Stream Subcircuits

223 In spite of the fact that no regions or connections were pre-
224 specified, our data-driven system identification methods were
225 highly successful in accurately identifying linked subcircuits
226 (Figure 3), each of which was consistent with independently
227 validated experiments in the rodent, macaque, and human.
228 Each subcircuit was uniquely specified with respect to its
229 showing the largest absolute **D** value in response to its re-
230 spective stimulus-type (See Methods 5).

231 As per the translationally-established *prefrontal-limbic*
232 *subcircuit* (PFLC) for processing of affective stimuli, in
233 which the thalamus provides a hub connecting a "low road"
234 pathway to the amygdala and a "high road" pathway to the or-
235 bitofrontal cortex (OFC) and ventromedial prefrontal cortex
236 (vmPFC)(23–26), our model recovered all key components
237 and their relationships (Figure 3, red).

238 Likewise, for processing object form and recognition of
239 non-affective stimuli, the *ventral stream subcircuit* has been
240 shown to originate in the thalamus, project to V1-V2-V4, ter-
241 minating in the inferior temporal gyrus (ITG), then progress
242 to the inferior frontal gyrus (ITG, specifically the ventro-
243 lateral prefrontal cortex) and then to the orbitofrontal cor-
244 tex/ventromedial prefrontal cortex (vmPFC)(27–29). Our
245 model was able to recover nearly all of ventral stream (V1-
246 V2-V4 were implicit in connecting the thalamus to the ITG),
247 including "top-down" feedback (30) from the inferior frontal
248 gyrus (IFG) to the inferior temporal gyrus (ITG) (Figure 3,
249 blue).

250 In the context of psychiatry it is important to note that,
251 while in human neuroimaging the ventral stream's (VS) input
252 to the IFG has most been most often studied in the context
253 of disambiguation of semantic meaning(31, 32), this same
254 subcircuit also disambiguates perceptual meaning in the con-
255 text of ambiguous threat in assessing risk ("threat general-
256 ization"). Indeed, it is dysregulation of the VS subcircuit—in
257 response to *non-affective*, rather than affective stimuli—that
258 we have previously shown (in four independent datasets, to-
259 taling N=226) to track the spectrum of trait to clinical anxiety
260 (33), which most closely relates to the UK Biobank variable
261 "neuroticism."

262 The two subcircuits were found to be mutually interacting
263 with each other, with the pivot point centered at the thalamus,
264 a known neuroanatomical hub shared by both circuits(29).
265 Our system identification methods identified the thalamus to
266 have two inputs providing negative feedback, one from hip-
267 pocampus for the prefrontal-limbic subcircuit and one from
268 inferior frontal gyrus for the ventral stream subcircuit; i.e.

$$Hippocampus \xrightarrow[-]{PFLC} Thalamus \xleftarrow[-]{VS} IFG^1$$

¹Notation: Region A $\xrightarrow[\text{circuit}]{\text{connection type}}$ Region B. Excitatory connections are denoted by + and inhibitory connections by -. Direction of the arrow head represents indicates directionality of causation.

