

# 1 The Challenges and Prospects of Brain-based Prediction of 2 Behavior

3 Jianxiao Wu<sup>1,2\*</sup>, Jingwei Li<sup>1,2</sup>, Simon B. Eickhoff<sup>1,2</sup>, Dustin Scheinost<sup>3,4,5,6,7</sup>, and  
4 Sarah Genon<sup>1,2\*</sup>

5 <sup>1</sup>Institute of Neuroscience and Medicine, Brain & Behaviour (INM-7), Research  
6 Center Jülich, Jülich, Germany

7 <sup>2</sup>Institute for Systems Neuroscience, Medical Faculty, Heinrich-Heine University  
8 Düsseldorf, Düsseldorf, Germany

9 <sup>3</sup>Department of Radiology and Biomedical Imaging, Yale School of Medicine, New  
10 Haven, USA

11 <sup>4</sup>Department of Statistics and Data Science, Yale University, New Haven, USA

12 <sup>5</sup>Child Study Center, Yale School of Medicine, New Haven, USA

13 <sup>6</sup>Interdepartmental Neuroscience Program, Yale School of Medicine, USA

14 <sup>7</sup>Department of Biomedical Engineering, Yale School of Engineering and Applied  
15 Sciences, New Haven, USA

16

17 \* Corresponding Author. Email address: [j.wu@fz-juelich.de](mailto:j.wu@fz-juelich.de); [s.genon@fz-juelich.de](mailto:s.genon@fz-juelich.de)

## 18 Abstract

19 Relating individual brain patterns to behavior is fundamental in system neuroscience. Recently,  
20 the predictive modeling approach has become increasingly popular, largely due to the recent  
21 availability of large open datasets and access to computational resources. This means that we can  
22 use machine learning models, and interindividual differences at the brain level represented by  
23 neuroimaging features to predict interindividual differences in behavioral measures. By doing so,  
24 we could identify biomarkers and neural correlates in a data-driven fashion. Nevertheless, this  
25 budding field of neuroimaging-based predictive modelling is facing issues that may limit its  
26 potential applications. Here, we review these existing challenges, as well as those that we  
27 anticipate as the field develops. We focus on the impact of these challenges on brain-based  
28 predictions. We suggest potential solutions to address the resolvable challenges, while keeping in  
29 mind that some general and conceptual limitations may also underlie the predictive modeling  
30 approach.

31

32 The study of the relationships between individual differences in brain phenotypes and individual  
33 behaviors is fundamental in neuroscience, both from a basic scientific perspective and an applied  
34 perspective. The term ‘predictive modeling’ refers to the use of machine learning techniques to  
35 build a statistical model for the estimation of behavioral variables from brain-based neuroimaging  
36 data, either structural or functional<sup>1,2</sup>. More precisely, a prediction model is trained to predict  
37 particular behavioral variables from brain-based data from a number of individuals (the training  
38 set), and its performance is then evaluated on unseen data (test set).

39 The potential practical applications promised by such prediction approaches in precision medicine,  
40 healthcare, human resources and education<sup>1,3-5</sup> are certainly exciting. Potential future applications  
41 may include prediction of individual treatment outcomes to guide treatment choices and dosage,  
42 classification of clinical subgroups with different brain pathology and thus different treatment  
43 requirements, as well as prediction of future cognitive abilities and mental health at developmental  
44 stage. As concrete examples that could be envisioned, brain-based predictions may provide  
45 objective biomarkers when evaluating the effect of cognitive training or cognitive-behavior  
46 therapies (e.g., for mild functional cognitive alterations and anxio-depressive phenotypes,  
47 respectively). While the effect of these interventions could be more readily investigated with  
48 standard cognitive tests and interview/questionnaires respectively, such approaches are prone to  
49 many biases (e.g., practice effects, subjectivity biases, expectations biases). As a recent working  
50 example, a prediction model of sustained attention provided a neuromarker of sustained attention<sup>6</sup>.  
51 This neuromarker can be used both for predicting attention deficit symptoms, and for localizing  
52 targets of potential brain-based treatments. Ultimately, brain-based prediction could be expected  
53 to provide objective biomarkers that can inform us about the brain mechanisms behind the effects  
54 under scientific investigations. Aided by the publicly available large neuroimaging datasets,  
55 accessible computational resources, as well as code sharing practices, predictive modeling has  
56 become a powerful tool towards these future outlooks.

57 Among the various types of neuroimaging data, functional data may be an intuitive choice for  
58 relating brain organization to behavioral functions. In particular, task-free resting-state functional  
59 Magnetic Resonance Imaging (rs-fMRI) scans can be readily collected for large groups of subjects<sup>7</sup>,  
60 making them popular choices for neuroimaging-based predictions. In the last ten years, RSFC has  
61 been the most popular input features to brain-behavior prediction models<sup>2</sup>, in predictions of various  
62 phenotypes including fluid intelligence<sup>8-10</sup>, attention<sup>6,11,12</sup>, and working memory<sup>13-15</sup>. Brain-based  
63 psychometric prediction using other features such as task-based functional connectivity, gray  
64 matter volume, cortical thickness, and structural connectivity has also been investigated in  
65 predictions of general cognitive abilities<sup>16-18</sup>, attentional control<sup>19</sup>, and working memory<sup>20,21</sup>.  
66 However, and although this may be expected to change in the future, as far, the majority of studies  
67 forming the scientific literature have used RSFC alone or in combination with other features, for  
68 psychometric prediction<sup>2</sup>.

69 As a budding and growing field, brain-based psychometric predictions remain to be improved and  
70 validated. Many reviews have analyzed methodological options based on the current state of the  
71 field and given guidance for future studies<sup>1,2,4,5,22-24</sup>. Practical tutorials have also been published

72 for guidance on specific implementation details<sup>22,25</sup>. Nevertheless, the field also faces general and  
73 conceptual issues that are likely to limit the future usefulness of predictive modeling.

74 In this review, we discuss the current and anticipated future challenges in psychometric prediction  
75 based on neuroimaging features. For each challenge, we identify both inherent limitations in brain-  
76 based psychometric predictions which may not be readily solved based on current resources and  
77 aspects that could be addressed with potential solutions. In the following sections, we discuss the  
78 general challenges of low prediction accuracies, followed by two core issues, generalizability and  
79 interpretability. Finally, we briefly discuss the potential vulnerability of brain-based prediction  
80 models to enhancement and adversarial attacks.

## 81 Low prediction accuracies

82 Low prediction accuracies limit any potential application of the model. The general procedure of  
83 prediction model development and validation is described in Fig 1. A prediction model is assessed  
84 by applying it in a validation sample separate from the training sample, and by measuring the  
85 similarity or dissimilarity between the values predicted for the subjects in this sample and the truly  
86 observed values of the psychometric variable for these subjects (Box 1). Fig. 2 shows three  
87 examples of the most commonly used measure of model accuracy (Pearson's correlation  
88 coefficient), and the predicted-observed relationships underlying the accuracies. This measure  
89 indicates the global linear trend between predicted and observed values, but cannot identify  
90 systematic biases and size of errors. Presently, prediction accuracies of various psychometric  
91 variables have been reported from as low as 0.06 to as high as 0.908<sup>1,2</sup>. This wide range of  
92 accuracies with both low and high values close to the value bound reflects the complexity of brain-  
93 based psychometric prediction study design, as model accuracy can be affected by methodological  
94 decision and data characteristics (e.g., the amount of relevant variance in behavioral and/or brain  
95 data). While many studies that showed high prediction accuracies also appear to have used very  
96 small samples, in studies using large samples, the prediction accuracies are usually reported in the  
97 range of 0.2 to 0.4<sup>26-29</sup>, implying a generally lower accuracy when evaluating brain-based  
98 predictions in population-representative samples. A recent literature survey have evidenced a  
99 correlation of  $r=-0.265$  between the sizes of the training sample and the reported prediction  
100 accuracies, demonstrating the generality of this trend.

