



# Relationship between prediction accuracy and feature importance reliability: An empirical and theoretical study

Jianzhong Chen<sup>a,b,c,d,1</sup>, Leon Qi Rong Ooi<sup>a,b,c,d,e,1</sup>, Trevor Wei Kiat Tan<sup>a,b,c,d,e</sup>,  
Shaoshi Zhang<sup>a,b,c,d,e</sup>, Jingwei Li<sup>f,g</sup>, Christopher L. Asplund<sup>a,b,d,h,i,j</sup>, Simon B Eickhoff<sup>f,g</sup>,  
Danilo Bzdok<sup>k,l</sup>, Avram J Holmes<sup>m</sup>, B.T. Thomas Yeo<sup>a,b,c,d,e,n,\*</sup>

<sup>a</sup> Centre for Sleep and Cognition, Yong Loo Lin School of Medicine, National University of Singapore, Singapore

<sup>b</sup> Centre for Translational MR Research, Yong Loo Lin School of Medicine, National University of Singapore, Singapore

<sup>c</sup> Department of Electrical and Computer Engineering, National University of Singapore, Singapore

<sup>d</sup> N.1 Institute for Health & Institute for Digital Medicine (WisDM), National University of Singapore, Singapore

<sup>e</sup> Integrative Sciences and Engineering Programme (ISEP), National University of Singapore, Singapore

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

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

<sup>h</sup> Division of Social Sciences, Yale-NUS College, Singapore

<sup>i</sup> Department of Psychology, National University of Singapore, Singapore

<sup>j</sup> Duke-NUS Medical School, Singapore

<sup>k</sup> Department of Biomedical Engineering, Montreal Neurological Institute, McGill University, Montreal, Quebec, Canada

<sup>l</sup> Mila - Quebec AI Institute, Montreal, Canada

<sup>m</sup> Departments of Psychology and Psychiatry, Yale University, New Haven, CT, USA

<sup>n</sup> Martinos Center for Biomedical Imaging, Massachusetts General Hospital, Charlestown, MA, USA

## A B S T R A C T

There is significant interest in using neuroimaging data to predict behavior. The predictive models are often interpreted by the computation of feature importance, which quantifies the predictive relevance of an imaging feature. Tian and Zalesky (2021) suggest that feature importance estimates exhibit low split-half reliability, as well as a trade-off between prediction accuracy and feature importance reliability across parcellation resolutions. However, it is unclear whether the trade-off between prediction accuracy and feature importance reliability is universal. Here, we demonstrate that, with a sufficient sample size, feature importance (operationalized as Haufe-transformed weights) can achieve fair to excellent split-half reliability. With a sample size of 2600 participants, Haufe-transformed weights achieve average intra-class correlation coefficients of 0.75, 0.57 and 0.53 for cognitive, personality and mental health measures respectively. Haufe-transformed weights are much more reliable than original regression weights and univariate FC-behavior correlations. Original regression weights are not reliable even with 2600 participants. Intriguingly, feature importance reliability is strongly positively correlated with prediction accuracy across phenotypes. Within a particular behavioral domain, there is no clear relationship between prediction performance and feature importance reliability across regression models. Furthermore, we show mathematically that feature importance reliability is necessary, but not sufficient, for low feature importance error. In the case of linear models, lower feature importance error is mathematically related to lower prediction error. Therefore, higher feature importance reliability might yield lower feature importance error and higher prediction accuracy. Finally, we discuss how our theoretical results relate with the reliability of imaging features and behavioral measures. Overall, the current study provides empirical and theoretical insights into the relationship between prediction accuracy and feature importance reliability.

## 1. Introduction

Neuroimaging provides a non-invasive means to study human brain structure and function. *In vivo* imaging features have been linked to many clinically relevant phenotypes when contrasting populations of patients and healthy controls (Greicius et al., 2004; Kennedy et al., 2006). However, these group-level studies ignore inter-individual differences within and across patient populations (Zhang et al., 2016;

Xia et al., 2018; Zabihi et al., 2019; Tang et al., 2020; Wolfers et al., 2020). As a result, there is an increasing interest in the field to shift from group differences to accurate individual-level predictions (Dosenbach et al., 2010; Finn et al., 2015; Hsu et al., 2018; Nostro et al., 2018; Kong et al., 2019).

One goal of neuroimaging-based behavioral prediction is clinical usage to forecast practically useful clinical endpoints (Gabrieli et al., 2015). This ambition requires users to have trust in the predictive mod-

\* Corresponding author at: Centre for Sleep and Cognition, Yong Loo Lin School of Medicine, National University of Singapore, Singapore.

E-mail address: [thomas.yeo@nus.edu.sg](mailto:thomas.yeo@nus.edu.sg) (B.T.T. Yeo).

<sup>1</sup> These authors contributed equally to this work.

els, which often rests on a given models' interpretability (Bussone et al., 2015; Price, 2018; Anderson and Anderson, 2019; Diprose et al., 2020; Hedderich and Eickhoff, 2020). Indeed, the recently enacted European Union Global Data Protection Regulation (GDPR) states that patients have a right to "meaningful information about the logic involved" when automated decision-making systems are used (Vasey et al., 2022a, 2022b). Furthermore, in many studies, the derived predictive models are often interpreted to gain insights into the predictive principles and inter-individual differences that underpin observed brain-behavior relationships (Finn et al., 2015; Greene et al., 2018; Chen et al., 2022). Therefore, while many studies in the neuroimaging literature have focused on prediction accuracy (Dadi et al., 2019; He et al., 2020; Pervaiz et al., 2020; Schulz et al., 2020; Abrol et al., 2021), enhancing model interpretability remains an important issue.

One approach to interpret predictive models is the computation of feature-level importance, which quantifies the relevance of an imaging feature in the predictive model. In the case of linear models, most previous studies have interpreted the regression weights (Jiang et al., 2020; Sripada et al., 2020; Cropley et al., 2021; Xiao et al., 2021) of predictive models. However, the covariance structure among predictive features can lead to incorrect interpretations (Haufe et al., 2014). Instead, Haufe and colleagues demonstrated that it is necessary to perform an inversion of the linear models to yield the correct interpretation (Haufe et al., 2014). We refer to this inversion as the Haufe transform. Further explanation of the Haufe transform can be found in Section 2.7 and in the original study (Haufe et al., 2014).

A recent study suggested that in the context of behavioral predictions from functional connectivity (FC), the reliability of feature-level importance (original regression weights and Haufe-transformed weights) across independent samples was poor (Tian and Zalesky, 2021). Because the study utilized a maximum sample size of 400 and predicted only a small selection of cognitive measures and sex, it remains unclear whether the results generalize to other sample sizes and behavioral domains. Tian and Zalesky also found that higher resolution parcellations led to better prediction accuracy but lower feature importance reliability. However, it is unclear whether the trade-off between prediction accuracy and feature importance reliability is universal. A universal trade-off would be counterintuitive given that both feature importance reliability and prediction accuracy should reflect the reliability of brain-behavior relationship across independent datasets. More specifically, if the brain-behavior relationships in two independent data samples are highly similar, then we would expect that a model trained on one dataset to generalize well to the other dataset (i.e., high prediction accuracy). We would also expect the models trained on both datasets to be highly similar, leading to high feature importance reliability. Therefore, we hypothesize that there is not a universal trade-off between prediction accuracy and feature importance reliability.

In the present study, we used the Adolescent Brain Cognitive Development (ABCD) study to investigate the relationship between prediction accuracy and feature importance reliability. Resting-state functional connectivity was used to predict a wide range of 36 behavioral measures across cognition, personality (related to impulsivity), and mental health. We considered four commonly used prediction models: kernel ridge regression (KRR), linear ridge regression (LRR), least absolute shrinkage and selection operator (LASSO), and Random Forest (RF) models. Consistent with Tian and Zalesky (2021), we found that Haufe-transformed weights were more reliable than regression weights and univariate FC-behavior correlations. However, for sufficiently large sample sizes, we found fair to excellent split-half reliability for the Haufe-transformed weights. On the other hand, the original regression weights were unreliable even with thousands of participants. Intriguingly, feature importance reliability was strongly correlated with prediction accuracy across behavioral measures. Within a particular behavioral domain, there was no clear relationship between prediction performance and feature importance reliability across regression algorithms. **We show mathematically that split-half feature importance reliability is necessary**, but not

**sufficient, for low feature importance error.** In the case of linear models, prediction error closely reflects feature importance error. Overall, the current study provides empirical and theoretical insights into the relationship between prediction accuracy and feature importance reliability.

## 2. Methods

### 2.1. Dataset

The Adolescent Brain Cognitive Development (ABCD) dataset (2.0.1 release) was used for its large sample size, as well as its rich imaging and behavioral measures. The Institutional Review Board (IRB) at the University of California, San Diego approved all aspects of the ABCD study (Auchter et al., 2018). Parents or guardians provided written consent while the child provided written assent (Clark et al., 2018).

After quality control and excluding siblings, the final sample consisted of 5260 unrelated participants. Consistent with our previous studies (Chen et al., 2022; Ooi et al., 2022), each participant had a  $419 \times 419$  FC matrix as the imaging features, which were used to predict 36 behavioral measures across the behavioral domains of cognition, personality, and mental health.

### 2.2. Image preprocessing

Images were acquired across 21 sites in the United States with harmonized imaging protocols for GE, Philips, and Siemens scanners (Casey et al., 2018). We used structural T1 and resting-fMRI. For each participant, there were four resting-fMRI runs. Each resting-fMRI run was 300 s long. Preprocessing followed our previously published study (Chen et al., 2022). For completeness, the key preprocessing steps are summarized here.

Minimally preprocessed T1 data were used (Hagler et al., 2019). The structural data were further processed using FreeSurfer 5.3.0 (Dale et al., 1999; Fischl et al., 1999a, 1999b; Fischl et al., 2001; Ségonne et al., 2004, 2007), which generated accurate cortical surface meshes for each individual. Individuals' cortical surface meshes were registered to a common spherical coordinate system (Fischl et al., 1999a, 1999b). Individuals who did not pass recon-all quality control (Hagler et al., 2019) were removed.

Minimally preprocessed fMRI data (Hagler et al., 2019) were further processed with the following steps: (1) removal of initial frames, with the number of frames removed depending on the type of scanner (Hagler et al., 2019); and (2) alignment with the T1 images using boundary-based registration (Greve and Fischl, 2009) with FsFast (<http://surfer.nmr.mgh.harvard.edu/fswiki/FsFast>). Functional runs with boundary-based registration (BBR) costs greater than 0.6 were excluded. Framewise displacement (FD) (Jenkinson et al., 2002) and voxel-wise differentiated signal variance (DVARs) (Power et al., 2012) were computed using `fsl_motion_outliers`. Respiratory pseudomotion was filtered out using a bandstop filter (0.31–0.43 Hz) before computing FD (Power et al., 2019; Fair et al., 2020; Gratton et al., 2020). Volumes with  $FD > 0.3$  mm or  $DVARs > 50$ , along with one volume before and two volumes after, were marked as outliers and subsequently censored. Uncensored segments of data containing fewer than five contiguous volumes were also censored (Gordon et al., 2016; Kong et al., 2019). Functional runs with over half of their volumes censored and/or  $max FD > 5$  mm were removed. Individuals who did not have at least 4 min of data were also excluded from further analysis.