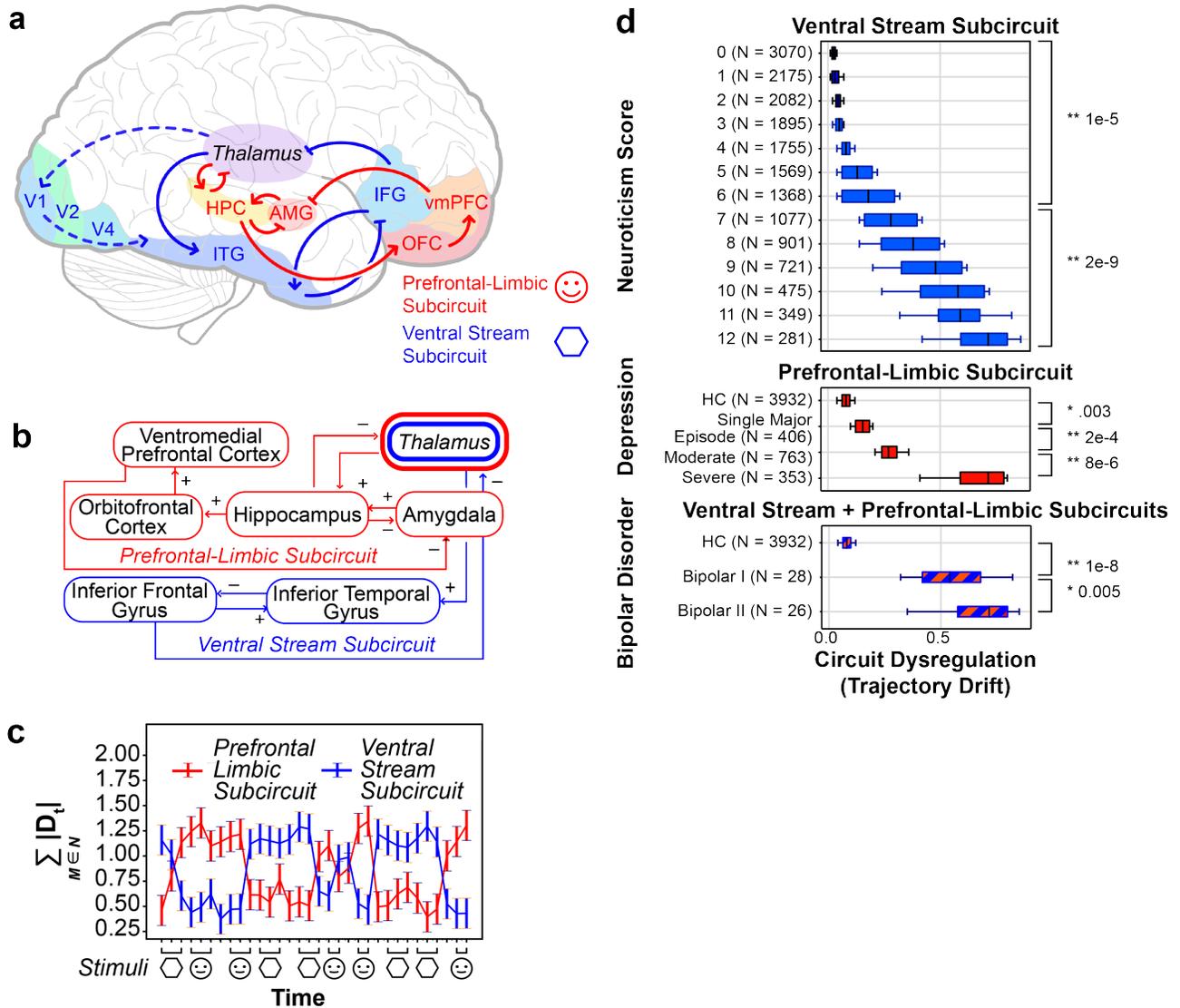


Fig. 3. Purely Data-Driven System Identification Accurately Recovers Known Subcircuits, Whose Degree of Dysregulation is Linked to Psychiatric Symptom Severity. (a-b) Without including any pre-specified information regarding regions or connections, our data-driven methods were highly successful in recovering two known subcircuits, interacting and linked via the thalamus, that have been independently validated in the rodent, macaque, and human (details in Discussion). (c) Results show selective dominance of each subcircuit with changing stimuli: the *prefrontal-limbic* subcircuit selectively engaged during processing of affective (faces) stimuli while the *ventral stream* subcircuit selectively engaged during processing of non-affective (shapes) stimuli (the y-axis is the relative dominance of one circuit versus the other, defined as absolute sum of relevant entries of the feedforward matrix D in our state space equation (Eqn. 3). For prefrontal-limbic subcircuit $M = \{thalamus, hippocampus, OFC, vmPFC\}$ and for ventral stream subcircuit $M = \{thalamus, ITF, IFG\}$) (d) We use trajectory drift (measured as the mean squared error between actual and predicted trajectories) as a measure of feedback control error, or dysregulation, of each subcircuit. Dysregulation of the prefrontal-limbic subcircuit was measured as the error in the thalamus trajectory as predicted from negative feedback by the hippocampus (hippocampus \rightarrow thalamus). Dysregulation of the ventral stream was measured as error in the thalamus trajectory as predicted from negative feedback by the inferior frontal gyrus (inferior frontal gyrus \rightarrow thalamus). Trajectory drift of the ventral stream subcircuit tracks severity of neuroticism, while trajectory drift of the prefrontal-limbic circuit tracks severity of depression. For bipolar disorder, the thalamic trajectory could not be predicted from either hippocampal or IFG trajectories, and thus reflects dysregulation of the full circuit. This could be due to either more systemic problems with feedback across both circuits, or that the full circuit is receiving dysregulated inputs from another, different, circuit not identified by these tasks. Bonferroni corrected * $P < 0.05$; ** $P < 0.01$

269 Similarly, the thalamus has two outputs, both excitatory
270 connections, one to hippocampus in the PFLC and the other
271 to ITG in VS i.e.

$$Hippocampus \xleftarrow{+}_{PFLC} Thalamus \xrightarrow{+}_{VS} ITG$$

272 In this way, the thalamus *completes* one of two competing
273 negative feedback loops, one for each of the two identified

subcircuits.

274 We measure the relative dominance of one subcircuit versus
275 the other at any point in time as the absolute sum of relevant
276 entries in the time-varying feedforward matrix D_t in our
277 state space equation (Methods, Equation 3). Note that each
278 entry in $D_t^{(a,b)}$ represents the causal dependence of trajectory
279 of b on trajectory of a (i.e. effective connectivity $a \rightarrow b$).
280

281 Our results show these two competing feedback loops to
282 be alternatively dominating in strength based on the stimuli
283 during the scan for each subject (Figure 3b). The prefrontal-
284 limbic subcircuit was found to be the dominant loop at points
285 in time when the subjects were tasked with matching facial
286 stimuli of angry or fearful affect. The opposite was observed
287 for ventral stream subcircuit, which was the dominant feed-
288 back loop when subjects were tasked with matching geomet-
289 ric shapes.

290 Trajectory Drift as Dysregulation

291 We use trajectory drift as measure of circuit-wide control er-
292 ror, and therefore dysregulation. This drift is measured as
293 the mean squared error between the actual trajectories and
294 the predicted trajectories. We further compare these varying
295 levels of dysregulation with the severity and type of psychi-
296 atric symptoms and diagnoses. These include scored degrees
297 of neuroticism (a measure of stress vulnerability, anxiety),
298 depression, and diagnosis of Type 1 and 2 bipolar disorders
299 (definitions in Methods).

300 To estimate dysregulation of the PFLC and VS subcircuits
301 identified for N=19,831 subjects, we use the same task fMRI
302 scans. However, unlike our identification of circuits in the
303 previous section, here we measure dysregulation across the
304 entire scan, independent of the design matrix.

305 Our results show marked association between greater dys-
306 regulation of specific subcircuits and the severity of the psy-
307 chiatric symptoms (Figure 3c).

308 Since the thalamus was identified as a pivot point for
309 switching between the two subcircuits, in determining which
310 of the two competing feedback loops dominates the system,
311 we specifically focused on regulation of the thalamus; i.e.
312 prediction of thalamus's trajectory as a function of negative
313 feedback by either the hippocampus (for PFLC) or the infe-
314 rior frontal gyrus (for VS).

315 Our results show more severe neuroticism to be associ-
316 ated with greater trajectory drift (control error) in the ven-
317 tral stream subcircuit (level 0 [N=3070] vs. 6 [N=1368]
318 $**p \leq 1e-5$; level 6 [N=1368] vs. 12 [N=281] $**p \leq 2e-9$)
319 (Figure 3c top), specific to weakened negative feedback from
320 the inferior frontal gyrus (IFG) to the thalamus (IFG $\xrightarrow{-}$ Tha-
321 lamus). This inhibitory connection is critical to stable regu-
322 lation of the ventral stream and was observed to in turn result
323 in greater dysregulation downstream with respect to thalamic
324 outputs to the inferior temporal gyrus (ITG) (Thalamus $\xrightarrow{+}$
325 ITG).

326 In contrast, our results show more severe depression to be
327 associated with greater trajectory drift (control error) in the
328 prefrontal-limbic subcircuit (HC [N=3932] vs. Single Major
329 Episode [N=406] $*p \leq 0.003$; Single Major Episode [N=406]
330 vs. Moderate [N=763] $**p \leq 2e-4$; Moderate [N=763] vs.
331 Severe [N=353] $**p \leq 8e-6$) (Figure 3c center), specific to
332 weakened negative feedback from the hippocampus to the
333 thalamus (Hippocampus $\xrightarrow{-}$ Thalamus). Note that this rela-

334 tionship is itself dependent on the excitatory inputs from the
335 thalamus to the hippocampus (Thalamus $\xrightarrow{+}$ Hippocampus)
336 completing the negative feedback loop.

337 In the case of bipolar disorder, the thalamus was observed
338 to be dysregulated with respect to both of its inhibitory in-
339 puts (Hippocampus $\xrightarrow[-PFLC]{-}$ Thalamus $\xleftarrow[-VS]{-}$ Inferior Tempo-
340 ral Gyrus) with greater dysregulation observed for subjects
341 with Bipolar I Disorder compared to subjects with Bipolar
342 II Disorder (Figure 3 bottom) (HC [N=3932] vs. Bipolar I
343 [N=28] $*p \leq 1e-8$; Bipolar I [N=28] vs. Bipolar II [N=26]
344 $*p \leq 0.0005$). In the case of Bipolar Disorder I & II, the
345 thalamic trajectory was observed to drift significantly from
346 its predicted trajectory, but the system was not dominated by
347 either of the two competing feedback loops. This could be
348 due either to more systemic problems with feedback across
349 both circuits, or that the full circuit is receiving dysregulated
350 inputs from another, different, circuit not identified by these
351 tasks.

