

# 1 Personalized Mood Prediction from Patterns 2 of Behavior Collected with Smartphones 3

4 Brunilda Balliu<sup>\*1,2,3</sup>, Chris Douglas<sup>4</sup>, Darsol Seok<sup>10</sup>, Liat Shenhav<sup>5</sup>, Yue Wu<sup>5</sup>, Doxa Chatzopoulou<sup>10</sup>, Bill  
5 Kaiser<sup>9</sup>, Victor Chen<sup>9</sup>, Jennifer Kim<sup>10</sup>, Sandeep Deverasetty<sup>10</sup>, Inna Arnaudova<sup>10</sup>, Robert Gibbons<sup>11</sup>, Eliza  
6 Congdon<sup>4</sup>, Michelle G. Craske<sup>4,6</sup>, Nelson Freimer<sup>2,7</sup>, Eran B. Halperin<sup>5,8</sup>, Sriram Sankararaman<sup>1,5,7</sup>,  
7 Jonathan Flint<sup>4,7\*</sup>

8  
9 Departments of <sup>1</sup>Computational Medicine, <sup>2</sup>Pathology and Laboratory Medicine, <sup>3</sup>Biostatistics, <sup>4</sup>Psychiatry and  
10 Biobehavioral Science, <sup>5</sup>Computer Science, <sup>6</sup>Psychology, <sup>7</sup>Human Genetics, <sup>8</sup>Anesthesiology, <sup>9</sup>Electrical  
11 Engineering and <sup>10</sup>Semel Institute for Neuroscience and Human Behavior, University of California Los Angeles,  
12 Los Angeles, USA, <sup>11</sup>Departments of Medicine, Public Health Sciences and Comparative Human Development,  
13 University of Chicago, USA

14

15 \* Corresponding authors (emails: [bballiu@ucla.edu](mailto:bballiu@ucla.edu) and [Jflint@mednet.ucla.edu](mailto:Jflint@mednet.ucla.edu))

16

17

## 18 Abstract

19 Over the last ten years, there has been considerable progress in using digital behavioral phenotypes,  
20 captured passively and continuously from smartphones and wearable devices, to infer depressive mood.  
21 However, most digital phenotype studies suffer from poor replicability, often fail to detect clinically  
22 relevant events, and use measures of depression that are not validated or suitable for collecting large and  
23 longitudinal data. Here, we report high-quality longitudinal validated assessments of depressive mood  
24 from computerized adaptive testing paired with continuous digital assessments of behavior from  
25 smartphone sensors for up to 40 weeks on 183 individuals experiencing mild to severe symptoms of  
26 depression. We apply a novel combination of cubic spline interpolation and idiographic models to  
27 generate individualized predictions of future mood from the digital behavioral phenotypes, achieving high  
28 prediction accuracy of depression severity up to three weeks in advance ( $R^2 \geq 80\%$ ) and a 65.7%  
29 reduction in the prediction error over a baseline model which predicts future mood based on past  
30 depression severity alone. Finally, our study verified the feasibility of obtaining high-quality longitudinal  
31 assessments of mood from a clinical population and predicting symptom severity weeks in advance using  
32 passively collected digital behavioral data. Our results indicate the possibility of expanding the repertoire  
33 of patient-specific behavioral measures to enable future psychiatric research.

## 34 Introduction

35 Major depressive disorder (MDD) affects almost one in five people<sup>1</sup> and is now the world's leading cause  
36 of disability<sup>2</sup>. However, it is often undiagnosed: only about half of those with MDD are identified and  
37 offered treatment<sup>3,4</sup>. In addition, for many people, MDD is a chronic condition characterized by periods of  
38 relapse and recovery that requires ongoing monitoring of symptoms. MDD diagnosis and symptom  
39 monitoring is typically dependent on clinical interview, a method that rarely exceeds an inter-rater

40 reliability of 0.7<sup>5,6</sup>. Furthermore, sufferers are unlikely to volunteer that they are depressed because of the  
41 reduced social contact associated with low mood and because of the stigma attached to admitting to being  
42 depressed. Developing new ways to quickly and accurately diagnose MDD or monitor depressive  
43 symptoms in real time would substantially alleviate the burden of this common and debilitating condition.

44 The advent of electronic methods of collecting information, e.g., smartphone sensors or wearable  
45 devices, means that behavioral measures can now be obtained as individuals go about their daily lives.  
46 Over the last ten years there has been considerable progress in using these digital behavioral phenotypes to  
47 infer mood and depression<sup>7-15</sup>. Yet, most digital mental health studies suffer from one or more of the  
48 following limitations<sup>16-18</sup>. First, many studies are likely underpowered to meet their analytic  
49 objectives<sup>10,12,19,20</sup>. Second, most studies do not follow up subjects long enough to adequately capture  
50 changes in signal within an individual over time<sup>10,11,19,21,22</sup>, even though such changes are highly  
51 informative for clinical care. The few studies with longitudinal assessments use ecological momentary  
52 assessments<sup>19,20,23</sup> to measure state mood, rather than a psychometrically validated symptom scale for  
53 depression. Furthermore, they examine associations between behavior and mood at a population level<sup>23</sup>.  
54 This nomothetic approach is limited by the fact that both mood and its relationship to behavior can vary  
55 substantially between individuals. Last, many of the existing studies focus on healthy subjects, thus  
56 prohibiting evaluation of how well digital phenotypes perform in predicting depression<sup>24</sup>.

57 Here, we overcome these limitations by using a validated measure of depression from  
58 computerized adaptive testing<sup>25</sup> to obtain high-quality longitudinal measures of mood. Computerized  
59 adaptive testing is a technology for interactive administration of tests that tailors the test to the examinee  
60 (or, in our application, to the patient)<sup>26</sup>. Tests are 'adaptive' in the sense that the testing is driven by an  
61 algorithm that selects questions in real time and in response to the on-going responses of the patient. By  
62 employing item response theory to select a small number of questions from a large bank, the test provides

63 a powerful and efficient way to detect psychiatric illness without suffering response fatigue. We also use  
64 smartphone sensing<sup>27</sup> to passively and continuously collect behavioral phenotypes for up to 40 weeks on  
65 183 individuals experiencing mild to severe symptoms of depression (3,005 days with mood assessment  
66 and 29,254 days with behavioral assessment). To account for inter-individual heterogeneity and provide  
67 individual-specific predictors of depression trajectories we use an idiographic (or, personalized) modeling  
68 approach. Ultimately, we expect that this approach will provide patient-specific predictors of depressive  
69 symptom severity to guide personalized intervention, as well as enable future psychiatric research, for  
70 example in genome and phenome-wide association studies.

## 71 Results

### 72 Study participants and treatment protocol

73 Participants (N = 437; 76.5% female, 26.5% white) are University of California Los Angeles (UCLA)  
74 students experiencing mild to severe symptoms of depression or anxiety enrolled as part of the Screening  
75 and Treatment for Anxiety and Depression<sup>28</sup> (STAND) study (Sup Figure 1). The STAND eligibility  
76 criteria and treatment protocol are described extensively elsewhere<sup>29</sup>. Briefly, participants are initially  
77 assessed using the Computerized Adaptive Testing Depression Inventory<sup>31</sup> (CAT-DI), an online adaptive  
78 tool that offers validated assessments of depression severity (measured on a 0-100 scale). After the initial  
79 assessment, participants are routed to appropriate treatment resources depending on depression severity:  
80 those with mild ( $35 \leq \text{CAT-DI} < 65$ ) to moderate ( $65 \leq \text{CAT-DI} < 75$ ) depression at baseline received  
81 online support with or without peer coaching<sup>30</sup> while those with severe depression ( $\text{CAT-DI} \geq 75$ )  
82 received in-person care from a clinician (Materials and Methods).

83 STAND enrolled participants in two waves, each with different inclusion criteria and CAT-DI  
84 assessment and treatment protocol (Sup Figure 2A). Wave 1 was limited to individuals with mild to

85 moderate symptoms at baseline (N=182) and treatment lasted for up to 20 weeks. Wave 2 included  
86 individuals with mild to moderate (N=142) and severe (N=124) symptoms and treatment lasted for up to  
87 40 weeks. Eleven individuals participated in both waves. Depression symptom severity was assessed up to  
88 every other week for the participants that received online support (both waves), i.e., those with mild to  
89 moderate symptoms, and every week for the participants that received in-person clinical care, i.e., those  
90 with severe symptoms (Materials and Methods).

91

## 92 [Adherence to CAT-DI assessment protocol](#)

93 Overall, participants provided a total of 4,507 CAT-DI assessments (out of 11,218 expected by the study  
94 protocols). Participant adherence to CAT-DI assessments varied across enrollment waves (Likelihood  
95 ratio test [LRT] P-value  $< 2.2 \times 10^{-16}$ ), treatment groups (LRT P-value  $< 2.2 \times 10^{-16}$ ), and during the follow-  
96 up period (LRT P-value =  $1.29 \times 10^{-6}$ ). Specifically, participants that received clinical care were more  
97 adherent than those which only received online support (Sup Figure 2B). Attrition for participants which  
98 received clinical care was linear over the follow-up period, with 1.7% of participants dropping out CAT-  
99 DI assessments within two weeks into the study. Attrition for participants that received online support was  
100 large two weeks into the study (33.5% of Wave 1 and 37.3% of Wave 2 participants) and linear for the  
101 remaining of the study. More information about factors impacting study adherence can be found in the  
102 Supplementary Material.

103 For building personalized mood prediction models, we focus on 183 individuals (49 from Wave 1  
104 and 134 from Wave 2) who had at least five mental health assessments during the study (Materials and  
105 Methods). For these individuals we obtained a total of 3,005 CAT-DI assessments with a median of 13  
106 assessments, 171 follow-up days, and 10 days between assessments per individual (Figure 1A-C).

## 107 Computerized adaptive testing captures treatment-related changes in depression severity

108 We assessed what factors contribute to variation in the CAT-DI severity scores (Figure 1E, Materials and  
109 Methods). Subjects are assigned to different treatments (online support or clinical care) depending on their  
110 CAT-DI severity scores, so not surprisingly we see a significant source of variation attributable to the  
111 treatment group (10.3% of variance explained, 95% CI: 8.37 - 12.68%). Once assigned to a treatment  
112 group, we expect to see changes over time as treatment is delivered to individuals with severe symptoms  
113 at baseline. This is reflected in a significant source of variation attributable to the interaction between the  
114 treatment group and the number of weeks spent in the study (8.54% of variance explained, 95% CI: 5.92 -  
115 10.4%) and the improved scores for individuals with severe symptoms at baseline as they spend more time  
116 in the study (Sup Figure 3). We found no statistically significant effect of the COVID pandemic, sex, and  
117 other study parameters. The largest source of variation in depression severity scores is attributable to  
118 between-individual differences (41.78% of variance explained, 95% CI: 38.31 - 42.02%), suggesting that  
119 accurate prediction of CAT-DI severity requires learning models tailored to each individual.

120

## 121 Digital behavioral phenotypes capture changes in behavior

122 We set out to examine how digital behavioral phenotypes change over time for each person and with  
123 CAT-DI severity scores. For example, we want to know how hours of sleep on a specific day for a specific  
124 individual differs from the average hours of sleep in the previous week, or month. To answer these  
125 questions, we extracted digital behavioral phenotypes (referred to hereinafter as features) captured from  
126 participants' smartphone sensors and investigated which features predicted the CAT-DI scores. STAND  
127 participants had the AWARE framework<sup>27</sup> installed on their smartphones, which queried phone sensors to  
128 obtain information about a participant's location, screen on/off behavior, and number of incoming and  
129 outgoing text messages and phone calls. We processed these measurements (Materials and Methods and

130 Supplementary Material) to obtain daily aggregate measures of activity (23 features), social interaction  
131 (18 features), sleep quality (13 features), and device usage (two features). In addition, we processed these  
132 features to capture relative changes in each measure for each individual, e.g., changes in average amount  
133 of sleep in the last week compared to what is typical over the last month. In total, we obtained 1,325  
134 features. Missing daily feature values (Sup Figure 4) were imputed using two different imputation  
135 methods, AutoComplete<sup>31</sup> and softImpute<sup>32</sup> (Materials and Methods), resulting in 29,254 days of logging  
136 events across all individuals.

137 Several of these features map onto the DSM-5 MDD criteria of anhedonia, sleep disturbance, and  
138 loss of energy (Supplementary Material; Sup Figure 5). We computed correlations between these features  
139 and an individuals' depression severity score and found that these features often correlate strongly with  
140 changes in depression (Figure 2A). For example, for one individual, the number of unique locations  
141 visited during the day shows a strong negative correlation with their depression severity scores during the  
142 study (Pearson's  $\rho = -0.65$ , Benjamini-Hochberg [BH]-adjusted P-value =  $2.50 \times 10^{-20}$ ). We observed a lot  
143 of heterogeneity in the strength and direction of the correlation of these features with depression severity  
144 across individuals. For example, features related to location entropy are positively (Pearson's  $\rho=0.40$ , BH-  
145 adjusted P-value =  $3.55 \times 10^{-11}$ ) correlated with depression severity for some individuals but negatively  
146 ( $\rho=-0.59$ , BH-adjusted P-value =  $1.11 \times 10^{-22}$ ) or not correlated ( $\rho=-3.92 \times 10^{-04}$ , BH-adjusted P-value =  
147 0.995) for others. Finally, as expected from the large heterogeneity in these correlation between  
148 individuals, the correlation of these features with depression severity scores across individuals was very  
149 poor, the strongest correlation was observed for the wake-up time ( $\rho=.07$ , BH-adjusted P-value =  $2.47 \times 10^{-03}$ ).  
150

151 Figure 2B illustrates an individual with severe depressive symptoms for whom we can identify a  
152 window of disrupted sleep that co-occurred with a clinically significant increase in symptom severity

153 (from mild to severe CAT-DI scores). Subsequently, a return to baseline patterns of sleep coincided with  
154 symptom reduction. Quantifying this relationship poses a number of issues, which we turn to next.

155

### 156 Predicting CAT-DI scores from digital phenotypes

157 To predict future depression severity scores using digital behavioral phenotypes, we considered three  
158 analytical approaches. First, we applied an idiographic approach, whereby we build a separate prediction  
159 model for each of the participants. Specifically, for each individual, we train an elastic net regression  
160 model using the first 70% of their depression scores and predict the remaining 30% of scores. Second, we  
161 applied a nomothetic approach that used data from all participants to build a single model for depression  
162 severity prediction using the same analytical steps: we train an elastic net regression model using the first  
163 70% of depression scores of each individual and predict the remaining 30% of scores (Materials and  
164 Methods). The result of this nomothetic approach was a single elastic net regression model that makes  
165 predictions in all participants.

166 The main difference between the nomothetic and idiographic approach is that the nomothetic  
167 model assumes that each feature has the same relationship with the CAT-DI scores across individuals, for  
168 example, that a phone interaction is always associated with an increase in depression score. However, it is  
169 possible, and we see this in our data, that an increase in phone interaction can be associated with an  
170 increase in symptom severity for one person, but a decrease in another (Figure 2A). The idiographic  
171 model allows for this possibility by using a different slope for each feature and individual. In addition, we  
172 know that large differences exist in average depression scores between individuals (Figure 1E). To  
173 understand the impact of accounting for these differences in a nomothetic approach, we also applied a  
174 third approach (referred to as nomothetic\*) which includes individual indicator variables in the elastic net



175 regression model in order to allow for potentially different intercepts for each individual. All three models  
176 include stay day as a covariate.