101 While big data and deep learning has enabled substantial successes in many fields, neither has  
102 been particularly helpful in improving the performance of brain-based prediction models. To begin  
103 with, even the easy-to-collect rs-fMRI data are considerably more difficult to collect than pictures  
104 or texts typically used in the field of computer vision and natural language processing, respectively.  
105 The lack of truly big data in cognitive neuroscience may explain why deep learning has often been  
106 reported to not outperform simpler models<sup>1,24,27</sup>. The potential of deep learning as more powerful  
107 models would thus depend on the possibility of collecting truly big neuroimaging datasets.

108 Alternatively, techniques such as few-shot learning could inspire new solutions to utilize deep  
109 learning without acquiring big data. From the data perspective, the few-shot learning strategy  
110 called data augmentation can be employed to artificially increase the sample size. Furthermore,  
111 simulated rs-fMRI and RSFC data have been used to generate additional datasets recently<sup>30-33</sup>.

112 Their applications for predictive modeling of behavior remain to be further investigated. From the  
113 parameter perspective, the meta-learning paradigm of few-shot learning can be useful by training  
114 a generalized model on a large dataset, which can be used for prediction of different targets in  
115 smaller datasets<sup>34</sup>. Nevertheless, both strategies impose some requirements and may not appear  
116 beneficial for all types of brain-based predictions. Augmented or simulated data are limited by the  
117 characteristics of the existing data used for augmentation or simulation. Accordingly, a  
118 nonrepresentative dataset (e.g., including only a certain age group or ethnicity) cannot become  
119 population representative through augmentation. As for the meta-learning strategy, its  
120 performance depends on the similarity of the prediction target in the large dataset and the  
121 prediction target in the smaller dataset<sup>34</sup>. This means that the meta-learning model would only be  
122 beneficial for smaller datasets which use the same or very similar instrument for behavioral  
123 measurement as those existing in the larger datasets in which the original model is developed.

124 It may instead be more feasible to capitalize on existing data, including neuroimaging features  
125 from multiple modalities, to boost prediction accuracies. Structural, functional, and diffusion MRI  
126 probe different neurobiological aspects, offering complementary information for psychometric  
127 prediction. In prediction studies based on functional MRI, resting-state and task functional MRI  
128 features are often combined<sup>13,35-37</sup>. However, the benefit of combining these features in terms of  
129 prediction performance has not been comprehensively investigated. Prediction studies using  
130 multimodal data have found different type of features to contribute to the prediction, including  
131 local connectome<sup>18</sup>, cortical area<sup>18</sup>, cortical thickness<sup>17</sup>, gray matter volume<sup>21</sup>, RSFC<sup>17,38-40</sup>, and  
132 task functional connectivity<sup>39,41</sup>. Some studies reported that integrating multimodal MRI data did  
133 not actually improve the prediction performance than using a single modality<sup>21,39</sup>. Furthermore,  
134 combining multimodal features inevitably increases the feature dimension and in turn the risk of  
135 overfitting, requiring feature selection or reduction techniques, such as stacking<sup>18,38,41</sup>. Generally,  
136 a systematic evaluation of multimodal psychometric prediction across multiple distinct cohorts,  
137 with an extensive set of neuroimaging features, psychometric measures, and model design, would  
138 be an important next step for validating this research direction.

139 Moreover, psychometric prediction accuracies are dependent on the target psychometric variable  
140 to predict. For behavioral traits in cognition and socioaffective domains, the definition of the  
141 abstract constructs measured by many behavioral variables and relatedly the construct validity of  
142 these variables are still debated<sup>42-44</sup>. The reliability and validity of these behavioral traits require  
143 improvement through both theoretical and experimental validations. Interestingly, many studies  
144 have reported higher prediction accuracies for cognitive measures compared to mental health  
145 traits<sup>5,34,38,45,46</sup>. It may be assumed that prediction of mental health would be particularly difficult  
146 in healthy population because the participants would show very limited variations in mental health  
147 measures. As the largest and highest neuroimaging quality datasets open to the research  
148 community include mainly healthy population, studies attempting to develop predictive models of  
149 mental health may be limited either by data availability and quality for clinical populations, or  
150 lower prediction accuracies when using easily accessible data. Relatedly, low test-retest reliability  
151 of functional MRI measures may be another source of poor prediction accuracies<sup>47,48</sup>. As the  
152 reliability of connectivity features computed may depend on data collection protocols<sup>49-51</sup>, the  
153 selection of reliable data would further restrict the available sample size.

154 One pessimistic view is that current modeling approaches may not be able to handle the  
155 heterogeneity of the population-representative samples, or that the brain-behavior relationships  
156 captured in neuroimaging datasets may simply be too weak<sup>52-55</sup>. Neuroimaging patterns may be a  
157 reduced summary of endogenous factors and the exposome that has a limited power to explain  
158 interindividual variability in behavior. Crucially, brain-based prediction models need to be  
159 justified based on the additional predictive power not already provided by non-neuroimaging  
160 features that can be easily collected by questionnaires and interviews, especially considering the  
161 high cost of MRI. From a practical standpoint, it may be useful to investigate the prediction  
162 performance of hybrid models making use of all types of data available in a realistic situation. For  
163 instance, RSFC patterns may vary with age due to developmental effects in younger population  
164 and due to aging in older population. Similarly, cognitive measures may be affected by age  
165 differently in different age subgroups. Allowing the prediction model to learn these interactions  
166 across a large age range would thus help the model to predict the target variable more accurately  
167 in general. Finally, the variability of brain-behavior association patterns across different subgroups  
168 brings forth another crucial challenge: model generalizability.

## 169 **Generalizability of prediction models**

170 The utility of a prediction model depends on its generalizability. That is, its ability to make accurate  
171 predictions on unseen data, firstly the test set data and ultimately data from the broader population.  
172 In the context of brain-based psychometric predictions, we discuss generalizability both in terms  
173 of generalizing to completely new cohorts and of generalizing to different subgroups of the  
174 population.

### 175 **Cross-cohort generalizability**

176 Cross-cohort generalizability can be defined as the prediction performance of a model in a different  
177 dataset from the training dataset (Fig. 1). Generalizable models are important both for discovering  
178 neurobiological insights general to the population and for deploying prediction models to broader  
179 settings. In most present brain-based prediction studies, the training and test sets are drawn from  
180 the same cohort under a cross-validation scheme<sup>2</sup>. While cross-validation helps to evaluate model  
181 performance without requiring additional datasets, to rigorously test the cross-cohort  
182 generalizability of a model, it is necessary to evaluate the model on completely unrelated datasets.  
183 Among studies which employed both internal and external validation, many studies found similar  
184 prediction accuracies in internal and external test sets<sup>9,13,56-59</sup>. Nevertheless, most of these studies  
185 had small external test samples ( $N < 200$ ), calling into question the representativeness of these test  
186 cohorts. In two studies with large test cohorts ( $N \sim 1000$ ), drops in prediction accuracies were  
187 observed when generalizing to new cohorts<sup>26,60</sup>. It has been suggested that reproducible brain-  
188 behavior association may only be found using samples with thousands of participants<sup>55,61</sup>.  
189 However, it has also been shown that generalizable associations and predictions can be achieved  
190 with much smaller samples in some specific cases<sup>62,63</sup>. Additionally, it should be noted that  
191 generalizability of the statistical model is not a direct indication of the generalizability of brain-  
192 behavior association derived from the model, the latter showing a low to moderate extent of  
193 generalizability across cohorts<sup>60</sup>.

