





## RESEARCH ARTICLE

# Longitudinally stable, brain-based predictive models mediate the relationships between childhood cognition and socio-demographic, psychological and genetic factors

Narun Pat<sup>1</sup>  | Yue Wang<sup>1</sup> | Richard Anney<sup>2</sup>  | Lucy Riglin<sup>2</sup>  |  
Anita Thapar<sup>2</sup>  | Argyris Stringaris<sup>3,4</sup> 

<sup>1</sup>Department of Psychology, University of Otago, Dunedin, New Zealand

<sup>2</sup>MRC Centre for Neuropsychiatric Genetics and Genomics, Division of Psychological Medicine and Clinical Neurosciences, School of Medicine and Wolfson Centre for Young People's Mental Health, Cardiff University, Cardiff, UK

<sup>3</sup>Division of Psychiatry, Department of Clinical, Educational and Health Psychology, University College London, London, UK

<sup>4</sup>Department of Psychiatry, National and Kapodistrian University of Athens, Athens, Greece

## Correspondence

Argyris Stringaris, Division of Psychiatry, Department of Clinical, Educational and Health Psychology, University College London, 1-19 Torrington Pl, London WC1E 7HB, UK.  
Email: [a.stringaris@ucl.ac.uk](mailto:a.stringaris@ucl.ac.uk)

Narun Pat, also known as Narun Pornpattananankul, Department of Psychology, University of Otago, William James Building, 275 Leith Walk, Dunedin 9016, New Zealand.  
Email: [narun.pat@otago.ac.nz](mailto:narun.pat@otago.ac.nz)

## Funding information

Health Research Council of New Zealand, Grant/Award Number: 21/618; University of Otago

## Abstract

Cognitive abilities are one of the major transdiagnostic domains in the National Institute of Mental Health's Research Domain Criteria (RDoC). Following RDoC's integrative approach, we aimed to develop brain-based predictive models for cognitive abilities that (a) are developmentally stable over years during adolescence and (b) account for the relationships between cognitive abilities and socio-demographic, psychological and genetic factors. For this, we leveraged the unique power of the large-scale, longitudinal data from the Adolescent Brain Cognitive Development (ABCD) study ( $n \sim 11$  k) and combined MRI data across modalities (task-fMRI from three tasks: resting-state fMRI, structural MRI and DTI) using machine-learning. Our brain-based, predictive models for cognitive abilities were stable across 2 years during young adolescence and generalisable to different sites, partially predicting childhood cognition at around 20% of the variance. Moreover, our use of 'opportunistic stacking' allowed the model to handle missing values, reducing the exclusion from around 80% to around 5% of the data. We found fronto-parietal networks during a working-memory task to drive childhood-cognition prediction. The brain-based, predictive models significantly, albeit partially, accounted for variance in childhood cognition due to (1) key socio-demographic and psychological factors (proportion mediated = 18.65% [17.29%–20.12%]) and (2) genetic variation, as reflected by the polygenic score of cognition (proportion mediated = 15.6% [11%–20.7%]). Thus, our brain-based predictive models for cognitive abilities facilitate the development of a robust, transdiagnostic research tool for cognition at the neural level in keeping with the RDoC's integrative framework.

## KEYWORDS

adolescent brain cognitive development, general cognition, longitudinal large-scale data, machine learning, multimodal MRI, polygenic score, research domain criteria

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2022 The Authors. *Human Brain Mapping* published by Wiley Periodicals LLC.

## 1 | INTRODUCTION

According to the Research Domain Criteria (RDoC), cognitive abilities are considered one of the major transdiagnostic domains, cutting across mental disorders (Morris & Cuthbert, 2012). In children and adults, cognitive abilities are related to various mental disorders, including but not limited to depression (Shilyansky et al., 2016), attention-deficit/hyperactivity disorder (ADHD) (Thaler et al., 2013) and psychotic disorders (Sheffield et al., 2018). Cognitive abilities that span across cognitive tasks, such as language, mental flexibility and memory, reflect a trait, known as general cognition or the *g*-factor (Flynn, 2009). Yet, we still do not have predictive models that can robustly capture the relationship between the *g*-factor and the brain. Having a brain-based predictive model for the *g*-factor is a key for us to adapt the RDoC's integrative approach—to understand cognitive abilities across units of analyses, from behaviours to brain and genes that reflect the influences of socio-demographical and psychological factors across the lifespan (Insel et al., 2010; National Institute of Mental Health [NIMH], n.d.-a).

Developing the brain-based predictive models for children's *g*-factor to be used in the RDoC framework faces several challenges. The first challenge is longitudinal stability, which is one of the requirements in the RDoC framework (Insel et al., 2010; NIMH, n.d.-a). Predictive models should not only be generalisable to out-of-sample data (i.e., be predictive of children's *g*-factor that were not part of the original sample) but also be developmentally stable (Sui et al., 2020) in order to capture the *g*-factor across the lifespan (Tucker-Drob, 2009). Here, we started to tackle this challenge by using—for the first time to the best of our knowledge—longitudinal, large-scale data in children, from the Adolescent Brain Cognitive Development (ABCD) study (Yang & Jernigan, n.d.), to demonstrate the longitudinal stability of the brain-based predictive models across 2 years during adolescence.

The second challenge is multimodal integration. So far, brain-based predictive models have been mainly built from a single MRI modality without integrating different sources of information from different MRI modalities. For instance, the *g*-factor is associated with activity during certain cognitive tasks, such as working memory (Gray et al., 2003; Waiter et al., 2009) (task-based functional MRI; task-fMRI), the intrinsic functional connectivity between different areas (Dubois et al., 2018; Pamplona et al., 2015; Sripada, Rutherford, et al., 2020) (resting-state fMRI [rs-fMRI]) and the anatomy of grey matter (Narr et al., 2007) (structural MRI [sMRI]) and white matter (Genç et al., 2018; Góngora et al., 2020) (diffusion tensor imaging [DTI]). However, recent findings, mainly in adults, have started to show the benefits of integrating data across modalities, rather than relying solely on a single modality (Jiang et al., 2020; Rasero et al., 2021; Sui et al., 2020). Here, we adapted a machine-learning framework, called stacking (Wolpert, 1992), to integrate information across MRI modalities into a 'stacked' model. Briefly, we separately built models to predict the *g*-factor based on each brain modality, resulting in one predicted value from each modality for each participant. We then built a 'stacked' model to predict the *g*-factor based on these predicted values. We tested if the stacked model indeed enhanced predictive performance over single modalities in predicting children's *g*-factor.