352 Finally, we compare trajectory drifts of discovered cir-
353 cuits with more conventional methods currently prevalent in
354 clinical neuroscience (Table 1). These standard methods in-
355 clude Stochastic DCM, correlation-based functional connec-
356 tivity (34), and activation based Generalized Linear Models
357 (GLM) (35). Note, however, that different methods answer
358 fundamentally different (if complementary) questions: DCM
359 and Trajectory Drift capture circuit-wide dynamics, func-
360 tional connectivity provides the strength of (undirected) sig-
361 naling across pairs of regions, and GLM provides activation
362 of individual regions. To allow for a more direct compari-
363 son of our circuit-wide measure to activation and networks,
364 in Table 1, we report results for the subcircuit regions and
365 connections that we identified as tracking symptom severi-
366 ty. Our results show that modeling psychiatric disorders in
367 terms of circuit dysregulation achieves markedly greater de-
368 tection sensitivity across all three sets of psychiatric symp-
369 toms. Beyond identification of differences, however, the most
370 important advantage of our method is that it uses data-driven
371 methods to construct generative computational neuroscience
372 models that explicitly consider homeostatic regulation across
373 negative feedback loops. This has the potential to allow hy-
374 potheses regarding dysregulation across psychiatrically rele-
375 vant circuits to be rigorously specified and empirically tested.

376 Discussion

377 In this work, we present a scalable fMRI data-driven techni-
378 que that allows for construction of generative circuits in the
379 human brain, and provides a quantitative measure—trajectory
380 drift—of their control error, or circuit-wide dysregulation.
381 We demonstrate the effectiveness of our technique in re-
382 covering artificially generated circuits of varied architectures
383 and transfer functions. To demonstrate the applicability of
384 the technique to computational psychiatry, we use large-
385 scale fMRI data to identify two subcircuits and demonstrate
386 that their dysregulation tracks with symptom severity with

		Neuroticism Score		Depression			Bipolar Disorder	
		0 (N=3070)	6 (N=1368)	HC (N=3932)	Single Major Episode (N=406)	Moderate (N=763)	HC (N=3932)	Bipolar I (N=28)
		vs.	vs.	vs.	vs.	vs.	vs.	vs.
		6 (N=1368)	12 (N=281)	Single Major Episode (N=406)	Moderate (N=763)	Severe (N=353)	Bipolar I (N=28)	Bipolar II (N=26)
DCM	<i>Failed to Converge</i>							
Trajectory Drift	Hippocampus \rightarrow Thalamus			*0.003	**2e-4	**8e-6	**1e-8	*0.005
	Inferior Frontal Gyrus \rightarrow Thalamus	**1e-5	**2e-9					
Functional Connectivity	Hippocampus \leftrightarrow Thalamus			0.06	0.1	*0.001	**1e-4	0.3
	Inferior Frontal Gyrus \leftrightarrow Thalamus	*0.002	**1e-4				0.12	0.26
GLM	Thalamus	0.08	*0.01	0.13	*0.012	0.16	*0.01	0.2

Table 1. Comparison of Trajectory Drift with fMRI Analytical Methods: Dynamic Causal Modeling (DCM), Functional Connectivity, and Activation Based Generalized Linear Models (GLM). We took a whole-brain purely data-driven approach in identifying circuits for the two circuit-based methods: DCM and Trajectory Drift. Of these, only Trajectory Drift was able to converge for the parcellation (139 regions of interest) and sample size (UK Biobank N=19,831). For the two subcircuits identified by Trajectory Drift: Prefrontal-Limbic and Ventral Stream, we then tested how the key regulatory components for Prefrontal-Limbic (negative feedback by the hippocampus) and Ventral Stream (negative feedback by the inferior frontal gyrus) were interpreted by non-circuit-based methods: GLM and Functional connectivity. For each comparison of psychiatric variables, we report p-values from statistical significance testing using Welch's t-test for unequal variances and sample sizes. Bonferroni corrected * $P < 0.05$; ** $P < 0.01$. Cells that were not applicable are greyed out.

387 markedly greater detection sensitivity than standard analytic
388 methods.

389 fMRI has conventionally been used to either compute
390 brain activation maps, as areas of differential hemodynamic
391 response, or to quantify pairwise connectivity between brain
392 regions using Pearson correlation(36). More recent devel-
393 opments in fMRI analyses consider graph-theoretic mea-
394 sures (37) and a shift towards dynamic patterns of connectiv-
395 ity using time varying connections (38, 39). What all of these
396 methods lack, however, is a conceptual and mathematical
397 framework for considering the implications of closed feed-
398 back loops. Without these, activation maps and connectivity-
399 derived networks can suggest the presence of neural circuits,
400 but can neither define nor simulate their behavior, which in-
401 cludes their regulation. Given the assumption that psychiatric
402 disorders are grounded in the failure of circuits to maintain
403 homeostatic regulation, the ability to identify trajectories–
404 including drift from normative trajectories–is thus an impor-
405 tant step in the development of computational psychiatry and
406 its characterization of dynamical disease(40, 41).

407 ACKNOWLEDGEMENTS

408 The research described in this article was funded by the Baszucki Brain Re-
409 search Fund (L.R.M.-P.)

410 DATA AVAILABILITY

411 The neuroimaging and phenotypic data used in this work was obtained from UK-
412 Biobank under Data Access Application 37462 and is available upon application at
413 <http://www.ukbiobank.ac.uk/register-apply/>.