194 At present, the main challenge from the perspective of cross-cohort generalizability is the lack of  
195 awareness from scientific investigators and hence the lack of assessment. The need for large  
196 external test cohorts for evaluating prediction models is often overlooked during the planning  
197 phase of a study, and later dismissed on the grounds that such large cohorts are not available for  
198 the specific psychometric measure investigated. More generally, cross-cohort generalizability of  
199 prediction models may be affected and limited by the similarity of data collection and processing  
200 protocols in the different cohorts<sup>60</sup>. The need for large datasets has led to researchers' reliance on  
201 whatever data is provided by the several publicly (or semi-publicly) shared datasets. Many studies  
202 have trained and evaluated prediction models using the Human Connectome Project Young Adult  
203 data, which were processed with a specific pipeline not always adopted or viable in other  
204 datasets<sup>64,65</sup>. Ideally, standardizing data collection protocols and processing pipelines would  
205 improve model generalizability in both research and practical situations. However, imaging  
206 conditions in samples involving children or older adults would often result in lower scan duration,  
207 making it difficult to achieve the same standards that can be set in healthy young samples<sup>66</sup>. The  
208 need for large cohorts and varied data specification may not be fully reconcilable. Partial solutions  
209 would be more robust preprocessing strategies and prediction models to harmonize data  
210 differences or to extract generalizable information despite the data differences.

### 211 **Generalizability across subgroups (within a single dataset)**

212 Typically, the test set for evaluating a prediction model is randomly selected from the cohort or  
213 taken from an external validation cohort. The composition of the test set may be completely  
214 random or stratified for balanced distributions of age, gender, and other variables of interest. While  
215 the model performances reflect the average performance in the test cohort population, they are not  
216 informative of potential prediction biases between test subjects. In both medical and non-medical  
217 applications, model bias has been reported for potential mistreatments of subgroups based on  
218 gender, ethnicity and socioeconomic status<sup>67-69</sup>. In connectivity-based prediction, ethnicity-based  
219 bias has been reported where prediction accuracies were lower in African American subjects in  
220 comparison to White American subject, even if models were trained on only African American  
221 subjects<sup>70</sup>. Moreover, models tend to predict lower cognitive scores and higher negative social  
222 behavior scores for African American subjects<sup>70</sup>, demonstrating the potential biases in applications  
223 of the prediction models. Such robust biases call for more balanced samples in scientific  
224 approaches, including not only more data collection in underrepresented population, but also the  
225 development of brain templates, atlases, and preprocessing tools based on balanced samples.

226 Common concepts used to define population subgroups like gender and ethnicity are complex  
227 notions themselves often entangled with socioeconomic factors. Relatedly, brain-based prediction  
228 models do not see the population divided into distinct gender-based or ethnic groups but have been  
229 shown to learn complex profiles relating brain measures, covariates, and psychometric variables<sup>70</sup>.  
230 It was recently demonstrated that individuals that do not follow the majority trend of brain-  
231 phenotype relationships in the training sample can cause consistent prediction failure<sup>71</sup>. For  
232 instance, if most older subjects in the training sample scored lower for a cognitive test, a few older  
233 subjects in the validation sample who scored high for the cognitive test would become outliers and  
234 lead to prediction failures. In other words, model bias may be caused by any form of stereotypical

235 brain-behavior relationships in the training sample, not specific to an ethnic or gender group. This  
236 could lead to further difficulty in collecting balanced samples since these stereotypical  
237 relationships can hardly be anticipated during data collection phase.

238 In the case where differences in brain-behavior relationships can be assumed across different  
239 subgroups in the sample, group-specific models have been used to improve prediction accuracies  
240 within certain subgroups or provide insights into the differences in brain-behavior association  
241 across subgroups<sup>17,37,72,73</sup>. Nevertheless, the validity and potential bias in subgroup definition, for  
242 instance ambiguity in ethnicity reporting, could limit the validity of any insights generated.  
243 Furthermore, brain-behavior relationships inferred from group-specific models should not be  
244 simplified in terms of causal relations with the subgroups, lest we fall into the trap of model bias  
245 and mistreatment again<sup>70</sup>. Alternatively, an ensemble learning technique called boosting may be  
246 useful for capturing different brain-behavior relationships without defining subgroups. In boosting,  
247 a sequence of models is trained where each model assigns more importance to subjects that were  
248 wrongly predicted by previous models, thereby automatically identifying the outlying subjects.

249 From a basic neuroscience perspective, the insights gained from a biased prediction model may  
250 lead to false conclusions regarding behavior and social identities, while from a practical  
251 perspective, a biased model deployed for social applications would easily lead to inequitable  
252 treatment of target populations. In order to develop a fair prediction model, both dedicated study  
253 design and model transparency are vital. This hence calls for more population-representative  
254 samples, clearly documented study and model parameters, as well as interpretable models.

## 255 Model interpretability

256 While accuracy and generalizability are requirements of any predictive model, interpretability is  
257 another crucial goal, if less easy to quantify. From a basic neuroscience perspective, prediction  
258 models need to be interpretable to contribute to our knowledge about brain-behavior relationships,  
259 while from a practical perspective, interpretability is required to evaluate the neurobiological  
260 validity of the model and, relatedly, its trustworthiness. A model with lower accuracy but higher  
261 interpretability may be preferred to a black-box model with higher accuracy, as the transparency  
262 of the former model allows assessments of the model trustworthiness. For instance, model bias  
263 against an ethnic minority could be identified earlier if the model can be interpreted easily.  
264 Nonetheless, achieving good model interpretability is not trivial and sometimes requires  
265 compromise in prediction performance<sup>22</sup>.

266 Many early studies provided an illusion of interpretability by treating regression weights from  
267 machine learning models as feature importance for neuroscientific interpretation. Later studies  
268 have demonstrated that these weights are neither stable across cross-validation folds<sup>46,74</sup>, nor  
269 conceptually valid as brain feature importance<sup>75</sup>. It may still be possible to interpret the regression  
270 weights after transforming them into corresponding forward model weights using the Haufe  
271 transform<sup>75</sup>. While stable predictive networks may be identified for cognition<sup>45</sup>, the stability in  
272 cross-validation and generalizability to new cohorts of the transformed weights were still reported  
273 to be low<sup>60,74</sup>. The reliability of transformed weights may improve with larger sample size<sup>76</sup>,  
274 making this technique potentially suitable in large cohorts. Nevertheless, when using functional

275 connectivity as features, it may be difficult to align the connectivity edges to brain mapping  
276 literature, or to summarize the feature importance values into practically useful information.  
277 Feature importance of connectivity edges may be more easily visualized and interpreted by  
278 grouping the connectivity edges in networks or finding the top connections. For instance, Fig. 3a  
279 shows groups of important connections for predicting cognition within the visual network, within  
280 the default mode network, as well as between the default mode network and other networks, while  
281 Fig. 3b shows that the most important connections for predicting fluid intelligence tend to be cross-  
282 hemispheric between medial regions or between temporal regions.

283 Many other solutions have been proposed for interpreting prediction models. Using a feature-  
284 dropping concept used in random forests<sup>77</sup>, feature importance for each feature can be quantified  
285 as the decrease in prediction performance when that feature is removed from the feature set<sup>62,78-80</sup>.  
286 This has been sometimes referred to as a ‘virtual lesion’ approach in the computational  
287 neuroimaging field. Such simple implementations may not, however, scale well to large feature  
288 sets as each feature is dealt with independently. Alternatively, using sparse regression models, only  
289 a small subset of features is selected by the regression algorithm for prediction. This leads to an  
290 inbuilt binary interpretation where only the small set of selected features is considered important.  
291 For instance, Fig. 3c shows the feature importance assignment for predicting novelty seeking by a  
292 sparse algorithm, helping the model interpretations to focus on frontal-subcortical, parietal-frontal,  
293 and within-frontal connections. This approach identifies predictive features in a data-driven  
294 manner, albeit limited to research questions where sparsity can be safely assumed. When using  
295 highly correlated features like functional connectivity, some algorithms may fail to include all  
296 important features that are correlated to each other<sup>81</sup>. Considering a large set of features without  
297 feature selection, it may still be possible to assess feature importance using Shapley Additive  
298 exPlanation (SHAP)<sup>82,83</sup>. SHAP determines each feature’s contribution similar to the ‘virtual  
299 lesion’ approach, but in all possible subsets of features, providing a distribution of feature  
300 importance for each feature. Finally, using a recently proposed region-wise framework, each brain  
301 region’s features set can be assessed instead of individual features. Concretely, a region-wise  
302 model is trained and tested to provide a model accuracy specific to the brain region<sup>60</sup>.  
303 Interpretations based on region-wise models are easy to illustrate (Fig. 3d) and to some extent align  
304 with the brain mapping literature. Nevertheless, the distributed aspect of brain organization is not  
305 modeled by the region-wise models, limiting strong interpretations to mostly region-specific  
306 properties.