The third challenge is missing data. Children's neuroimaging data are notoriously affected by movement artefacts (Fassbender et al., 2017). For example, the ABCD study recommended a set of quality control variables for detecting noisy data from each modality (Hagler et al., 2019; Yang & Jernigan, n.d.), resulting in a listwise exclusion of 17% to over 50% of data depending on a modality. If we were to exclude children who have noisy data from any single modality, we would have to exclude almost 80% of the data, strictly limiting the generalisability of our model to children with highly clean data (who are unlikely to be representative of the rest of the sample). We overcame this problem by using a recently developed framework, built on top of the stacking framework, called 'opportunistic stacking' (Engemann et al., 2020). Briefly, we first duplicated predicted values from each modality-specific model, and imputed the missing value in each duplicate either with an arbitrarily high or low value. We then used Random Forest (Breiman, 2001) to create a final prediction from the imputed, predicted values. Accordingly, opportunistic stacking allows us to keep the data as long as there is at least one modality available, leaving more data in the model-building process and reducing the risk of missing-data bias.

Beyond demonstrating a robust out-of-sample relationship between the brain and the *g*-factor, the brain-based predictive models have to demonstrate the construct validity, especially for them to be used according to the RDoC framework (Insel et al., 2010). For instance, RDoC stipulates that cognitive abilities are affected by socio-demographic and psychological factors (Morris & Cuthbert, 2012; NIMH, n.d.-b). This is in line with recent studies showing that cognitive abilities are related to factors such as socio-economic status (Farah et al., 2006), mental health (Biederman et al., 2004; Goodall et al., 2018) and extracurricular activities (Kirlic et al., 2021). Accordingly, for the brain-based predictive models to demonstrate RDoC's construct validity, the brain-based predictive models should be able to explain the associations between the *g*-factor and these socio-demographic and psychological factors.

Likewise, RDoC stipulates that cognitive abilities should not be studied as a unitary construct, but should rather be studied through different units of analysis, from behaviours to the brain and genes (Insel et al., 2010; Morris & Cuthbert, 2012; NIMH, n.d.-a, n.d.-c). Thus, the brain-based predictive models for cognitive abilities should be related to the 'gene-based' predictive models for cognitive abilities, given that they both reflect different units of analysis of the same RDoC's domain. A polygenic score (PGS), a composite measure of common gene variants, can be considered a predicted value from the gene-based predictive models (Bogdan et al., 2018). For cognitive abilities (Plomin & Deary, 2015), a PGS is based on the associations between several single nucleotide polymorphisms (SNPs) and cognitive abilities in a separate Genome-Wide Association Study (GWAS) (Davies et al., 2011), such as in a recent GWAS among 257,841 adults (Lee et al., 2018). Accordingly, for the brain-based predictive models to demonstrate RDoC's construct validity, the brain-based predictive models should also be able to explain the associations between the *g*-factor and the PGS of cognitive abilities (Lee et al., 2018).

To develop brain-based predictive models for the *g*-factor, we (i) used behavioural performance from cognitive tasks to derive the *g*-

factor and (ii) built brain-based predictive models to predict this behaviourally derived *g*-factor from multimodal MRI data. We used the ABCD Release 3.0 (Yang & Jernigan, *n.d.*), including baseline data (age 9–10 years old) from over 11,000 children and follow-up data (age 11–12 years old) from roughly half of the participants. We first derived children's *g*-factor from their behavioural performance on six cognitive tasks using confirmatory factor analysis (CFA). We then built brain-based predictive models by treating multimodal MRI data as the features and the children's *g*-factor derived from behavioural performance as the target. More specifically, in our models, we implemented opportunistic stacking (Engemann et al., 2020) to integrate MRI data across modalities and to deal with missing values from each modality. There were six modalities in total: three task-based fMRI (working-memory 'N-Back', reward 'Monetary Incentive Delay [MID]' and inhibitory control 'Stop Signal'), rs-fMRI, sMRI and DTI. To determine the robustness and longitudinal stability of the brain-based predictive models, we tested how well the models predicted the *g*-factor of unseen children at the same ages and at 2 years older as well as at different data-collection sites. Next, to demonstrate whether multimodal integration led to better predictive performance, we applied bootstrapping to compare the stacked model with the best-performing modality-specific model. To explain the feature importance of the final models (i.e., determining brain features that contributed highly to the prediction of the *g*-factor), we applied several 'explainers', including eNetXplorer (Candia & Tsang, 2019a), conditional permutation importance (CPI) (Strobl et al., 2008) and SHapley Additive exPlanations (SHAP) (Lundberg & Lee, 2017).

We then conducted mediation analyses to ensure that the brain-based predictive models for the *g*-factor demonstrated RDoC's construct validity. In these analyses, we tested the extent to which our brain-based predictive models could account for the relationships between the behaviourally derived *g*-factor and key socio-demographic, psychological and genetic factors. For this purpose, in addition to the brain-based predictive models, we also computed two additional predictive models that predicted the behaviourally derived *g*-factor, either from (a) 70 socio-demographic and psychological variables (Kircic et al., 2021) or (b) genes via a PGS of cognitive abilities (Lee et al., 2018). The 70 socio-demographic and psychological variables covered children's and/or their parents' socio-demographics, mental health, personality, sleep, physical activity, screen use, drug use, developmental adversity and social interaction. This resulted in three predicted values of the *g*-factor, based on features of the predictive models: 'brain-based *g*-factor', 'socio-demography-and-psychology-based *g*-factor' and 'gene-based *g*-factor'. We then computed these predicted values on unseen children at each hold-out data collection site and applied the mediation analyses. Here, we treated (i) the socio-demography-and-psychology-based and gene-based *g*-factors as the independent variables, (ii) the brain-based *g*-factor as the mediator and (iii) the behaviourally derived *g*-factor as the dependent variable. Through these mediation analyses, we quantified the extent to which the brain-based predictive models for cognitive abilities developed in this study mediated the relationships between the behaviourally derived *g*-factor and socio-demographic, psychological and genetic factors.