The following nuisance covariates were regressed out of the fMRI time series: global signal, six motion correction parameters, averaged ventricular signal, averaged white matter signal, and their temporal derivatives (18 regressors in total). Regression coefficients were estimated from the non-censored volumes. We chose to regress the global signal because we were interested in behavioral prediction, and global

signal regression has been shown to improve behavioral prediction performance (Greene et al., 2018; Li et al., 2019). The brain scans were interpolated across censored frames using least squares spectral estimation (Power et al., 2014), band-pass filtered ( $0.009 \text{ Hz} \leq f \leq 0.08 \text{ Hz}$ ), projected onto FreeSurfer fsaverage6 surface space and smoothed using a 6 mm full-width half maximum kernel.

### 2.3. Functional connectivity

We used a whole-brain parcellation comprising 400 cortical regions of interest (ROIs) (Schaefer et al., 2018) and 19 subcortical ROIs (Fischl et al., 2002). For each participant and each fMRI run, functional connectivity (FC) was computed as the Pearson's correlations between the average time series of each pair of ROIs. FC matrices were then averaged across runs, yielding a  $419 \times 419$  FC matrix for each participant. Correlation values were converted to z-scores using Fisher's r-to-z transformation prior to averaging and converted back to correlation values after averaging. Censored frames were ignored when computing FC.

### 2.4. Behavioral data

Following our previous study (Chen et al., 2022), we considered 16 cognitive, 11 mental health, and 9 impulsivity-related personality measures. The cognitive measures were vocabulary, attention, working memory, executive function, processing speed, episodic memory, reading, fluid cognition, crystallized cognition, overall cognition, short delay recall, long delay recall, fluid intelligence, visuospatial accuracy, visuospatial reaction time, and visuospatial efficiency. The mental health measures were anxious depressed, withdrawn depressed, somatic complaints, social problems, thought problems, attention problems, rule-breaking behavior, aggressive behavior, total psychosis symptoms, psychosis severity, and mania. The impulsivity-related personality measures were negative urgency, lack of planning, sensation seeking, positive urgency, lack of perseverance, behavioral inhibition, reward responsiveness, drive, and fun seeking.

Participants who did not have all behavioral measures were excluded from further analysis. As recommended by the ABCD consortium, individuals from Philips scanners were also excluded due to incorrect preprocessing. Finally, by excluding siblings, the main analysis utilized data from 5260 unrelated children.

### 2.5. Split-half cross-validation

ABCD is a multi-site dataset. To reduce sample size variability across sites, smaller sites were combined to create 10 "site-clusters", each containing at least 300 individuals (Table S1). Thus, participants within a site were in the same site-cluster.

A split-half cross-validation procedure was utilized to evaluate the prediction performance and the split-half reliability of feature importance. For each split, 5 site-clusters were selected as the training set and the remaining 5 were selected as the test set. Prediction models were trained on the training set to predict the behavioral measures from the FC matrices. The prediction models were then evaluated on the test set.

Here, we considered kernel ridge regression (KRR), linear ridge regression (LRR), and least absolute shrinkage and selection operator (LASSO) models for prediction. Hyperparameters were tuned using cross-validation within the training set (Chen et al., 2022). We also explored the use of random forests (RF) for prediction (Breiman, 2001). Because of the large number of FC features in the current study, the RF is much slower than KRR: a single RF model required 2 hours of training compared with 10 seconds for KRR. Therefore, tuning the hyperparameters of the RFs would be computationally infeasible. Consequently, the hyperparameters of the RF models were fixed with the number of trees set to 100 and the depth of each tree set to be 4.

Prediction accuracy was defined as the Pearson's correlation between the predicted and observed behavior of test participants. Feature importance of the regression models was computed in the training set (see Section 2.7). After the prediction model was trained and evaluated, the training and test sets were swapped. The model training and evaluation procedure were then repeated. Thus, for a given regression approach and interpretation method, each data split yielded two prediction accuracies and two sets of feature importance.

For each data split, the two accuracy numbers were averaged yielding an overall prediction accuracy for the split. On the other hand, the two sets of feature importance ( $f_1$  and  $f_2$ ) were used to compute split-half reliability (Tian and Zalesky, 2021), which we refer to as split-half intra-class correlation coefficient (ICC). Note that  $f_1$  and  $f_2$  are vectors of length  $K$ , where  $K$  is the total number of features. The  $k$ -th element of  $f_1$  (or  $f_2$ ) is the feature importance of the  $k$ -th feature in the first (or second) set of feature importance.

$$\text{ICC} = \frac{1}{Ks^2} \sum_{k=1}^K (f_{1,k} - \bar{f})(f_{2,k} - \bar{f}), \quad (1)$$

$f_{1,k}$  and  $f_{2,k}$  are the  $k$ -th element of  $f_1$  and  $f_2$  respectively.  $\bar{f}$  is the pooled mean computed from both  $f_1$  and  $f_2$ , given by  $\frac{1}{2K} \sum_{k=1}^K (f_{1,k} + f_{2,k})$ .

$s^2$  is the pooled variance computed from both  $f_1$  and  $f_2$ , given by  $\frac{1}{2K} \sum_{k=1}^K ((f_{1,k} - \bar{f})^2 + (f_{2,k} - \bar{f})^2)$ .

To ensure stability, the data split was repeated 126 (the number of unique ways to split ten site-clusters into two halves, which is 10 choose 5 divided by 2) times.

### 2.6. Reliability across different sample sizes

The procedure in the previous section utilized the full sample size. To evaluate feature importance reliability across different sample sizes, the previous procedure (Section 2.5) was repeated, but the participants were subsampled for each split-half cross-validation to achieve a desired sample size  $N$ . More specifically, we considered sample sizes of 200, 400, 1000, and 1500. For each sample size  $N$ , we first split the 10 site-clusters into two halves, each containing 5 site-clusters (Section 2.5).  $N/10$  samples were then randomly sampled from each site-cluster. The procedure was repeated 126 (the number of unique ways to split ten site-clusters into two halves, which is 10 choose 5 divided by 2) times.

### 2.7. Original and Haufe-transformed weights

We used KRR, LRR, LASSO and RF to predict 36 behavioral measures from FC features. In particular, the lower triangular entries of the FC matrix were used as input for the regression models. LRR, LASSO and RF are commonly used in the literature. We have previously demonstrated that KRR is a powerful approach for resting-FC behavioral prediction (He et al., 2020).

Since KRR is less commonly used in the literature, we will provide a high-level explanation here. Briefly, let  $y_i$  and  $FC_i$  be the behavioral measure and FC of training individual  $i$ . Let  $y_t$  and  $FC_t$  be the behavioral measure and FC of a test individual. Then, kernel regression would predict the test individual's behavior as the weighted average of the training individuals' behavior, i.e.  $y_t \approx \sum_{i \in \text{training set}} \text{Similarity}(FC_i, FC_t) y_i$ , where  $\text{Similarity}(FC_i, FC_t)$  was defined as the Pearson's correlation between  $FC_i$  and  $FC_t$ . Thus, kernel regression assumed that individuals with more similar FC exhibit more similar behavior. To reduce overfitting, an  $\ell_2$ -regularization term was included, which was tuned in the training set (Kong et al., 2019; Li et al., 2019; He et al., 2020).

To interpret the trained models, we considered both the regression weights and Haufe-transformed weights. Since LRR and LASSO are linear models, the regression weights were straightforward to obtain. In the

case of KRR, the kernel regression model was converted to an equivalent linear regression model, yielding one regression weight for each feature (Liu et al., 2007; Chen et al., 2022). We note that this conversion was possible because we used the correlation kernel, which is linear when the input features are pre-normalized. In the case of RF models, feature importance was extracted through calculating the out-of-bag error using a conditional permutation procedure, which reduced selection bias of correlated variables (Strobl et al., 2008). We refer to this approach as conditional variable importance.

Each prediction model was also inverted using the Haufe transform (Haufe et al., 2014). To motivate the Haufe transform (Haufe et al., 2014; Chen et al., 2022), suppose we seek to predict target variable  $y$  (e.g., fluid intelligence) from the FC of two edges ( $FC1$  and  $FC2$ ). In this example, let us assume that  $FC1 = y - motion$ , and  $FC2 = motion$ . Then a prediction model with 100% accuracy would be  $1 \times FC1 + 1 \times FC2$ . The regression weights of this model are both one for  $FC1$  and  $FC2$ . Based on the weights of the regression model, we would conclude that both  $FC1$  and  $FC2$  are strongly related to the target variable  $y$ . The Haufe transform resolves this issue by computing the covariance between the predicted target variable and each FC feature in the training set. In this toy example,  $FC2$  will be assigned a weight of zero by the Haufe transform, consistent with the intuition that  $FC2$  is not related to the target variable even though it is helpful for predicting the target variable.

Further informative examples can be found in Haufe et al. (2014). More generally, Haufe and colleagues demonstrated that for a linear predictive model, the appropriate transformation can be obtained by computing the covariance of each feature and the predicted target variable in the training set. When applied to nonlinear models, the Haufe transform recovers the “best” linear interpretation of the nonlinear models in the least square sense (Haufe et al., 2014). Therefore, our application of the Haufe transform to the random forests will only yield a partial (linear) interpretation of the random forests.

## 2.8. Mass univariate associations

Besides predictive models, we also examined the split-half reliability of mass univariate associations between FC and behavioral measures, which is sometimes referred to as brain-wide association analysis (Marek et al., 2022). We note that mass univariate associations are often used for feature selection in neuroimaging predictive models (Finn et al., 2015). The selected features are then used to interpret the model (Finn et al., 2015; Shen et al., 2017). Therefore, mass univariate associations are a good proxy for such approaches. Here, univariate association is defined as the correlation between each FC feature and each behavioral measure. To study the split-half reliability of univariate associations, we performed the same split-half procedure (Sections 2.5 and 2.6). However, instead of training a predictive model in the training set, we correlated the FC features and the behavioral measures of the training participants to obtain one t-statistic for each feature and each behavioral measure. This procedure was repeated for the test participants. Split-half reliability was defined as the split-half ICC of the t-statistic values between the two halves of the dataset (i.e., training and test sets).

## 2.9. Data and code availability

The ABCD data are publicly available via the NIMH Data Archive (NDA). Processed data from this study have been uploaded to the NDA. Researchers with access to the ABCD data will be able to download the data: <https://dx.doi.org/10.15154/1528762>. Analysis code specific to this study was can be found on GitHub: [https://github.com/ThomasYeoLab/CBIG/tree/master/stable\\_projects/predict\\_phenotypes/ChenOoi2023\\_ICCW](https://github.com/ThomasYeoLab/CBIG/tree/master/stable_projects/predict_phenotypes/ChenOoi2023_ICCW). Co-authors TWKT and SZ

reviewed the code before merging it into the GitHub repository to reduce the chance of coding errors.

## 3. Results

### 3.1. Haufe-transformed weights exhibit fair to excellent split-half reliability with large sample sizes

We computed resting-state functional connectivity (RSFC) among 400 cortical (Schaefer et al., 2018) and 19 subcortical (Fischl et al., 2002) regions for 5260 participants from the ABCD dataset (Casey et al., 2018). The lower triangular entries of the  $419 \times 419$  RSFC matrix were then vectorized to predict 36 behavioral scores that span across 3 domains: cognition, personality, and mental health.

Feature importance of KRR predictive models was interpreted using two approaches: regression weights and Haufe-transformed weights. For comparison, t-statistics from mass univariate associations were also computed. We used a split-half procedure to compute the split-half reliability of feature importance. For each split, we fitted the KRR model on each half and obtained the feature importance. The split-half reliability was defined as the split-half ICC of the feature importance values between the two halves.