414 CODE AVAILABILITY

415 All source code is available at
416 <https://github.com/fahadsultan/QuantifyingDysregulation>.

Bibliography

1. Lilianne Mujica-Parodi and Helmut Strey. Making sense of computational psychiatry. *The international journal of neuropsychopharmacology*, 23, 03 2020. doi: 10.1093/ijnp/pyaa013. 418
2. E. J. Hinch. *Perturbation Methods*. Cambridge Texts in Applied Mathematics. Cambridge University Press, 1991. doi: 10.1017/CBO9781139172189. 419
3. Michael CK Khoo. *Physiological control systems: analysis, simulation, and estimation*. John Wiley & Sons, 2018. 420
4. Petar Kokotović, Hassan K. Khalil, and John O'Reilly. *Singular Perturbation Methods in Control: Analysis and Design*. Society for Industrial and Applied Mathematics, 1999. doi: 10.1137/1.9781611971118. 421
5. Karl J Friston, Lee Harrison, and Will Penny. Dynamic causal modelling. *Neuroimage*, 19 (4):1273–1302, 2003. 422
6. Karl J Friston, Joshua Kahan, Bharat Biswal, and Adeel Razi. A dcm for resting state fmri. *Neuroimage*, 94:396–407, 2014. 423
7. Stefan Frässle, Ekaterina I Lomakina, Lars Kasper, Zina M Manjaly, Alex Leff, Klaas P Pruessmann, Joachim M Buhmann, and Klaas E Stephan. A generative model of whole-brain effective connectivity. *Neuroimage*, 179:505–529, 2018. 424
8. Stefan Frässle, Ekaterina I Lomakina, Adeel Razi, Karl J Friston, Joachim M Buhmann, and Klaas E Stephan. Regression dcm for fmri. *Neuroimage*, 155:406–421, 2017. 425
9. Daniel A Handwerker, John M Ollinger, and Mark D'Esposito. Variation of bold hemodynamic responses across subjects and brain regions and their effects on statistical analyses. *Neuroimage*, 21(4):1639–1651, 2004. 426
10. Geoffrey Karl Aguirre, Eric Zarahn, and Mark D'Esposito. The variability of human, bold hemodynamic responses. *Neuroimage*, 8(4):360–369, 1998. 427
11. Olivier David, Isabelle Guillemain, Sandrine Saillet, Sebastien Reyt, Colin Deransart, Christoph Segebarth, and Antoine Depaulis. Identifying neural drivers with functional mri: an electrophysiological validation. *PLoS biology*, 6(12):e315, 2008. 428
12. Pedro A Valdes-Sosa, Alard Roebroeck, Jean Daunizeau, and Karl Friston. Effective connectivity: influence, causality and biophysical modeling. *Neuroimage*, 58(2):339–361, 2011. 429
13. Clare Bycroft, Colin Freeman, Desislava Petkova, Gavin Band, Lloyd T Elliott, Kevin Sharp, Allan Motyer, Damjan Vukcevic, Olivier Delaneau, Jared O'Connell, et al. The uk biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726):203–209, 2018. 430
14. Ahmad R Hariri, Alessandro Tessitore, Venkata S Mattay, Francesco Fera, and Daniel R Weinberger. The amygdala response to emotional stimuli: a comparison of faces and scenes. *Neuroimage*, 17(1):317–323, 2002. 431
15. Deanna M Barch, Gregory C Burgess, Michael P Harms, Steven E Petersen, Bradley L Schlaggar, Maurizio Corbetta, Matthew F Glasser, Sandra Curtiss, Sachin Dixit, Cindy Feldt, et al. Function in the human connectome: task-fmri and individual differences in behavior. *Neuroimage*, 80:169–189, 2013. 432
16. Richard B Buxton, Kámil Uludağ, David J Dubowitz, and Thomas T Liu. Modeling the hemodynamic response to brain activation. *Neuroimage*, 23:S220–S233, 2004. 433
17. Seiji Ogawa, RS Menon, David W Tank, SG Kim, H Merkle, JM Ellermann, and K Ugurbil. Functional brain mapping by blood oxygenation level-dependent contrast magnetic resonance imaging. a comparison of signal characteristics with a biophysical model. *Biophysical journal*, 64(3):803–812, 1993. 434
18. David C Van Essen, Stephen M Smith, Deanna M Barch, Timothy EJ Behrens, Essa Yacoub, Kamil Ugurbil, Wu-Minn HCP Consortium, et al. The wu-minn human connectome project: an overview. *Neuroimage*, 80:62–79, 2013. 435
19. Paul Ekman and Wallace V Friesen. Measuring facial movement. *Environmental psychology and nonverbal behavior*, 1(1):56–75, 1976. 436
20. Rahul S Desikan, Florent Ségonne, Bruce Fischl, Brian T Quinn, Bradford C Dickerson, 437

468 Deborah Blacker, Randy L Buckner, Anders M Dale, R Paul Maguire, Bradley T Hyman, 554
469 et al. An automated labeling system for subdividing the human cerebral cortex on mri scans 555
470 into gyral based regions of interest. *Neuroimage*, 31(3):968–980, 2006. 556

471 21. Jean A Frazier, Sufen Chiu, Janis L Breeze, Nikos Makris, Nicholas Lange, David N 557
472 Kennedy, Martha R Herbert, Eileen K Bent, Vamsi K Koneru, Megan E Dieterich, et al. 558
473 Structural brain magnetic resonance imaging of limbic and thalamic volumes in pediatric 559
474 bipolar disorder. *American Journal of Psychiatry*, 162(7):1256–1265, 2005. 560

475 22. Jörn Diedrichsen, Joshua H Balsters, Jonathan Flavell, Emma Cussans, and Narender 561
476 Ramnani. A probabilistic mr atlas of the human cerebellum. *Neuroimage*, 46(1):39–46, 562
477 2009. 563

478 23. Elizabeth A. Phelps and Joseph E. LeDoux. Contributions of the amygdala to emotion 564
479 processing: From animal models to human behavior. *Neuron*, 48, 2005. ISSN 08966273. 565
480 doi: 10.1016/j.neuron.2005.09.025. 566

481 24. M. Alexandra Kredlow, Robert J. Fenster, Emma S. Laurent, Kerry J. Ressler, and Eliza- 567
482 beth A. Phelps. Prefrontal cortex, amygdala, and threat processing: implications for ptsd. 568
483 *Neuropsychopharmacology*, 47, 2022. ISSN 1740634X. doi: 10.1038/s41386-021-01155-7. 569

484 25. Joseph E. LeDoux. Emotion circuits in the brain. *Annual Review of Neuroscience*, 23(1): 570
485 155–184, 2000. 571

486 26. MR Bennett. The prefrontal–limbic network in depression: Modulation by hypothalamus, 572
487 basal ganglia and midbrain. *Progress in neurobiology*, 93(4):468–487, 2011. 573

488 27. Dwight J. Kravitz, Kadharbatcha S. Saleem, Chris I. Baker, Leslie G. Ungerleider, and Mor- 574
489 timer Mishkin. The ventral visual pathway: An expanded neural framework for the pro- 575
490 cessing of object quality. *Trends in Cognitive Sciences*, 17, 2013. ISSN 13646613. doi: 576
491 10.1016/j.tics.2012.10.011. 577

492 28. Cornelius Weiller, Marco Reiser, Ivo Peto, Jrgen Hennig, Nikos Makris, Michael Petrides, 578
493 Michel Rijntjes, and Karl Egger. The ventral pathway of the human brain: A continuous as- 579
494 sociation tract system. *NeuroImage*, 234, 2021. ISSN 10959572. doi: 10.1016/j.neuroimage. 580
495 2021.117977. 581