## 2 | MATERIALS AND METHODS

We employed the ABCD Study Curated Annual Release 3.0 (Yang & Jernigan, *n.d.*), which included 3 T MRI data and cognitive tests from 11,758 children (female = 5631) at the baseline (9–10 years old) and 5693 children (female = 2617) at the 2-year follow-up (11–12 years old). The study recruited the children from 21 sites across the United States (Garavan et al., 2018). We further excluded 54 children based on Snellen Vision Screener (Luciana et al., 2018; Snellen, 1862). These children either could not read any line, could only read the first (biggest) line, or could read up to the fourth line but indicated difficulty in reading stimuli on the iPad used for administering cognitive tasks (see below). The ethical considerations of the ABCD study, such as informed consent, confidentiality and communication with participants about assessment results, have been detailed elsewhere (Clark et al., 2018). Institutional Review Boards where the data were collected approved the study's protocols.

### 2.1 | The *g*-factor

We derived the *g*-factor using children's behavioural performance from six cognitive tasks. These six tasks, collected on an iPad during a 70-min in-session visit outside of MRI (Luciana et al., 2018; Thompson et al., 2019), were available in both baseline and follow-up datasets. First, the Picture Vocabulary measured vocabulary comprehension and language (Gershon et al., 2014). Second, the Oral Reading Recognition measured reading and language decoding (Bleck et al., 2013). Third, the Flanker measured conflict monitoring and inhibitory control (Eriksen & Eriksen, 1974). Fourth, the Pattern Comparison Processing measured the speed of processing (Carlozzi et al., 2013). Fifth, the Picture Sequence Memory measured episodic memory (Bauer et al., 2013). Sixth, the Rey-Auditory Verbal Learning measured memory recall after distraction and a short delay (Daniel & Wahlstrom, 2014).

Similar to the previous work (Ang et al., 2020; Pat et al., 2021; Thompson et al., 2019), we applied the second-order model of the *g*-factor using CFA to encapsulate the *g*-factor as the higher-order latent variable underlying performance across cognitive tasks. More specifically, our input data were standardised performance from each cognitive task. In our second-order model, we had the *g*-factor as the second-order latent variable. We also had three first-order latent variables in the model: language (underlying the Picture Vocabulary and Oral Reading Recognition), mental flexibility (underlying the Flanker and Pattern Comparison Processing), and memory recall (underlying the Picture Sequence Memory and Rey-Auditory Verbal Learning).

We fixed latent factor variances to one and applied Maximum Likelihood with Robust standard errors (MLR) using Huber-White standard errors and scaled test statistics. To demonstrate model fit, we used scaled and/or robust indices, including comparative fit index (CFI), Tucker-Lewis index (TLI), root mean squared error of approximation (RMSEA) and standardized root mean square residual (SRMR) as well as used internal consistency, OmegaL2 (Jorgensen et al., 2018), of the *g*-factor. To implement the CFA, we used lavaan

(Rosseel, 2012) (version = .6-6) and semTools (Jorgensen et al., 2018) along with semPlot (Epskamp, 2015) for visualisation. Note to ensure the robustness of the chosen *g*-factor model, we also examined the similarity in factor scores of the *g*-factor based on three different CFA models: the second-order model, the single-factor model, and the mixture between exploratory factor analysis (EFA) and CFA models (Appendix S1).

## 2.2 | Multimodal MRI

We used MRI data from six modalities: three task-based fMRI, rs-fMRI, sMRI and DTI. Note ‘modalities’ here referred to sets of features in our predictive models, as such we treated three task-based fMRI as separate modalities even though they were task-based fMRI. The ABCD study provided detailed procedures on data acquisition and MRI image processing elsewhere (Casey et al., 2018; Hagler et al., 2019; Yang & Jernigan, n.d.). We strictly followed their recommended exclusion criteria based on automated and manual QC review of each modality, listed under the *abcd\_imgincl01* table (Yang & Jernigan, n.d.). The ABCD created an exclusion flag for each modality (with a prefix ‘imgincl’) based on several criteria, involving image quality, MR neurological screening, behavioural performance, number of repetition time (TRs) among others. We removed participants with an exclusion flag at any MRI indices, separately for each modality. We also applied the three interquartile range ( $3 \times$  IQR) rule (i.e., datapoint with a value over 3 IQRs away from the nearest quartile) with listwise deletion to remove observations with outliers in any indices within each modality. Additionally, to adjust for between-site variability, we used an Empirical Bayes method, ComBat (Fortin et al., 2017; Nielson et al., 2018). We applied ComBat to all modalities except for task-based fMRI, given that between-site variability was found to be negligible for task-based contrasts (Nielson et al., 2018). See below for our approach to mitigate data leakage due to  $3 \times$  IQR and ComBat.

### 2.2.1 | Three task-based fMRI

We used task-based fMRI from three tasks. First, in the working-memory ‘N-Back’ task (Barch et al., 2013; Casey et al., 2018), children saw pictures of houses and emotional faces. Depending on the block, children reported if a picture matched either: (a) a picture that was shown 2 trials earlier (2-back), or (b) a picture that was shown at the beginning of the block (0-back). To focus on working-memory-related activity, we used the (2-back vs. 0-back) linear contrast (i.e., high vs. low working memory load). Second, in the MID task (Casey et al., 2018; Knutson et al., 2000), children needed to respond before the target disappeared. And doing so would provide them with a reward, if and only if the target followed the ‘reward cue’ (but not the ‘neural cue’). To focus on reward anticipation-related activity, we used the (Reward Cue vs. Neutral Cue) linear contrast. Third, in the Stop-Signal Task (SST) (Casey et al., 2018; Whelan et al., 2012), children needed to withhold or interrupt their motor response to a ‘Go’ stimulus when it was followed unpredictably by a Stop signal. To focus on inhibitory control-related activity, we used the (Any

Stop vs. Correct Go) linear contrast. Note that, for the SST, we used two additional exclusion criteria, *tfmri\_sst\_beh\_glitchflag*, and *tfmri\_sst\_beh\_violatorflag*, to address glitches in the task as recommended by the study (Bissett et al., 2020; Garavan et al., 2020). For all tasks, we used the average contrast values across two runs. More specifically, these contrasts were unthresholded, similar to previous work (Bolt et al., 2017). These values were embedded in the brain parcels based on FreeSurfer’s (Dale et al., 1999) Destrieux (Destrieux et al., 2010) and ASEG (Fischl et al., 2002) atlases (148 cortical surface and 19 subcortical volumetric regions, resulting in 167 features for each task-based fMRI task).