177 To assess whether digital behavioral phenotypes predict mood, we have to deal with the problem  
178 that digital phenotypes are acquired daily, while CAT-DI are usually administered every week (and often  
179 much less frequently, on average every 10 days). We assume that the CAT-DI indexes a continuously  
180 variable trait, but what can we use as the target for our digital predictions when we have such sparsely  
181 distributed measures? We can treat this as a problem of imputation, in which case the difficulty reduces to  
182 knowing the likely distribution of missing values. However, we also assume that both CAT-DI and digital  
183 features only imperfectly reflect a fluctuating latent trait of depression. Thus, our imputation is used not  
184 only to fill in missing data points but also to be a closer reflection of the underlying trait that we are trying  
185 to predict, namely, depressive severity.

186 We interpolate the unmeasured estimates of depression by modeling the latent trait as a cubic  
187 spline with different degrees of freedom (Figure 3A). For many individuals, CAT-DI values fluctuate  
188 considerably during the study, while for others less so. To accommodate this variation, we alter the  
189 degrees of freedom of the cubic spline: the more degrees of freedom, the greater the allowed variation. For  
190 each individual, we used cubic splines with four degrees of freedom, denoted by CS(4df), degrees of  
191 freedom corresponding to the number of observed CAT-DI categories in the training set, denoted by  
192 CS(2-4df), and degrees of freedom identified by leave-one-out cross-validation in the training set, denoted  
193 by CS(cv). For comparison purposes, we also used a last-observation-carried-forward (LOCF) approach,  
194 a naive interpolation method which does not apply any smoothness to the observed trait. In addition, we  
195 also include results from analyses done without interpolating CAT-DI but rather modeling the (bi)weekly  
196 measurements. Because spline interpolation will cause data leakage across the training-testing split and  
197 upwardly bias prediction accuracy, we train our prediction models using cubic spline interpolation on only

198 the training data (first 70% of time series of each individual) and assess prediction accuracy performance  
199 in the testing set (last 30%) using the time series generated by applying cubic splines to the entire time  
200 series (Figure 3B).

201 We evaluated the prediction performance of each model and for each latent trait across and within  
202 participants. We refer to the former as group level prediction and the later as individual level prediction.  
203 Looking at group level prediction performance, compared to within each participant separately, allows us  
204 to compute prediction accuracy metrics, e.g.,  $R^2$ , as a function of the number of days ahead we are  
205 predicting and test for their statistical significance across all predicted observations.

206 We first evaluated group level prediction accuracy. Figure 4 shows group level prediction  
207 performance for each latent trait using the nomothetic, nomothetic\*, and idiographic model when the  
208 features were imputed with Autocomplete and CAT-DI was modeled using a logistic elastic net  
209 regression. We observed that across all latent traits the nomothetic model shows lower prediction accuracy  
210 (mean absolute percentage error [MAPE] = 25-28% and  $R^2 < 5\%$  for all latent traits), compared to the  
211 nomothetic\* (MAPE = 18-25% and  $R^2 = 30-46\%$ ) or idiographic (MAPE = 16-23% and  $R^2 = 37-66\%$ )  
212 models (Figure 4A-B). This is in line with the large proportion of depression scores variance explained by  
213 between-individual differences (Figure 1E) which get best captured by the nomothetic\* and idiographic  
214 models. The idiographic model also showed higher prediction accuracy than the nomothetic(\*) model  
215 when the features were imputed using softImpute or when CAT-DI was modeled using a linear elastic net  
216 regression (Sup Figure 6A-B) as well as when CAT-DI was modeled at the (bi)weekly level without  
217 interpolation to get daily level data (Sup Figure 8A-B).

218 We also compared the prediction performance for each of the different latent traits. As expected,  
219 we achieve a higher prediction accuracy for the more highly penalized cubic spline latent traits compared  
220 to the LOCF latent trait, as the latter has, by default, a larger amount of variation left to be explained by

221 the features. For example, for the idiographic models, we obtained an  $R^2= 66.4\%$  for CS(2-4df) versus  
222 36.9% for LOCF, implying that weekly patterns of depression severity, which are more likely to be  
223 captured by the LOCF latent trait, are harder to predict than depression severity patterns over a couple of  
224 weeks or months, which are more likely to be captured by the cubic spline latent traits with smallest  
225 degrees of freedom.

226 To understand the effect of time on prediction accuracy, we assessed prediction performance as a  
227 function of the number of weeks ahead we are predicting from the last observation in the training set  
228 (Figure 4C). The idiographic models achieved high prediction accuracy for depression scores up to three  
229 weeks from the last observation in the training set, e.g.,  $R^2= 84.2\%$  and  $73.2\%$  for the CS(2-4df) latent  
230 trait to predict observations one week and four weeks ahead, respectively. Prediction accuracy falls below  
231 80% after four weeks.

232 To quantify the contribution of features on group-level prediction accuracy, we assessed to what  
233 extent the features improve the prediction of each model above that achieved by a baseline model that  
234 includes just the intercept and study day. Figure 4D-E shows the log2 fold change in CAT-DI prediction  
235 accuracy, as measured by MAPE and  $R^2$ , of the feature-based model over the baseline model. The  
236 baseline nomothetic model often predicts the same value, i.e., training set intercept, so we cannot compute  
237  $R^2$ . The feature-based idiographic model achieved the greatest improvement in prediction accuracy over  
238 the corresponding baseline model, resulting in 65.7% reduction in the MAPE and 7.1% increase in  $R^2$   
239 over the baseline model for the CS(2-4df) latent trait. The idiographic model also showed higher  
240 prediction accuracy than the corresponding baseline model when the features were imputed using  
241 softImpute or when CAT-DI was modeled using a linear elastic net regression (Sup Figure 6C-D) as well  
242 as when CAT-DI was modeled at the (bi)weekly level (Sup Figure 8C-D). These results suggest that the

243 passive phone features enhance prediction, over and above past CAT-DI and study day, for most  
244 individuals in our study.

245 We next evaluated individual level prediction accuracy (Figure 5). For this analysis, in order to be  
246 able to assess the statistical significance of our prediction accuracy within each individual, we only keep  
247 individuals with at least five mental health assessments in the test set (N=143). In line with the group level  
248 prediction performance, the idiographic model outperformed the other models at the individual level  
249 (Figure 5A; median MAPE across individuals for all latent traits = 13.3 - 18.9% versus 20.1-23% for the  
250 nomothetic and 14.5-20.4% for the nomothetic\* model). Using an idiographic modeling approach, we  
251 significantly predicted the future mood for 79.0% of individuals (113 out of 143 with  $R > 0$  and  $FDR < 5\%$   
252 across individuals) for at least one of the latent traits (Figure 5B), compared to 58.7% and 65.7% of  
253 individuals for the nomothetic and nomothetic\* model, respectively. The median  $R^2$  value across  
254 significantly predicted individuals for the idiographic models was 47.0% (Figure 5C), compared to 23.7.%  
255 and 28.4% for the nomothetic and nomothetic\* model, respectively. In addition, for 41.3% of these  
256 individuals, the idiographic model had prediction accuracy greater than 70%, demonstrating high  
257 predictive power in inferring mood from digital behavioral phenotypes for these individuals, compared to  
258 6.2% and 9.7% for the nomothetic and nomothetic\* model, respectively (Figure 5C). The idiographic  
259 model also outperformed the nomothetic(\*) model when the features were imputed using softImpute or  
260 when CAT-DI was modeled using a linear elastic net regression (Sup Figure 7A) as well as when CAT-DI  
261 was modeled at the (bi)weekly level (Sup Figure 8E).

262 Next, we compared individual-level prediction accuracy of each model against the corresponding  
263 baseline model that includes just the intercept and study day. Figure 5D and Sup Figure 7C-D show the  
264 distribution across individuals of the log2 fold change in CAT-DI prediction accuracy of the feature-based  
265 model over the baseline model. In accordance with the group level prediction performance, the feature-

266 based idiographic model achieved the greatest improvement in prediction accuracy over the corresponding  
267 baseline model, resulting in a median of over two-fold reduction in the MAPE (Figure 5D; median MAPE  
268 of feature-based model across individuals for all latent traits = 13.3 - 18.9% versus 40.1-41.4% for the  
269 baseline model). The idiographic model also showed greatest improvement in prediction accuracy over the  
270 corresponding baseline model than the nomothetic(\*) model when CAT-DI was modeled at the (bi)weekly  
271 level (Sup Figure 8F).

272 To identify the features that most robustly predict depression in each person we extracted top-  
273 feature predictors for each individual's best-fit idiographic model. We limit this analysis to the 113  
274 individuals which showed significant prediction accuracy for at least one of the latent traits. As expected,  
275 the study day was predictive of the mood for 63% of individuals and was mainly associated with a  
276 decrease in symptom severity (median odds ratio [OR] = 0.86 across individuals). Although no behavioral  
277 feature uniformly stood out, as expected by the high correlation between features and heterogeneity in  
278 correlation between features and CAT-DI across individuals (Figure 2A), the variation within the last 30  
279 days in the proportion of unique contacts for outgoing texts and messages (a proxy for erratic social  
280 behavior), the time of last (first) interaction with the phone after midnight (in the morning) (a proxy for  
281 erratic bedtime [wake up time] and sleep quality), and the proportion of time spent at home during the day  
282 (a proxy for erratic activity level) were among the top predictors of future mood and were often associated  
283 with an increase in symptom severity (OR = 1.05 - 1.23 across features and individuals). The heatmap  
284 display of predictor importance in Figure 6 highlights the heterogeneity of passive features for predicting  
285 the future across individuals. For example, poor mental health, as indicated by high CAT-DI depression  
286 severity scores, was associated with decreased variation in location entropy in the evenings (a proxy for  
287 erratic activity level) in the past 30 days for one individual (OR = 0.94) while for another individual it was  
288 associated with increased variation (OR = 1.20).

## 289 Factors associated with prediction performance

290 Using digital behavioral features to predict future mood was useful for 74-77% of our cohort and the  
291 contribution of the features to the prediction performance varies across these individuals. What might  
292 contribute to this variation? Identifying the factors involved might allow us to develop additional models  
293 with higher prediction accuracy. To identify factors that are associated with prediction performance, we  
294 computed the correlation between accuracy metrics (prediction  $R^2$  and MAPE of feature-based model and  
295 difference in MAPE between feature-based and baseline models) with different study parameters e.g.,  
296 treatment group, sex, etc. (Figure 7).

297       Increased variability in depression scores during the study, as measured by the number of unique  
298 CAT-DI categories for each individual, were correlated with poorer prediction performance of the feature-  
299 based model, as measured by MAPE (Spearman's  $\rho=0.49$  and  $0.23$ ,  $p\text{-value} = 2.25 \times 10^{-2}$  and  $9.79 \times 10^{-4}$   
300 for LOCF and CS(4df) latent traits, respectively). In addition, larger differences in median depression  
301 scores between the training and test set for each individual were correlated with poorer prediction  
302 performance, as measured by MAPE (Spearman's  $\rho=0.32$ ,  $p\text{-value} = 9.11 \times 10^{-4}$  for the CS(4df) latent  
303 trait). This suggests that, for some individuals in the study, the training depression scores are higher/lower  
304 than the test depression scores (as expected by Sup Figure 4) and that adding the study day or digital  
305 phenotypes as a predictor does not completely mediate this issue. The size of the training and test set as  
306 well as demographic variables were not strongly correlated to prediction performance.

307       While we had poorer prediction performance for individuals whose mood shows greater variability  
308 during the course of the study, these are also the individuals for which using a feature-based model  
309 improves prediction accuracy compared to a baseline model that predicts based on past depression  
310 severity and study day alone. Specifically, larger variability in depression scores for each individual was  
311 correlated with better prediction performance of a feature-based model than a baseline model, as measured

312 by difference in MAPE between the two models (Spearman's  $\rho=-0.54$  and  $-0.49$ ,  $p\text{-value} = 5.96 \times 10^{-4}$  and  
313  $p\text{-value} = 4.46 \times 10^{-3}$  for the CS(4df) and CS(2-4df) latent traits, respectively).

314

## 315 Discussion

316 In this paper, we showed the feasibility of longitudinally measuring depressive symptoms over 183  
317 individuals for up to 10 months using computerized adaptive testing and passively and continuously  
318 measuring behavioral data captured from the sensors built into smartphones. Using a novel combination of  
319 cubic spline interpolation and idiographic prediction models, we were able to impute and predict a latent  
320 depression trait on a hold-out set of each individual several weeks in advance.

321 Our ability to longitudinally assess depressive symptoms and behavior within many individuals  
322 and over a long period of time enabled us to assess how far out we can predict depressive symptoms, how  
323 variable prediction accuracy can be across different individuals, and what factors contribute to this  
324 variability. In addition, it enabled us to assess the contribution of behavioral features to prediction  
325 accuracy above and beyond that of prior symptom severity or study day alone. We observed that  
326 prediction accuracy dropped below 70% after four weeks. In addition, prediction accuracy varied  
327 considerably across individuals as did the contribution of the features to this accuracy. Individuals with  
328 large variability in symptom severity during the course of the study (such as those in clinical care) were  
329 harder to predict but benefited the most from using behavioral features. We expect that pairing digital  
330 phenotypes from smartphones with behavioral phenotypes from wearable devices, which are worn  
331 continuously and might measure behavior with less error, as well as addition of phenotypes, like those  
332 from electronic health records, could help address some of these challenges.

333 Our results are consistent with other studies that predict daily mood as measured by ecological  
334 momentary assessments or a short screener (i.e., PHQ2<sup>17</sup>) and confirm the superior prediction

335 performance of idiographic models over nomothetic ones. Our study goes further, by exploring if the  
336 superior prediction accuracy of idiographic models is a result of better modeling the relationship between  
337 features and mood or simply of better modeling the baseline mood of each individual. We show that a  
338 large part of the increase in prediction performance of idiographic models is due to the latter, as indicated  
339 by the increase in prediction performance between the nomothetic and modified nomothetic models.

340 High-burden studies over long time periods may result in drop-out, particularly for depressed  
341 individuals<sup>33</sup>. In our case, we observed that attrition for CAT-DI assessment was linear over the follow-up  
342 period, except for the first two weeks during which a large proportion of individuals which received online  
343 support dropped out (typical of online mental health studies<sup>34</sup>). In addition, participants which received  
344 clinical care were more adherent than those which received online support, despite endorsing more severe  
345 depressive symptoms. These participants had regular in-person treatment sessions during which they were  
346 instructed to complete any missing assessments emphasizing the importance of using reminders or  
347 incentives for online mental health studies.

348 There are several limitations in the current study. First, the idiographic models that we use here are  
349 fit separately for each individual and might not thus maximize statistical power. In addition, they assume a  
350 (log-)linear relationship between behavioral features and depression severity and will fit poorly if this  
351 assumption is violated. One potential alternative is to employ mixed models that jointly model data from  
352 all individuals using individual-specific slopes and low degree polynomials. However, due to the high  
353 dimensionality of our data, such models are hard to implement. Second, while it is well established that  
354 Computerized Adaptive Testing can be repeatedly administered to the same person over time without  
355 response set bias due to adaptive question sets<sup>25</sup>, extended use over months might still lead to limited  
356 response bias<sup>35</sup>. Third, the adaptive nature of CAT-DI, which might assess different symptoms for  
357 different individuals, frustrates joint analyses. Fourth, the imputation method used for imputing digital



358 behavioral features assumes data to be missing at random (MAR), meaning missingness depended on  
359 observed data<sup>36</sup>. While this assumption is hard to test, MAR seems quite plausible in our study given that  
360 the data is missing more often for participants that did not receive regular reminders. In addition, research  
361 has shown that violation of the MAR assumption does not seriously distort parameter estimates<sup>37</sup>. Finally,  
362 the age and gender distribution in our participants may limit the generalizability of our findings to the  
363 wider population.