307 Ultimately, useful model interpretations are reliant on the prediction accuracy and generalizability  
308 of the model. With very low accuracies, the interpretations generated from the models may be  
309 arbitrary at best, while with low generalizability the interpretations may be valid only for the  
310 training sample. The challenge hence lies in designing models where interpretability can be  
311 achieved with minimal or no compromise in accuracy. Potential directions may include more  
312 powerful generative models, more informative priors, and interpretable deep neural networks.  
313 Generative models and deep neural network models may be combined into deep generative models  
314 to bring forth the benefits of both interpretability and accuracy, with specific approaches including  
315 variational autoencoders<sup>84</sup>, generative adversarial networks<sup>85</sup>, and autoregressive models<sup>86,87</sup>. With  
316 traditional machine learning models, feature importance based on existing models can help to



317 reduce feature dimensionality in new models in new cohorts, which offers new interpretations to  
318 validate against the existing model's interpretation. In this way, a positive reinforcement loop may  
319 exist between boosting prediction accuracy and interpretability, reducing the need to sacrifice one  
320 for the other.

## 321 Enhancement and adversarial attacks

322 Enhancement and adversarial attacks can threaten the trustworthiness of neuroimaging-based  
323 predictive models. Enhancement attacks are those where purposeful data alterations can lead to  
324 falsely enhanced model performance, while adversarial attacks are those where specifically  
325 designed noise are added to the data to cause a model to fail<sup>88</sup>. An artificially enhanced model may  
326 be the result of scientific malpractice or fraud which, if not discovered, could lead to large amount  
327 of time and resources wasted in the wrong research direction. Successful adversarial attacks on  
328 deployed models mean that prediction outcomes would become unreliable. In biomedical  
329 development for example, the effect of a treatment or of a drug could be faked or exaggerated with  
330 data manipulation of the machine learning model to mislead financial investors. . Similar to the  
331 issue of generalizability, the main challenge of these attacks in the field of neuroimaging-based  
332 psychometric prediction is the lack of awareness. As practical applications of neuroimaging-based  
333 prediction models are still far-fetched at present, there is a lack of motivation for researchers to  
334 anticipate models that are robust to these attacks. Furthermore, replication studies that might detect  
335 enhancement attacks are still rather lacking in the field. While there is no evidence of existing  
336 enhancement or adversarial attacks in the field and no practical solution proposed against them  
337 currently, these are crucial issues to address in the perspective of the deployment of brain-based  
338 prediction models for applications in the society.

339 Simple data enhancement can be done by biased subject selection. Subjects may be retroactively  
340 selected based on their individual prediction outcome, or only chosen if they follow certain brain-  
341 phenotype stereotypes. Such manipulations can be detected if data characteristics and exclusion  
342 criteria are reported faithfully, especially when outliers are excluded based on a threshold. A more  
343 advanced approach involves adding patterns correlated to the behavioral variable of interest to the  
344 imaging features, boosting the prediction accuracies to almost perfect accuracies without causing  
345 the features to become significantly different from the original features<sup>88</sup>. Furthermore, it is  
346 possible to design data enhancements to cause machine learning models to learn brain-behavior  
347 relationships not existing in the original data. This means that enhancement attacks may also be  
348 detrimental from a basic neuroscience perspective as conclusions drawn would not be valid. This  
349 type of attack may be detected when a replication study fails to generalize the model to new cohorts  
350 but can only really be confirmed if the raw data and data processing code can be openly examined.

351 The effects of adversarial attacks in machine learning models for clinical applications have been  
352 investigated<sup>89,90</sup>. For brain-based prediction models in healthy population, it has also been shown  
353 that very minor data manipulations can cause the classification accuracy to drop to 0%<sup>88</sup>. To design  
354 this type of attack, the model parameters must be known, hence bringing forth additional  
355 challenges in achieving both open science and practical utility. Data validation to identify  
356 manipulated data, if possible, may become paramount in the future of adversarial attacks.

357 Potentially, machine learning models employed in practical applications can make use of online  
358 learning where a trained model continues to receive new batches of data for additional training,  
359 while only sharing the model at baseline for scientific purposes.

360 In face of potential enhancement and adversarial attacks, model and study reproducibility enabled  
361 by open science is necessary to detect and address these data manipulations. With transparent study  
362 design and provenance tracking, the field can benefit from multiple aspects including easier  
363 replication, enhancement attack monitoring, comparison across studies, and results pooling<sup>22,88</sup>.