### 2.2.2 | Resting-state fMRI

During rs-fMRI collection, the children viewed a crosshair for 20 min. The ABCD’s preprocessing strategy has been published elsewhere (Hagler et al., 2019). Briefly, the study parcellated regions into 333 cortical-surface regions (Gordon et al., 2016) and correlated their time-series (Hagler et al., 2019). They then grouped these correlations based on 13 predefined large-scale networks (Gordon et al., 2016): auditory, cingulo-opercular, cingulo-parietal, default-mode, dorsal-attention, frontoparietal, none, retrosplenial-temporal, salience, sensorimotor-hand, sensorimotor-mouth, ventral-attention and visual networks. Note that ‘none’ refers to regions that do not belong to any networks. After applying the Fisher’s *r*-to-*z* transformation, the study computed mean correlations between pairs of regions within each large-scale network ( $n = 13$ ) and between large-scale networks ( $n = 78$ ) and provided these mean correlations in their Releases (Yang & Jernigan, n.d.). This resulted in 91 features for the rs-fMRI. Given that the correlations between (not within) large-scale networks were highly collinear with each other (e.g., the correlation between auditory and cingulo-opercular was collinear with that between auditory and default-mode), we further decorrelated them using partial correlation. We first applied the inverse Fisher’s *r*-to-*z* transformation, then partial correlation transformation, and then reapplied the Fisher *r*-to-*z* transformation.

### 2.2.3 | Structural MRI

The ABCD study processed sMRI, including cortical reconstruction and subcortical volumetric segmentation, using FreeSurfer (Dale et al., 1999). Here, we considered FreeSurfer-derived Destrieux (Destrieux et al., 2010) regional cortical thickness measures ( $n = 148$  cortical surface) and ASEG (Fischl et al., 2002) regional subcortical volume measures ( $n = 19$ ), resulting in 167 features for sMRI. We also adjusted regional cortical thickness and volumetric measures using mean cortical thickness and total intracranial volume, respectively.

### 2.2.4 | Diffusion tensor imaging

Here, we focused on fractional anisotropy (FA) (Alexander et al., 2007) of DTI. FA characterises the directionality of the

distribution of diffusion within white matter tracts, which can indicate the density of fibre packing (Alexander et al., 2007). The ABCD study segmented major white matter tracts using AtlasTrack (Hagler et al., 2009, 2019). Here, we considered FA of 23 major tracts, 10 of which were separately labelled for each hemisphere. These tracts included corpus callosum, forceps major, forceps minor, cingulate and parahippocampal portions of cingulum, fornix, inferior frontal occipital fasciculus, inferior longitudinal fasciculus, pyramidal/corticospinal tract, superior longitudinal fasciculus, temporal lobe portion of superior longitudinal fasciculus, anterior thalamic radiations and uncinate. This left 23 features for DTI.

### 2.3 | Predictive models of multimodal MRI: opportunistic stacking

To integrate multimodal MRI into one predictive model and to control for missing values across modalities, we applied opportunity stacking (Engemann et al., 2020) (Figure 1). We started with the first-layer training set. Here, we used standardised features from each modality to separately predict the  $g$ -factor via a penalised regression. The main advantage of a penalised regression is its ease of interpretation given that the prediction is made based on a weighted sum of features. Moreover, predictive performance of penalised regressions for capturing brain-and-behaviour relationships in MRI appeared good, often on-par with other more black-box algorithms (Dadi et al., 2019; Dubois et al., 2018; Engemann et al., 2020; Niu et al., 2020; Rasero et al., 2021). Following previous research (Dubois et al., 2018), we used Elastic Net (Zou & Hastie, 2005), a general form of penalised regression via the glmnet package (Friedman et al., 2010). Elastic Net requires two hyperparameters. First, the 'penalty' determines how strong the feature's slopes are regularised. Second, the 'mixture' determines the degree to which the regularisation is applied to the sum of squared coefficients (known as Ridge) versus to the sum of absolute values of the coefficients (known as LASSO). We tuned these two hyperparameters using a 10-fold cross-validation grid search and selected the model with the lowest mean absolute error (MAE). In the grid, we used 200 levels of the penalty from  $10^{-10}$  to 10, equally spaced on the logarithmic-10 scale and 11 levels of the mixture from 0 to 1 on the linear scale.

Once we obtained the final modality-specific models from the first-layer training set, we fit these models to data in the second-layer training set. This gave us six predicted values of the  $g$ -factor from six modalities, and these are the features to predict the  $g$ -factor in the second-layer training set. To handle missing observations when combining these modality-specific features, we applied the opportunistic stacking approach (Engemann et al., 2020) by creating duplicates of each modality-specific feature. After standardisation, we coded missing observations in one as an arbitrarily large value of 1000 and in the other as an arbitrarily small value of  $-1000$ , resulting in 12 features. That is, as long as a child had at least one modality available, we would be able to include this child in stacked modelling.

Previous research (Engemann et al., 2020) advocated for a more flexible algorithm that can capture non-linear and interactive

relationships at the second-layer training set. Here, we used the Random Forests algorithm (Breiman, 2001) from the ranger package (Wright & Ziegler, 2017) to predict the  $g$ -factor from the 12 features (Engemann et al., 2020; Josse et al., 2020). Random Forests use a multitude of decision trees on various sub-samples of the data and implement averaging to enhance prediction and to control over-fitting. We used 1000 trees and turned two hyperparameters. First 'mtry' is the number of features randomly sampled at each split. Second 'min\_n' is the minimum number of observations in a node needed for the node to be split further. We implemented a 10-fold cross-validation grid search and selected the model with the lowest root mean squared error (RMSE). In the grid, we used 12 levels of the mtry from 1 to 12, and 101 levels of the min\_n from 1 to 1000, both on the linear scale. This resulted in the 'stacked' model that incorporated data across modalities.

To prevent data leakage, we fit the CFA model to the observations in the first-layer training data and then computed factor scores of the  $g$ -factor on all training and test data. Note that to demonstrate the stability of the factor scores of the  $g$ -factor when applied to unseen data (i.e., not part of the modelling process), we also compared the factor scores of the  $g$ -factor estimated from the first-layer training data and the scores estimated from the whole baseline data (Appendix S2). Similarly, we also applied the  $3 \times$  IQR rule and Combat separately for first-layer training, second-layer training, baseline test and follow-up test data. For the machine learning workflow, we used 'tidymodels' ([www.tidymodels.org](http://www.tidymodels.org)).