Fig. 1 shows the split-half reliability of the two interpretation methods and mass univariate associations across 126 splits for different sample sizes and behavioral domains. Consistent with previous studies, split-half reliability of feature importance increases with larger sample sizes across all behavioral domains and interpretation methods (Tian and Zalesky, 2021; Marek et al., 2022). The Haufe-transformed weights were consistently more reliable than univariate associations (t-statistics), which were in turn more reliable than the regression weights. Haufe-transformed weights at a sample size of 200 were more reliable than the original regression weights at a sample size of 2630.

At the largest sample size of 2630, an average split-half ICC of 0.75 was achieved for Haufe-transformed weights of models predicting cognitive measures, which is considered “excellent” split-half reliability (Cicchetti, 1994). On the other hand, an average split-half ICC of 0.57 and 0.53 were achieved for personality and mental health at the full sample size, which are considered “fair” split-half reliabilities (Cicchetti, 1994). Under the same sample size and interpretation method, the split-half reliability of feature importance for mental health and personality was consistently lower than that of cognition.

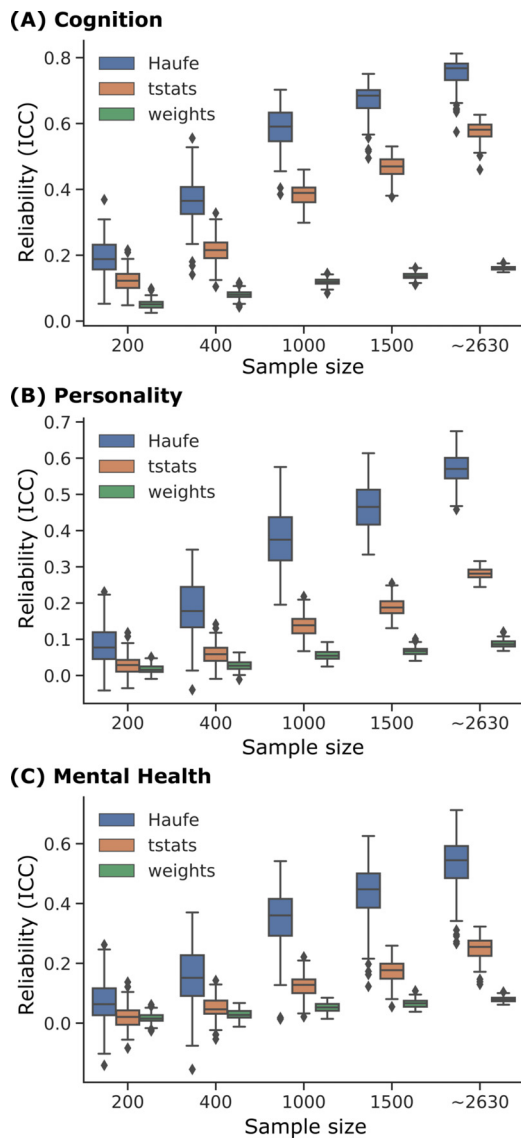
Similar conclusions were obtained with linear ridge regression (Fig. 2), LASSO (Fig. S1) and RF (Fig. S2). In the case of RF models, conditional variable importance was computed (instead of weights). Note that univariate associations (tstats) were computed independent of regression models and are therefore the same across Figs. 1, 2, S1 and S2. Overall, we found that Haufe-transformed weights achieved fair to excellent split-half reliability with sufficiently large samples.

### 3.2. Haufe-transformed weights are highly consistent across prediction models

The previous section investigated the reliability of feature importance across different data samples. Here, we seek to examine the reliability of feature importance across different prediction models in the full sample of 5260 participants. For each split-half of the 5260 participants, we computed the similarity (Pearson's correlation) of feature importance across the prediction models.

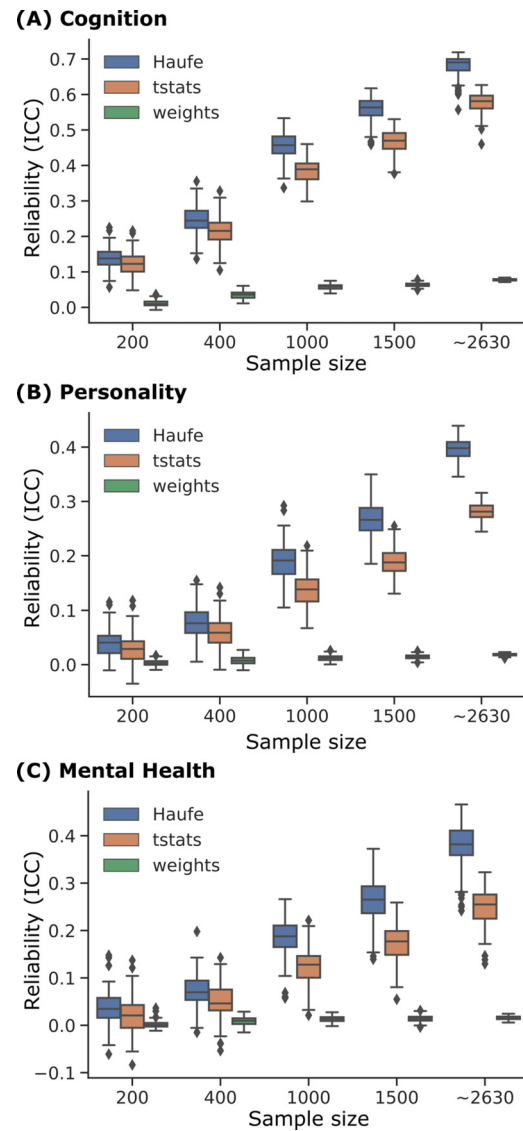
Fig. 3 shows the similarity of feature importance across prediction models. Consistent with Tian and Zalesky (2021), we found that Haufe-transformed weights showed better consistency than the original regression weights. Unlike Tian and Zalesky (2021), because of our significantly larger sample size, excellent consistency was observed for the Haufe-transformed weights (max = 0.97, min = 0.63). Interestingly, although random forests have rather different inductive biases from linear





**Fig. 1.** Split-half reliability of feature importance of kernel ridge regression (KRR) models across different sample sizes, interpretation methods, and behavioral domains: (A) cognition, (B) personality, and (C) mental health. Split-half reliability was computed as split-half interclass correlation coefficients (ICC) of feature importance obtained from two non-overlapping split-halves of the ABCD participants. After splitting, participants were randomly subsampled to show the effect of sample size on feature importance reliability. Full data without subsampling was reported as a sample size of ~2630. “~” was used because the two halves have similar (but not exactly the same) sample sizes that summed to 5260 (total number of participants). Split-half ICC values were reported for Haufe-transformed model weights (Haufe), mass univariate associations (tstats), and original regression weights (weights). Boxplots show the distribution of average split-half ICC within each behavioral domain across 126 split-half pairs. For each boxplot, the box extends from the lower to upper quartile values of the data, with a line at the median. The whiskers extend from the box to show the data range (excluding outliers). Outliers are defined as data points beyond 1.5 times the interquartile range and shown as flier points past the whiskers. Overall, across different sample sizes and behavioral domains, Haufe-transformed weights were more reliable than mass univariate associations (tstats), which were in turn more reliable than regression weights. Similar conclusions were obtained with linear ridge regression (Fig. 2), LASSO (Fig. S1) and random forests (Fig. S2).

models, the Haufe-transformed weights of the random forests still exhibited strong similarity with the linear models, especially when predicting cognition.

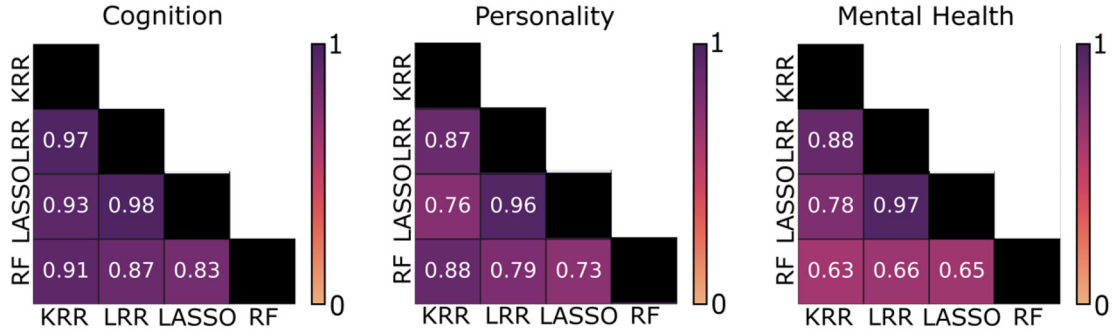


**Fig. 2.** Split-half reliability of feature importance of linear ridge regression (LRR) models across different sample sizes, interpretation methods, and behavioral domains: (A) cognition, (B) personality, and (C) mental health. Same as Fig. 1, except using LRR as the prediction model. Split-half reliability was computed as split-half interclass correlation coefficients (ICC) of feature importance obtained from two non-overlapping split-halves of the dataset. After splitting, data were randomly subsampled to show the effect of sample size on feature importance reliability. Full data without subsampling was reported as a sample size of ~2630. “~” was used because the two halves have similar (but not exactly the same) sample sizes that summed to 5260 (total number of participants). Note that mass univariate associations (tstats) were computed independent of regression models and are therefore the same across Figs. 1, 2, S1 and S2. Overall, across different sample sizes and behavioral domains, Haufe-transformed weights were more reliable than mass univariate associations (tstats), which were in turn more reliable than original regression weights.

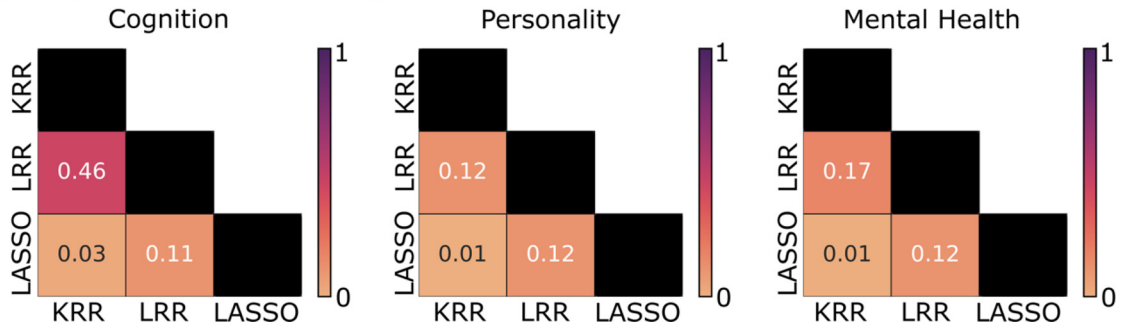
### 3.3. Feature importance reliability is strongly positively correlated with prediction accuracy across behavioral measures

So far, our results have been largely consistent with Tian and Zalesky (2021), except our larger sample sizes led to better split-half reliability of the Haufe-transformed weights. Next, we investigated the relationship between prediction accuracy and split-half reliability of feature importance using the full sample of 5260 participants.

## (A) Haufe-transformed weights



## (B) Original regression weights



**Fig. 3.** Similarity of feature importance across three predictive models in the full sample of 5260 participants. (A) Consistency of feature importance for Haufe-transformed weights. (B) Consistency of feature importance for original regression weights. We note that RF did not have regression weights, so did not appear in panel B. Similarity was computed as the Pearson's correlation between feature importance values across different predictive models (KRR, LRR, LASSO and RF). Similarity was computed for each split-half and then averaged across the 126 data splits. Excellent consistency was observed for the Haufe-transformed weights.