364 In conclusion, our study verified the feasibility of using passively collected digital behavioral  
365 phenotypes from smartphones to predict depressive symptoms weeks in advance. Its key novelty lies in  
366 the use of computerized adaptive testing, which enabled us to obtain high-quality longitudinal assessments  
367 of mood on 183 individuals over many months, and in the use of personalized prediction models, which  
368 offer a much higher predictive power compared to nomothetic models. Ultimately, we expect that the  
369 method will lead to a screening and detection system that will alert clinicians in real-time to initiate or  
370 adapt treatment as required. Moreover, as passive phenotyping becomes more scalable for hundreds of  
371 thousands of individuals, we expected that this method will enable large genome and phenome-wide  
372 association studies for psychiatric genetic research.

373

## 374 **Materials and Methods**

### 375 **Study participants and treatment protocol**

376 Participants are University of California Los Angeles (UCLA) students experiencing mild to severe  
377 symptoms of depression or anxiety enrolled as part of the STAND program<sup>29</sup> developed under the UCLA  
378 Depression Grand Challenge<sup>38</sup> treatment arm. All UCLA students aged 18 or older who had internet  
379 access and were fluent in English were eligible to participate. STAND enrolled participants in two waves.

380 The first wave enrolled participants from April 2017 to June 2018. The second wave of enrollment began  
381 at the start of the academic year in 2018 and continued for three years, during which time, from March  
382 2020, a Safer-At-Home order was imposed in Los Angeles to control the spread of COVID-19. All  
383 participants are offered behavioral health tracking through the AWARE<sup>27</sup> framework and had to install the  
384 app in order to be included in the study. All participants provided written informed consent for the study  
385 protocol approved by the UCLA institutional review board (IRB #16-001395 for those receiving online  
386 support and #17-001365 for those receiving clinical support).

387 Depression symptom severity at baseline and during the course of the study was assessed using the  
388 Computerized Adaptive Testing Depression Inventory<sup>25</sup> (CAT-DI), a validated online mental health  
389 tracker. Computerized adaptive testing is a technology for interactive administration of tests that tailors  
390 the test to the patient<sup>26</sup>. Tests are 'adaptive' in the sense that the testing is driven by an algorithm that  
391 selects questions in real-time and in response to the ongoing responses of the patient. CAT-DI uses item  
392 response theory to select a small number of questions from a large bank, thus providing a powerful and  
393 efficient way to detect psychiatric illness without suffering response fatigue.

394 Participants were classified into treatment groups based on their depression and anxiety scores at  
395 baseline, which indicated the severity of symptoms in those domains. Individuals who are not currently  
396 experiencing symptoms of depression (CAT-DI score < 35) or anxiety are offered the opportunity to  
397 participate in the study without an active treatment component by contributing CAT-DI assessment. These  
398 individuals are excluded from our analyses as they do not show any variation in CAT-DI. Participants that  
399 exhibited scores below the moderate depression range (CAT-DI score < 74) were offered internet-based  
400 cognitive behavioral therapy, which includes adjunctive support provided by trained peers or clinical  
401 psychology graduate students via video chat or in person. Eligible participants with symptoms in this  
402 range were excluded if they were currently receiving cognitive behavioral therapy, refused to install the

403 AWARE phone sensor app, or were planning an extended absence during the intervention period.  
404 Participants that exhibited scores in the range of severe depression symptoms (CAT-DI score 75-100) or  
405 who endorsed current suicidality were offered in-person clinical care which included evidence-based  
406 psychological treatment with option for medication management. Additional exclusion criteria were  
407 applied to participants with symptoms in this range, which included clinically-assessed severe  
408 psychopathology requiring intensive treatment, multiple recent suicide attempts resulting in  
409 hospitalization, or significant psychotic symptoms unrelated to major depressive or bipolar manic  
410 episodes. These criteria were determined through further clinical assessment. Participants with symptoms  
411 in this range were also excluded if they were unwilling to provide a blood sample or transfer care to the  
412 study team while receiving treatment in the STAND program.

413 Depression symptom severity was assessed up to every other week for the participants that  
414 received online support (both waves), i.e., those with mild to moderate symptoms, and every week for the  
415 participants that received in-person clinical care, i.e., those with severe symptoms (Sup Figure 2A).  
416 Participants that received in-person care had also four in-person assessment events, at weeks 8, 16, 28, and  
417 40, prior to the COVID-19 pandemic. Thus, Wave 1 participants can have a maximum of 13 CAT-DI  
418 assessments while Wave 2 participants can have a maximum of 21 (online support) or 44 assessments,  
419 depending on severity and excluding initial assessments prior to treatment assignment.

420 CAT-DI was assessed at least one time for 437 individuals that installed the AWARE app. Here,  
421 we limit our prediction analyses to individuals that have at least five CAT-DI assessments (N=238; since  
422 we need at least four points to interpolate CAT-DI in the training set), have at least 60 days of sensor data  
423 in the same period for which CAT-DI data is also available (N=189), and show variation in their CAT-DI  
424 scores in the training set (N=183), which is necessary in order to build prediction models.

425

## 426 Adherence to CAT-DI assessment protocol and factors affecting adherence

427 To assess if participant adherence to CAT-DI assessments varied across enrollment waves and treatment  
428 groups, we used a logistic regression with the proportion of CAT-DI assessments a participant completed  
429 as the dependent variable and the enrollment waves or treatment groups as independent variables. A  
430 similar model was used to assess impact of sex and age on participant adherence (results presented in the  
431 Supplement). To assess if participant adherence varied with time in the study, we used a logistic  
432 regression random effect model, as implemented in the lmerTest<sup>39</sup> R package, with an indicator variable  
433 for the individual remaining in the study for each required assessment as the dependent variable and a  
434 continuous study week as an independent variable. An individual-specific random effect was used to  
435 account for repeated measurement of each individual during the study. A likelihood ratio test was used to  
436 test for the significance of the effect of each independent variable against the appropriate null model.

## 437 438 Variance partition of CAT-DI metrics

439 We calculate the proportion of CAT-DI severity variance explained by different study parameters using a  
440 linear mixed model as implemented in the R package variancePartition<sup>40</sup> with the subject id, study id,  
441 season, sex, and year modeled as random variables while the day of the study, the age of the subject, and a  
442 binary variable indicating the dates before or after the safer at home order was issued in California  
443 modeled as fixed, i.e.,

$$444 \quad y = \sum_j X_j \beta_j + \sum_k Z_k a_k + \epsilon$$

445 where  $y$  is the vector of the CAT-DI values across all subjects and time points,  $X_j$  is the matrix of  $j^{\text{th}}$  fixed  
446 effect with coefficients  $\beta_j$ ,  $Z_k$  is the matrix corresponding to the  $k^{\text{th}}$  random effect with coefficients  $a_k$   
447 drawn from a normal distribution with variance  $\sigma_{a_k}^2$ . The noise term,  $\epsilon$ , is drawn from a normal  
448 distribution with variance  $\sigma_\epsilon^2$ . All parameters are estimated with maximum likelihood<sup>42</sup>. Variance terms

449 for the fixed effects are computed using the post hoc calculation  $\hat{\sigma}_{\beta_j}^2 = \text{var}(X_j\beta_j)$ . The total variance is  
450  $\hat{\sigma}_{Total}^2 = \hat{\sigma}_{\beta_j}^2 + \hat{\sigma}_{a_k}^2 + \hat{\sigma}_{\epsilon}^2$  so that the fraction of variance explained by the  $j^{\text{th}}$  fixed effect is  $\hat{\sigma}_{\beta_j}^2/\hat{\sigma}_{Total}^2$ , by  
451 the  $k^{\text{th}}$  random effect is  $\hat{\sigma}_{a_k}^2/\hat{\sigma}_{Total}^2$ , and the residual variance is  $\hat{\sigma}_{\epsilon}^2/\hat{\sigma}_{Total}^2$ . Confidence intervals for  
452 variance explained were calculated using parametric bootstrap sampling as implemented in the R package  
453 `variancePartition`<sup>41</sup>.

#### 454 Feature extraction from smartphone sensors

455 We describe feature extraction in detail in the Supplementary Material. Broadly, we  
456 extracted 23 features related to mobility, e.g., location entropy, 13 related to sleep and circadian  
457 rhythm, e.g., hours of uninterrupted sleep, 18 related to social interaction, e.g., duration of  
458 outgoing calls, and two related to mobile device usage, e.g., number of interactions with phone  
459 per day. Each of these features was calculated on a daily basis. Furthermore, each of these  
460 features was computed over three daily non-overlapping time windows of equal duration (night  
461 00:00-08:00, day 08:00-16:00, evening 16:00-00:00), under the hypothesis that participant  
462 behavior may be more or less variable based on external constraints such as a regular class  
463 schedule during daytime hours.

464 In addition, considering a participant's current mental state may be influenced by patterns  
465 of behavior from days prior, sliding window averages of each of the daily features were  
466 calculated over multiple sliding windows ranging from three days to one month prior to the  
467 current day, i.e., windows of length three, seven, 14, and 30 days. The variance of each feature  
468 was also calculated over these same windows, to estimate whether behavior had been stable or  
469 variable during that time, e.g., were there large fluctuations in sleep time over the past week?

470 Finally, under the hypothesis that recent *changes* in behavior may be more indicative of  
471 changes in mental state than absolute measures, a final set of transformations were applied to

472 each feature. These transformations compared the sliding window means of two different  
473 durations against each other, to estimate the change in behavior during one window over that of a  
474 longer duration window (the longer window serving as a local baseline for the participant). This  
475 allowed estimates from the raw features of whether, e.g., the participant had slept less last night  
476 than typical over the past week or slept less on average in the last week than typical over the last  
477 month. All of these transformations were applied to the base features extracted from sensor data  
478 and included as separate features fed into subsequent regression approaches.

479 In total, 1,325 raw and transformed features were extracted and included in the final  
480 analysis.

481

## 482 [Imputation of smartphone-based features](#)

483 To address the missing features problem (Sup Figure 4), we considered two different imputation  
484 methods: matrix completion via iterative soft-thresholder SVD, as implemented in the R package  
485 softImpute, and AutoComplete, a deep-learning imputation method that employs copy-masking to  
486 propagate missingness patterns present in the data. Both approaches were applied separately to each  
487 individual as follows. First, we removed features that exhibited  $> 90\%$  missingness for that individual.  
488 Next, we trained the imputation model on the training split alone. Finally, each imputation model was  
489 applied to the training and test dataset to impute the features for that individual. Before prediction, we  
490 normalize all features to have zero mean and unit standard deviation using mean and standard deviation  
491 estimates from the training set alone.

492

## 493 Imputation and of CAT-DI severity scores for prediction models

494 To get daily-level CAT-DI severity scores, we interpolate the scores for each individual across the  
495 whole time series (ground truth) or only the time series corresponding to the training set (70% of the time  
496 series) by moving the last CAT-DI score forward, denoted by LOCF, or by smoothing the CAT-DI scores  
497 using cubic splines with different degrees of freedom (Figure 3A). Cubic smoothing spline fitting was  
498 done using the *smooth.spline* function from the *stats* package in R. We consider cubic splines with four  
499 degrees of freedom (denoted by CS(4df) and corresponding to the number of possible CAT-DI severity  
500 categories, i.e. normal, mild, moderate, and severe), cubic splines with degrees of freedom equal to the  
501 number of observed CAT-DI categories for each individual in the training set (ranging from two to four  
502 and denoted by CS(2-4df)), and degrees of freedom identified by ordinary leave-one-out cross-validation  
503 in the training set (denoted by CS(cv)).

504

## 505 Nomothetic and idiographic prediction models of future mood

506 We split the data for each individual into a training (70% of trajectory) and a test set (remaining  
507 30% of trajectory). To predict the future mood of each individual in the test set from smartphone-based  
508 features in the test set, we train an elastic net logistic or linear regression model<sup>41</sup> in the train set. We set  
509  $\alpha$ , i.e., the mixing parameter between ridge regression and lasso, to 0.5 and use 10-fold cross-validation to  
510 find the value for parameter  $\lambda$ , i.e., the shrinkage parameter. For the idiographic models, we train separate  
511 elastic net models for each individual while for the nomothetic and modified nomothetic models we train  
512 one model across all individuals. To account for individual differences in the average CAT-DI severity  
513 scores in the training set, the modified nomothetic model fits individual-specific intercepts by including  
514 individual indicator variables in the regression model. This is similar in nature to a random intercept  
515 mixed model where each individual has their own intercept. Note that the test data are the same for all of

516 these models, i.e., the remaining 30% of each individual's trajectories. Predictions outside the CAT-DI  
517 severity range, i.e., [0,100], are set to NA and not considered for model evaluation. We compute  
518 prediction accuracy metrics by computing the Pearson's product-moment correlation coefficient (R)  
519 between observed and predicted depression scores in the test set across and within individuals as well as  
520 the squared Pearson coefficient ( $R^2$ ). To assess the significance of the prediction accuracy we use a one-  
521 sided paired test for Pearson's product-moment correlation coefficient, as implemented in the *cor.test*  
522 function of the stats<sup>42</sup> R package, and a likelihood ratio test for the significance of  $R^2$ . We use the  
523 Benjamini-Hochberg procedure<sup>43</sup> to control the false discovery rate across individuals at 5%.

## 524 Acknowledgments

525 The authors gratefully acknowledge all study participants. S.S was funded in part by NIH grant  
526 R35GM125055 and NSF grants III-1705121 and CAREER-1943497. DS was supported by NSF-NRT  
527 #1829071.

## 528 Author Contributions

529 BB and JF conceived of the project. D.C., V.C., and J.K. participated in the subject recruitment and data  
530 collection. BB lead the data analysis with contributions from CD, LS, AW, and DS. BB and JF wrote the  
531 first draft of the manuscript with contribution from CD and SS. All authors contributed to subsequent edits  
532 of the manuscript and approved the final manuscript.

## 533 Competing Interests

534 The authors declare no competing interests.



## 535 Data and Code Availability

536 The datasets generated and analyzed during the current study are available from the  
537 corresponding author upon reasonable request. The code that supports the findings of this study  
538 is available online at [https://github.com/BrunildaBalliu/stand\\_mood\\_prediction](https://github.com/BrunildaBalliu/stand_mood_prediction).