### 2.4 | Testing the robustness of the predictive models of multimodal MRI

We examined the predictive ability of the models based on multimodal MRI between predicted versus observed  $g$ -factor, using Pearson's correlation ( $r$ ), coefficient of determination ( $R^2$ , calculated using the sum of square definition), MAE, and RMSE. To investigate the predictive ability of the modality-specific models, we used the models tuned from the first-layer training set. To investigate the predictive ability of the stacked model, we used the model tuned from both the first-layer and second-layer training sets.

#### 2.4.1 | Out-of-sample predictive ability of multimodal MRI: Baseline and follow-up samples

We first split the data into four parts (Figure 1): (1) first-layer training set ( $n = 3041$ ), (2) second-layer training set ( $n = 3042$ ), (3) baseline test set ( $n = 5622$ ) and (4) follow-up test set ( $n = 5656$ ). Especially noteworthy is that children who were in the baseline test set were also in the follow-up test set. In other words, none of the children in the first-layer and second-layer training sets was in either of the test sets. We used the baseline test set for out-of-sample, same-age predictive abilities, while we used the follow-up test sets for out-of-sample, longitudinal predictive abilities.



### Data Splitting

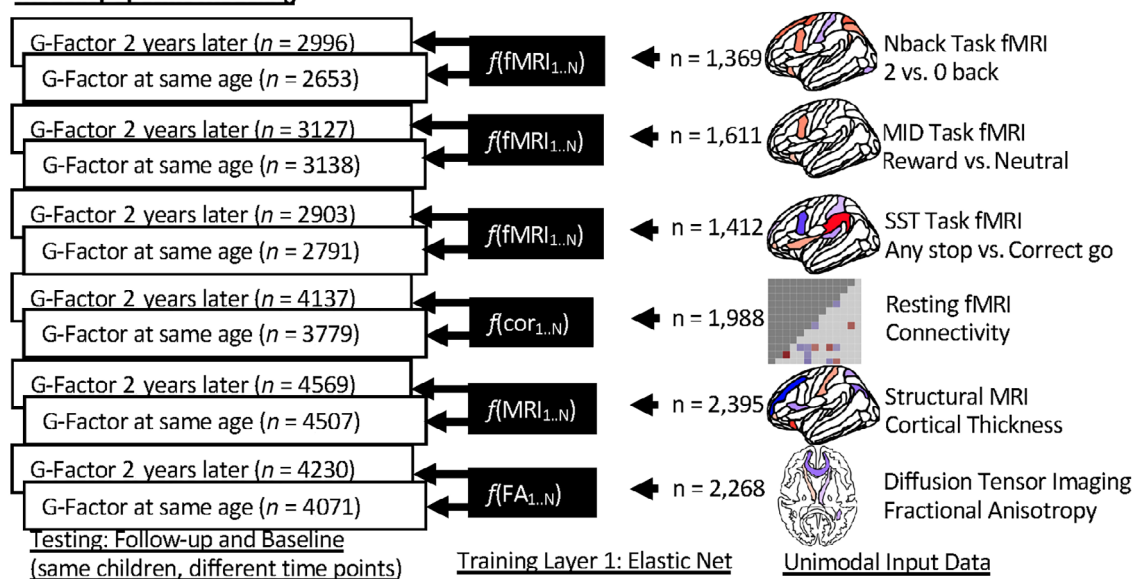
Baseline data (age 9–10 years old)

1st-layer training set - CFA for G-Factor - 10-fold CV for Elastic Net tuning	2nd-layer training set - 10-fold CV for Random Forests tuning	Baseline test set
--	---	-------------------

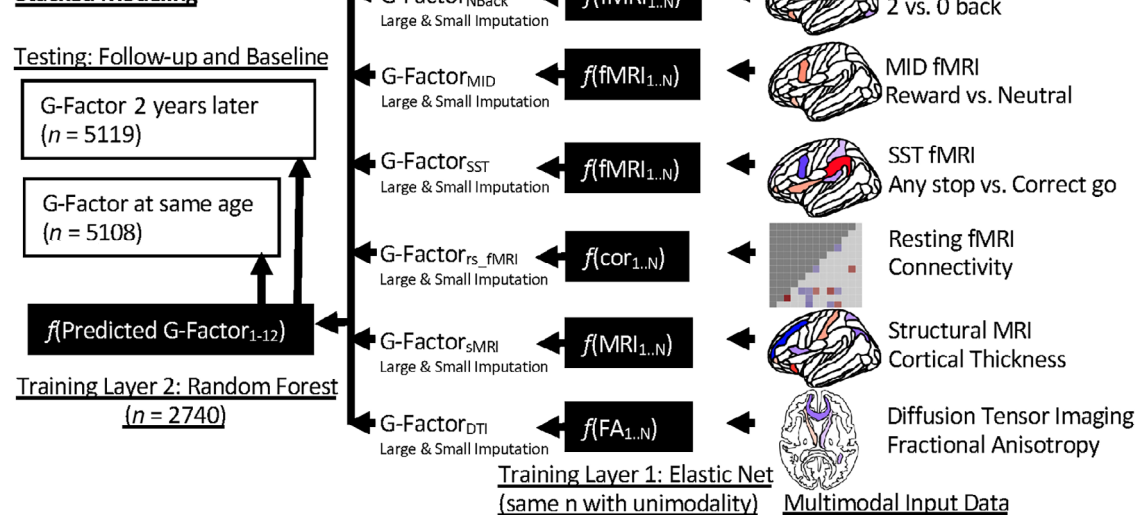
Follow-Up data (age 11–12 years old)

Data not yet released	Follow-up test set
-----------------------	--------------------

### Modality-specific modeling



### Stacked Modeling



**FIGURE 1** Longitudinal predictive modelling approach used for out-of-sample predictive ability of multimodal MRI. We split the data into four sets: First-layer training, second-layer training, baseline test, and follow-up test. We used the same participants in the baseline test and follow-up test sets. Modality-specific modelling only used the first-layer training set, while stacked modelling used both training sets to combine predicted values across modalities. At the first training layer, using elastic net, we separately predicted the  $g$ -factor based on each of the six modalities, resulting in six predicted values. At the second training layer, we applied opportunistic stacking by duplicating these six predicted values, and then imputed missing observations in one as an arbitrarily large value of 1000 and in the other as an arbitrarily small value of  $-1000$ , resulting in 12 predicted values. We then used Random Forest to predict the  $g$ -factor based on these 12 predicted values. The number of observations was different depending on the quality control of data from each modality. “Data not yet released” reflects the fact that ABCD release 3.0 (Yang & Jernigan, *n.d.*) only provided half of the follow-up data (age 11–12 years old), while providing the full baseline data (age 9–10 years old). CFA, confirmatory factor analysis; cor, correlation; CV, cross-validation; FA, fractional anisotropy

To examine the performance of opportunistic stacking as a function of missing values, we further split the test sets based on the presence of each modality. First, Stacked All required data with at least one modality present. This allowed us to examine the stacked model's performance when the missing values were all arbitrarily coded. Second, Stacked Complete required data with all modalities present. This represents the situation when the data were as clean as possible. Third, Stacked Best had the same missing values as the modality with the best prediction. This allowed us to make a fair comparison in performance between the stacked model and the model with the best modality, given their same noise level from missing value. Fourth, Stacked No Best did not have any data from the modality with the best prediction and had at least one modality present. This represents the highest level of noise possible.