Split-half reliability and prediction accuracy of each behavioral score were computed for each split-half of the dataset, followed by averaging across the 126 data splits. Fig. 4A shows the correlation between feature importance reliability and prediction accuracy across the 36 behavioral measures for KRR. Prediction accuracy was highly correlated with split-half reliability of Haufe-transformed model weights ( $r = 0.78$ ), t-statistics ( $r = 0.94$ ) and original regression weights ( $r = 0.97$ ). This suggests that a behavioral measure that was predicted with higher accuracy also enjoyed better feature importance reliability.

Similar conclusions were obtained with linear ridge regression (Fig. 4B), LASSO (Fig. 4C) and RF (Fig. 4D). Overall, we found a strong positive relationship between feature importance reliability and prediction accuracy. We repeated the analysis with the behavioral component scores from Ooi et al. (2022) and found similar positive relationships (Fig. S3).

Furthermore, in the case of Haufe transform and univariate associations (t-stats), there appears to be a nonlinear relationship between prediction accuracies and split-half ICC (Fig. 4). More specifically, higher accuracies led to greater split-half ICC, but with diminishing returns for behavioral measures with higher accuracies.

### 3.4. No clear relationship between prediction accuracy and feature importance reliability across predictive models

Table 1 summarizes average prediction accuracies for different behavioral domains, as well as split-half ICC of feature importance using the full sample of 5260 participants.

Among the linear models, KRR exhibited the highest split-half ICC, but not necessarily the best prediction performance. LASSO generally

had the worse prediction performance and the worst split-half ICC. LRR exhibited the best prediction performance, but an intermediate level of split-half ICC. On the other hand, the RF models exhibited good prediction performance and Haufe-transform split-half ICC for cognition and personality, but not for mental health. Overall, there was no clear relationship between prediction performance and feature importance reliability.

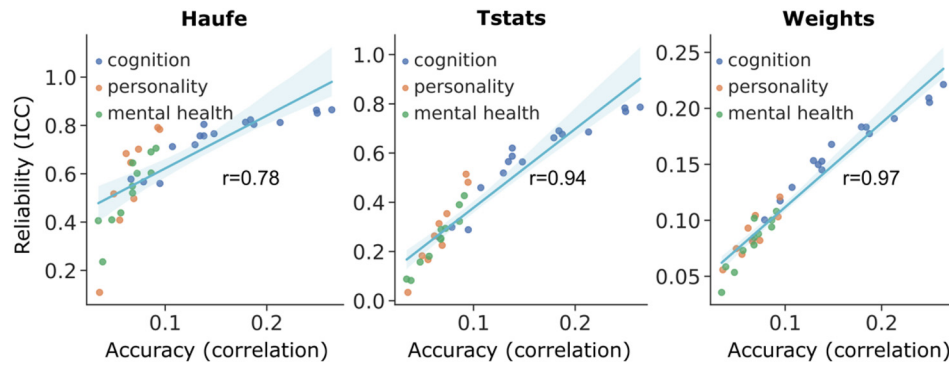
Note that in our other studies (Chen et al., 2022; Ooi et al., 2022), the prediction performance of KRR was similar to (or slightly better) than LRR, suggesting that depending on the dataset (or even across different samples within the same dataset), prediction accuracies can vary across prediction approaches.

### 3.5. Split-half reliability is necessary, but not sufficient, for correct feature importance

We have shown a strong positive correlation between feature importance reliability and prediction accuracy (Fig. 4). There is also a lack of relationship between prediction accuracy across prediction models and feature importance reliability (Table 1). In the remaining sections of this study, we will delve more deeply into the mathematical relationships among feature importance reliability, feature importance error and prediction error.

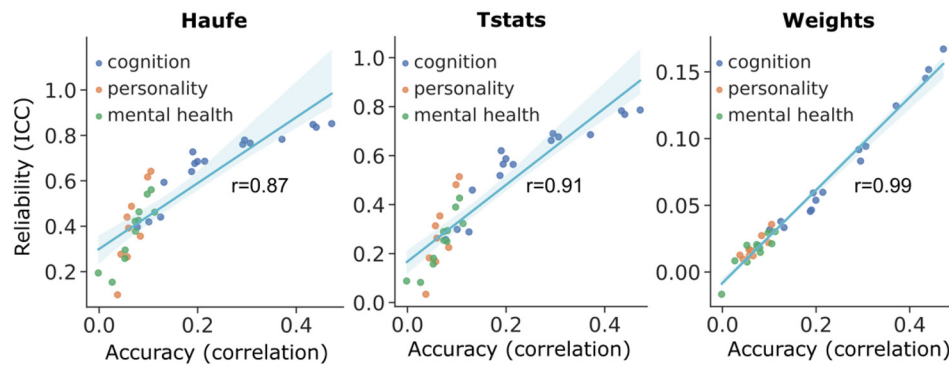
We begin by showing that split-half feature importance reliability is necessary but not sufficient for obtaining the “correct” feature importance. Let  $f_G$  be the hypothetical ground-truth feature importance that might be derived assuming the correct generative process relating brain features and behavioral measures is known. However, in the following analysis, we do not assume the ground truth generative process is known

## (A) Correlation between reliability and accuracy for KRR

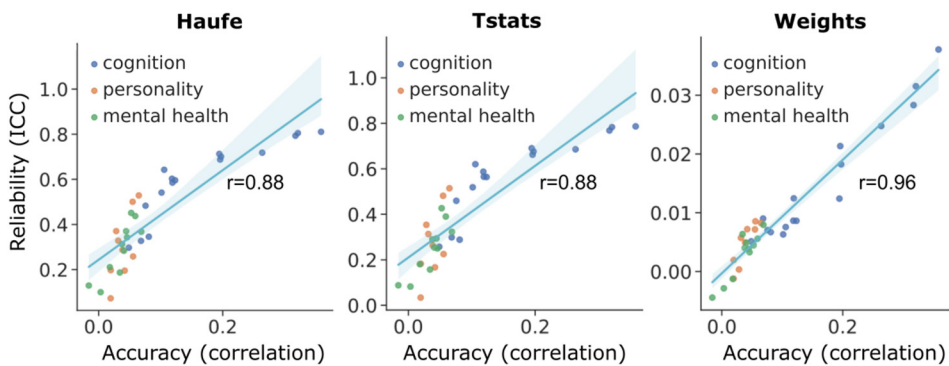


**Fig. 4.** Split-half reliability of feature importance is positively correlated with prediction accuracy across 36 behavioral measures for (A) kernel ridge regression (KRR), (B) linear ridge regression (LRR), (C) LASSO and (D) random forest (RF). Split-half reliability and prediction accuracy of each behavioral score were computed for each split-half of the dataset, followed by averaging across the 126 data splits.

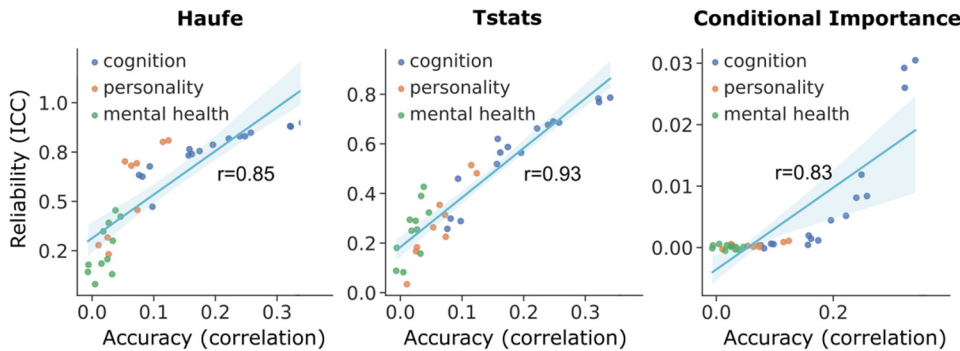
## (B) Correlation between reliability and accuracy for LRR



## (C) Correlation between reliability and accuracy for LASSO



## (D) Correlation between reliability and accuracy for RF



**Table 1**

Summary of average prediction performance for cognitive, personality and mental health measures, as well as split-half ICC of Haufe-transformed weights, original weights, conditional variable importance and univariate associations (t-statistics). In general, within a behavioral domain (e.g., cognition), lower (or higher) prediction performance for a given predictive model was not necessarily associated with lower (or higher) split-half ICC.

Cognition	Corr	ICC (Haufe)	ICC (Weights)	ICC (Conditional variable importance)	ICC (Univariate association)
KRR	0.16	0.75	0.16	N.A.	0.58
LRR	0.25	0.68	0.08	N.A.	
LASSO	0.17	0.60	0.02	N.A.	
RF	0.20	0.76	N.A.	0.01	
Personality	Corr	ICC (Haufe)	ICC (Weights)	ICC (Conditional variable importance)	ICC (Univariate association)
KRR	0.07	0.57	0.09	N.A.	0.28
LRR	0.07	0.40	0.02	N.A.	
LASSO	0.04	0.30	0.01	N.A.	
RF	0.06	0.55	N.A.	0.00	
Mental Health	Corr	ICC (Haufe)	ICC (Weights)	ICC (Conditional variable importance)	ICC (Univariate association)
KRR	0.07	0.53	0.08	N.A.	0.25
LRR	0.07	0.38	0.02	N.A.	
LASSO	0.04	0.29	0.01	N.A.	
RF	0.02	0.26	N.A.	0.00	

and we make no assumption about how  $f_G$  can be computed even if the ground truth generative process is known.

Let  $f_S$  be the feature importance estimated from data sample  $S$ . Both  $f_G$  and  $f_S$  are  $D \times 1$ , where  $D$  is the number of features. The expected feature importance error can be defined as the expectation of the squared error across different data samples  $S$ :  $E_S[(f_G - f_S)^T(f_G - f_S)]$ . Let  $\bar{f}_S = E_S[f_S]$  be the feature importance averaged across all possible data samples  $S$ . The feature importance error can then be decomposed into two terms:

$$E_S[(f_G - f_S)^T(f_G - f_S)] = (f_G - \bar{f}_S)^T(f_G - \bar{f}_S) + E_S[(f_S - \bar{f}_S)^T(f_S - \bar{f}_S)] \quad (2)$$

The proof is provided in the Appendix A. The decomposition of feature importance error as in Eq. (2) is similar in spirit (and derivation) to the classical bias-variance decomposition of prediction error.

The first term  $(f_G - \bar{f}_S)^T(f_G - \bar{f}_S)$  in Eq. (2) measures the bias of the feature importance estimation procedure. The second term  $E_S[(f_S - \bar{f}_S)^T(f_S - \bar{f}_S)]$  measures the variance of the estimated feature importance across different samples, which is the opposite of reliability. In other words, higher variance in feature importance estimation is the same as lower reliability. Therefore, from Eq. (2), we note that low feature importance variance (i.e., high feature importance reliability) is necessary but not sufficient for low feature importance error. Low feature importance variance must be coupled with low feature importance bias to achieve a small feature importance estimation error.

### 3.6. Prediction error reflects feature importance error for linear models

The previous section shows that the reliability of feature importance is not sufficient for low feature importance error. In this section, we show that when the ground truth data generation model is linear and feature importance is defined as regression weights (or Haufe-transformed weights), then the prediction error is directly related to the feature importance error.