## 539 References

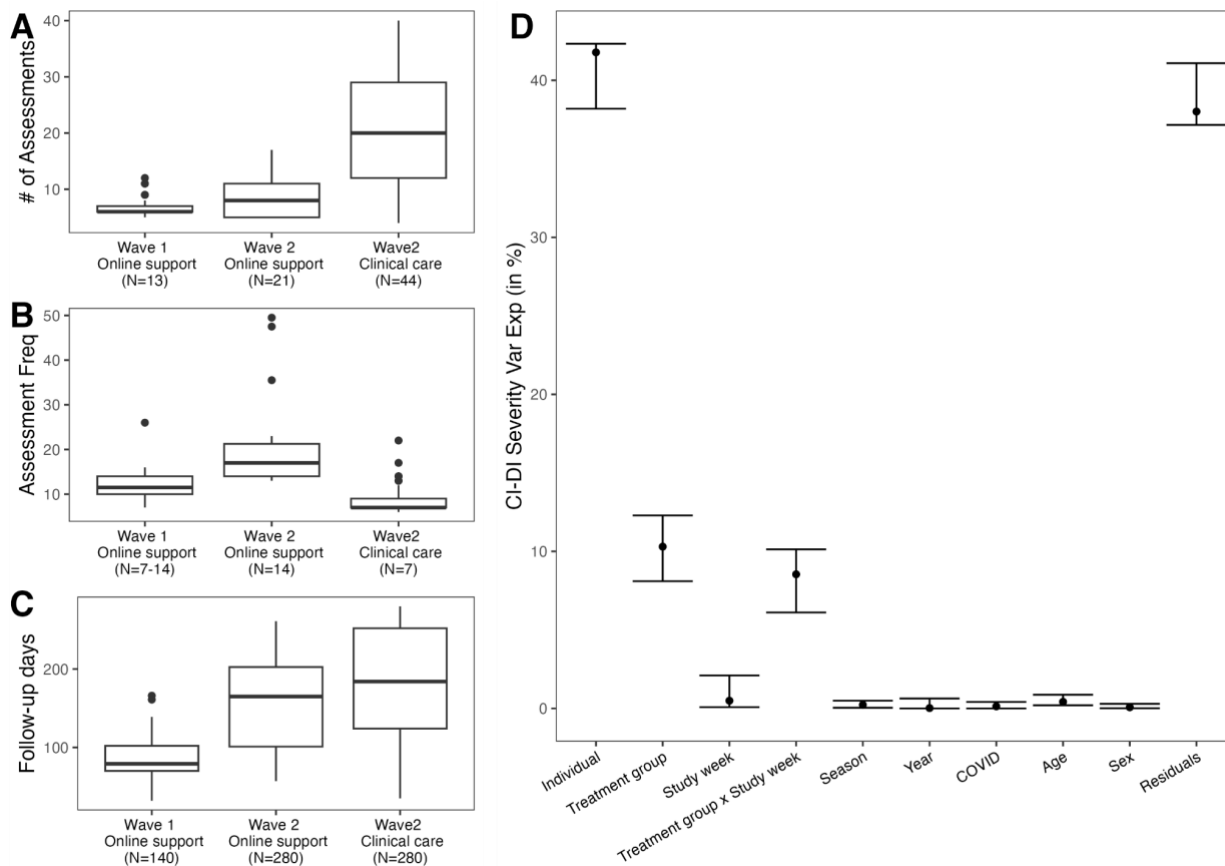
- 540 1. Hasin, D. S. *et al.* Epidemiology of Adult DSM-5 Major Depressive Disorder and Its  
541 Specifiers in the United States. *JAMA Psychiatry* **75**, 336–346 (2018).
- 542 2. World Health Organization. Depression and Other Common Mental Disorders: Global  
543 Health Estimates. Preprint at (2017).
- 544 3. Goldberg, D. Epidemiology of mental disorders in primary care settings. *Epidemiol. Rev.*  
545 **17**, 182–190 (1995).
- 546 4. Wells, K. B. *et al.* Detection of depressive disorder for patients receiving prepaid or fee-  
547 for-service care. Results from the Medical Outcomes Study. *JAMA* **262**, 3298–3302 (1989).
- 548 5. Spitzer, R. L., Forman, J. B. & Nee, J. DSM-III field trials: I. Initial interrater diagnostic  
549 reliability. *Am. J. Psychiatry* **136**, 815–817 (1979).
- 550 6. Regier, D. A. *et al.* DSM-5 field trials in the United States and Canada, Part II: test-retest  
551 reliability of selected categorical diagnoses. *Am. J. Psychiatry* **170**, 59–70 (2013).
- 552 7. Madan, A., Cebrian, M., Lazer, D. & Pentland, A. Social sensing for epidemiological  
553 behavior change. *MIT Web Domain* (2010).
- 554 8. Ma, Y., Xu, B., Bai, Y., Sun, G. & Zhu, R. Daily Mood Assessment Based on Mobile  
555 Phone Sensing. in *2012 Ninth International Conference on Wearable and Implantable Body*  
556 *Sensor Networks* 142–147 (2012). doi:10.1109/BSN.2012.3.
- 557 9. Chen, Z. *et al.* Unobtrusive sleep monitoring using smartphones. in *2013 7th*  
558 *International Conference on Pervasive Computing Technologies for Healthcare and*  
559 *Workshops* 145–152 (2013).
- 560 10. Likamwa, R., Liu, Y., Lane, N. & Zhong, L. MoodScope: Building a Mood Sensor from  
561 Smartphone Usage Patterns. in (2013). doi:10.1145/2462456.2464449.

- 562 11. Frost, M., Doryab, A., Faurholt-Jepsen, M., Kessing, L. & Bardram, J. Supporting  
563 disease insight through data analysis: Refinements of the MONARCA self-assessment  
564 system. *UbiComp 2013 - Proc. 2013 ACM Int. Jt. Conf. Pervasive Ubiquitous Comput.*  
565 (2013) doi:10.1145/2493432.2493507.
- 566 12. Doryab, A., Min, J.-K., Wiese, J., Zimmerman, J. & Hong, J. I. Detection of Behavior  
567 Change in People with Depression. in (2019).
- 568 13. Marsch, L. A. Digital health data-driven approaches to understand human behavior.  
569 *Neuropsychopharmacology* **46**, 191–196 (2021).
- 570 14. Matcham, F. *et al.* Remote Assessment of Disease and Relapse in Major Depressive  
571 Disorder (RADAR-MDD): recruitment, retention, and data availability in a longitudinal  
572 remote measurement study. *BMC Psychiatry* **22**, 136 (2022).
- 573 15. Pratap, A. *et al.* Real-world behavioral dataset from two fully remote smartphone-based  
574 randomized clinical trials for depression. *Sci. Data* **9**, 522 (2022).
- 575 16. Aledavood, T. *et al.* Smartphone-Based Tracking of Sleep in Depression, Anxiety, and  
576 Psychotic Disorders. *Curr. Psychiatry Rep.* **21**, 49 (2019).
- 577 17. De Angel, V. *et al.* Digital health tools for the passive monitoring of depression: a  
578 systematic review of methods. *NPJ Digit. Med.* **5**, 3 (2022).
- 579 18. Zarate, D., Stavropoulos, V., Ball, M., de Sena Collier, G. & Jacobson, N. C. Exploring  
580 the digital footprint of depression: a PRISMA systematic literature review of the empirical  
581 evidence. *BMC Psychiatry* **22**, 421 (2022).
- 582 19. Shah, R. V. *et al.* Personalized machine learning of depressed mood using wearables.  
583 *Transl. Psychiatry* **11**, 1–18 (2021).
- 584 20. Jacobson, N. C. & Bhattacharya, S. Digital biomarkers of anxiety disorder symptom

- 585 changes: Personalized deep learning models using smartphone sensors accurately predict  
586 anxiety symptoms from ecological momentary assessments. *Behav. Res. Ther.* **149**, 104013  
587 (2022).
- 588 21. Burns, R. A., Anstey, K. J. & Windsor, T. D. Subjective well-being mediates the effects  
589 of resilience and mastery on depression and anxiety in a large community sample of young  
590 and middle-aged adults. *Aust. N. Z. J. Psychiatry* **45**, 240–248 (2011).
- 591 22. Melcher, J. *et al.* Digital phenotyping of student mental health during COVID-19: an  
592 observational study of 100 college students. *J. Am. Coll. Health JACH* **71**, 736–748 (2023).
- 593 23. Servia-Rodríguez, S. *et al.* Mobile Sensing at the Service of Mental Well-being: a Large-  
594 scale Longitudinal Study. in *Proceedings of the 26th International Conference on World*  
595 *Wide Web* 103–112 (International World Wide Web Conferences Steering Committee,  
596 2017). doi:10.1145/3038912.3052618.
- 597 24. Stachl, C. *et al.* Predicting personality from patterns of behavior collected with  
598 smartphones. *Proc. Natl. Acad. Sci.* **117**, 17680–17687 (2020).
- 599 25. Gibbons, R. D., Weiss, D. J., Frank, E. & Kupfer, D. Computerized Adaptive Diagnosis  
600 and Testing of Mental Health Disorders. *Annu. Rev. Clin. Psychol.* **12**, 83–104 (2016).
- 601 26. Wainer, H., Dorans, N. J., Flaugher, R., Green, B. F. & Mislevy, R. J. *Computerized*  
602 *Adaptive Testing: A Primer*. (Routledge, 2000).
- 603 27. Ferreira, D., Kostakos, V. & Dey, A. K. AWARE: Mobile Context Instrumentation  
604 Framework. *Front. ICT* **2**, (2015).
- 605 28. UCLA Depression Grand Challenge | Screening and Treatment for Anxiety & Depression  
606 (STAND) Program. <https://www.stand.ucla.edu/>.
- 607 29. A Novel and Integrated Digitally Supported System of Care for Depression and Anxiety:

- 608 Findings From an Open Trial. *JMIR Preprints* <https://preprints.jmir.org/preprint/46200>.
- 609 30. Rosenberg, B. M., Kodish, T., Cohen, Z. D., Gong-Guy, E. & Craske, M. G. A Novel  
610 Peer-to-Peer Coaching Program to Support Digital Mental Health: Design and  
611 Implementation. *JMIR Ment. Health* **9**, e32430 (2022).
- 612 31. An, U. *et al.* Deep Learning-based Phenotype Imputation on Population-scale Biobank  
613 Data Increases Genetic Discoveries. 2022.08.15.503991 Preprint at  
614 <https://doi.org/10.1101/2022.08.15.503991> (2022).
- 615 32. Hastie, T., Mazumder, R., Lee, J. D. & Zadeh, R. Matrix Completion and Low-Rank  
616 SVD via Fast Alternating Least Squares. *J. Mach. Learn. Res.* **16**, 3367–3402 (2015).
- 617 33. DiMatteo, M. R., Lepper, H. S. & Croghan, T. W. Depression Is a Risk Factor for  
618 Noncompliance With Medical Treatment: Meta-analysis of the Effects of Anxiety and  
619 Depression on Patient Adherence. *Arch. Intern. Med.* **160**, 2101–2107 (2000).
- 620 34. Egilsson, E., Bjarnason, R. & Njardvik, U. Usage and Weekly Attrition in a Smartphone-  
621 Based Health Behavior Intervention for Adolescents: Pilot Randomized Controlled Trial.  
622 *JMIR Form. Res.* **5**, e21432 (2021).
- 623 35. Devine, J. *et al.* Evaluation of Computerized Adaptive Tests (CATs) for longitudinal  
624 monitoring of depression, anxiety, and stress reactions. *J. Affect. Disord.* **190**, 846–853  
625 (2016).
- 626 36. Applied Missing Data Analysis: Second Edition. *Guilford Press*  
627 [https://www.guilford.com/books/Applied-Missing-Data-Analysis/Craig-](https://www.guilford.com/books/Applied-Missing-Data-Analysis/Craig-Enders/9781462549863)  
628 [Enders/9781462549863](https://www.guilford.com/books/Applied-Missing-Data-Analysis/Craig-Enders/9781462549863).
- 629 37. Collins, L. M., Schafer, J. L. & Kam, C. M. A comparison of inclusive and restrictive  
630 strategies in modern missing data procedures. *Psychol. Methods* **6**, 330–351 (2001).

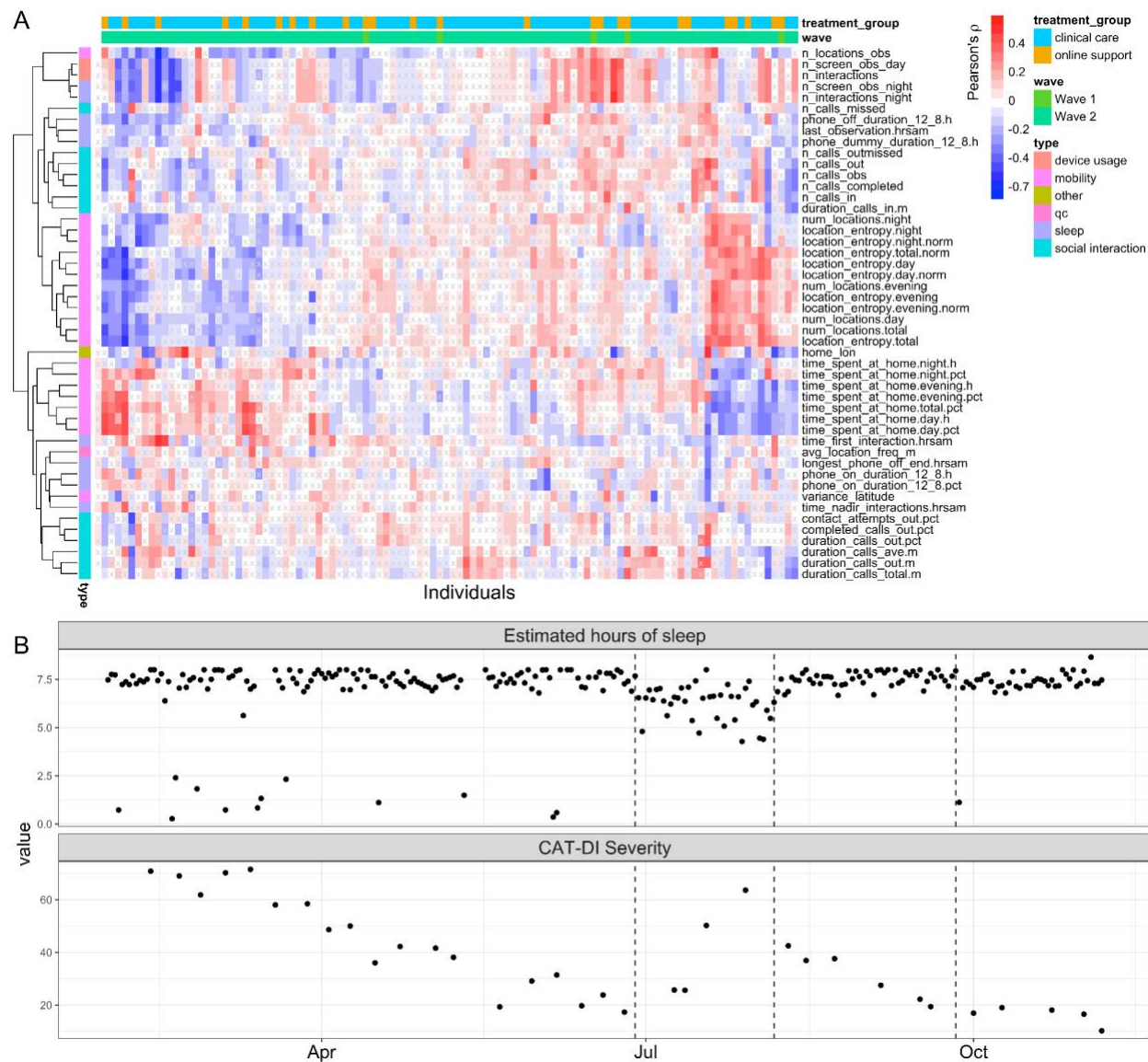
- 631 38. UCLA Depression Grand Challenge. [https://depression.semel.ucla.edu/studies\\_landing](https://depression.semel.ucla.edu/studies_landing).
- 632 39. Kuznetsova, A., Brockhoff, P. B., Christensen, R. H. B. & Jensen, S. P. lmerTest: Tests  
633 in Linear Mixed Effects Models. (2020).
- 634 40. Hoffman, G. E. & Schadt, E. E. variancePartition: interpreting drivers of variation in  
635 complex gene expression studies. *BMC Bioinformatics* **17**, 483 (2016).
- 636 41. Zou, H. & Hastie, T. Regularization and Variable Selection via the Elastic Net. *J. R. Stat.*  
637 *Soc. Ser. B Stat. Methodol.* **67**, 301–320 (2005).
- 638 42. stats-package: The R Stats Package. <https://rdr.io/r/stats/stats-package.html>.
- 639 43. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and  
640 Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300  
641 (1995).
- 642