#### 2.4.2 | Comparing out-of-sample predictive ability of multimodal MRI between the stacked model and the model with the best modality: baseline and follow-up samples

Here, we made a statistical comparison in the out-of-sample predictive ability between Stacked Best and the modality-specific model with the highest predictive performance, two of which had the same number of missing values in the test sets. We applied bootstrapping with 5000 iterations to examine the differences in performance indices (including,  $r$ ,  $R^2$ , MAE and RMSE) on both baseline and follow-up test sets. If stacking truly led to enhanced predictive performance, then we should see 95% CI of the bootstrapped differences to be different from 0.

#### 2.4.3 | Out-of-site predictive ability of multimodal MRI

To examine out-of-site predictive ability, we applied leave-one-site-out cross-validation to the baseline data. This enabled us to extract predicted values of the  $g$ -factor based on multimodal MRI data at each hold-out site, and in turn, to examine the generalisability of different models on different data collection sites. Different sites involved different MRI machines, experimenters as well as demographics across the United States (Garavan et al., 2018). Moreover, using leave-one-site-out cross-validation also prevented having the participants from the same family in the training and test sets. Here, we first removed data from one site that only recruited 34 participants and removed participants from six families who were separately scanned at different sites. We then held out data from one site as a test set and divided the rest to be first- and second-layer training sets. We cross-validated predictive ability across these hold-out sites. We applied the same modelling approach with the out-of-sample predictive models, except for two configurations to reduce the amount of ram used and computational time. Specifically, in our grid search, we used 100 levels of penalty (as opposed to 200) for Elastic Net and limited the maximal

min\_n to 500 (as opposed to 1000) for Random Forests. For the stacked model, we tested its predictive ability on children with at least one modality (i.e., stacked all). We examined the out-of-site prediction between predicted versus observed  $g$ -factor at each hold-out site.

### 2.5 | Feature importance of multimodal MRI models

To understand which features contribute to the prediction of the modality-specific (i.e., Elastic Net) models, we applied permutation from the eNetXplorer (Candia & Tsang, 2019b) package to the first-layer training set of the out-of-sample predictive ability splits (Figure 1). We first chose the best mixture from the previously run grid and fit two sets of several Elastic Net models. The first 'target' models used the true  $g$ -factor as the target, while the second 'null' models used the randomly permuted  $g$ -factor as the target. eNetXplorer split the data into 10 folds 100 times/runs. For each run, eNetXplorer performed cross-validation by repeatedly training the target models on nine folds and tested on the leftover fold. Also, in each cross-validation run, eNetXplorer trained the null models 25 times. eNetXplorer then used the mean of non-zero model coefficients across all folds in a given run as a coefficient for each run,  $k'$ . Across runs, eNetXplorer weighted the mean of a model coefficient by the frequency of obtaining a non-zero model coefficient per run. Formally, we defined an empirical  $p$ -value as:

$$p_{val} = \frac{1}{1 + n_{run} * n_{per}} \left\{ 1 + \sum_{run=1}^{n_{run}} \sum_{per=1}^{n_{per}} \Theta \left( \left| \beta_{null}^{run,per} \right| - \left| \beta_{target}^{run} \right| \right) \right\}, \quad (1)$$

where  $p_{val}$  is an empirical  $p$ -value,  $run$  is a run index,  $n_{run}$  is the number of runs,  $per$  is a permutation index,  $n_{per}$  is the number of permutation,  $\Theta$  is the right-continuous Heaviside step function and  $|\beta|$  is the magnitude of feature coefficient. That is, to establish statistical significance for each feature, we used the proportion of runs in which the null models performed better than the target models. We plotted the target models' coefficients with  $p_{val} < .05$  on the brain images using the ggseg (Mowinckel & Vidal-Piñero, 2020) package.

To identify which modalities contributed strongly to the prediction of the stacked (i.e., Random Forests) model, we applied two methods: (1) CPI (Debeer & Strobl, 2020) and (2) SHAP (Lundberg & Lee, 2017) to the second-layer training set. CPI is an explainer, designed specifically for Random Forest. We implemented CPI using the 'permimp' package, as detailed elsewhere (Debeer & Strobl, 2020). Briefly, the original permutation importance (Breiman, 2001) shuffled the observations of one feature at a time while holding the target and other features in the same order. Researchers then examined decreases in predictive accuracy in the out-of-bag observations due to the permutation of some features. Stronger decreases are then assumed to reflect the importance of such features. However, this method has shown to be biased when there are correlated features (Strobl et al., 2007). CPI corrected for this bias by constraining the

feature permutation to be within partitions of other features, which was controlled by the threshold 's' value. We used the default s value at 0.95, which assumed dependencies among features (Debeer & Strobl, 2020).

SHAP (Lundberg & Lee, 2017) is a model-agnostic explainer, designed to explain the contribution of each feature to the prediction from any machine learning models via Shapley values (Roth, 1988). We implemented SHAP using the 'fastshap' package (<https://bgreenwell.github.io/fastshap/>). Based on the cooperative game theory, a Shapley value (Roth, 1988) quantifies a fair distribution of a payout to each player based on his/her contribution in all possible coalitions where each coalition includes a different subset of players. When applying Shapley values to machine learning, researchers treat each feature as a player in a game, a model output as a pay out and subsets of features as coalitions. Shapley values reflect the weighted differences in a model output when each feature is included versus not included in all possible subsets of features. SHAP (Lundberg & Lee, 2017) offers a computationally efficient approach to estimate Shapley values.