496 29. Luiz Pessoa and Ralph Adolphs. Emotion processing and the amygdala: From a ‘low road’ 582
497 to ‘many roads’ of evaluating biological significance. *Nature Reviews Neuroscience*, 11, 583
498 2010. ISSN 1471003X. doi: 10.1038/nrn2920. 584

499 30. Charles D. Gilbert and Mariano Sigman. Brain states: Top-down influences in sensory 585
500 processing. *Neuron*, 54, 2007. ISSN 08966273. doi: 10.1016/j.neuron.2007.05.019. 586

501 31. Mirjana Bozic, Lorraine K. Tyler, David T. Ives, Billi Randall, and William D. Marslen-Wilson. 587
502 Bihemispheric foundations for human speech comprehension. *Proceedings of the National 588
503 Academy of Sciences of the United States of America*, 107, 2010. ISSN 00278424. doi: 589
504 10.1073/pnas.1000531107. 590

505 32. Jennifer M. Rodd, Ingrid S. Johnsrude, and Matthew H. Davis. Dissociating frontotemporal 591
506 contributions to semantic ambiguity resolution in spoken sentences. *Cerebral Cortex*, 22, 592
507 2012. ISSN 10473211. doi: 10.1093/cercor/bhr252. 593

508 33. Lilliane R Mujica-Parodi, Jiook Cha, and Jonathan Gao. From anxious to reckless: a control 594
509 systems approach unifies prefrontal-limbic regulation across the spectrum of threat detec- 595
510 tion. *Frontiers in systems neuroscience*, 11:18, 2017. 596

511 34. Karl J Friston. Functional and effective connectivity in neuroimaging: a synthesis. *Human 597
512 brain mapping*, 2(1-2):56–78, 1994. 598

513 35. Karl J Friston, Andrew P Holmes, Keith J Worsley, J-P Poline, Chris D Frith, and Richard SJ 599
514 Frackowiak. Statistical parametric maps in functional imaging: a general linear approach. 600
515 *Human brain mapping*, 2(4):189–210, 1994. 601

516 36. Jean-Baptiste Poline and Matthew Brett. The general linear model and fmri: does love last 602
517 forever? *Neuroimage*, 62(2):871–880, 2012. 603

518 37. Danielle Smith Bassett and ED Bullmore. Small-world brain networks. *The neuroscientist*, 604
519 12(6):512–523, 2006. 605

520 38. Daniel J Lurie, Daniel Kessler, Danielle S Bassett, Richard F Betzel, Michael Breakspear, 606
521 Shella Kheihholz, Aaron Kucyi, Raphaël Liégeois, Martin A Lindquist, Anthony Randal McIn- 607
522 tosh, et al. Questions and controversies in the study of time-varying functional connectivity 608
523 in resting fmri. *Network Neuroscience*, 4(1):30–69, 2020. 609

524 39. Leonardo L Gollo and Michael Breakspear. The frustrated brain: from dynamics on motifs 610
525 to communities and networks. *Philosophical Transactions of the Royal Society B: Biological 611
526 Sciences*, 369(1653):20130532, 2014. 612

527 40. Quentin JM Huys, Michael Browning, Martin P Paulus, and Michael J Frank. Advances in 613
528 the computational understanding of mental illness. *Neuropsychopharmacology*, 46(1):3–19, 614
529 2021. 615

530 41. Daniel Durstewitz, Quentin JM Huys, and Georgia Koppe. Psychiatric illnesses as disorders 616
531 of network dynamics. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 6 617
532 (9):865–876, 2021. 618

533 42. Lionel Barnett, Adam B Barrett, and Anil K Seth. Granger causality and transfer entropy are 619
534 equivalent for gaussian variables. *Physical review letters*, 103(23):238701, 2009. 620

535 43. Steen Moeller, Essa Yacoub, Cheryl A Olman, Edward Auerbach, John Strupp, Noam Harel, 621
536 and Kâmil Uğurbil. Multiband multislice ge-epi at 7 tesla, with 16-fold acceleration using partial 622
537 parallel imaging with application to high spatial and temporal whole-brain fmri. *Magnetic 623
538 resonance in medicine*, 63(5):1144–1153, 2010. 624

539 44. Christian F Beckmann and Stephen M Smith. Probabilistic independent component analysis 625
540 for functional magnetic resonance imaging. *IEEE transactions on medical imaging*, 23(2): 626
541 137–152, 2004. 627

542 45. Gholamreza Salimi-Khorshidi, Gwenaëlle Douaud, Christian F Beckmann, Matthew F 628
543 Glasser, Ludovica Griffanti, and Stephen M Smith. Automatic denoising of functional mri 629
544 data: combining independent component analysis and hierarchical fusion of classifiers. 630
545 *Neuroimage*, 90:449–468, 2014. 631

546 46. Mark W Woolrich, Brian D Ripley, Michael Brady, and Stephen M Smith. Temporal autocor- 632
547 relation in univariate linear modeling of fmri data. *Neuroimage*, 14(6):1370–1386, 2001. 633

548 47. Fidel Alfaro-Almagro, Mark Jenkinson, Neal K Bangerter, Jesper LR Andersson, Ludovica 634
549 Griffanti, Gwenaëlle Douaud, Stamatios N Sotiropoulos, Saad Jbabdi, Moises Hernandez- 635
550 Fernandez, Emmanuel Vallee, et al. Image processing and quality control for the first 10,000 636
551 brain imaging datasets from uk biobank. *Neuroimage*, 166:400–424, 2018. 637

552 48. Stephen B Manuck, Sarah M Brown, Erika E Forbes, and Ahmad R Hariri. Temporal stability 638
553 of individual differences in amygdala reactivity. *American Journal of Psychiatry*, 164(10): 639
1613–1614, 2007. 640

49. Daniel J Smith, Barbara I Nicholl, Breda Cullen, Daniel Martin, Zia Ul-Haq, Jonathan Evans, 641
Jason MR Gill, Beverly Roberts, John Gallacher, Daniel Mackay, et al. Prevalence and 642
characteristics of probable major depression and bipolar disorder within uk biobank: cross- 643
sectional study of 172,751 participants. *PLoS one*, 8(11):e75362, 2013. 644

Methods

Time Varying AutoRegressive eXogenous (TVARX) Model

Functional connectivity network for N regions-of-interest is traditionally defined as an $N \times N$ adjacency matrix A where

$$A_{x,y} = \frac{\text{cov}(X,Y)}{\sigma_x \sigma_y}, A \in \mathbb{R}^{(N \times N)}, x, y \leq N \quad (1)$$

When time series are normalized (zero mean and unit standard deviation), a common fMRI preprocessing step, the Pearson correlation coefficient is equal to the slope of the regression. Thus resulting in a linear regression model of the form

$$y_t = A_{x,y} x + b_t \quad (2)$$

Most dynamic variants simply extend this definition by adding an additional temporal dimension resulting in a time-varying adjacency matrix $A \in \mathbb{R}^{N \times N \times T}$ where T is either the length of the time series or the number of sliding time-windows.