643

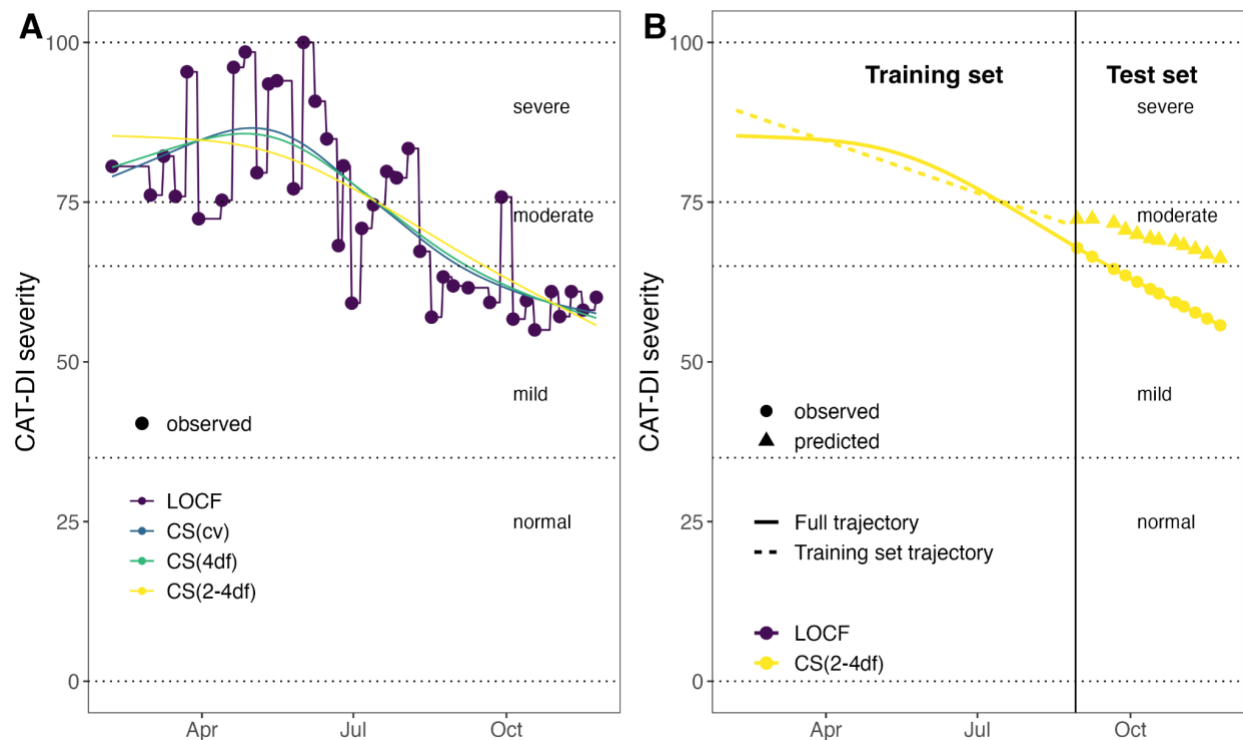
644 *Figure 1: Overview of CAT-DI assessment frequency and source of variation in CAT-DI. (A-C) Boxplot of the observed*  
645 *number of CAT-DI assessments (A), median number of days between assessments (B), and follow-up time in days (C) for each*  
646 *wave and treatment group. The numbers in the parentheses indicate the expected values for each of these metrics according to*  
647 *study design (Sup Figure 2). (D) Proportion of CAT-DI severity variance explained (VE) by inter-individual differences and*  
648 *other study parameters with 95% confidence intervals. The proportion of variance attributable to each source was computed*  
649 *using a linear mixed model with the individual id and season (two multilevel categorical variables) modeled as random variables*  
650 *and all other variables modeled as fixed (see Materials and Methods).*

651

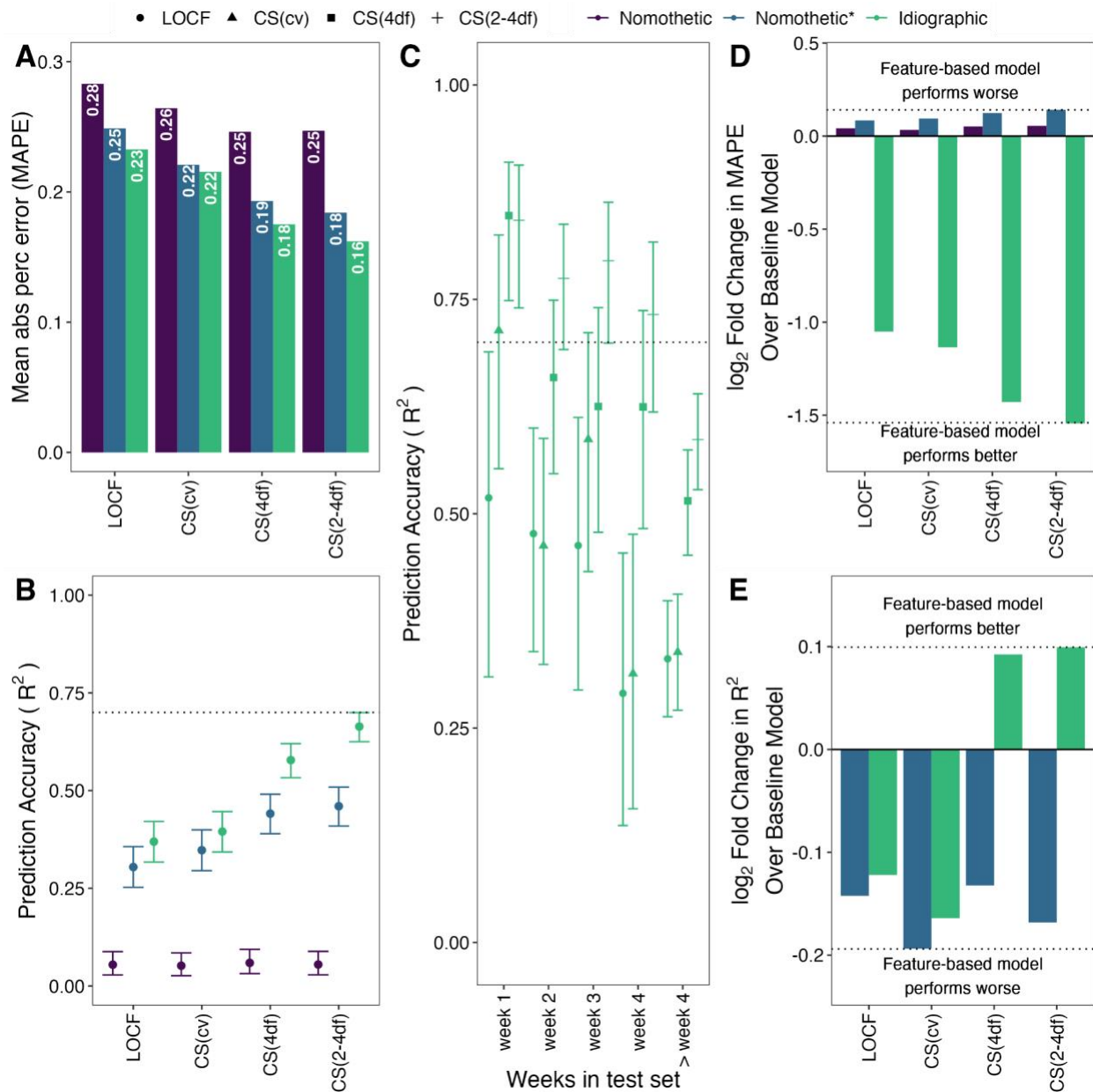


652  
 653 *Figure 2: Overview of correlation between depression severity scores and features. (A) Heatmap for Pearson's correlation*  
 654 *coefficient (color of cell) between CAT-DI scores and behavioral features (y-axis) across individuals (first column) and within*  
 655 *each individual (x-axis). Correlation coefficients with BH-adjusted p-values > 0.05 are indicated by x. For plotting ease, we limit*  
 656 *to untransformed features (N=50, see Materials and Methods). Rows and columns are annotated by feature type and by each*  
 657 *individual's wave and treatment group. Rows and columns are ordered using hierarchical clustering with Euclidean distance.*  
 658 *(B) Example of identifying window of potential sleep disruption using sensor data related to phone usage and screen on/off*  
 659 *status. The top panel shows estimated hours of sleep for an individual during the study while the bottom panel shows the*  
 660 *depression severity scores during the same period. The dotted lines indicate the dates at which a change point is estimated to*  
 661 *have occurred in the estimated hours of sleep as estimated using a change point model framework for sequential change*  
 662 *detection (Materials and Methods). BH: Benjamini Hochberg.*

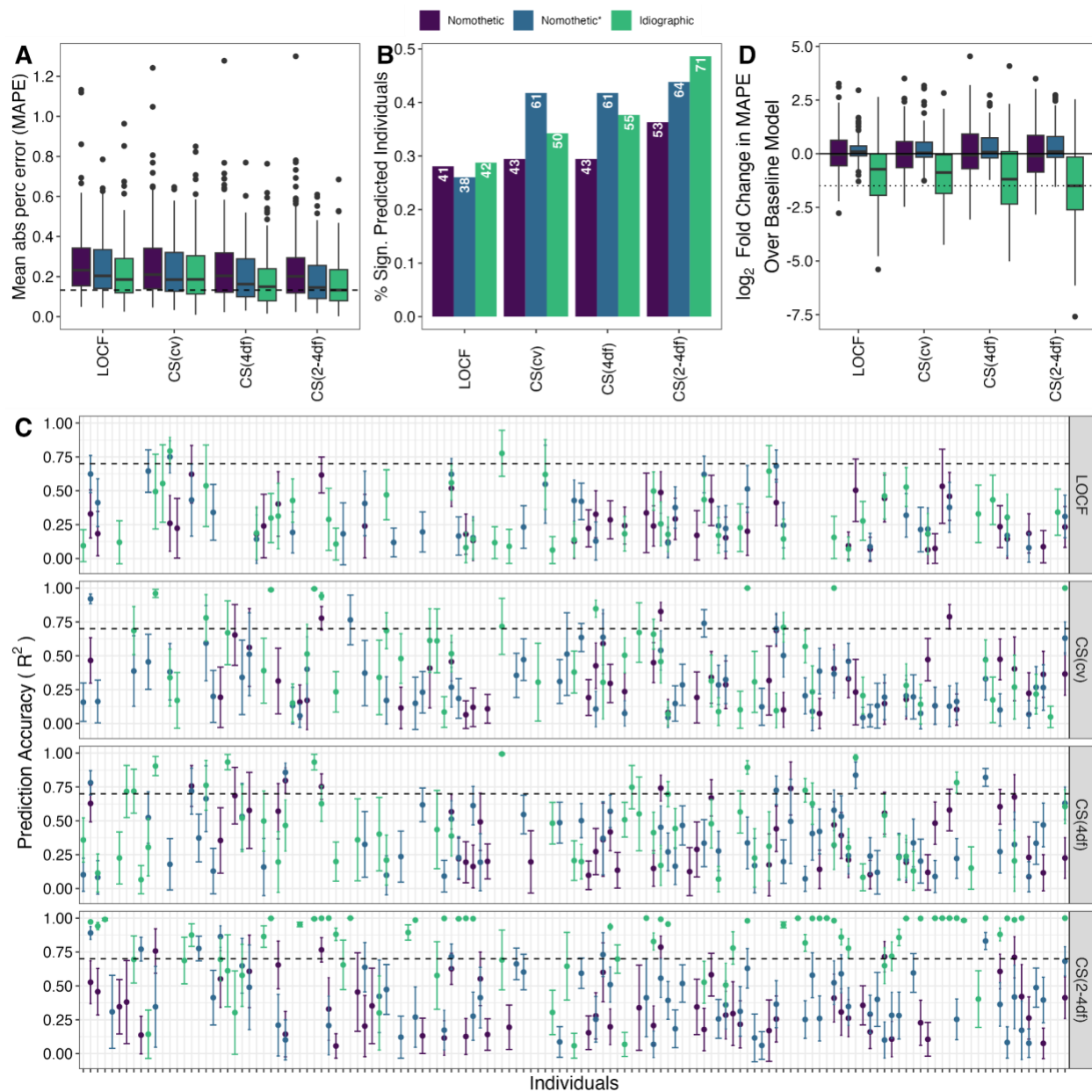




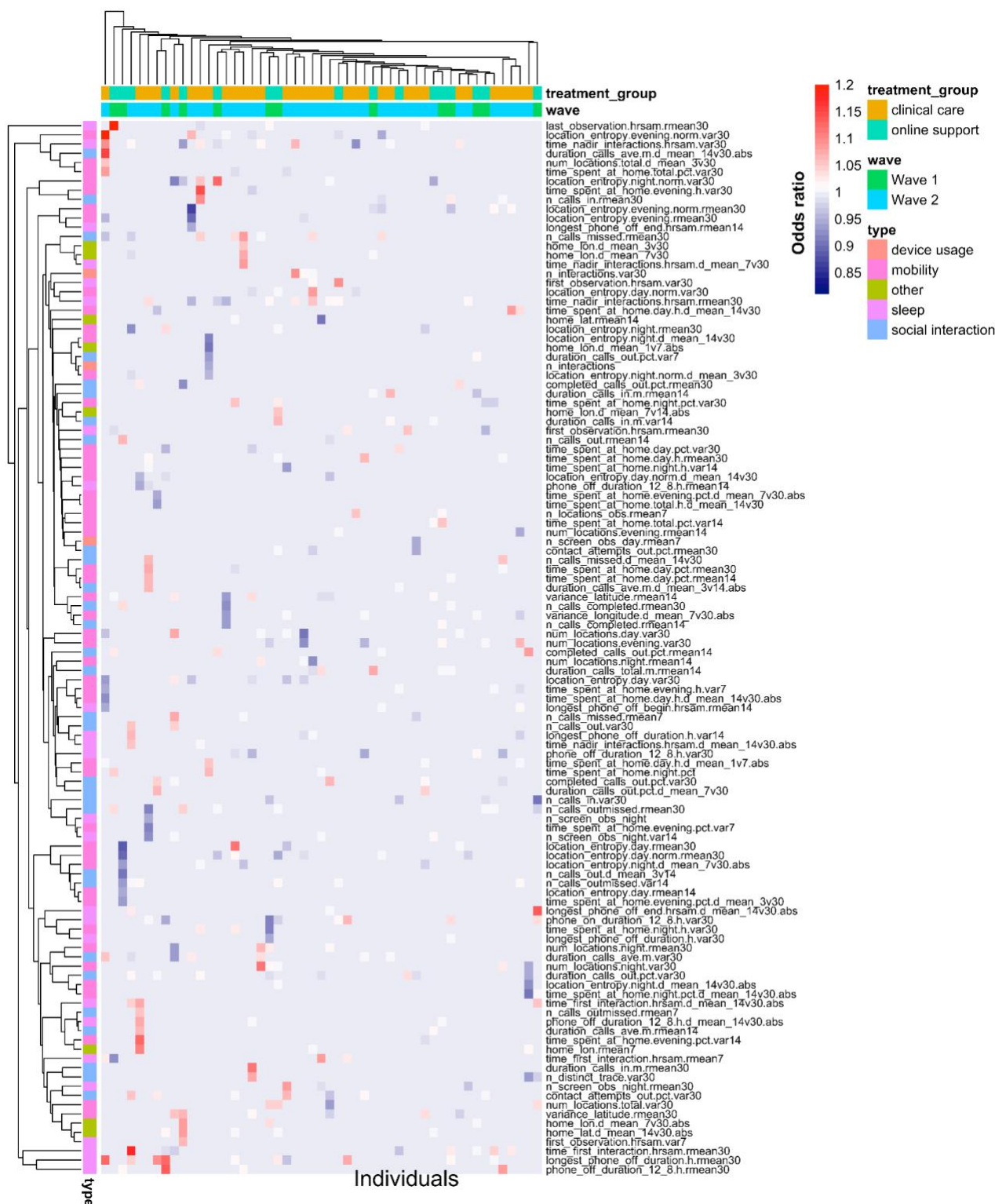
663  
 664 **Figure 3: Interpolation of depression severity scores and latent trait inference.** (A) Illustration of different interpolation  
 665 methods considered for imputing the depression severity scores and inferring the latent depression traits. The dotted horizontal  
 666 lines indicate the depression severity score thresholds for the normal ( $0 \leq \text{CAT-DI} < 35$ ), mild ( $35 \leq \text{CAT-DI} < 65$ ), moderate ( $65$   
 667  $\leq \text{CAT-DI} < 75$ ), and severe ( $75 \leq \text{CAT-DI} \leq 100$ ) depression severity categories. (B) Illustration of the prediction method for the  
 668 CS(2-4df) interpolation method. We first infer the latent trait on the full CAT-DI trajectory of an individual (continuous yellow  
 669 line). We then split the trajectory into a training set (days 1 until  $t$ ) and a test set (days  $t+1$  until  $T$ ), infer the latent trait on the  
 670 training set (dashed yellow line), and predict the trajectory in the test set (yellow triangles). Finally, we compute prediction  
 671 accuracy metrics by comparing the observed (yellow circles) and predicted (yellow triangles) depression scores in the test set.  
 672 We follow a similar approach for the other interpolation methods. The vertical line indicates the first date of the test set  
 673 trajectory, i.e., the last 30% of the trajectory. LOCF: last observation carried forward. CS ( $xd$ ): cubic spline with  $x$  degrees of  
 674 freedom. CS (cv): best-fitting cubic spline according to leave-one-out cross-validation.