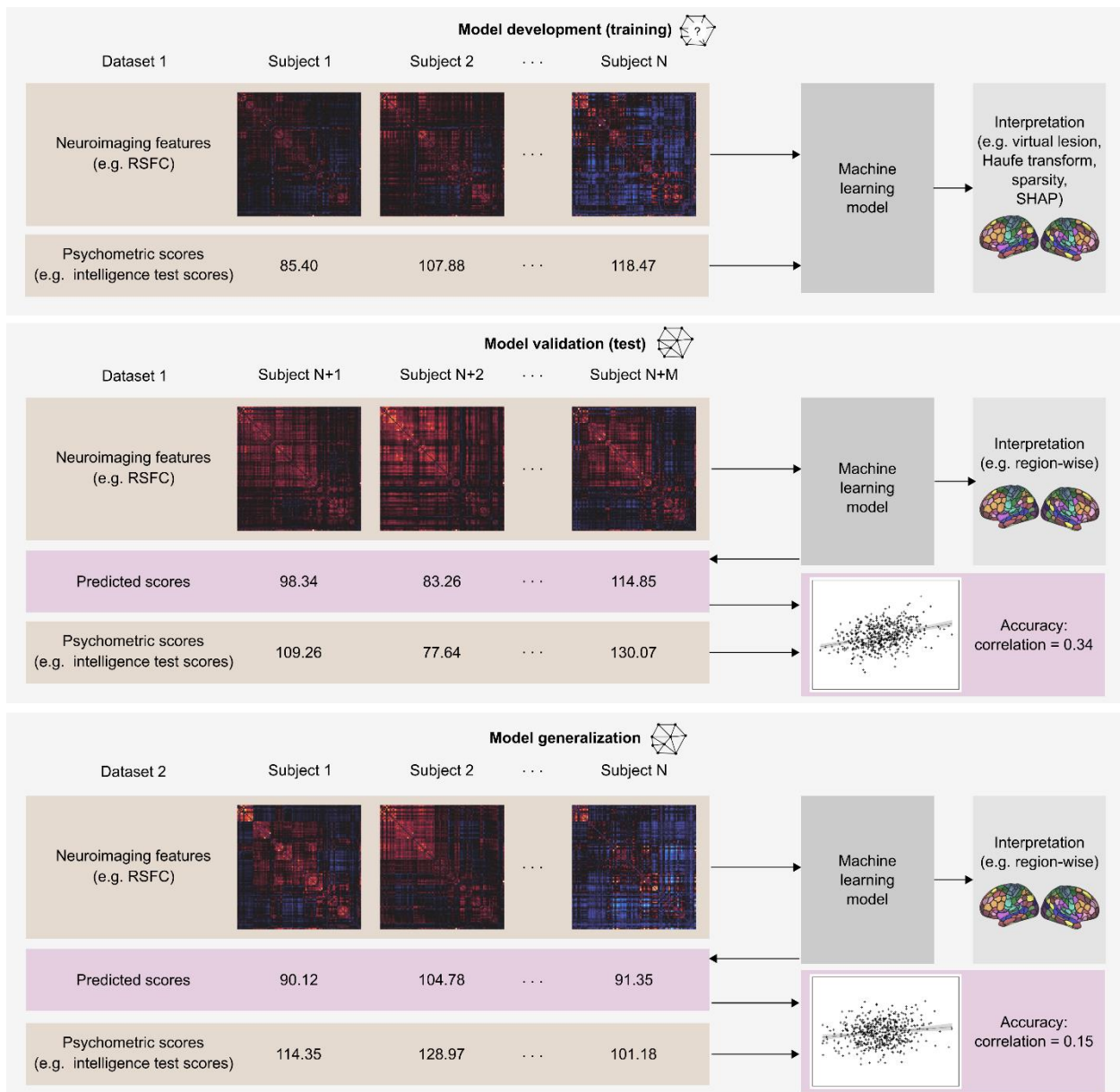
## 364 Conclusions

365 Many challenges lie in the way of brain-based predictive modeling of behavior before it can be  
366 substantially useful for understanding complex brain-behavior relationships or for practical  
367 applications. While some limitations are inherent, such as smaller sample sizes in studies interested  
368 in phenotypic measure uncommon in large open datasets, others are solvable, such as assessment  
369 and improvement of generalizability. By acknowledging this and addressing the solvable issues,  
370 brain-based psychometric predictions can steadily progress towards scientific and practical utility.  
371 We encourage more comprehensive study design, comprising multiple cohorts to cover more  
372 population-representative samples, and ensuring model validity with careful confound handling.  
373 Furthermore, we advocate for model evaluation based on both accuracies and generalizability.  
374 Predictive modeling in neuroscience is a necessarily interdisciplinary field, which requires  
375 combinations of neuroscientific knowledge, statistical concepts, and machine learning techniques  
376 to achieve its potential. Beyond this interdisciplinarity, transparent models, diverse data, and  
377 rigorous study designs are the keys to move forward.

## 378 Acknowledgements

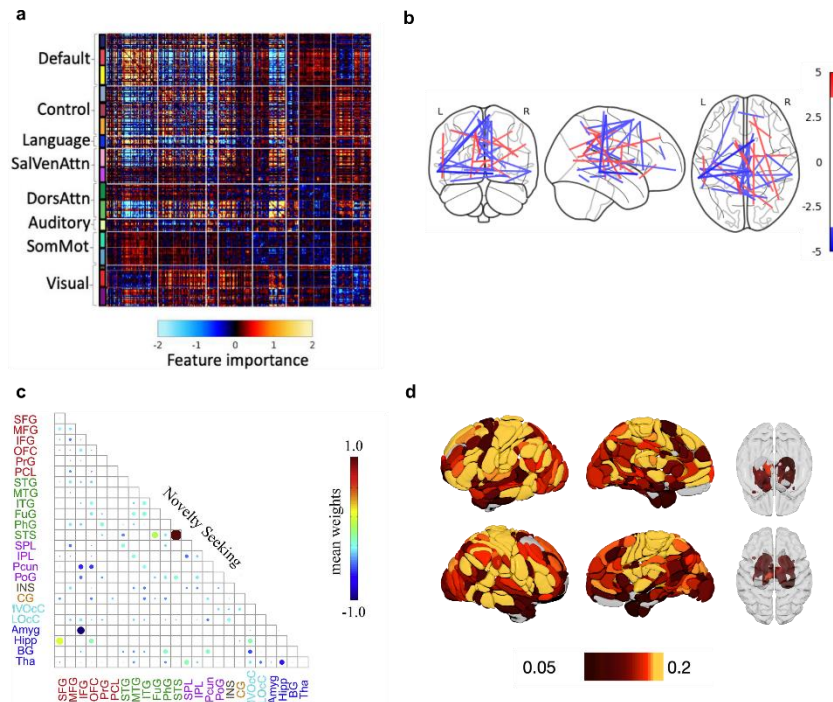
379 This work was supported by the Deutsche Forschungsgemeinschaft (DFG, GE 2835/2–1, EI 816/  
380 4–1), the Helmholtz Portfolio Theme ‘Supercomputing and Modelling for the Human Brain’ and  
381 the European Union’s Horizon 2020 Research and Innovation Programme under Grant Agreement  
382 No. 720270 (HBP SGA1) and Grant Agreement No. 785907 (HBP SGA2).

383



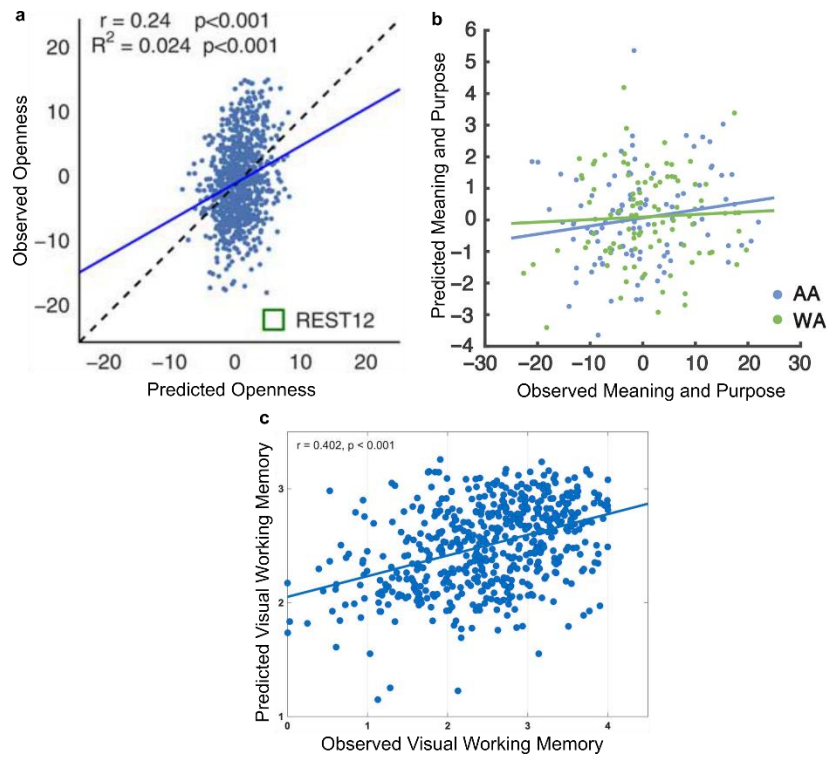
384

385 **Fig. 1 | Model development and validation for neuroimaging-based psychometric predictions.** A machine learning  
 386 model is first trained using neuroimaging features and psychometric scores from subjects 1 to N (from the training  
 387 set). The model learns a relationship between the neuroimaging features and the psychometric scores. For validation,  
 388 the model takes in neuroimaging features from subjects N+1 to N+M (from the test set), and outputs predicted values  
 389 for the psychometric scores. The predicted scores can then be compared to the actual scores using various accuracy  
 390 measures (see Box 1) to evaluate the performance of the model. To assess the generalizability of the model,  
 391 the model needs to be applied to a new dataset in a similar way to its application in the test set.



392

393 **Fig. 2 | Prediction accuracies measured by Pearson's correlation.** **a**, Scatter plot of observed and predicted openness  
 394 trait, with a Pearson's correlation accuracy of 0.24<sup>91</sup>. Blue line shows the fitted line between observed and predicted  
 395 values, while black dashed line marks the line with unit slope and zero intercept. It can be noted that, while  
 396 (standardized) observed values have a wide range of variations (roughly between -20 and 15, predicted values remain  
 397 tightly scattered around zero. **b**, Scatter plot of observed and predicted scores of meaning and purpose, with a  
 398 Pearson's correlation accuracy of 0.17 in African American subjects and 0.049 in White American subjects<sup>70</sup>. Blue and  
 399 green lines show the fitted line between observed and predicted values in African American and White American  
 400 subjects respectively. The correlation appears slightly higher in African American than White American, while the  
 401 prediction errors may actually be greater in the former group. **c**, Scatter plot of observed and predicted visual working  
 402 memory performance, with a Pearson's correlation accuracy of 0.402<sup>21</sup>. Blue line shows the fitted line between  
 403 observed and predicted values. Overall, from all three plots, it can be observed that the Pearson's correlation  
 404 coefficient is higher when the fitted line has a slope closer to one. It is also noteworthy that predicted values in all  
 405 cases tend to have smaller variances compared to the observed values. This reflects the tendency of machine learning  
 406 algorithms to generate predictions closer to the sample mean. Finally, outliers or prediction failures can be observed  
 407 in all plots even when correlation accuracies are moderate. As the correlation accuracies measure the relative  
 408 goodness-of-fit, they are less affected by (or reflective of) outliers compared to accuracy measures based on absolute  
 409 errors.