In this work, we extend this simple prevalent linear model by modeling BOLD time series observed for a brain region using a state space model of the form:

$$\begin{aligned} y_t &= Z_t \alpha_t + D_t u_t + d_t + \epsilon \\ \alpha_{t+1} &= T_t \alpha_t + B_t u_t + c_t + R_t \eta_t \end{aligned} \quad (3)$$

where y_t refers to the observation vector at time t , u_t refers to the input (or control) vector from other regions of the brain, α_t refers to the (unobserved) state vector at time t , and where the irregular components are defined as $\epsilon_t \sim N(0, H_t)$ and $\eta_t \sim N(0, Q_t)$.

The remaining variables in the equations are matrices describing the process. The total length of the time series being T , the number of ROIs being N and K being the number of states, their variable names and dimensions are as follows: design $Z \in \mathbb{R}^{N \times K \times T}$, input $B \in \mathbb{R}^{N \times K \times T}$, observation intercept $d \in \mathbb{R}^{N \times T}$, observation covariance $H \in \mathbb{R}^{N \times N \times T}$, transition $T \in \mathbb{R}^{K \times K \times T}$, state intercept $c \in \mathbb{R}^{K \times T}$, selection $R \in \mathbb{R}^{K \times K \times T}$, state covariance $Q \in \mathbb{R}^{K \times K \times T}$

Note that this formulation is a generalized framework with prevailing definitions of static functional connectivity as correlations ($D_t = D_{t+1} \forall t, Z_t = O \forall t$) and dynamic functional connectivity as time-varying correlations ($Z_t = O \forall t$) as special cases. Table 2 breaks down existing models and presents

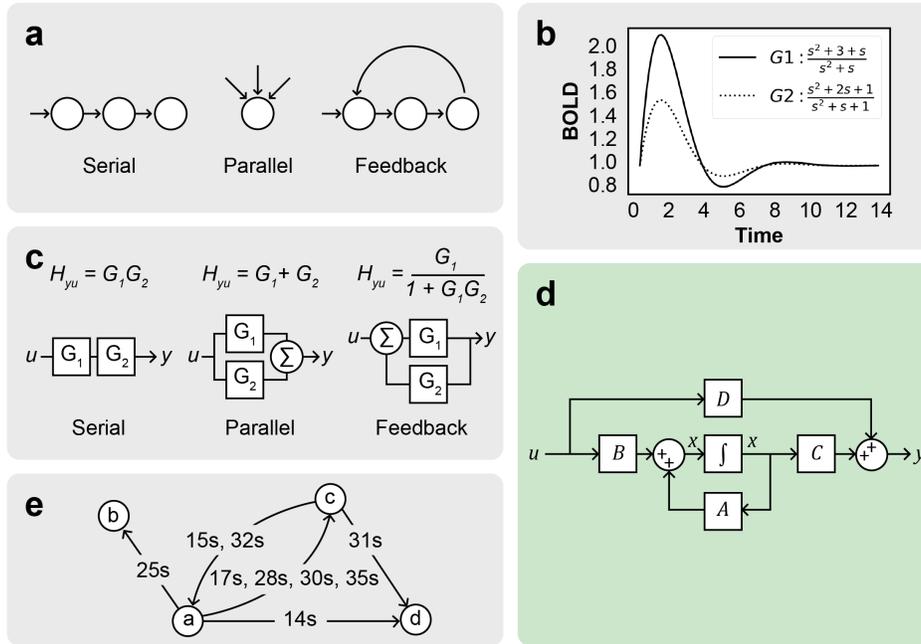


Fig. 4. Inferring Closed-Loop Circuits(a) Standard connectivity ("network") analyses depend upon linear regressions, which are only capable of modeling a very specific topology: parallel inputs. In contrast to parallel inputs A , most neurobiological circuits of relevance to psychiatry also require serial B and feedback C components, structures that could lead to an explosion of error propagation using standard statistical methods. (b) Impulse response for two hemodynamic response functions (HRFs) with different relaxation times and transfer functions. (c) To illustrate how transfer function structure changes with different circuit topologies, we show three transfer functions, each of which corresponds to a different kind of "motif," with series, parallel and feedback connections. By using pairs of inputs u and outputs y to obtain their transfer function, we systematically infer circuit topology. (d) Block diagram representation of state-space equations (e) Dynamic effective connectivity as a time-stamped temporal graph.

Model	D (feedforward)	T (transition)	B (input)	Z (design)
Correlations (sFCN / X)	Time invariant	O	O	O
Autoregressive (AR)	O	Time invariant	O	I
Time-varying Correlations (dFCN / TVX)	Time varying	O	O	O
Autoregressive w/ eXogenous inputs (ARX)	Time invariant	Time invariant	O	O
Time-varying Autoregressive w/ eXogenous inputs (TVARX)	Time varying	Time invariant	O	I

Table 2. Bridging the gap between network theory and control theory: extending existing correlation based models to our TVARX model. The breakdown and comparison in terms of state space parameters elucidates how our model is a generalized version of existing definitions of static (sFCN) and dynamic (dFCN) functional connectivity and extends networks to circuits. O : zero matrix, I : identity matrix

595 a comparison with our extended model (TVARX) in terms of
596 parameters in our state space equations.

The time-varying feedthrough matrix D is used as dynamic effective connectivity between nodes. Note that effective connectivity defined this way is akin to a general form of Granger causality or transfer entropy (42). U granger-causes (\dot{Y}) if

$$Y = D\dot{U} + T\dot{Y}$$

597 where \dot{U} and \dot{Y} represent lagged values of U and V .

$$y_t = \sum_{i=1}^m d_i u_{t-i} + \sum_{j=1}^n \tau_j y_{t-c} + c$$

null hypothesis: $D = 0$ (lagged values of U do not explain variance in Y)

This dynamic effective connectivity graph is a temporal graph as shown in 4e and can be formally defined as a set of time-stamped edges, each with its own connectivity strength $\{(a, b, t_1, D_{a,b}), (c, d, t_2, D_{c,d}), \dots, (x, y, T, D_{x,y})\}$.