## 2.6 | Testing whether the brain-based predictive models mediated the relationships of the behaviourally derived *g*-factor with socio-demographic, psychological and genetic factors

Using leave-one-site-out cross-validation, we built three predictive models for the *g*-factor from (1) multimodal MRI (see above under 'Out-of-site Predictive Ability of Multimodal MRI'), (2) key socio-demographic and psychological factors and (3) a PGS. This resulted in three types of predicted values of the *g*-factor of unseen children at each hold-out data collection site: the brain-based *g*-factor, the socio-demography-and-psychology-based *g*-factor and the gene-based *g*-factor, respectively. We then test if the brain-based *g*-factor mediated the relationship that the behaviourally derived *g*-factor had with the socio-demography-and-psychology-based and gene-based *g*-factors.

## 2.7 | Key socio-demographic and psychological factors

We performed leave-one-site-out cross-validation to build 'socio-demographic-and-psychological-based' predictive models. These models predicted the behaviourally derived *g*-factor from key socio-demographic and psychological factors on the baseline data, similar to using leave-one-site-out cross-validation to create the 'brain-based' predictive models above. This enabled us to extract predicted values of the *g*-factor based on key socio-demographic and psychological factors at each hold-out site, called socio-demography-and-psychology-based *g*-factor. Here, we applied a similar modelling approach with leave-one-site-out cross-validation for multimodal MRI, except that we used only one layer of Elastic Net tuned with 200 levels of the penalty (from  $10^{-10}$  to 10) and 11 levels of the mixture (from 0 to 1). For pre-processing, we first imputed missing values of the

categorical features via mode replacement and then converted them to dummy variables. We next normalised these dummy variables and all numerical features and the behaviourally derived *g*-factor. At the last pre-processing step, we used k-nearest neighbour with five neighbours to impute the missing values of the normalised, numerical features.

Key socio-demographic and psychological factors included 70 features (Kirlic et al., 2021) collected at the baseline (9–10 years old): child's mental health based on symptom scales in Child Behavioral Checklist (Achenbach et al., 2017) (eight features), primary caretaker's mental health based on personal strengths and symptom scales in Aseba Adult Self Report (Achenbach et al., 2017) and General Behavior Inventory-Mania (Youngstrom et al., 2008) (nine features), child's personality based on Behavioral Inhibition System/Behavioral Activation System (Carver & White, 1994) and the UPPS-P Impulsive Behavior Scale (Zapolski et al., 2010) (nine features), child's sleep problems based on Sleep Disturbance Scale (Bruni et al., 1996) (eight features), child's physical activities based on Youth Risk Behavior Survey (Adolescent and School Health | CDC, 2020) (four features), child screen use (Bagot et al., 2018) (four features), parental use of alcohol, tobacco and marijuana after pregnancy based on Developmental History Questionnaire (Kessler et al., 2009; Merikangas et al., 2009) (three features), child developmental adversity (prematurity, birth complications and pregnancy complications) based on Developmental History Questionnaire (Kessler et al., 2009; Merikangas et al., 2009) (three features), child socio-demographics (Zucker et al., 2018) including sex, race, bilingual use (Dick et al., 2019), parental marital status, parental education, parental income, household size, economic insecurities, area deprivation index (Kind et al., 2014), lead risk (Frostenson & Kliff, 2016), crime reports (United States Department of Justice. Office of Justice Programs. Federal Bureau of Investigation, 2012), neighbourhood safety (Echeverria et al., 2004) and school environment, involvement and disengagement (Stover et al., 2010) (17 features) and child social interactions based on Parental Monitoring Scale (Chilcoat & Anthony, 1996), Child Report of Behavior Inventory (Schaefer, 1965), Strengths and Difficulties Questionnaire (Goodman et al., 2003) and Moos Family Environment Scale (Moos & Humphrey, 1974) (five features).

## 2.8 | Polygenic scores

To extract predicted values of the *g*-factor based on genetics, we used PGSs for adult cognitive ability (Lee et al., 2018). The ABCD study provided details on genotyping elsewhere (Uban et al., 2018). Briefly, the study took saliva and whole blood samples and genotyped them using Smokescreen™ Array. The ABCD applied quality control based on calling signals and variant call rates, ran the Ricopili pipeline and imputed the data with TOPMED. We excluded data from problematic plates and with a subject-matching issue, identified by the ABCD. We further quality controlled the data as follows. First, we removed individuals with minimal or excessive heterozygosity. We also excluded SNPs based on minor allele frequency (<5%) and violations of Hardy-Weinberg equilibrium ( $P < 1E-10$ ). We limited the analysis to 'unrelated individuals' as defined



by individuals with low genetic relatedness (more than third-degree relative pairs; identical by descent [IBD]  $\geq 0.0422$ ).

We defined alleles associated with the  $g$ -factor as those related to cognitive abilities in a large-scale discovery GWAS sample of European ancestry ( $N = 257,841$ ) (Lee et al., 2018). Given the lower predictive performance of PGS when the ancestry of a sample does not match with that of the discovery GWAS sample (Duncan et al., 2019), we restricted all analyses related to PGS to children of European ancestry (Duncan et al., 2019). We considered children to be genetically similar to the ancestry reference if they were within four standard deviations of the mean of the top four principal components (PCs) of the super-population individuals in the 1000 genomes Phase 3 reference genotypes (1000 Genomes Project Consortium, 2015).

We used the Pthreshold PGS approach where we defined risk alleles as those associated with cognitive abilities within the discovery GWAS sample (Lee et al., 2018) at 10 different thresholds from  $p < .5-.00000001$  (referred to as PGS thresholds). The final sample for PGS included 4,814 children (2,263 females;  $M_{\text{age}} = 9.94$  [SD = .61] years). We computed PGS as the Z-scored, weighted mean number of linkage independent risk alleles. While the  $g$ -factor was significantly related to the PGS of cognitive ability across thresholds (Figure 7), the relationship at the  $p < .01$  PGS threshold was the numerically strongest ( $r = 0.21$ ,  $p < .001$  [95%CI = 0.18–0.24]). Accordingly, we focused our analyses using the  $p < .01$  PGS threshold and treated the PGS at this threshold as our gene-based  $g$ -factors.