410

411 **Fig. 3 | Visualizations of model interpretations.** **a**, Feature importance of all RSFC edges for predicting general  
 412 cognition in a young adult cohort with parcels grouped under networks<sup>38</sup>. Colors correspond to the Haufe transformed  
 413 weight values. Important connections can be found within the visual network, within the default mode network,  
 414 between the default mode network and the control network, as well as between the default mode network and the  
 415 attention networks. **b**, Feature importance of top RSFC edges for predicting fluid intelligence in a young adult cohort  
 416 shown in their corresponding positions in the brain<sup>60</sup>. Colors correspond to the Haufe transformed weight values.  
 417 Most top connections can be found between medial regions or temporal regions across the hemispheres. **c**, Feature  
 418 importance of all RSFC edges for predicting novelty seeking in a young adult cohort when a sparse algorithm was used.  
 419 Colors correspond to the mean weight values across cross-validation splits<sup>92</sup>. The sparse set of selected features  
 420 mostly include frontal-subcortical, parietal-frontal, and within-frontal connections. **d**, Brain region importance for  
 421 predicting fluid cognition in an aging cohort based on the RSFC features using the region-wise approach<sup>60</sup>. Colors  
 422 correspond to the prediction accuracies achieved using brain regional connectivity profiles. The relatively more  
 423 predictive regions can be identified in the cingulate cortex, the peripheral visual area, the right supramarginal gyrus,  
 424 the right anterior insula, the central sulcus, and the right lateral frontal cortex.

425

## Box 1 | Measures of model accuracy

In order to evaluate a model in a validation sample, its predictions need to be compared against the actual values of the psychometric variable. The closer the predicted values are to the actual values, the more accurate the model is. This degree of closeness can be represented either by correlation metrics examining the linear trend between all predicted and observed values, or by error metrics examining the absolute differences between each pair of predicted and observed values.

The most common metric in the literature is the Pearson's correlation coefficient ( $r$ ) between predicted and observed values<sup>Error! Bookmark not defined.</sup>, measuring the normalized covariance between the two variables ( $r = \frac{cov(pred,obs)}{\sigma_{pred}\sigma_{obs}}$ ). This correlation coefficient is an indication of the extent to which a given increase or decrease in one variable is associated with a similar increase or decrease in the other variable. Similarly, Spearman's correlation can be used to measure the ranked correlation between predicted and observed values, providing an indication of how well the two groups of values are monotonically related.

Common error metrics include mean absolute error, mean squared error (MSE), and root mean squared error, measuring the average difference between predicted and observed values in the validation sample in slightly different manners. In general, the error values should be normalized by the standard deviation (or the range of predicted values for absolute errors) of the validation sample, so that they are comparable to standardized measures from other samples<sup>Error! Bookmark not defined.</sup>.

While a high correlation suggests that predicted values are generally higher when observed values are higher, it does not mean that predicted values are numerically close to the observed values. As a result, the correlation metrics cannot detect systematic biases where the predicted values are consistently higher (or lower) than the observed value. It may be recommended that high correlation accuracies should be validated with error-based accuracies to check for systematic bias. On the other hand, correlation metrics might be more useful when generalizing a model to new data where the psychometric variables are similar to but not the same as those in the training sample and numerical closeness between predicted and observed values may not be required.

Finally, a useful metric for model evaluation is the coefficient of determination (or  $R^2$ ), providing a measure of goodness-of-fit of the model. A simple form of  $R^2$  is  $r^2$  may also be computed as the square of the correlation coefficient from the correlation metric. It should be noted that this  $r^2$  measures the goodness-of-fit between the predicted-observed relation and its fitted line, and hence is not a direct measure of model fit itself. Using error metrics such as MSE, the more general  $R^2$  can be computed as  $R^2 = 1 - \frac{MSE}{\sigma^2}$ , measuring the goodness-of-fit of the regression equation estimated by the prediction model to the validation data. The  $R^2$  values can also be interpreted as the ratio of explained variance by the model to the total variance in the sample, offering an intuitive way to explain the accuracies measured.

## 427 References

- 428 1. Sui, J., Jiang, R., Bustillo, J. & Calhoun, V. Neuroimaging-based Individualized Prediction of  
429 Cognition and Behavior for Mental Disorders and Health: Methods and Promises. *Biol. Psychiatry* **88**,  
430 818-828 (2020).
- 431 2. Yeung, A.W.K., More, S., Wu, J. & Eickhoff, S.B. Reporting details of neuroimaging studies on  
432 individual traits prediction: a literature survey. *NeuroImage* **256**, 119275 (2022).
- 433 3. Cirillo, D. & Valencia, A. Big data analytics for personalized medicine. *Curr. Opin. Biotechnol.* **58**,  
434 161-167 (2019).
- 435 4. Dadi, K. et al. Benchmarking functional connectome-based predictive models for resting-state fMRI.  
436 *NeuroImage* **192**, 115-134 (2019).
- 437 5. Dhamala, E., Yeo, B.T.T. & Holmes, A.J. One Size Does Not Fit All: Methodological Considerations  
438 for Brain-Based Predictive Modelling in Psychiatry. *Biol. Psychiatry* **93**, 717-728 (2023).
- 439 6. Rosenberg, M.D. et al. A neuromarker of sustained attention from whole-brain functional connectivity.  
440 *Nat. Neurosci.* **19**, 165-171 (2016).
- 441 7. Lee, M.H., Smyser, C.D. & Shimony, J.S. Resting-state fMRI: a review of methods and clinical  
442 applications. *Am. J. Neuroradiol.* **34**, 1866-1872 (2013).
- 443 8. Finn, E.S. et al. Functional connectome fingerprinting: identifying individuals using patterns of brain  
444 connectivity. *Nat. Neurosci.* **18**, 1664-1671 (2015).
- 445 9. Ferguson, M.A., Anderson, J.S., Spreng, R.N. Fluid and flexible minds: Intelligence reflects synchrony  
446 in the brain's intrinsic network architecture. *Netw. Neurosci.* **1**, 192-207 (2017).
- 447 10. Li, J. et al. A neuromarker of individual general fluid intelligence from the white-matter functional  
448 connectome. *Transl. Psychiatry* **10**, 147 (2020).
- 449 11. Kumar, S. et al. An information network flow approach for measuring functional connectivity and  
450 predicting behavior. *Brain Behav.* **9**, e01346 (2019).
- 451 12. Rosenberg, M.D. et al. Functional connectivity predicts changes in attention observed across minutes,  
452 days, and months. *Proc. Natl. Acad. Sci. USA* **117**, 3797-3807 (2020).
- 453 13. Avery, E.W. et al. Distributed Patterns of Functional Connectivity Predict Working Memory  
454 Performance in Novel Healthy and Memory-impaired Individuals. *J. Cogn. Neurosci.* **32**, 241-255  
455 (2020).
- 456 14. Pläschke, R.N. et al. Age differences in predicting working memory performance from network-based  
457 functional connectivity. *Cortex* **132**, 441-459 (2020).
- 458 15. Zhang, H. et al. Do intrinsic brain functional networks predict working memory from childhood to  
459 adulthood? *Hum. Brain Mapp.* **41**, 4574-4586 (2020).
- 460 16. Girault, J.B. et al. White matter connectomes at birth accurately predict cognitive abilities at age 2.  
461 *NeuroImage* **192**, 145-155 (2019).
- 462 17. Jiang, R. et al. Multimodal data revealed different neurobiological correlates of intelligence between  
463 males and females. *Brain Imaging Behav.* **14**, 1979-1993 (2020).
- 464 18. Rasero, J., Sentis, A.I., Yeh, F.C. & Verstynen, T. Integrating across neuroimaging modalities boosts  
465 prediction accuracy of cognitive ability. *PLoS Comput. Biol.* **17**, e1008347 (2021).
- 466 19. Wei, L., et al. Grey matter volume in the executive attention system predict individual differences in  
467 effortful control in young adults. *Brain Topogr.* **32**, 111-117 (2019).
- 468 20. Kaufmann, T. et al. Task modulations and clinical manifestations in the brain functional connectome  
469 in 1615 fMRI datasets. *NeuroImage* **147**, 243-252 (2017).
- 470 21. Xiao, Y. et al. Predicting visual working memory with multimodal magnetic resonance imaging. *Hum.*  
471 *Brain Mapp.* **42**, 1446-1462 (2021).