Linear time invariant systems represented in state space form can be converted into input/output transfer functions by

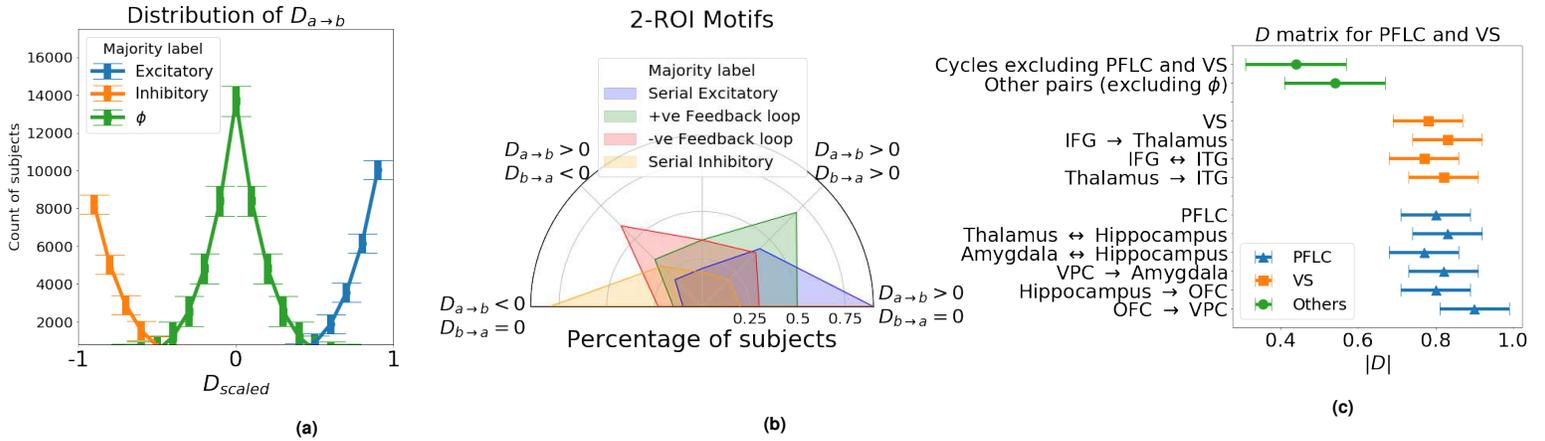


Fig. 5. Three steps of Circuit Discovery (a) As step 1, based on D in equation 3 after fitting on task fMRI, directed relationship between pairs of ROIs is labeled as either Excitatory, Inhibitory or as having no relationship (ϕ). Label for $a \rightarrow b$ is assigned through majority count between subjects. (b) In step 2, based on pairs of directed relationship labels $a \rightarrow b$ from step 1, an undirected label is assigned to each pair of regions, based on individual directed relationships in terms of corresponding values in the D matrix: $D_{a \rightarrow b}$ and $D_{b \rightarrow a}$. These labels are assigned also using majority count across subjects (c) In step 3, once labels of pairs of regions are determined in step 2, cycles are identified via depth-first traversal and the cycle with the largest D value is picked as the active circuit. In our experiments, Prefrontal Limbic Circuit (PFLC) and Ventral Stream (VS) were found to be two cycles with largest total D values for the two stimuli (faces and shapes) respectively.

applying Laplace transform

$$G(s) = \frac{\text{num}(s)}{\text{den}(s)} = \frac{a_0 s^m + a_1 s^{m-1} + \dots + a_m}{b_0 s^n + b_1 s^{n-1} + \dots + b_n} \quad (4)$$

where n is generally greater than or equal to m (for a proper transfer function).

State space systems can be manipulated using standard arithmetic operations as well as the feedback, parallel, and series. Vice versa, each of the connection types: feedback, parallel, and series represent arithmetic operations over state space systems and/or transfer functions as given in Figure 4 Panel c.

The parameters of the TVARX model are learned by maximizing loglikelihood via Kalman filter. The method for calculating the covariance matrix of parameter estimates uses outer product of gradient estimator using Broyden-Fletcher-Goldfarb-Shanno (BFGS) solver. The method by which the Hessian is numerically approximated is outer product of gradients.

The model is fit on seventy-five percent of the time series for each subject and dysregulation is measured as the error in prediction of the remaining twenty-five percent of the time series from the actual signal.

After fitting the TVARX model to data, trajectory drift is measured as error between actual trajectories and trajectories predicted by fitted model. This error is measured using mean squared error. Trajectory drift is used as a measure of feedback control error, or dysregulation, of each subcircuit.

The statistical significance testing between error distributions between different cohort of individuals is carried out using Welsch's t-test to account for skewed distributions between healthy and diseased populations.

Circuit Discovery

A circuit in our framework is defined as a set of connections such that there exists an elementary/Eulerian circuit (simple cycles) of length > 2 . Each connection is defined for a pair of regions and of the following four types: excitatory, inhibitory, negative feedback loop and positive feedback loop defined over elements of the feedthrough matrix D . Excitatory and inhibitory connections between x and y are defined simply as effective connections where $D_{x,y} > 0$ and $D_{x,y} < 0$. Feedback loops are defined as Eulerian/elementary circuits of length = 2. Positive feedback loops are ones where both connections are excitatory. Inversely, negative feedback loops are ones with at least one inhibitory connection i.e.

$$\text{Serial} : Y = f(U), |D| > 0$$

$$\text{Inhibitory} : D < 0$$

$$\text{Excitatory} : D > 0$$

$$\text{Parallel} : Y = f(U, V), |DU| > 0 \wedge |DV| > 0$$

$$\text{Feedback} : Y = f(U) \text{ and } U = f(Y)$$

where $f(U)$ is represented by Eq. (3)

Parallel connections/inputs in the circuit are implicit as BOLD signal for a region y at time t is fitted against multiple inputs. All excitatory and inhibitory connections are series by default.

Discovery of circuits in our experiments is done in the following three steps:

1. Once TVARX model is fitted on human fMRI data, based on fitted values of D matrix in equation 3, directed relationships between pairs of ROIs is labeled

as either excitatory, inhibitory or as having no relationship (ϕ). The fitted value of $D_{a \rightarrow b}$ varied for different subjects. Since, we use a 139-ROI parcellation, the number of directed relationships equal 139^2 . For each of these directed relationships $D_{a \rightarrow b}$, we have a count for each of $count_{D < 0}$, $count_{D = 0}$ and $count_{D > 0}$, where $count_{D > 0} + count_{D = 0} + count_{D < 0} = N(19,831)$. We assign a final label for $a \rightarrow b$ through majority count between subjects i.e. the label that was observed for most subjects was used for all subjects. In our results (Figure 5 panel a), we observed exponential curves, with a majority of labels being ϕ (no relationship).

2. Based on pairs of directed relationship labels $a \rightarrow b$ from step 1, an undirected label is assigned to each pair of ROIs based on individual directed relationships as determined by corresponding values in the D matrix: $D_{a \rightarrow b}$ and $D_{b \rightarrow a}$. Just as in step 1, different subjects have different $D_{a \rightarrow b}$ and $D_{b \rightarrow a}$ values. These conflicts are resolved by assigning a final label based on majority count across subjects (Figure 5 panel b).
3. Once labels of pairs of regions are determined in step 2, cycles are identified via depth-first traversal and the cycle with the largest D value is picked as the active circuit. In our experiments, Prefrontal Limbic Circuit (PFLC) and Ventral Stream (VS) subcircuit were found to be two cycles with largest absolute cumulative D values for the two stimuli (faces and shapes) respectively. The absolute D values for PFLC and VS were found to be significantly larger than for other circuits and relationships between regions not part of PFLC and VS (Figure 5 panel c).