A linear regression model assumes that the data is generated through a linear combination of features. For example, assume that a given data point  $(x_i, y_i)$  is generated by a linear model  $y_i = x_i^T w_G + \epsilon$ . Here,  $y_i$  is a scalar,  $x_i$  is a  $D \times 1$  vector,  $w_G$  is the groundtruth  $D \times 1$  regression weights, and  $D$  is the number of features.  $\epsilon$  is an independent noise term with zero mean. Without loss of generality, we assume that the expectation of  $y$  across data samples is 0 and the expectation of  $x$  across data samples is 0 for every feature. In the case of FC prediction of behavioral traits, each data sample is a participant.

Suppose data sample  $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$  is drawn as the training set. We can then train a linear regression model (e.g., LRR or LASSO) on  $S$  and obtain the regression weights  $w_S$ . The resulting prediction model will be  $\hat{y} = x^T w_S$ . Let the difference between the ground truth and estimated weights be  $\Delta_w(S) = w_G - w_S$ . Thus, the regression weights error (on average across different training sets  $S$ ) can be defined as  $E_S[(w_G - w_S)^T(w_G - w_S)] = E_S[\Delta_w(S)^T \Delta_w(S)]$ .

On the other hand, the expected prediction error of the prediction algorithm can be defined as  $E_S E_{x,y}[(y - x^T w_S)^2]$ . Here,  $E_{x,y}$  is the expectation of the squared prediction error over out-of-sample test data points sampled from the distribution of  $(x, y)$ . We note that the test data points are sampled independently from the sampling of the training dataset  $S$ . Then, the expected test error can be decomposed into:

$$E_S E_{x,y}[(y - x^T w_S)^2] = \text{Var}(\epsilon) + E_S[\Delta_w(S)^T * \text{COV}(X) * \Delta_w(S)] \quad (3)$$

The proof is found in the Appendix B. In Eq. (3) the first term is the irreducible error  $\text{Var}(\epsilon)$ , which is the variance of the noise. The second term  $E_S[\Delta_w(S)^T * \text{COV}(X) * \Delta_w(S)]$  is determined by both the regression weights error  $\Delta_w(S)$  and the covariance of features  $\text{COV}(X)$ .

We can consider three different scenarios for the covariance matrix  $\text{COV}(X)$ . First, suppose  $\text{COV}(X)$  is an identity matrix, which implies the features are independent and of unit variance. Then, the prediction error (Eq. (3)) can be written as  $\text{Var}(\epsilon) + E_S[\Delta_w(S)^T \Delta_w(S)]$ . Therefore, the prediction error is simply the sum of the regression weights error and the irreducible error.

Second, suppose  $\text{COV}(X)$  is a diagonal matrix, i.e.,  $\text{COV}(X) = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_d)$ , which implies the features are independent. In this case, the prediction error (Eq. (3)) can be written as  $\text{Var}(\epsilon) + E_S[\sum_{d=1}^D \sigma_d \Delta_{w(d)}(S)^2]$ . Here,  $\Delta_{w(d)}(S)$  is the regression weight error of the  $d$ -th feature based on the training dataset  $S$ . In this scenario, a bigger regression weights error still leads to a bigger prediction error, but the weights error of features with larger variance results in a larger prediction error than features with small variance.

Third, suppose we do not make any independence assumptions about the features. Since  $\text{COV}(X)$  is a symmetric matrix, we can decompose  $\text{COV}(X)$  as  $\text{COV}(X) = R^T D R$ . Here,  $R$  is a rotation matrix where  $R^T R$  is equal to an identity matrix and  $D$  is a diagonal matrix. Then, we can rewrite the prediction error (Eq. (3)) as:

$$\text{Var}(\epsilon) + E_S[\Delta_w(S)^T * R^T * D * R * \Delta_w(S)] \quad (4)$$

To summarize the three scenarios for  $\text{COV}(X)$ , regression weights errors of all features are related to prediction error, but features with a



larger variance (up to a rotation) have a stronger relationship with the prediction error.

We can also establish a similar relationship between the Haufe-transformed weights error and the prediction error. Note that the Haufe-transformed weights can be computed as  $(X_S)^T * w_S$ . Here the  $w_S$  is the original regression weights and  $COV(X_S)$  is the feature covariance of training sample equation S. Assuming that the sample covariance is close to the true covariance, i.e.,  $COV(X_S) \approx COV(X)$ , then the Haufe-transformed weights error can be written as:

$$\begin{aligned} E_S [\Delta_w(S)^T * COV(X) * COV(X) * \Delta_w(S)] \\ = E_S [\Delta_w(S)^T * R^T D^2 R * \Delta_w(S)] \end{aligned} \quad (5)$$

Comparing the Haufe-transformed weights error (Eq. (5)) with the prediction error (Eq. (4)), we see that the Haufe-transformed weights error is closely related to the prediction error, given that Eqs. (4) and (5) only differ by the square of the diagonal matrix  $D$ .

Overall, we conclude that higher original regression weights errors (Eq. (4)) and higher Haufe-transformed errors (Eq. (5)) are related to greater prediction error up to a scaling by the feature covariance matrix.

#### 4. Discussion

We have provided empirical and theoretical evidence on the relationship between prediction accuracy and feature importance reliability.

##### 4.1. Haufe-transformed model weights are more reliable than original regression weights and univariate FC-behavior correlations

Consistent with Tian and Zalesky (2021), we found that Haufe-transformed weights were much more reliable than original regression weights. In our experiments, we note that even with a sample size of ~2630 participants, the original kernel regression weights achieved a split-half ICC of less than 0.2 when predicting cognitive measures, which is less than the split-half ICC of Haufe-transformed weights with a sample size of 200 (Fig. 1A). This is perhaps not surprising since it has been empirically shown that regression weights contain more noise than the Haufe-transformed weights (Haufe et al., 2014). Furthermore, for predictive models with sparse regularization (e.g., LASSO), it is well-known that noise in the features can lead to very different features being selected, which will lead to low split-half reliability in the regression weights. In the case of random forests, the poor split-half reliability was probably due to the large number of functional connectivity features (87,571 features). Using a random forest with 100 trees and a depth of 4 would mean that a maximum of 1500 unique features being chosen, so the conditional variable importance was susceptible to the random choice of features included in the random forest.

Also consistent with Tian and Zalesky (2021), we found that Haufe-transformed weights were more reliable than univariate brain-behavior correlations. In our experiments, we note that with a sample size of ~2630 participants, the univariate FC-behavior correlations achieved a split-half ICC of less than 0.6 for cognitive measures, which is less than the split-half ICC of Haufe-transformed weights with a sample size of 1000 (Fig. 1A). The higher split-half ICC of Haufe-transformed weights over univariate associations is somewhat surprising. A previous study has suggested that the predicted outcomes of predictive models is substantially more reliable than the functional connectivity features themselves (Taxali et al., 2021). Here, we speculate that the predicted behavioral measures might even be more reliable than the raw behavioral measures themselves. The reason is that the regularization of many predictive models serves to “shrink” the predicted outcomes towards the population mean, which should increase reliability. If predicted behavioral measures are more reliable than raw behavioral measures, then the covariance of the predicted behavioral measures with FC (i.e., haufe-transformed weights) should be more reliable than the correlation between raw behavioral measures and FC (i.e., univariate associations).

It is also worth mentioning that Tian and Zalesky (2021) found that the split-half ICC of Haufe-transformed weights remained lower than 0.4 across split-half of 800 participants (i.e., two groups of 400 participants), which is consistent with our results (see sample size of 400 in our Figs. 1 and 2). Not surprisingly, we obtained higher reliability with larger sample sizes. More specifically, with a sample size of about 2600 participants, Haufe-transformed weights achieve average intra-class correlation coefficients of 0.75, 0.57 and 0.53 for cognitive, personality and mental health measures respectively (Fig. 1). Overall, the use of Haufe-transformed weights might help to alleviate reliability issues highlighted in previous neuroimaging studies (Kharabian Masouleh et al. 2019; Marek et al., 2022). On the other hand, we recommend that regression weights should not be used for model interpretation given their low split-half reliability even in the large sample regime of a few thousand participants.

##### 4.2. There is not always an empirical trade-off between feature importance reliability and prediction accuracy

Tian and Zalesky (2021) found that FC-based prediction using lower resolution atlases (compared with higher resolution atlases) had higher feature importance reliability but lower prediction accuracy. Our study suggests that this trade-off between prediction accuracy and feature importance reliability is not universal. For example, we found that behavioral measures that are predicted better also enjoy better feature importance reliability (Fig. 4).

Furthermore, in our current study, within a behavioral domain, there was no clear relationship between prediction performance and feature importance reliability across regression algorithms (Table 1). Similarly, as can be seen in Figure 2 of Tian and Zalesky (2021), higher prediction accuracy does not necessitate lower split-half reliabilities, e.g., kernel ridge regression enjoyed better prediction accuracy and feature importance reliability than connectome-based predictive modeling.

Overall, these empirical results show that it is possible to achieve high prediction accuracy and high feature importance reliability, suggesting that there is not always a trade-off between prediction accuracy and feature importance reliability.

##### 4.3. There is not a theoretical trade-off between feature importance reliability and prediction accuracy

Eq. (2) shows that feature importance reliability is necessary but not sufficient for obtaining the “correct” feature importance (or low feature importance error). More specifically, feature importance error can be decomposed into a bias term and a variance term, where the variance term is the opposite of feature importance reliability. Consequently, low feature importance variance (i.e., high feature importance reliability) is necessary but not sufficient for low feature importance error.

This result echoes previous studies in neuroimaging (Noble et al., 2017), as well as other areas of quantitative research (Kirk and Miller, 1986), demonstrating that reliability is not the same as validity. To give an extreme example, if we utilized an extremely strong regularization in our regression models, the regression weights would be driven to zero. In this scenario, the feature importance (regression weights) would be highly reliable across data samples, but the feature importance would not be valid or close to the ground truth values (derived from the ground truth generative process).

In the case of linear models, we further showed in Eq. (3) that higher feature importance error (operationalized by original regression weights) is related to worse prediction accuracy, up to a rotation and scaling by the feature covariance matrix. In Eq. (5), we showed that higher feature importance error (operationalized by Haufe-transformed weights) is related to worse prediction accuracy, up to a scaling of the eigenvalues of the feature covariance matrix.

Overall, these theoretical results suggest that at least in the case of linear models, there is not necessarily a trade-off between feature impor-

tance reliability and prediction accuracy. In fact, improving prediction performance might even reduce feature importance error and potentially improve feature importance reliability.

Given the link between feature importance error and prediction performance, in some sense, feature importance error might be a more meaningful metric than feature importance reliability. However, we cannot directly measure feature importance error, so a useful proxy might be to consider both feature importance reliability and prediction performance.

#### 4.4. Reliability of functional connectivity and behavioral measures

There is a significant literature on the reliability of FC (Noble et al., 2019). Recent studies have also emphasized that the reliability of behavioral measures (in addition to FC reliability) is important for good prediction performance (Nikolaidis et al., 2022; Gell et al., 2023). How do FC and behavioral reliability relate to our theoretical results?

Recall that assuming  $y = x^T w_G + \epsilon$ , then as shown in Eq. (3), the prediction error can be written as  $Var(\epsilon) + E_S[\Delta_w(S)^T * COV(X) * \Delta_w(S)]$ , where  $COV(X)$  is the  $D \times D$  feature covariance matrix and  $\Delta_w(S)$  is the  $D \times 1$  regression weights error.  $Var(\epsilon)$  is the variance of the irreducible noise  $\epsilon$ , which can be thought of as the variance of the behavioral measure unrelated to the features.