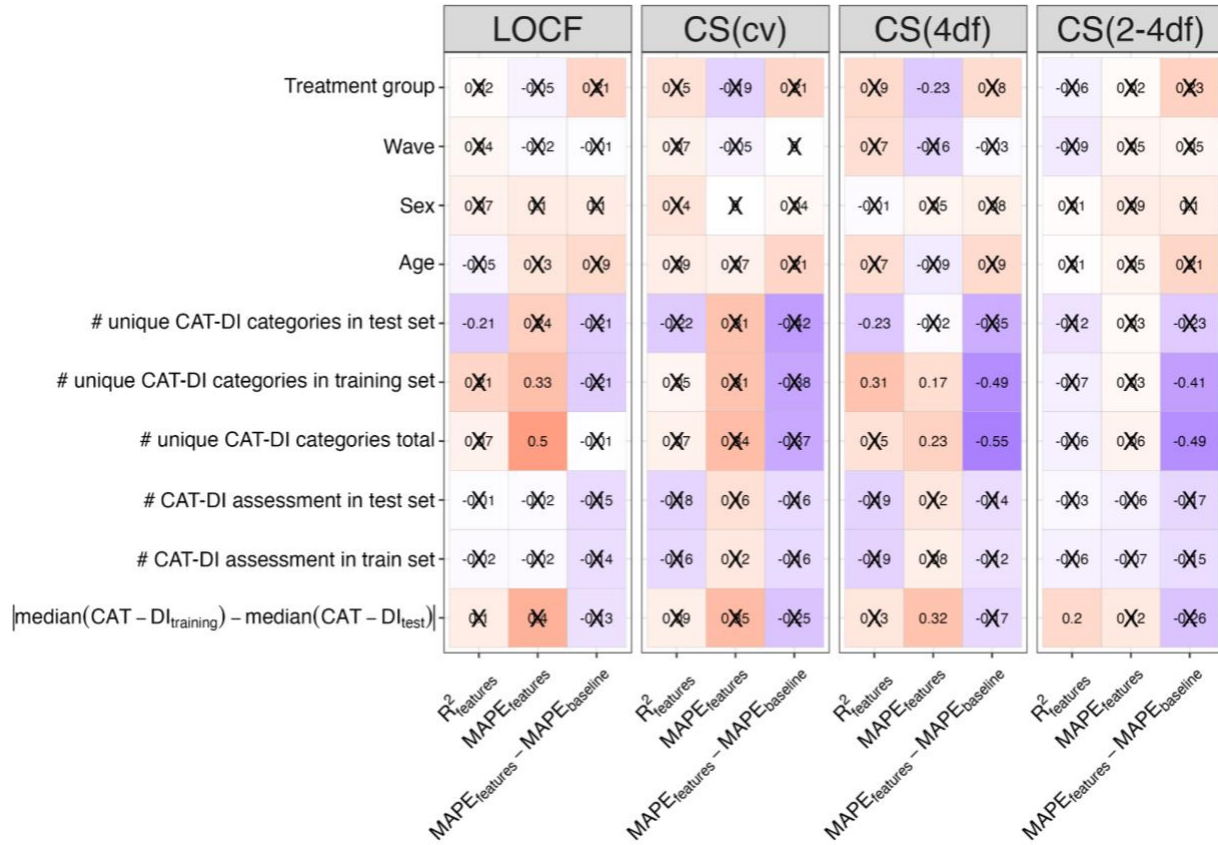
675  
 676 **Figure 4: Idiographic models achieve higher group level prediction accuracy than nomothetic models.** (A-B) CAT-DI  
 677 prediction accuracy across all individuals in the test set as measured by MAPE (A) and  $R^2$  (B) across all individuals for different  
 678 models and latent depression traits. The dotted line in B indicates 70% prediction accuracy and bars indicate 95% confidence  
 679 intervals of  $R^2$ . (C)  $R^2$  versus the number of weeks ahead we are predicting from the last observation in the training set. The  
 680 dotted line indicates 70% prediction accuracy. Bars indicate 95% confidence intervals of  $R^2$ . (D-E)  $\log_2$  fold change in CAT-DI  
 681 prediction accuracy, as measured by MAPE (D) and  $R^2$  (E), of feature-based model over the baseline model. Negative  $\log_2$  fold  
 682 change in MAPE and positive  $\log_2$  fold change in  $R^2$  mean that the feature-based model performs better than the baseline model.  
 683 A  $\log_2$  fold change in MAPE of -1 means that the prediction error of the baseline model is twice as large as that of the feature-  
 684 based model. The dotted line indicates the  $\log_2$  fold change for the best and worse performing model/latent trait combination.  
 685 Features were imputed with Autocomplete and CAT-DI was modeled using a logistic elastic net regression. MAPE: mean  
 686 absolute percent error. LOCF: last observation carried forward. CS(xdf): cubic spline with x degrees of freedom. CS(cv): best-  
 687 fitting cubic spline according to leave-one-out cross-validation.



688  
 689 **Figure 5: Idiographic models achieve higher individual level prediction accuracy than nomothetic models.** (A) Box plots of  
 690 distribution of MAPE across individuals for different models and latent depression traits. The dashed line indicates the median  
 691 MAPE of the best performing model/latent trait combination, i.e., idiographic model and CS(2-4df) spline. (B) Bar plots of the  
 692 proportion of individuals with significantly predicted mood ( $R > 0$  at  $FDR < 5\%$  across individuals) for each latent trait and  
 693 prediction model. (C) Prediction accuracy ( $R^2$ ) with 95% CI across all individuals and latent traits. (D)  $\log_2$  fold change in  
 694 CAT-DI prediction accuracy, as measured by MAPE, of feature-based model over the baseline model. Negative  $\log_2$  fold change  
 695 in MAPE mean that the feature-based model performs better than the baseline model. All plots are based on individuals with at  
 696 least five assessments in the test set ( $N=143$ ). Features were imputed with Autocomplete and CAT-DI was modeled using a  
 697 logistic elastic net regression. LOCF: last observation carried forward. CS(xdf): cubic spline with x degrees of freedom. CS(cv):  
 698 best-fitting cubic spline according to leave-one-out cross-validation.



699  
700 *Figure 6: Most predictive behavioral features according to idiographic models. Heatmap of idiographic elastic net regression*  
701 *coefficients for significantly predicted individuals (N=113 with R>0 and FDR<5%). Columns indicate individuals and rows*  
702 *indicate features. For visualization ease, we limit plot to features that have an odds ratio coefficient value above 1.05 or below*  
703 *0.95 in at least one individual and individuals with at least one feature passing this threshold. The heatmap color indicates the*  
704 *elastic net regularized odds ratio for each feature and individual.*



705  
706  
707  
708  
709  
710

Figure 7: **Factors associated with prediction performance of CAT-DI severity scores.** Correlation between prediction accuracy of an individual (metrics on the y-axis) and the number of CAT-DI assessment available in the training and test set, the difference in median CAT-DI severity between the training and test set, the number of the unique CAT-DI categories (normal to severe) observed (total and in training and test sets), age, sex, wave, and treatment group (a proxy for depression severity). MAPE: mean absolute percentage error.

## 711 Supplementary Material

### 712 Factors impacting study adherence

713 Participant adherence to CAT-DI assessments varied with sex and age. Among participants that  
714 received online support, men were less likely to complete all CAT-DI assessments in wave 1 (OR= 0.86,  
715 LRT P-value =  $2.9 \times 10^{-4}$ ) but more likely to complete them in wave 2 (OR= 1.31, LRT P-value =  $3.1 \times 10^{-11}$ ).  
716 Participant adherence did not vary with sex for those receiving clinical support. In addition, among  
717 participants that received online support in wave 2, older participants were more likely to complete all  
718 CAT-DI assessments than younger participants (OR=1.13, LRT P-value  $< 2.2 \times 10^{-16}$ ). Participant  
719 adherence did not vary with age for participants in wave 1 or those receiving clinical support in wave 2.

720

### 721 Feature extraction from smartphone sensors

#### 722 Preprocessing features

723 Each sensor collected through the AWARE framework is stored separately with a common set of  
724 data items (device identifier, timestamp, etc.) as well as a set of items unique to each sensor  
725 (sensor-specific items such as GPS coordinates, screen state, etc.). Data from each sensor was  
726 preprocessed to convert Unix UTC timestamps into local time, remove duplicate logging entries,  
727 and remove entries with missing sensor data. Additionally, some data labels that are numerically  
728 coded during data collection (e.g., screen state) were converted to human-readable labels for ease  
729 of interpretation.

730

## 731 Mobility features

732 Location data was divided into 24-hour windows starting and ending at midnight each  
733 day. To identify locations where participants spent time, GPS data were filtered to identify  
734 observations where the participants were stationary since the previous observation. Stationary  
735 observations were those defined as having an average speed of  $<0.7$  meters per second  
736 (approximately half the average walking speed of the average adult). These stationary  
737 observations were then clustered using hierarchical clustering to identify unique locations in  
738 which participants spent time during each day. Hierarchical clustering was chosen over k-means  
739 and density-based approaches such as DBSCAN due to its ability to deterministically assign  
740 clusters to locations with a precisely defined and consistent radius, independent of occasional  
741 data missingness.

742 Locations were defined to have a maximum radius of 400 m, a sufficient radius to  
743 account for noise in GPS observations. Clusters were then filtered to exclude any location in  
744 which the participant spent less than 15 minutes over the day to exclude location artifacts, e.g., a  
745 participant being stuck in traffic during daily commute, or passing through the same area of  
746 campus multiple times in a day. To address data missingness in situations where GPS  
747 observations were not received at regular intervals, locations were linearly interpolated to  
748 provide an estimated location every 3 minutes.

749 For each day, a home location was assigned based on the location each participant spent  
750 the most time in between the hours of midnight to eight am. This approach allowed for better  
751 interpretation of behavior for participants who split time between multiple living situations, for  
752 example, students who return home for the weekend or a vacation. Next, multiple features were  
753 extracted from this location data, including total time spent at home each day, total number of

754 locations visited, overall location entropy, and normalized location entropy. Each of these  
755 features was additionally computed over three daily non-overlapping time windows of equal  
756 duration (night 00:00-08:00, day 08:00-16:00, evening 16:00-00:00), under the hypothesis that  
757 participant behavior may be more or less variable based on external constraints such as a regular  
758 class schedule during daytime hours. In total, 28 mobility features were extracted.

759

### 760 [Sleep and circadian rhythm features](#)

761 Sleep and circadian rhythm features were extracted from logs of participant interactions  
762 with their phone, following prior work showing that last interaction with the phone at night can  
763 serve as a reasonable proxy for bedtime, and first interaction in the morning for waketime<sup>46</sup>. The  
764 longest phone-off period (or assumed uninterrupted sleep duration) was tracked each night, as  
765 well as the beginning and end time of that window as estimates of bedtime and waketime. To  
766 account for participants who may have interrupted sleep, the time spent using the phone between  
767 the hours of midnight and 8 am was also tracked to account for participants who may use their  
768 phone briefly in the middle of the night but are otherwise asleep for the majority of that window.  
769 Finally, time-varying kernel density estimates were derived using the total set of phone  
770 interactions, to estimate the daily time nadir of interactions, as an additional proxy for the time of  
771 overall circadian digital activity nadir. In total, 12 sleep and circadian rhythm features were  
772 extracted.

773

### 774 [Social interaction and other device usage features](#)

775 Additional social interaction features were extracted from anonymized logs of participant  
776 calls and text messages sent and received from their smartphone device. Features extracted from



777 this data include, for example, the total number of phone calls made, total time spent on the  
778 phone, and percentage of calls connected that were outgoing (i.e., dialed by the participant)  
779 versus incoming. In total, 18 social interaction and device usage features were extracted. Due to  
780 OS restrictions, sensors needed to extract text message features are not available on iOS devices  
781 and were only computed for the 15 participants with Android devices.

## 782 Mapping of behavioral features to DSM-5 Major Depressive Disorder criteria

783 The set of features described above map onto only a subset of DSM criteria that are closely  
784 associated with externally observable behaviors (Sup Figure 5) - sleep, loss of energy, and anhedonia (to  
785 the extent it is severe enough to globally reduce self-initiated activity). Other DSM criteria such as weight  
786 change, appetite disturbance, and psychomotor agitation/retardation are in theory also directly observable,  
787 but less so with the set of sensors available on a standard smartphone. For these criteria, other device  
788 sensors - for instance, smartwatch sensors - may be more applicable in the detection of e.g., fidgeting  
789 associated with psychomotor agitation. A final set of DSM criteria include those primarily subjective  
790 findings - depressed mood, feelings of worthlessness, suicidal ideation - which inherently require self-  
791 report to directly assess. Given that only 5 of 9 criteria are required for the diagnosis of MDD, an  
792 individual patient's set of symptoms may overlap minimally with those symptoms we expect to measure  
793 with the features described above. However, for others, the above features may cover a more significant  
794 portion of their symptom presentation and do a better job directly quantifying fluctuations in DSM-5  
795 criteria for that individual.