- 472 22. Scheinost, D. et al. Ten simple rules for predictive modeling of individual differences in neuroimaging.  
473 *NeuroImage* **193**, 35-45 (2019).
- 474 23. Gabrieli, J.D.E., Ghosh, S.S., Whitfield-Gabrieli, S. Prediction as a humanitarian and pragmatic  
475 contribution from human cognitive neuroscience. *Neuron* **85**, 11-26 (2015).
- 476 24. Pervaiz, U., Vidaurre, D., Woolrich, M.W., Smith, S.M. Optimising network modelling methods for  
477 fMRI. *NeuroImage* **221**, 116604 (2020).
- 478 25. Poldrak, R.A., Huckins, G. & Varoquax, G. Establishment of Best Practices for Evidence for Prediction:  
479 A Review. *J. Am. Med. Assoc. Psychiatry* **77**, 534-540 (2020).
- 480 26. Sripada, C. et al. Prediction of neurocognition in youth from resting state fMRI. *Mol. Psychiatry* **25**,  
481 3413-3421 (2019).
- 482 27. He, T. et al. Deep neural networks and kernel regression achieve comparable accuracies for functional  
483 connectivity prediction of behavior and demographics. *NeuroImage* **206**, 116276 (2020).
- 484 28. He, L. et al. Functional Connectome Prediction of Anxiety Related to the COVID-19 Pandemic. *Am. J.*  
485 *Psychiatry* **178**, 530-540 (2021).
- 486 29. Gao, S., Greene, A.S., Constable, R.T. & Scheinost, D. Combining multiple connectomes improves  
487 predictive modeling of phenotypic measures. *NeuroImage* **201**, 116038 (2019).
- 488 30. Zalesky, A., Fornito, A., Cocchi, L., Gollo, L.L. & Breakspear, M. Time-resolved resting-state brain  
489 networks. *Proc. Natl. Acad. Sci. USA* **111**, 10341-10346 (2014).
- 490 31. Bahg, G., Evans, D.G., Galdo, M. & Turner, B.M. Gaussian process linking functions for mind, brain,  
491 and behavior. *Proc. Natl. Acad. Sci. USA* **117**, 29398-29406 (2020).
- 492 32. Mihalik, A. et al. Canonical Correlation Analysis and Partial Least Squares for identifying brain-  
493 behaviour associations: a tutorial and a comparative study. *Biol. Psychiatry* **7**, 1055-1067 (2022).
- 494 33. Gal, S., Tik, N., Bernstein-Eliav, M. & Tavor, I. Predicting individual traits from unperformed tasks.  
495 *NeuroImage* **249**, 118920 (2022).
- 496 34. He, T. et al. Meta-matching as a simple framework to translate phenotypic predictive models  
497 from big to small data. *Nat. Neurosci* **25**, 795-804 (2022).
- 498 35. Takagi, Y., Hirayama, J.I., Tanaka, S.C. State-unspecific patterns of whole-brain functional  
499 connectivity from resting and multiple task states predict stable individual traits. *NeuroImage*  
500 **201**, 116036 (2019).
- 501 36. Burr, D.A. et al. Functional connectivity predicts the dispositional use of expressive  
502 suppression but not cognitive reappraisal. *Brain Behav.* **10**, e01493 (2020).
- 503 37. Jiang, R. et al. Task-induced brain connectivity promotes the detection of individual  
504 differences in brain-behavior relationships. *NeuroImage* **207**, 116370 (2020).
- 505 38. Ooi, L.Q.R. et al. Comparison of individualized behavioral predictions across anatomical,  
506 diffusion and functional connectivity MRI. *NeuroImage* **263**, 119636 (2022).
- 507 39. Dhamala, E., Jamison, K.W., Jaywant, A., Dennis, S. & Kuceyeski, A. Distinct functional and  
508 structural connections predict crystallised and fluid cognition in healthy adults. *Hum. Brain*  
509 *Mapp.* **42**, 3102-3118 (2021).
- 510 40. Mansour, L.S, Tian, Y., Yeo, B.T.T., Croyley, V. & Zalesky, A. High-resolution connectomic  
511 fingerprints: Mapping neural identity and behavior. *NeuroImage* **229**, 117695 (2021).
- 512 41. Pat, N. et al. Longitudinally stable, brain-based predictive models mediate the relationships  
513 between childhood cognition and socio-demographic, psychological and genetic factors. *Hum.*  
514 *Brain Mapp.* **43**, 5520-5542 (2022).
- 515 42. Hurtz, G.M. & Donovan, J.J. Personality and job performance: The Big Five revisited. *J. Appl.*  
516 *Psychol.* **85**, 869-879 (2000).



- 517 43. Kane, M.J., Conway, A.R.A., Miura, T.K. & Colflesh, G.J.H. Working memory, attention  
518 control, and the n-back task: A question of construct validity. *J. Exp. Psychol.* **33**, 615-622  
519 (2007).
- 520 44. Sanchez-Cubillo, I. et al. Construct validity of the Trail Making Test: Role of task-switching,  
521 working memory, inhibition/interference control, and visuospatial abilities. *J. Int.*  
522 *Neuropsychol. Soc.* **15**, 438-450 (2009).
- 523 45. Chen, J. et al. Shared and unique brain network features predict cognitive, personality, and  
524 mental health scores in the ABCD study. *Nat. Commun.* **13**, 2217 (2022).
- 525 46. Wu, J. et al. A Connectivity-Based Psychometric Prediction Framework for Brain-Behavior  
526 Relationship Studies. *Cereb. Cortex* **31**, 3732-3751 (2021).
- 527 47. Noble, S., Scheinost, D. & Constable, R.T. A decade of test-retest reliability of functional  
528 connectivity: A systematic review and meta-analysis. *NeuroImage* **203**, 116157 (2019).
- 529 48. Elliott, M.L. et al. What Is the Test-Retest Reliability of Common Task-Functional MRI  
530 Measures? New Empirical Evidence and a Meta-Analysis. *Psychol. Sci.* **31**, 792-806 (2020).
- 531 49. Patriat, R. et al. The effect of resting condition on resting-state fMRI reliability and consistency:  
532 A comparison between resting with eyes open, closed, and fixated. *NeuroImage* **78**, 463-473  
533 (2013).
- 534 50. Birn, R.M. et al. The effect of scan length on the reliability of resting-state fMRI connectivity  
535 estimates. *NeuroImage* **83**, 550-558 (2013).
- 536 51. Bennett, C.M. & Miller, M.B. fMRI reliability: Influences of task and experimental design.  
537 *Cogn. Affect. Behav. Neurosci.* **13**, 690-702 (2013).
- 538 52. Cremers, H.R., Wager, T.D., Yarkoni, T. The relation between statistical power and inference  
539 in fMRI. *PLoS One* **12**, e0184923 (2017).
- 540 53. Kharabian Masouleh, S., Eickhoff, S.B., Hoffstaedter, F., Genon, S. & Alzheimer's Disease  
541 Neuroimaging Initiative. Empirical examination of the replicability of associations between  
542 brain structure and psychological variables. *eLife* **8**, e43464 (2019).
- 543 54. Genon, S., Eickhoff, S.B., Kharabian, S. Linking interindividual variability in brain structure  
544 to behaviour. *Nat. Rev. Neurosci.* **23**, 307-318 (2022).
- 545 55. Marek, S. et al. Reproducible brain-wide association studies require thousands of individuals.  
546 *Nature* **603**, 654-660 (2022).
- 547 56. Beaty, R.E. et al. Robust prediction of individual creative ability from brain functional  
548 connectivity. *Proc. Natl. Acad. Sci. USA* **115**, 1087-1092 (2018).
- 549 57. Liu, P. et al. The functional connectome predicts feeling of stress on regular days and during  
550 the COVID-19 pandemic. *Neurobiol. Stress* **14**, 100285 (2021).
- 551 58. Ren, Z. et al. Connectome-based Predictive Modeling of Creativity Anxiety. *NeuroImage* **225**,  
552 117469 (2021).
- 553 59. Fong, A.H.C. et al. Dynamic functional connectivity during task performance and rest predicts  
554 individual differences in attention across studies. *NeuroImage* **188**, 14-25 (2019).
- 555 60. Wu, J. et al. Cross-cohort replicability and generalizability of connectivity-based psychometric  
556 prediction patterns. *NeuroImage* **262**, 119569 (2022).
- 557 61. Tervo-Clemmens, B. et al. Reply to: Multivariate BWAS can be replicable with moderate  
558 sample sizes. *Nature* **615**, E8-E12 (2023).