To think about the effect of the reliability of behavioral measure  $y$ , suppose we add more noise to the behavioral measure  $y$ , so that  $y = x^T w_G + \epsilon_2$ , where  $Var(\epsilon_2) > Var(\epsilon)$ . In this scenario, the prediction error can now be written as  $Var(\epsilon_2) + E_S[\Delta_w(S)^T * COV(X) * \Delta_w(S)]$ . Therefore, the equation is basically the same as before except for the larger noise variance  $Var(\epsilon_2)$ . In addition, the larger noise in  $y$  would also lead to greater regression weights error  $\Delta_w(S)$ . Overall, worse behavioral reliability leads to larger  $Var(\epsilon)$  and  $\Delta_w(S)$ , and thus worse prediction error. The worse regression weights error  $\Delta_w(S)$  is in turn associated with worse feature importance bias and/or reliability (via Eq. (2)).

On the other hand, suppose we add more noise to the FC features  $x$  to reduce FC reliability, so that the new features  $x_2 = x + \epsilon_2$ , where  $\epsilon_2$  is a  $D \times 1$  noise vector with zero mean. Then  $y = x^T w_G + \epsilon = x_2^T w_G - \epsilon_2^T w_G + \epsilon$ . Using the same derivation as Appendix B, the prediction error can now be decomposed as  $Var(\epsilon) + Var(\epsilon_2^T w_G) + E_S[\Delta_w(S)^T * COV(X_2) * \Delta_w(S)]$ . Therefore, worse FC reliability leads to an additional noise term  $Var(\epsilon_2^T w_G)$ , larger feature covariance matrix  $COV(X_2)$  and greater regression weight error  $\Delta_w(S)$ , thus leading to worse prediction performance. The worse regression weights error  $\Delta_w(S)$  is in turn associated with worse feature importance bias and/or reliability (via Eq. (2)).

It is important to note that perfect feature and behavioral reliability is not a panacea. For example, if features and behavioral measure both have perfect reliability, but they are not related to each other, e.g.,  $x$  and  $y$  both follow white Gaussian distributions. Then,  $y = x^T 0 + \epsilon$  (i.e.,  $w_G = 0$ ), and the prediction error of  $y$  still follows Eq. (3):  $Var(\epsilon) + E_S[\Delta_w(S)^T * COV(X) * \Delta_w(S)]$ . In this case,  $Var(\epsilon)$  is simply the variance of the behavioral measure. We note that if we simply predict the mean of the behavioral measure (i.e., completely ignore  $x$ ), then the prediction error of  $y$  will be  $Var(\epsilon)$ . On the other hand, if we are fitting a model to predict  $y$  from  $x$ , then because of the finite sample size, the regression weights will not be equal to the ground truth, so the regression weight error  $\Delta_w(S)$  is actually non-zero. Therefore, the overall prediction will be worse than simply predicting the mean of the behavioral measure, which would lead to a negative coefficient of determinant (a measure of prediction performance).

#### 4.5. Reconciling theoretical and empirical results

Our theoretical results suggest a link between feature importance reliability and prediction performance.

Consistent with the theoretical results, there was empirically a strong correlation between feature importance reliability and prediction performance across behavioral measures (Fig. 4). Similar to our previous studies (Kong et al., 2021; Chen et al., 2022; Ooi et al., 2022), cognitive measures were predicted better than other behavioral measures (Figs. 1, 2 and 4). One possible explanation for the variation in prediction performance across behavioral measures might be the reliability of the behavioral measures, as discussed in Section 4.4 and previous studies (Nikolaidis et al., 2022; Gell et al., 2023). Another possible explanation is the strength of the relationship between FC features and target behavioral measures (again discussed in Section 4.4). Therefore, behavioral measures with higher reliability and/or stronger relationship with FC features might be predicted better, as well as enjoyed better feature importance error and reliability.

On the other hand, there was empirically not a clear relationship between prediction performance and feature importance reliability across predictive models (Table 1). For example, when predicting cognition, KRR exhibited worse prediction performance than LRR (0.16 versus 0.25), but better Haufe-transformed feature importance reliability (0.75 vs 0.68). There are several ways these empirical and theoretical results can be reconciled.

First, although KRR exhibited better Haufe-transformed feature importance reliability than LRR, it is possible that KRR had worse feature importance bias than LRR, so that the overall feature importance error is worse than LRR, resulting in worse prediction performance.

Second, recall that according to Eq. (5), the prediction error can be expressed as  $E_S[\Delta_w(S)^T * COV(X) * COV(X) * \Delta_w(S)]$ , where  $\Delta_w(S)$  is the  $D \times 1$  feature importance error (where  $D$  is the number of features) and  $COV(X)$  is the  $D \times D$  covariance matrix of the features. Because of the middle covariance term, not all feature importance errors are equally important. It is possible that KRR has lower feature importance errors on average across all features ( $\Delta_w(S)^T \Delta_w(S)$ ) than LRR, but the feature importance error is greater for certain features that are more intrinsically linked to the prediction error via the covariance term  $COV(X) * COV(X)$ .

Third, Eq. (5) assumes that the true data generation process is linear, i.e., there is a linear relationship between FC features and the target variable. Therefore, Eq. (5) might not hold if the true relationship between FC features and target variable is nonlinear.

Finally, the prediction error in Eq. (5) is the average across infinite instances of training set  $S$  and an infinite test set. Therefore, the equation can be violated in the finite sample scenario (Table 1).

## 5. Conclusion

In this study, we show that Haufe-transformed weights are much more reliable than original regression weights when computing feature importance. Furthermore, feature importance reliability is strongly positively correlated with prediction accuracy across phenotypes. However, within a particular behavioral domain, there is no clear relationship between prediction performance and feature importance reliability across regression models. We also show mathematically that feature importance reliability is necessary, but not sufficient, for low feature importance error. In the case of linear models, lower feature importance error is mathematically related to lower prediction error. Therefore, higher feature importance reliability might yield lower feature importance error and higher prediction accuracy. Overall, our study provides theoretical and empirical insights into the relationships among imaging feature reliability, behavioral measure reliability, feature importance reliability and behavioral prediction accuracies.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Credit authorship contribution statement

**Jianzhong Chen:** Conceptualization, Methodology, Software, Validation, Project administration, Writing – original draft. **Leon Qi Rong Ooi:** Conceptualization, Methodology, Software, Validation, Writing – review & editing. **Trevor Wei Kiat Tan:** Software, Validation. **Shaoshi Zhang:** Software, Validation. **Jingwei Li:** Methodology, Writing – review & editing. **Christopher L. Asplund:** Methodology, Writing – review & editing. **Simon B Eickhoff:** Methodology, Writing – review & editing. **Danilo Bzdok:** Methodology, Writing – review & editing. **Avram J Holmes:** Methodology, Writing – review & editing. **B.T. Thomas Yeo:** Conceptualization, Methodology, Supervision, Writing – review & editing, Funding acquisition.

## Data availability

The authors do not have permission to share data.

## Acknowledgements

Our research is currently supported by the Singapore National Research Foundation (NRF) Fellowship (Class of 2017), the NUS Yong Loo Lin School of Medicine (NUHSRO/2020/124/TMR/LOA), the Singapore National Medical Research Council (NMRC) LCG (OFLCG19May-0035), NMRC STaR (STaR20nov-0003), and the USA NIH (R01MH120080). Our computational work was partially performed on resources of the National Supercomputing centre, Singapore (<https://www.nscg.sg>). Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not reflect the views of the Singapore NRF or the Singapore NMRC.

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development<sup>SM</sup> (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9–10 and follow them over 10 years into early adulthood. The ABCD Study<sup>®</sup> is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at [https://abcdstudy.org/consortium\\_members/](https://abcdstudy.org/consortium_members/). ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators. The ABCD data repository grows and changes over time. The ABCD data used in this report came from <http://dx.doi.org/10.15154/1504041>.

## Supplementary materials

Supplementary material associated with this article can be found, in the online version, at [doi:10.1016/j.neuroimage.2023.120115](https://doi.org/10.1016/j.neuroimage.2023.120115).

## Appendix A

In this appendix, we will provide proof of Eq. (2), which decomposes the feature importance error  $E_S[(f_G - f_S)^T(f_G - f_S)]$  into a bias term

$$\begin{aligned} & (f_G - \bar{f}_S)^T(f_G - \bar{f}_S) \text{ and a variance term } E_S[(f_S - \bar{f}_S)^T(f_S - \bar{f}_S)]. \\ & E_S[(f_G - f_S)^T(f_G - f_S)] \\ & = E_S[(f_G - \bar{f}_S) - (f_S - \bar{f}_S)]^T[(f_G - \bar{f}_S) - (f_S - \bar{f}_S)] \\ & = E_S[(f_G - \bar{f}_S)^T(f_G - \bar{f}_S)] + E_S[(f_S - \bar{f}_S)^T(f_S - \bar{f}_S)] \\ & \quad - 2E_S[(f_G - \bar{f}_S)^T(f_S - \bar{f}_S)] \\ & = (f_G - \bar{f}_S)^T(f_G - \bar{f}_S) + E_S[(f_S - \bar{f}_S)^T(f_S - \bar{f}_S)], \end{aligned}$$

where the last equality is true because  $E_S[(f_G - \bar{f}_S)^T(f_S - \bar{f}_S)] = (f_G - \bar{f}_S)^T(\bar{f}_S - \bar{f}_S) = 0$ .

## Appendix B

In this appendix, we will provide proof of Eq. (3), which establishes the relationship between the prediction error  $E_S E_{x,y}[(y - x^T w_S)^2]$  and regression weights error  $\Delta_w(S)$ , assuming an underlying linear model  $y_i = x_i^T w_G + \epsilon$ :

$$\begin{aligned} & E_S E_{x,y}[(y - x^T w_S)^2] \\ & = E_S E_{x,y}[(x^T w_G + \epsilon - x^T w_S)^2] \\ & = E_S E_{x,y}[(x^T \Delta_w(S) + \epsilon)^2], \text{ where } \Delta_w(S) = w_G - w_S \\ & = E_S E_{x,y}[\epsilon^2] + E_S E_{x,y}[(x^T \Delta_w(S))^2] + 2 * E_S E_{x,y}[\epsilon * x^T \Delta_w(S)] \\ & = Var(\epsilon) + E_S E_{x,y}[\Delta_w(S)^T x x^T \Delta_w(S)], \text{ because } E_{x,y}(\epsilon * x^T) \\ & = E_{x,y}(\epsilon) E_{x,y}(x^T) = 0 \\ & = Var(\epsilon) + E_S[\Delta_w(S)^T * COV(X) * \Delta_w(S)] \end{aligned}$$