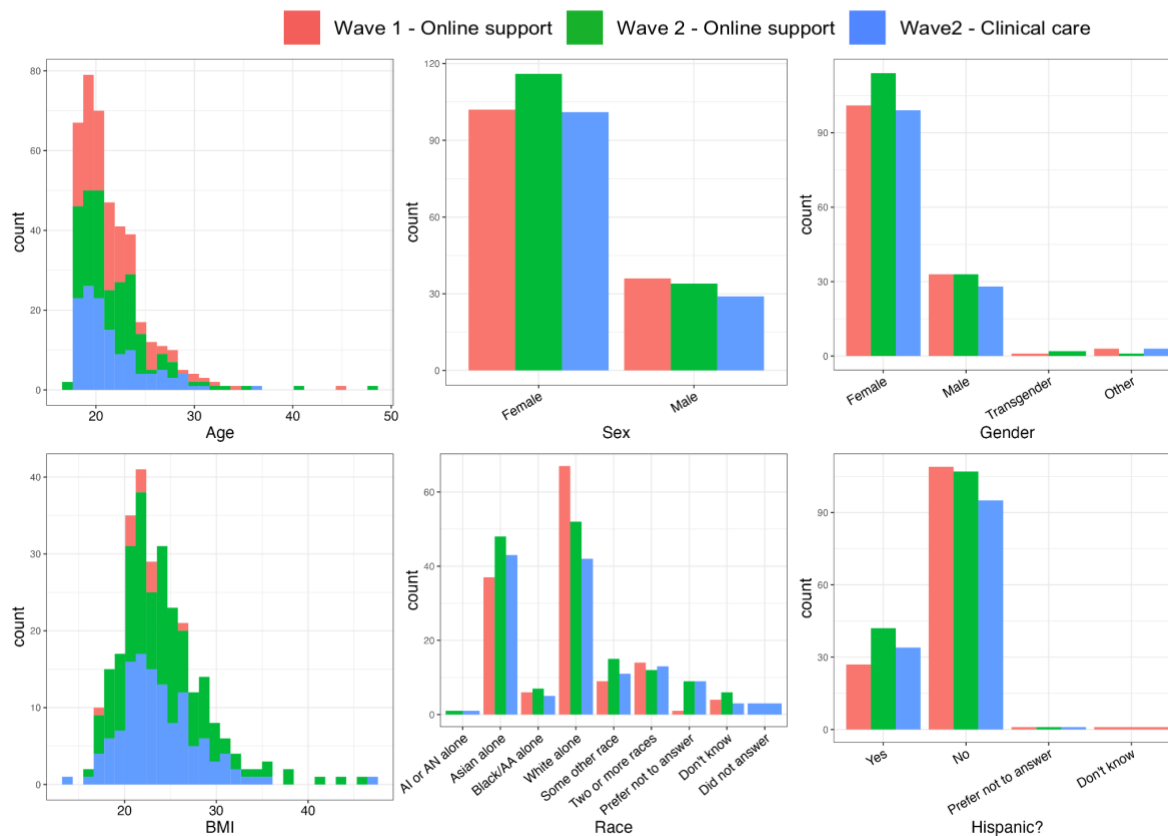
796

797

798

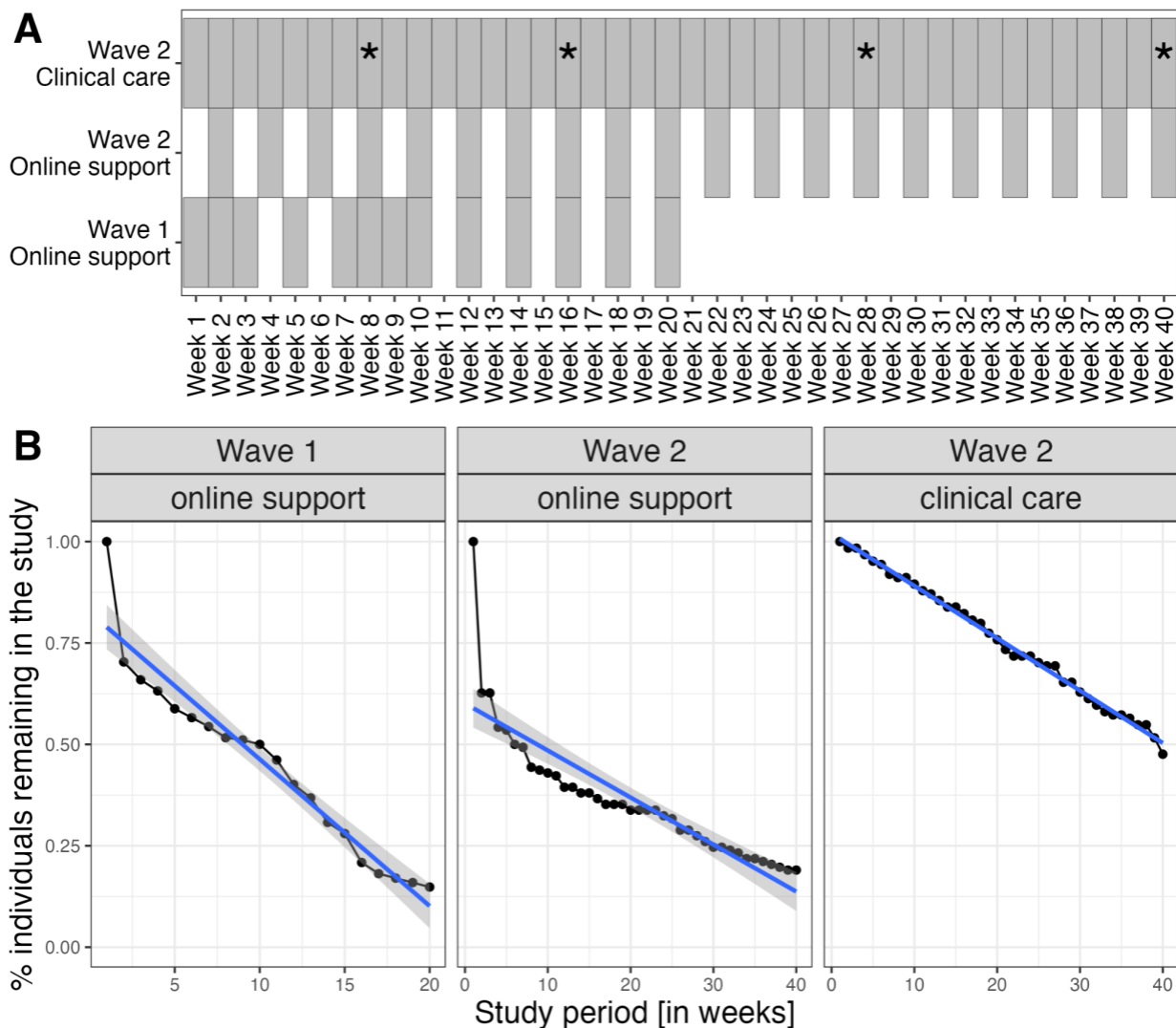
799

800



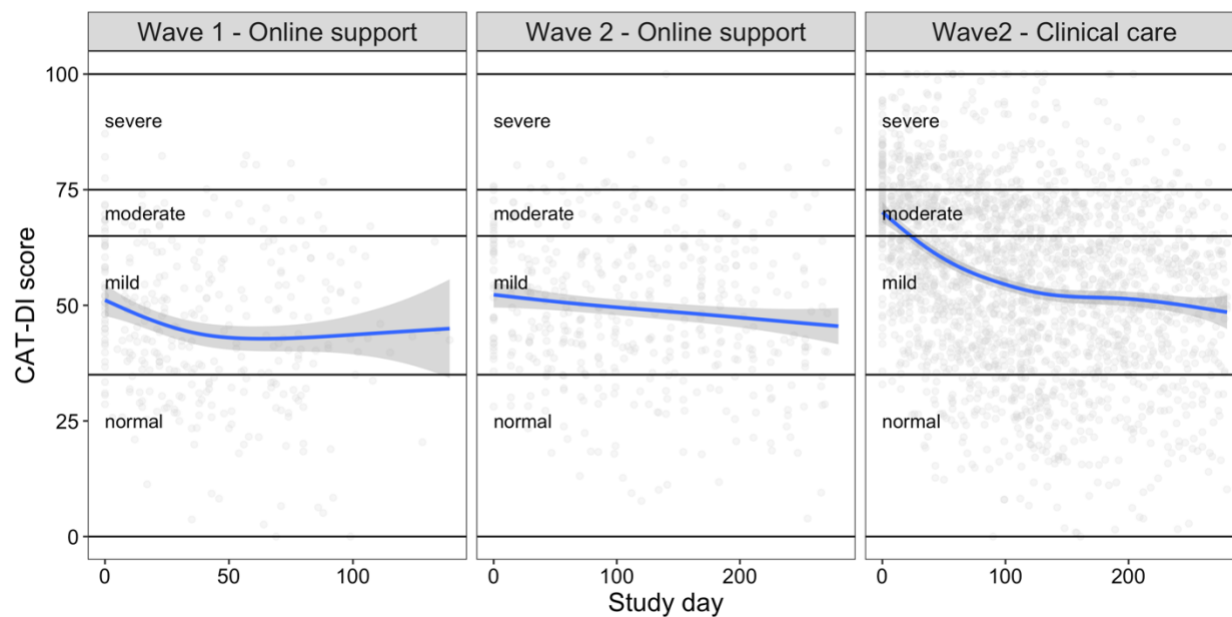
801  
802  
803  
804

*Sup Figure 1: Demographic information for participants in each wave. First row: histogram of age and bar plot of sex and gender. Second row: histogram of BMI and bar plot of race and ethnicity. Color indicates wave and treatment protocol combination. AI or AN: American Indian or Alaska Native. AA: African American.*

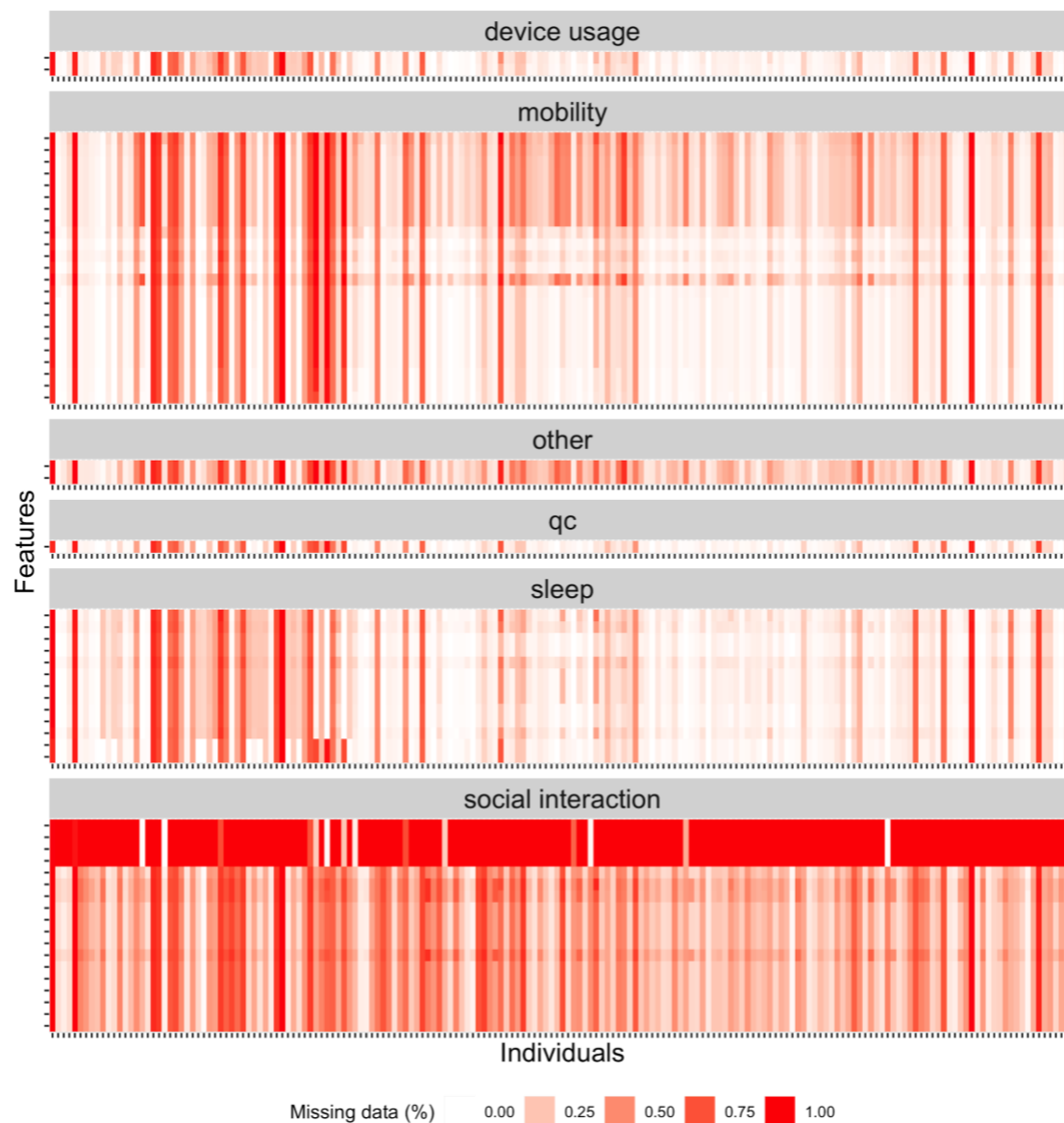


805  
 806 *Sup Figure 2: CAT-DI administration protocol and compliance with CAT-DI assessment protocol for each wave and*  
 807 *treatment group. (A) CAT-DI administration schedule. Each box indicates a week during which participants in each group were*  
 808 *expected to complete the CAT-DI. Asterisks indicate weeks with additional in-person administrations of CAT-DI for Wave 2*  
 809 *participants which received clinical care. (B) Participant CAT-DI retention rate for each enrollment wave and treatment group.*  
 810 *The x-axis shows weeks from the beginning of the study for each participant while the y axis shows the proportion of individuals*  
 811 *that were still completing the CAT-DI at that week. The continuous lines show the linear regression fit with 95% confidence*  
 812 *intervals (gray shading).*

813



814  
815 *Sup Figure 3: Effect of therapy per wave and treatment group. The x-axis shows the study day with zero indicating the first day*  
816 *of CAT-DI assessment for each individual. The y-axis indicates the CAT-DI severity score for each individual / day in the study.*  
817 *The blue line indicates the fit of a generalized additive model with  $y \sim s(\text{day} + \text{wave: treatment group}, \text{bs} = "cs")$  and gaussian*  
818 *family.*

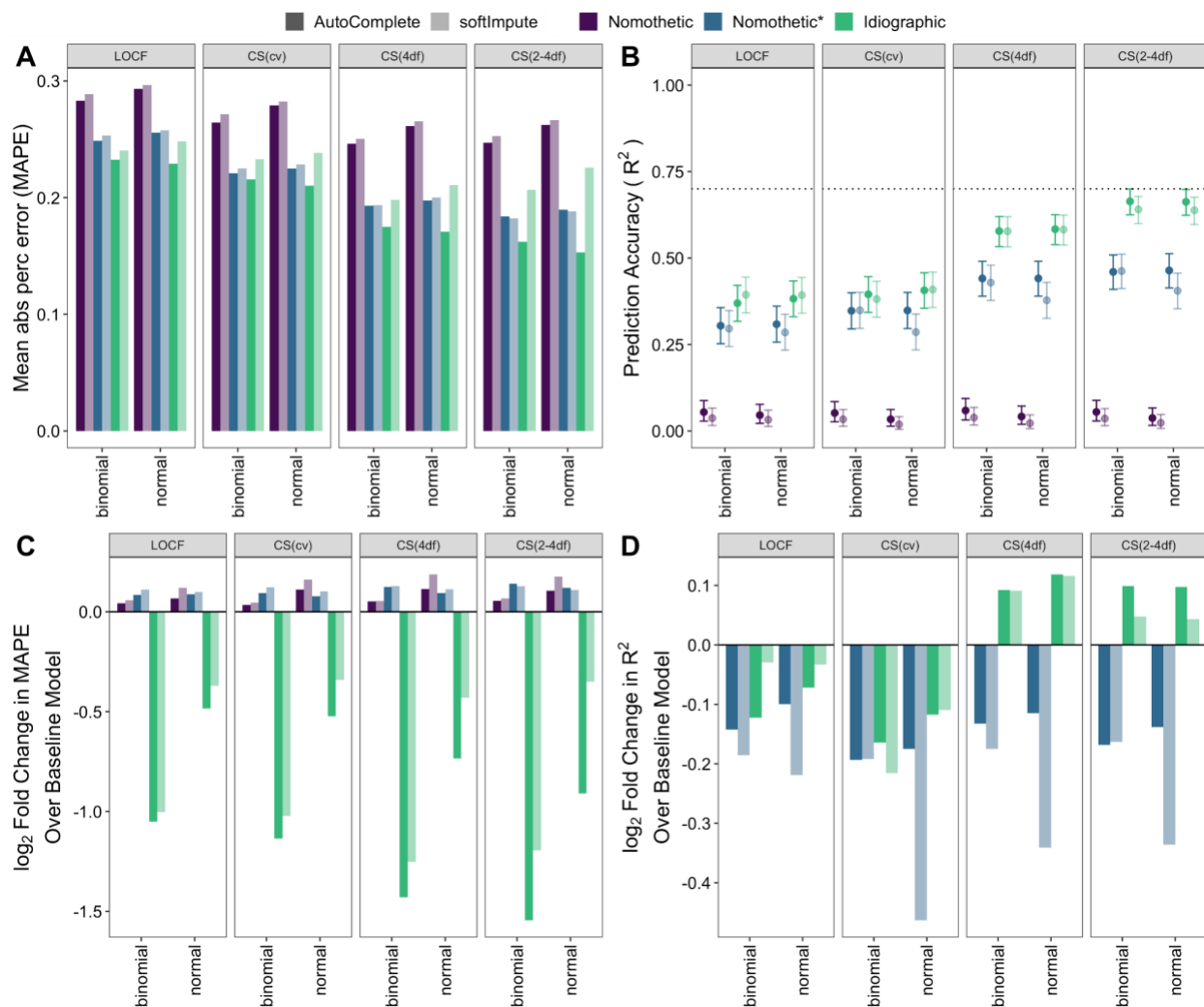


819  
820 *Sup Figure 4: **Missing feature data summary.** Heat map showing missing data percentage in each of the four types of features*  
821 *extracted from smartphone data for all individuals. Each tick on the x-axis (y-axis) represents an individual (feature). For ease of*  
822 *plotting, we have excluded transformation-based features. For participants with iOS devices (majority of individuals), we did not*  
823 *have any information on social interaction features related to text message information due to permission. These features are*  
824 *excluded from analyses when considering individuals with iOS devices.*