- 559 62. Rosenberg, M.D. & Finn, E.S. How to establish robust brain-behavior relationships without  
560 thousands of individuals. *Nat. Neurosci.* **25**, 835-837 (2022).
- 561 63. Spisak, T., Bingel, U. & Wager, T. Replicable multivariate BWAAS with moderate sample sizes.  
562 Preprint at <https://www.biorxiv.org/content/10.1101/2022.06.22.497072v1> (2022).
- 563 64. Van Essen, D.C. et al. The WU-Minh Human Connectome Project: an overview. *NeuroImage*  
564 **80**, 62-79 (2013).
- 565 65. Glasser, M.F. et al. The Minimal Preprocessing Pipelines for the Human Connectome Project.  
566 *NeuroImage* **80**, 105-124 (2013).
- 567 66. Harms, M.P. et al. Extending the Human Connectome Project across ages: imaging protocols  
568 for the lifespan development and aging projects. *NeuroImage* **183**, 972-984 (2018).
- 569 67. Chouldechova, A., Benavides-Prado, D., Fialko, O. & Vaithianathan, R. A case study of  
570 algorithm-assisted decision making in child maltreatment hotline screening decisions. *Proc.*  
571 *Mach. Learn. Res.* **81**, 134-148 (2018).
- 572 68. Martin, A.R. et al. Clinical use of current polygenic risk scores may exacerbate healthy  
573 disparities. *Nat. Genet.* **51**, 584-591 (2019).
- 574 69. Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. Dissecting racial bias in an  
575 algorithm used to manage the health of populations. *Science* **366**, 447-453 (2019).
- 576 70. Li, J. et al. Cross-ethnicity/race generalization failure of behavioral prediction from resting-  
577 state functional connectivity. *Sci. Adv.* **8**, eabj1812 (2022).
- 578 71. Greene, A.S. et al. Brain-phenotype models fail for individuals who defy sample stereotypes.  
579 *Nature* **609**, 109-118 (2022).
- 580 72. Greene, A.S., Gao, S., Scheinost, D. & Constable, R.T. Task-induced brain state manipulation  
581 improves prediction of individual traits. *Nat. Commun.* **9**, 2807 (2018).
- 582 73. Nostro, A.D. et al. Predicting personality from network-based resting-state functional  
583 connectivity. *Brain Struct. Funct.* **223**, 2699-2719 (2018).
- 584 74. Tian, Y. & Zalesky, A. Machine learning prediction of cognition from functional connectivity:  
585 Are feature weights reliable? *NeuroImage* **245**, 118648 (2021).
- 586 75. Haufe, S. et al. On the interpretation of weight vectors of linear models in multivariate  
587 neuroimaging. *NeuroImage* **87**, 96-110 (2014).
- 588 76. Chen J. et al. Relationship between prediction accuracy and feature importance reliability: An  
589 empirical and theoretical study. *NeuroImage* **274**, 120115 (2023).
- 590 77. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5-32 (2001).
- 591 78. Yip, S.W., Kiluk, B., & Scheinost, D. Towards Addiction Prediction: An Overview of Cross-  
592 Validated Predictive Modeling Findings and Considerations for Future Neuroimaging  
593 Research. *Biol. Psychiatry Cogn Neuroscie. Neuroimaging* **5**, 748-758 (2020).
- 594 79. Jiang, R., Woo, C.W., Qi, S., Wu, J. & Sui, J. Interpreting Brain Biomarkers: Challenges and  
595 solutions in interpreting machine learning-based predictive neuroimaging. *IEEE Signal*  
596 *Process. Mag.* **39**, 107-118 (2022).
- 597 80. Chormai, P., Pu, Y., Hu, H., Fisher, S.E., Francks, C. & Kong, X.Z. Machine learning of large-  
598 scale multimodal brain imaging data reveals neural correlates of hand preference. *NeuroImage*,  
599 **262**, 119534 (2022).
- 600 81. Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *J. Royal*  
601 *Statistical Soc. B* **67**, 301-320 (2005).

- 602 82. Lundberg, S.M. & Lee, S.-I. A unified approach to interpreting model prediction. *Adv. Neural*  
603 *Inf. Process. Syst.* **30** (2017).
- 604 83. Pat, N., Wang, Y., Bartonicek, A., Candia, J. & Stringaris, A. Explainable machine learning  
605 approach to predict and explain the relationship between task-based fMRI and individual  
606 differences in cognition. *Cereb. Cortex* **33**, 2682-2703 (2023).
- 607 84. Kingma, D.P. & Welling, M. Auto-Encoding Variational Bayes. Preprint at  
608 <https://arxiv.org/abs/1312.6114> (2013).
- 609 85. Goodfellow, I.J. et al. Generative Adversarial Networks. Preprint at  
610 <https://arxiv.org/abs/1406.2661> (2014).
- 611 86. van den Oord, A., Kalchbrenner, N. & Kavukcuoglu, K. Pixel Recurrent Neural Networks.  
612 Preprint at <https://arxiv.org/abs/1601.06759> (2016).
- 613 87. Fried, D. et al. Speaker-Follower Models for Vision-and-Language Navigation. In *NeuroIPS*  
614 (2018).
- 615 88. Rosenblatt, M. et al. Connectome-based machine learning models are vulnerable to subtle data  
616 manipulations. *Patterns* (In press).
- 617 89. Finlayson, S.G. et al. Adversarial attacks on medical machine learning. *Science* **363**, 1287-  
618 1289 (2019).
- 619 90. Finlayson, S.G., Chung, H.W., Kohane, I.S. & Beam, A.L. Adversarial Attacks Against  
620 Medical Deep Learning Systems. Preprint at <https://doi.org/10.48550/arXiv.1804.05296>  
621 (2019).
- 622 91. Dubois, J., Galdi, P., Han, Y., Paul, L.K. & Adolphs, R. Resting-state functional brain  
623 connectivity best predicts the personality dimension of openness to experience. *Personal.*  
624 *Neurosci.* **1**, E6 (2018).
- 625 92. Jiang, R. et al. Connectome-based individualized prediction of temperament trait scores.  
626 *NeuroImage* **183**, 366-374 (2018).

## 627 Competing interests

628 The authors declare no competing interests.