## References

- Abrol, A., Fu, Z., Salman, M., Silva, R., Du, Y., Plis, S., Calhoun, V., 2021. Deep learning encodes robust discriminative neuroimaging representations to outperform standard machine learning. *Nat. Commun.* 12 (1), 353.
- Anderson, M., Anderson, S.L., 2019. How should AI be developed, validated, and implemented in patient care? *AMA J. Ethics* 21 (2), E125–E130.
- Auchter, A.M., Hernandez Mejia, M., Heyser, C.J., Shilling, P.D., Jernigan, T.L., Brown, S.A., Tapert, S.F., Dowling, G.J., 2018. A description of the ABCD organizational structure and communication framework. *Dev. Cogn. Neurosci.* 32, 8–15.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- Bussone, A., Stumpf, S., O'Sullivan, D., 2015. The role of explanations on trust and reliance in clinical decision support systems. In: *Proceedings of the International Conference on Healthcare Informatics*, pp. 160–169. [ieeexplore.ieee.org](https://ieeexplore.ieee.org).
- Casey, B.J., Cannonier, T., Conley, M.L., Cohen, A.O., Barch, D.M., Heitzeg, M.M., Soules, M.E., Teslovich, T., Dellarco, D.V., Garavan, H., Orr, C.A., Wager, T.D., Banich, M.T., Speer, N.K., Sutherland, M.T., Riedel, M.C., Dick, A.S., Bjork, J.M., Thomas, K.M., Chaarani, B., Mejia, M.H., Hagler, D.J., Daniela Cornejo Jr, M., Scat, C.S., Harms, M.P., Dosenbach, N.U.F., Rosenberg, M., Earl, E., Bartsch, H., Watts, R., Polimeni, J.R., Kuperman, J.M., Fair, D.A., Dale, A.M. ABCD Imaging Acquisition Workgroup, 2018. The Adolescent Brain Cognitive Development (ABCD) study: imaging acquisition across 21 sites. *Dev. Cogn. Neurosci.* 32, 43–54.
- Chen, J., Tam, A., Kebets, V., Orban, C., Ooi, L.Q.R., Asplund, C.L., Marek, S., Dosenbach, N.U.F., Eickhoff, S.B., Bzdok, D., Holmes, A.J., Yeo, B.T.T., 2022. Shared and unique brain network features predict cognitive, personality, and mental health scores in the ABCD study. *Nat. Commun.* 13 (1), 2217.
- Cicchetti, D.V., 1994. Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. *Psychol. Assess.* 6 (4), 284–290.
- Clark, D.B., Fisher, C.B., Bookheimer, S., Brown, S.A., Evans, J.H., Hopfer, C., Hudziak, J., Montoya, I., Murray, M., Pfefferbaum, A., Yurgelun-Todd, D., 2018. Biomedical ethics and clinical oversight in multisite observational neuroimaging studies with children and adolescents: the ABCD experience. *Dev. Cogn. Neurosci.* 32, 143–154.
- Cropley, V.L., Tian, Y., Fernando, K., Mansour, L.S., Pantelis, C., Cocchi, L., Zalesky, A., 2021. Brain-Predicted Age Associates With Psychopathology Dimensions in Youths. *Biol. Psychiatry Cogn. Neurosci. Neuroimaging* 6 (4), 410–419.
- Dadi, K., Rahim, M., Abraham, A., Chyzyk, D., Milham, M., Thirion, B., Varoquaux, G. Alzheimer's Disease Neuroimaging Initiative, 2019. Benchmarking functional connectome-based predictive models for resting-state fMRI. *Neuroimage* 192, 115–134.



- Dale, A.M., Fischl, B., Sereno, M.I., 1999. Cortical surface-based analysis. I. Segmentation and surface reconstruction. *Neuroimage* 9 (2), 179–194.
- Diprose, W.K., Buist, N., Hua, N., Thuermer, Q., Shand, G., Robinson, R., 2020. Physician understanding, explainability, and trust in a hypothetical machine learning risk calculator. *J. Am. Med. Inform. Assoc. JAMIA* 27 (4), 592–600.
- Dosenbach, N.U.F., Nardos, B., Cohen, A.L., Fair, D.A., Power, J.D., Church, J.A., Nelson, S.M., Wig, G.S., Vogel, A.C., Lessov-Schlaggar, C.N., Barnes, K.A., Dubis, J.W., Feczko, E., Coalson, R.S., Pruett, J.R., Barch Jr, D.M., Petersen, S.E., Schlaggar, B.L., 2010. Prediction of individual brain maturity using fMRI. *Science* 329 (5997), 1358–1361 (New York, N.Y.).
- Fair, D.A., Miranda-Dominguez, O., Snyder, A.Z., Perrone, A., Earl, E.A., Van, A.N., Koller, J.M., Feczko, E., Tisdall, M.D., van der Kouwe, A., Klein, R.L., Mirro, A.E., Hampton, J.M., Adeyemo, B., Laumann, T.O., Gratton, C., Greene, D.J., Schlaggar, B.L., Hagler, D.J., Watts Jr, R., Garavan, H., Barch, D.M., Nigg, J.T., Petersen, S.E., Dale, A.M., Feldstein-Ewing, S.W., Nagel, B.J., Dosenbach, N.U.F., 2020. Correction of respiratory artifacts in MRI head motion estimates. *Neuroimage* 208, 116400.
- Finn, E.S., Shen, X., Scheinost, D., Rosenberg, M.D., Huang, J., Chun, M.M., Papademetris, X., Constable, R.T., 2015. Functional connectome fingerprinting: identifying individuals using patterns of brain connectivity. *Nat. Neurosci.* 18 (11), 1664–1671.
- Fischl, B., Liu, A., Dale, A.M., 2001. Automated manifold surgery: constructing geometrically accurate and topologically correct models of the human cerebral cortex. *IEEE Trans. Med. Imaging* 20 (1), 70–80.
- Fischl, B., Salat, D.H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., Montillo, A., Makris, N., Rosen, B., Dale, A.M., 2002. Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain. *Neuron* 33 (3), 341–355.
- Fischl, B., Sereno, M.I., Dale, A.M., 1999a. II: inflation, flattening, and a surface-based coordinate system. *Neuroimage* 9, 195–207.
- Fischl, B., Sereno, M.I., Tootell, R.B., Dale, A.M., 1999b. High-resolution intersubject averaging and a coordinate system for the cortical surface. *Hum. Brain Mapp.* 8 (4), 272–284.
- Gabrieli, J.D.E., Ghosh, S.S., Whitfield-Gabrieli, S., 2015. Prediction as a humanitarian and pragmatic contribution from human cognitive neuroscience. *Neuron* 85 (1), 11–26.
- Gell, M., Eickhoff, S.B., Omidvarnia, A., Küppers, V., Patil, K.R., Satterthwaite, T.D., Müller, V.I., and Langner, R., 2023. The Burden of Reliability: how Measurement Noise Limits Brain-Behaviour Predictions. *bioRxiv*.
- Gordon, E.M., Laumann, T.O., Adeyemo, B., Huckins, J.F., Kelley, W.M., Petersen, S.E., 2016. Generation and evaluation of a cortical area parcellation from resting-state correlations. *Cereb. Cortex* 26 (1), 288–303.
- Gratton, C., Dvoretzky, A., Coalson, R.S., Adeyemo, B., Laumann, T.O., Wig, G.S., Kong, T.S., Gratton, G., Fabiani, M., Barch, D.M., Tranel, D., Miranda-Dominguez, O., Fair, D.A., Dosenbach, N.U.F., Snyder, A.Z., Perlmuter, J.S., Petersen, S.E., Campbell, M.C., 2020. Removal of high frequency contamination from motion estimates in single-band fMRI saves data without biasing functional connectivity. *Neuroimage* 217, 116866.
- Greene, A.S., Gao, S., Scheinost, D., Constable, R.T., 2018. Task-induced brain state manipulation improves prediction of individual traits. *Nat. Commun.* 9 (1), 2807.
- Greicius, M.D., Srivastava, G., Reiss, A.L., Menon, V., 2004. Default-mode network activity distinguishes Alzheimer's disease from healthy aging: evidence from functional MRI. *Proc. Natl. Acad. Sci. USA* 101 (13), 4637–4642.
- Greve, D.N., Fischl, B., 2009. Accurate and robust brain image alignment using boundary-based registration. *Neuroimage* 48 (1), 63–72.
- Hagler Jr, D.J., Hatton, S., Cornejo, M.D., Makowski, C., Fair, D.A., Dick, A.S., Sutherland, M.T., Casey, B.J., Barch, D.M., Harms, M.P., Watts, R., Bjork, J.M., Garavan, H.P., Hilmer, L., Pung, C.J., Scat, C.S., Kuperman, J., Bartsch, H., Xue, F., Heitzeg, M.M., Laird, A.R., Trinh, T.T., Gonzalez, R., Tapert, S.F., Riedel, M.C., Squeglia, L.M., Hyde, L.W., Rosenberg, M.D., Earl, E.A., Howlett, K.D., Baker, F.C., Soules, M., Diaz, J., de Leon, O.R., Thompson, W.K., Neale, M.C., Herting, M., Sowell, E.R., Alvarez, R.P., Hawes, S.W., Sanchez, M., Bodurka, J., Breslin, F.J., Morris, A.S., Paulus, M.P., Simmons, W.K., Polimeni, J.R., van der Kouwe, A., Nencka, A.S., Gray, K.M., Pierpaoli, C., Matochik, J.A., Noronha, A., Aklin, W.M., Conway, K., Glantz, M., Hoffman, E., Little, R., Lopez, M., Pariyadath, V., Weiss, S.R., Wolff-Hughes, D.L., DelCarmen-Wiggins, R., Feldstein Ewing, S.W., Miranda-Dominguez, O., Nagel, B.J., Perrone, A.J., Sturgeon, D.T., Goldstone, A., Pfefferbaum, A., Pohl, K.M., Prouty, D., Urban, K., Bookheimer, S.Y., Dapretto, M., Galvan, A., Bagot, K., Giedd, J., Infante, M.A., Jacobus, J., Patrick, K., Shilling, P.D., Desikan, R., Li, Y., Sugrue, L., Banich, M.T., Friedman, N., Hewitt, J.K., Hopfer, C., Sakai, J., Tanabe, J., Cottler, L.B., Nixon, S.J., Chang, L., Cloak, C., Ernst, T., Reeves, G., Kennedy, D.N., Heeringa, S., Peltier, S., Schulenberg, J., Sripada, C., Zucker, R.A., Iacono, W.G., Luciana, M., Calabro, F.J., Clark, D.B., Lewis, D.A., Luna, B., Schirda, C., Brima, T., Foxe, J.J., Freedman, E.G., Muzek, D.W., Mason, M.J., Huber, R., McGlade, E., Prescott, A., Renshaw, P.F., Yurgelun-Todd, D.A., Allgaier, N.A., Dumas, J.A., Ivanova, M., Potter, A., Florsheim, P., Larson, C., Lisdahl, K., Charness, M.E., Fuemmeler, B., Hettema, J.M., Maes, H.H., Steinberg, J., Anokhin, A.P., Glaser, P., Heath, A.C., Madden, P.A., Baskin-Sommers, A., Constable, R.T., Grant, S.J., Dowling, G.J., Brown, S.A., Jernigan, T.L., Dale, A.M., 2019. Image processing and analysis methods for the adolescent brain cognitive development study. *Neuroimage* 202, 116091.
- Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.D., Blankertz, B., Bießmann, F., 2014. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage* 87, 96–110.
- He, T., Kong, R., Holmes, A.J., Nguyen, M., Sabuncu, M.R., Eickhoff, S.B., Bzdok, D., Feng, J., Yeo, B.T.T., 2020. Deep neural networks and kernel regression achieve comparable accuracies for functional connectivity prediction of behavior and demographics. *Neuroimage* 206, 116276.
- Hedderich, D.M., Eickhoff, S.B., 2020. Machine learning for psychiatry: getting doctors at the black box? *Mol. Psychiatry* 26 (1), 23–25.
- Hsu, W.T., Rosenberg, M.D., Scheinost, D., Constable, R.T., Chun, M.M., 2018. Resting-state functional connectivity predicts neuroticism and extraversion in novel individuals. *Soc. Cogn. Affect. Neurosci.* 13 (2), 224–232.
- Jenkinson, M., Bannister, P., Brady, M., Smith, S., 2002. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *Neuroimage* 17 (2), 825–841.
- Jiang, R., Calhoun, V.D., Fan, L., Zuo, N., Jung, R., Qi, S., Lin, D., Li, J., Zhuo, C., Song, M., Fu, Z., Jiang, T., Sui, J., 2020. Gender differences in connectome-based predictions of individualized intelligence quotient and sub-domain scores. *Cereb. Cortex* 30 (3), 888–900.
- Kennedy, D.P., Redcay, E., Courchesne, E., 2006. Failing to deactivate: resting functional abnormalities in autism. *Proc. Natl. Acad. Sci. USA* 103 (21), 8275–8280.
- Kharabian Masouleh, S., Eickhoff, S.B., Hoffstaedter, F., Genon, S., Alzheimer's Disease Neuroimaging Initiative, 2019. Empirical examination of the replicability of associations between brain structure and psychological variables. *eLife* 8, e43464.
- Kirk, J., Miller, M.J., 1986. Reliability and Validity in Qualitative Research. SAGE Publications, Inc.
- Kong, R., Li, J., Orban, C., Sabuncu, M.R., Liu, H., Schaefer, A., Sun, N., Zuo, X.N., Holmes, A.J., Eickhoff, S.B., Yeo, B.T.T., 2019. Spatial topography of individual-specific cortical networks predicts human cognition, personality, and emotion. *Cereb. Cortex* 29 (6), 2533–2551.
- Kong, R., Yang, Q., Gordon, E., Xue, A., Yan, X., Orban, C., Zuo, X.N., Spreng, N., Ge, T., Holmes, A., Eickhoff, S., Yeo, B.T.T., 2021. Individual-specific areal-level parcellations improve functional connectivity prediction of behavior. *Cereb. Cortex* 31 (10), 4477–4500.
- Li, J., Kong, R., Liégeois, R., Orban, C., Tan, Y., Sun, N., Holmes, A.J., Sabuncu, M.R., Ge, T., Yeo, B.T.T., 2019. Global signal regression strengthens association between resting-state functional connectivity and behavior. *Neuroimage* 196, 126–141.
- Liu, D., Lin, X., Ghosh, D., 2007. Semiparametric regression of multidimensional genetic pathway data: least-squares kernel machines and linear mixed models. *Biometrics* 63 (4), 1079–1088.
- Marek, S., Tervo-Clemmens, B., Calabro, F.J., Montez, D.F., Kay, B.P., Hatom, A.S., Donohue, M.R., Foran, W., Miller, R.L., Hendrickson, T.J., Malone, S.M., Kandal, S., Feczko, E., Miranda-Dominguez, O., Graham, A.M., Earl, E.A., Perrone, A.J., Cordova, M., Doyle, O., Moore, L.A., Conan, G.M., Uriarte, J., Snider, K., Lynch, B.J., Wilgenbusch, J.C., Pengo, T., Tam, A., Chen, J., Newbold, D.J., Zheng, A., Seider, N.A., Van, A.N., Metoki, A., Chauvin, R.J., Laumann, T.O., Greene, D.J., Petersen, S.E., Garavan, H., Thompson, W.K., Nichols, T.E., Yeo, B.T.T., Barch, D.M., Luna, B., Fair, D.A., Dosenbach, N.U.F., 2022. Reproducible brain-wide association studies require thousands of individuals. *Nature* 603 (7902), 654–660.
- Nikolaïdis, A., Chen, A.A., He, X., Shinohara, R., Vogelstein, J., Milham, M., and Shou, H., 2022. Suboptimal phenotypic reliability impedes reproducible human neuroscience. *bioRxiv*.
- Noble, S., Scheinost, D., Constable, R.T., 2019. A decade of test-retest reliability of functional connectivity: a systematic review and meta-analysis. *Neuroimage* 203, 116157.
- Noble, S., Spann, M.N., Tokoglu, F., Shen, X., Constable, R.T., Scheinost, D., 2017. Influences on the test-retest reliability of functional connectivity MRI and its relationship with behavioral utility. *Cereb. Cortex* 27 (11), 5415–5429.
- Nostro, A.D., Müller, V.I., Varikuti, D.P., Pläschke, R.N., Hoffstaedter, F., Langner, R., Patil, K.R., Eickhoff, S.B., 2018. Predicting personality from network-based resting-state functional connectivity. *Brain Struct. Funct.* 223 (6), 2699–2719.
- Ooi, L.Q.R., Chen, J., Shaoshi, Z., Kong, R., Tam, A., Li, J., Dhamala, E., Zhou, L.H., Holmes, A.J., and Thomas Yeo, B.T., 2022. Comparison of individualized behavioral predictions across anatomical, diffusion and functional connectivity MRI. *bioRxiv*.
- Pervaz, U., Vidaurre, D., Woolrich, M.W., Smith, S.M., 2020. Optimising network modelling methods for fMRI. *Neuroimage* 211, 116604.
- Power, J.D., Barnes, K.A., Snyder, A.Z., Schlaggar, B.L., Petersen, S.E., 2012. Spurious but systematic correlations in functional connectivity MRI networks arise from subject motion. *Neuroimage* 59 (3), 2142–2154.
- Power, J.D., Lynch, C.J., Silver, B.M., Dubin, M.J., Martin, A., Jones, R.M., 2019. Distinctions among real and apparent respiratory motions in human fMRI data. *Neuroimage* 201, 116041.
- Power, J.D., Mitra, A., Laumann, T.O., Snyder, A.Z., Schlaggar, B.L., Petersen, S.E., 2014. Methods to detect, characterize, and remove motion artifact in resting state fMRI. *Neuroimage* 84, 320–341.
- Price, W.N., 2018. Medical malpractice and black-box medicine. In: *Big Data, Health Law, and Bioethics*. Cambridge University Press, pp. 295–306.
- Schaefer, A., Kong, R., Gordon, E.M., Laumann, T.O., Zuo, X.N., Holmes, A.J., Eickhoff, S.B., Yeo, B.T.T., 2018. Local-global parcellation of the human cerebral cortex from intrinsic functional connectivity MRI. *Cereb. Cortex* 28 (9), 3095–3114.
- Schulz, M.A., Yeo, B.T.T., Vogelstein, J.T., Mourao-Miranada, J., Kather, J.N., Kording, K., Richards, B., Bzdok, D., 2020. Different scaling of linear models and deep learning in UKBiobank brain images versus machine-learning datasets. *Nat. Commun.* 11 (1), 4238.
- Ségonne, F., Dale, A.M., Busa, E., Glessner, M., Salat, D., Hahn, H.K., Fischl, B., 2004. A hybrid approach to the skull stripping problem in MRI. *Neuroimage* 22 (3), 1060–1075.
- Ségonne, F., Pacheco, J., Fischl, B., 2007. Geometrically accurate topology-correction of cortical surfaces using nonseparating loops. *IEEE Trans. Med. Imaging* 26 (4), 518–529.
- Shen, X., Finn, E.S., Scheinost, D., Rosenberg, M.D., Chun, M.M., Papademetris, X., Constable, R.T., 2017. Using connectome-based predictive modeling to predict individual behavior from brain connectivity. *Nat. Protoc.* 12 (3), 506–518.
- Sripada, C., Rutherford, S., Angstadt, M., Thompson, W.K., Luciana, M., Weigard, A.,