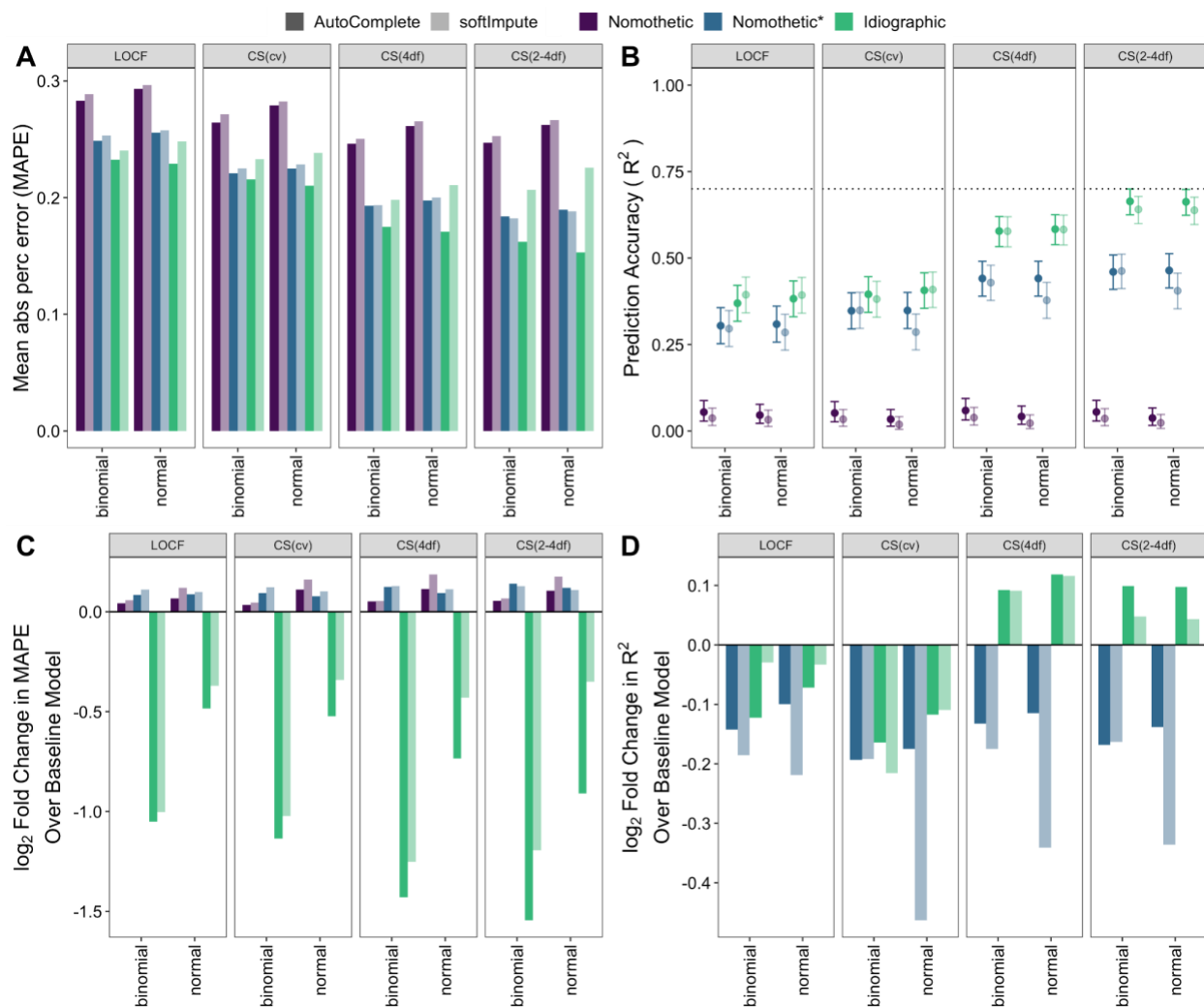
		Depressed mood	Diminished interest or loss of pleasure activities (anhedonia)	Weight change or appetite disturbance	Sleep disturbance	Psychomotor agitation or retardation	Fatigue or loss of energy	Feelings of worthlessness	Diminished concentration; indecisiveness	Suicidal ideation/intent
<b>Mobility features</b>	Variance in latitude	•								
	Variance in longitude	•								
	Number of locations visited in total	•								
	Number of locations visited at night	•								
	Number of locations visited during the day	•								
	Number of locations visited in the evening	•								
	Location entropy over full day	•								
	Location entropy over full day (normalized)	•								
	Location entropy at night	•								
	Location entropy at night (normalized)	•								
	Location entropy during the day	•								
	Location entropy during the day (normalized)	•								
	Location entropy in the evening	•								
	Location entropy in the evening (normalized)	•								
	Time spent at home total (hours)	•					•			
	Time spent at home at night (hours)	•					•			
	Time spent at home during the day (hours)	•					•			
	Time spent at home in the evening (hours)	•					•			
	Percentage of time spent at home in total (%)	•					•			
	Percentage of time spent at home at night (%)	•					•			
Percentage of time spent at home during the day (%)	•					•				
Percentage of time spent at home in the evening (%)	•					•				
<b>Phone interactions</b>	Number of phone interactions during day							•		
	Number of phone interactions at night									
<b>Sleep</b>	Duration of longest phone off period (hours)				•					
	Beginning of longest phone off period				•					
	End of longest phone off period				•					
	Phone on duration midnight to 8am (hours)				•					
	Phone off duration midnight to 8am (hours)				•					
	Percentage of time phone on midnight to 8am (%)				•					
	Time of nadir of phone interactions				•					
<b>Social interaction</b>	Total number of calls	•								
	Number of incoming calls attempted	•								
	Number of outgoing calls attempted	•								
	Number of completed calls	•								
	Number of missed calls	•								
	Number of unanswered outgoing calls	•								
	Percentage of completed outgoing calls (%)	•								
	Percentage of unanswered outgoing calls (%)	•								
	Duration of incoming calls (minutes)	•								
	Duration of outgoing calls (minutes)	•								
	Total time spent on phone (minutes)	•								
	Average duration per call (minutes)	•								
	Duration of outgoing calls (%)	•								
	Number of distinct message contacts	•								
	Total number of messages	•								
	Number of incoming messages	•								
	Number of outgoing messages	•								
	Percent of messages outgoing	•								

825  
826  
827  
828  
829

Sup Figure 5: **Mapping of sensor-derived behavioral features to DSM5 Major Depressive Disorder criteria.** The individual behavioral features derived from phone sensors map primarily to the DSM criteria of disrupted sleep, loss of energy, and anhedonia. Each of these base features is further transformed to look for deviations from individual baseline over varying time scales (e.g., last day's deviation from the weekly average) to arrive at the final set of behavioral features.

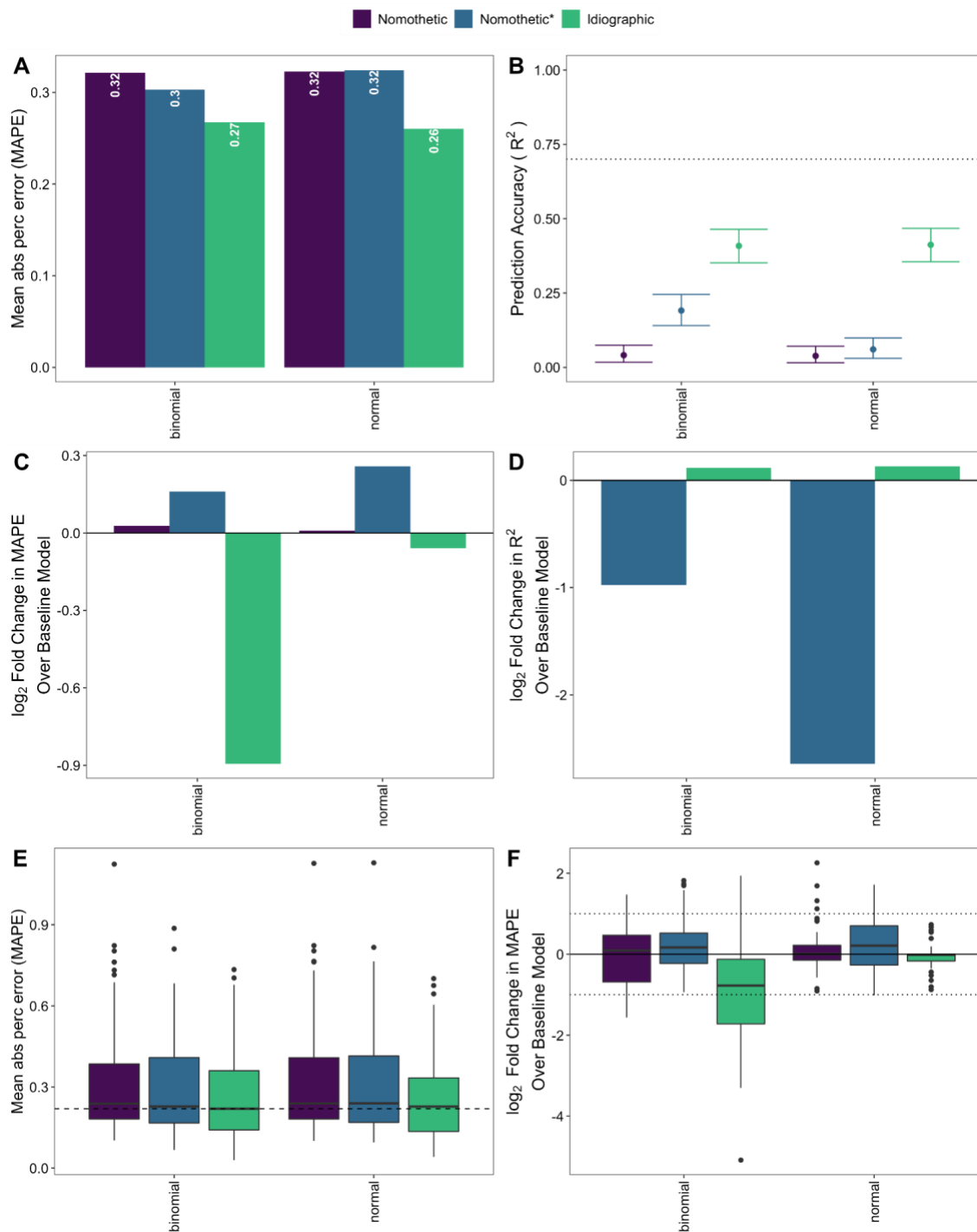


830  
 831 *Sup Figure 6: Idiographic models achieve higher group level prediction accuracy than nomothetic models. (A-B) CAT-DI*  
 832 *prediction accuracy across all individuals in the test set as measured by MAPE (A) and R<sup>2</sup> (B) across all individuals for different*  
 833 *latent depression traits (panel), modeling approaches (color), CAT-DI regression model (x-axis), and feature imputation methods*  
 834 *(transparency). The dotted line in B indicates 70% prediction accuracy and bars indicate 95% confidence intervals of R<sup>2</sup>. (C-D)*  
 835 *log<sub>2</sub> fold change in CAT-DI prediction accuracy, as measured by MAPE (C) and R<sup>2</sup> (D), of feature-based model over the*  
 836 *baseline model for different latent depression traits, modeling approaches, CAT-DI regression models, and feature imputation*  
 837 *methods. Negative log<sub>2</sub> fold change in MAPE and positive log<sub>2</sub> fold change in R<sup>2</sup> mean that the feature-based model performs*  
 838 *better than the baseline model. A log<sub>2</sub> fold change in MAPE of -1 means that the prediction error of the baseline model is twice*  
 839 *as large as that of the feature based model. MAPE: mean absolute percent error. LOCF: last observation carried forward.*  
 840 *CS(xdf): cubic spline with x degrees of freedom. CS(cv): best-fitting cubic spline according to leave-one-out cross-validation.*



841 **Sup Figure 7: Idiographic models achieve higher individual-level prediction accuracy than nomothetic models.** Box plots of  
 842 MAPE distribution of feature-based model (A) and log<sub>2</sub> fold change in CAT-DI prediction accuracy of feature-based model over  
 843 the baseline model (B) across individuals for different latent depression traits (panel), modeling approaches (color), CAT-DI  
 844 regression model (x-axis), and feature imputation methods (transparency). All plots are based on individuals with at least five  
 845 assessments in the test set (N=143). LOCF: last observation carried forward. CS(xdf): cubic spline with x degrees of freedom.  
 846 CS(cv): best-fitting cubic spline according to leave-one-out cross-validation.  
 847





848  
 849 **Sup Figure 8: Idiographic models achieve higher group and individual level prediction accuracy than nomothetic models**  
 850 **when CAT-DI is modeled at the (bi)weekly level.** (A-D) Population-level CAT-DI prediction accuracy of the feature-based  
 851 model (A-B) and  $\log_2$  fold change in population-level CAT-DI prediction accuracy of the feature-based model over the baseline  
 852 model (C-D) as measured by MAPE (A and C) and  $R^2$  (B and D) across all individuals in the test set for different modeling  
 853 approaches (color) and CAT-DI regression model (x-axis). The dotted line in B indicates 70% prediction accuracy and the bars  
 854 indicate 95% confidence intervals of  $R^2$ . (E-F) Box plots of MAPE distribution of feature-based model (E) and  $\log_2$  fold change  
 855 in CAT-DI prediction accuracy of feature-based model over the baseline model (F) in the test set across individuals for different  
 856 modeling approaches (color) and CAT-DI regression model (x-axis). Plots E-F are based on individuals with at least five  
 857 assessments in the test set ( $N=143$ ). For all six plots, features were imputed using AutoComplete. MAPE: mean absolute  
 858 percentage error.