Image Acquisition

Task fMRI data (tfMRI) were acquired on harmonized Siemens 3T Skyra scanners at four UK Biobank imaging centres (Cheadle, Manchester, Newcastle, and Reading). The scans were 2.4mm isotropic with TR of 0.735s and 332 frames per run (4 mins). Each subject had one run. The resolution of the images is 2.4x2.4x2.4 mm with a field-of-view of 88x88x64 matrix. The duration was four minutes (332 timepoints) with TR of 0.735 s and TE of 39ms, GE-EPI with x8 multislice acceleration, no iPAT, flip angle 52 degrees, and fat saturation.

A separate "single-band reference scan" was also acquired, as implemented in the Center for Magnetic Resonance Research (CMRR) multiband acquisition (43). This has the same geometry (including echo-planar imaging distortion) as the timeseries data, but has higher between-tissue contrast to noise, and is used as the reference scan in head motion correction and alignment to other modalities.

Data Preprocessing

Spatial smoothing, using a Gaussian kernel of FWHM 5 mm, was applied before the intensity normalisation, and neither Independent Component Analysis (ICA) (44) nor FMRIB's ICA-based X-noiseifier (FIX) (45) artefact removal was performed, both decisions being largely driven by the shorter timeseries in the tfMRI and because of the greater general reliance in tfMRI analysis on voxelwise timeseries modeling. All time series signal are standardized to z-scores (shifted to zero mean and scaled to unit variance) and the global signal is regressed out.

Pre-processing and task-induced activation modeling was carried out using FEAT (fMRI Expert Analysis Tool); time-series statistical analysis was carried out using FMRIB's Improved Linear Model (FILM) with local autocorrelation correction (46). The timings of the blocks of the two task conditions (shapes and faces) are defined in 2 text files. Display of the task video and logging of participant responses is carried out by ePrime software. The timings of the task blocks are fixed and already known as well as the correctness of subject responses. For more details on data collection, processing of collected images and quality control, please see (47).

Regions of Interest

The parcellation used in our experiments was provided by UK Biobank and included 139 regions of interest (ROIs). These ROIs are defined in MNI152 space, combining parcellations from the following atlases: the Harvard-Oxford cortical and subcortical atlases (20, 21) and the Diedrichsen cerebellar atlas (22).

Task

This task was adapted from the one developed by Hariri and colleagues which had shown evidence as a functional localizer (14) with moderate reliability across time (48). Participants are presented with blocks of trials that either ask them to decide which of two faces presented on the bottom of the screen match the face at the top of the screen, or which of two shapes presented at the bottom of the screen match the shape at the top of the screen. The faces have either angry or fearful expressions. Trials are presented in blocks of 6 trials of the same task (face or shape), with the stimulus presented for 2 s and a 1 s inter trial interval. Each block is preceded by a 3 s task cue ("shape" or "face"), so that each block is 21 s including the cue. Each of the two runs include 3 face blocks and 3 shape blocks.

For facial expressions, 12 different images were used, 6 of each gender and affect (angry or afraid), all derived from a standard set of pictures of facial affect (19). Simple geometric shapes (circles, vertical, and horizontal ellipses) were used as control stimuli.

Subjects were asked to match one of two simultaneously presented images with an identical target image. As a sensorimotor control task, the subjects were asked to match ge-

768 ometric shapes. For each face block, three images of each
769 gender and target affect (angry or fearful) were presented.
770 For each control block, six different geometric shapes were
771 presented as targets. During imaging, subjects responded by
772 pressing one of two buttons with their dominant hand, allow-
773 ing for the determination of accuracy and reaction time.

774 **Clinical Variables**

775 **Neuroticism**

776 Participants were assessed for twelve domains of neurotic
777 behaviours via the touchscreen questionnaire. Neuroticism
778 summarises the number of Yes answers across these twelve
779 questions into a single integer score for each participant. Par-
780 ticipants could answer Yes, No, Do not know or Prefer not to
781 answer. Questions included:

- 782 1. Does your mood often go up and down?
- 783 2. Do you ever feel 'just miserable' for no reason?
- 784 3. Are you an irritable person?
- 785 4. Are your feelings easily hurt?
- 786 5. Do you often feel 'fed-up'?
- 787 6. Would you call yourself a nervous person?
- 788 7. Are you a worrier?
- 789 8. Would you call yourself tense or 'highly strung'?
- 790 9. Do you worry too long after an embarrassing experi-
791 ence?
- 792 10. Do you suffer from 'nerves'?
- 793 11. Do you often feel lonely?
- 794 12. Are you often troubled by feelings of guilt?

795 This derived data field has come from Professor Jill Pell
796 from the Institute of Health and Wellbeing, University of
797 Glasgow (49).

798 **Depression**

799 Depression status of participants is defined from the touch-
800 screen questionnaire at baseline. Each of the three depression
801 states were defined based on a number of criteria:

- 802 1. Ever felt depressed for a whole week
- 803 2. Ever disinterested or unenthusiastic for a whole week
- 804 3. Only 1 episode
- 805 4. ≥ 2 episodes
- 806 5. Episode lasted ≥ 2 weeks
- 807 6. Ever seen a GP for nerves, anxiety, tension or depres-
808 sion
- 809 7. Ever seen a psychiatrist for nerves, anxiety, tension or
810 depression

811 Definitions Single Probable Major Depressive Episode: (1)
812 AND (3) AND (5) AND [(6) OR (7)] OR (1) AND (3) AND
813 (5) AND [(6) OR (7)]

Probable Recurrent Major Depression (Moderate): [(1) OR 814
(2)] AND (4) AND (5) AND (6) 815
Probable Recurrent Major Depression (Severe): [(1) OR (2)] 816
AND (4) AND (5) AND (7) 817

818 **Bipolar Disorder**

UKB data-fields from the touchscreen (which were based on 819
the Structured Clinical Interview for DSM IV Axis I Disor- 820
ders1) were classified into criteria groups to define a probable 821
case of Bipolar I or II. 822

Bipolar I (probable mania) was classified as (1) ever manic 823
or hyper for ≥ 2 days OR ever irritable or argumentative for 824
 ≥ 2 days AND (2) manic episodes characterised by at least 3 825
of 'more talkative', 'more active', 'needed less sleep', 'more 826
creative/more ideas' AND (3) longest manic episode \geq one 827
week duration AND (4) episode needed treatment or caused 828
problems at work. 829

Bipolar II (probable hypomania) classified as fulfilling criteria 830
(1), (2) and (3) of the Bipolar I definition, NOT criteria 831
(4). 832