- Hyde, L.H., Heitzeg, M., 2020. Prediction of neurocognition in youth from resting state fMRI. *Mol. Psychiatry* 25 (12), 3413–3421.
- Strobl, C., Boulesteix, A.L., Kneib, T., Augustin, T., Zeileis, A., 2008. Conditional variable importance for random forests. *BMC Bioinformatics* 9, 307.
- Tang, S., Sun, N., Floris, D.L., Zhang, X., Di Martino, A., Yeo, B.T.T., 2020. Reconciling dimensional and categorical models of autism heterogeneity: a brain connectomics and behavioral study. *Biol. Psychiatry* 87 (12), 1071–1082.
- Taxali, A., Angstadt, M., Rutherford, S., Sripada, C., 2021. Boost in test-retest reliability in resting state fMRI with predictive modeling. *Cereb. Cortex* 31 (6), 2822–2833.
- Tian, Y. and Zalesky, A., 2021. Machine learning prediction of cognition from functional connectivity: are feature weights reliable? *bioRxiv*.
- Vasey, B., Nagendran, M., Campbell, B., Clifton, D.A., Collins, G.S., Denaxas, S., Denniston, A.K., Faes, L., Geerts, B., Ibrahim, M., Liu, X., Mateen, B.A., Mathur, P., McCradden, M.D., Morgan, L., Ordish, J., Rogers, C., Saria, S., Ting, D.S.W., Watkinson, P., Weber, W., Wheatstone, P., McCulloch, P., 2022a. Reporting guideline for the early stage clinical evaluation of decision support systems driven by artificial intelligence: DECIDE-AI. *BMJ* 377, e070904.
- Vasey, B., Nagendran, M., Campbell, B., Clifton, D.A., Collins, G.S., Denaxas, S., Denniston, A.K., Faes, L., Geerts, B., Ibrahim, M., Liu, X., Mateen, B.A., Mathur, P., McCradden, M.D., Morgan, L., Ordish, J., Rogers, C., Saria, S., Ting, D.S.W., Watkinson, P., Weber, W., Wheatstone, P., McCulloch, P., 2022b. Reporting guideline for the early-stage clinical evaluation of decision support systems driven by artificial intelligence: DECIDE-AI. *Nat. Med.* 28 (5), 924–933.
- Wolters, T., Beckmann, C.F., Hoogman, M., Buitelaar, J.K., Franke, B., Marquand, A.F., 2020. Individual differences v. the average patient: mapping the heterogeneity in ADHD using normative models. *Psychol. Med.* 50 (2), 314–323.
- Xia, C.H., Ma, Z., Ciric, R., Gu, S., Betzel, R.F., Kaczkurkin, A.N., Calkins, M.E., Cook, P.A., García de la Garza, A., Vandekar, S.N., Cui, Z., Moore, T.M., Roalf, D.R., Ruparel, K., Wolf, D.H., Davatzikos, C., Gur, R.C., Gur, R.E., Shinohara, R.T., Bassett, D.S., Satterthwaite, T.D., 2018. Linked dimensions of psychopathology and connectivity in functional brain networks. *Nat. Commun.* 9 (1), 3003.
- Xiao, Y., Lin, Y., Ma, J., Qian, J., Ke, Z., Li, L., Yi, Y., Zhang, J., Cam-CAN, Dai, Z., 2021. Predicting visual working memory with multimodal magnetic resonance imaging. *Hum. Brain Mapp.* 42 (5), 1446–1462.
- Zabihi, M., Oldehinkel, M., Wolters, T., Frouin, V., Goyard, D., Loth, E., Charman, T., Tillmann, J., Banaschewski, T., Dumas, G., Holt, R., Baron-Cohen, S., Durston, S., Bölte, S., Murphy, D., Ecker, C., Buitelaar, J.K., Beckmann, C.F., Marquand, A.F., 2019. Dissecting the heterogeneous cortical anatomy of autism spectrum disorder using normative models. *Biol. Psychiatry Cogn. Neurosci. Neuroimaging* 4 (6), 567–578.
- Zhang, X., Mormino, E.C., Sun, N., Sperling, R.A., Sabuncu, M.R., Yeo, B.T.T. Alzheimer's Disease Neuroimaging Initiative, 2016. Bayesian model reveals latent atrophy factors with dissociable cognitive trajectories in Alzheimer's disease. *Proc. Natl. Acad. Sci. USA* 113 (42), E6535–E6544.