## 2.9 | Mediation analyses

To examine the extent to which brain-based, stacked predictive models of the  $g$ -factor accounted for the relationship between the behaviourally derived  $g$ -factor and the socio-demographic, psychological and genetic  $g$ -factors, we applied mediation analyses (MacKinnon et al., 2007). In these mediation analyses, we treated (i) the brain-based  $g$ -factor as the mediator, (ii) the socio-demography-and-psychology-based and gene-based  $g$ -factors as the independent variables and (iii) the behaviourally derived  $g$ -factor as the dependent variable. Note the behaviourally derived  $g$ -factor was computed based on the CFA models in the training data, which were later applied to each hold-out site. While the behaviourally derived  $g$ -factor was a latent variable, it represented the only 'observed' value here since the other three  $g$ -factors (brain-based, socio-demography-and-psychology-based and gene-based) were 'predicted' values from predictive models.

We conducted three mediation analyses. The first analysis only used the socio-demography-and-psychology-based  $g$ -factor as the independent variable. The second analysis only used the gene-based  $g$ -factor as the independent variable. The third analysis used both the socio-demography-and-psychology-based and gene-based  $g$ -factors as the independent variables, simultaneously in the same model. To control for population stratification in genetics, we also included four PCs as control variables in the mediation analyses involving the gene-based  $g$ -factor.

To implement the mediation analyses, we used structural equation modelling (SEM) with 5000 bootstrapping iterations via lavaan (Rosseel, 2012). We specifically calculated the *indirect effects* to show

whether the relationships between the behaviourally derived  $g$ -factor and the socio-demography-and-psychology-based and gene-based  $g$ -factors were significantly explained by the brain-based  $g$ -factor. Along with the indirect effects, we also computed the *proportion mediated* to demonstrate the proportion of variance accounted for by the brain-based  $g$ -factor.

## 2.10 | Data and code availability

We used publicly available data provided by the ABCD study (<https://abcdstudy.org>), held in the NIMH Data Archive (<https://nda.nih.gov/abcd/>).

We uploaded the R analysis script and detailed outputs for predictive modelling: <https://narunpat.github.io/GFactorModelingABCD3/GFactorModelingABCD3.html> and mediation analyses: <https://narunpat.github.io/GFactorModelingABCD3/MediationSocDemPsycPGSBrainABCD3.html>.

## 3 | RESULTS

### 3.1 | How robust are the factor scores of the $g$ -factor based on the second-order model?

Based on our CFA, the second-order model of the  $g$ -factor showed a good fit: (a) scaled, robust CFI = 0.995, (b) scaled, robust TLI = 0.988, (c) scaled, robust root mean square error of approximation (RMSEA) = 0.029 (90%CI = 0.015–0.043) and (d) robust SRMR = 0.014. The  $g$ -factor latent variable of the second-order model also had high internal consistency: OmegaL2 = 0.78.

See Appendix S1 and S2 for a more detailed CFA of the  $g$ -factor. In brief, firstly, the second-order model had better fit indices than the single-factor model. Additionally, factor scores of the  $g$ -factor from the second-order model, the single-factor model, and the mixture between EFA and CFA models were similar to each other at high magnitude (Pearson's  $r_s \geq 0.987$ ). Accordingly, the choice of  $g$ -factor models had only minimal effects on the estimation of the factor scores for the  $g$ -factor, and thus our brain-based predictive models should be generalisable to the factor scores of different  $g$ -factor CFA models beyond the second-order model. Lastly, the factor scores estimated from the first-layer training data were similar to the factor scores estimated from the full baseline data at high magnitude (Pearson's  $r_s > 0.997$ ), indicating the stability of the factor scores used.

### 3.2 | How robust are the brain-based predictive models?

#### 3.2.1 | Out-of-sample predictive ability of multimodal MRI

For hyperparameter-tuning results, see Appendix S3. Table 1 and Figure 2 summarise the out-of-sample predictive ability of

**TABLE 1** Out-of-sample and out-of-site predictive ability of multimodal MRI

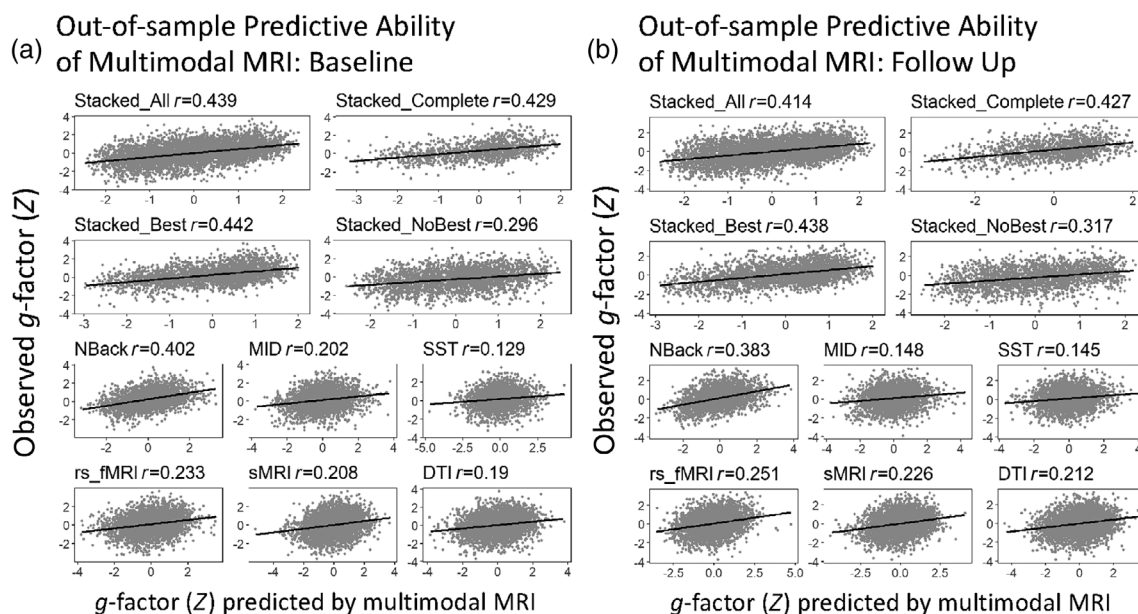
Out-of-sample predictive ability of multimodal MRI: Baseline samples				
Models	$r$	$R^2$	MAE	RMSE
Stacked_All	0.439	0.191	0.699	0.895
Stacked_Complete	0.429	0.183	0.61	0.78
Stacked_Best	0.442	0.195	0.62	0.798
Stacked_NoBest	0.296	0.085	0.783	0.987
N-Back	0.402	0.072	0.664	0.857
SST	0.129	-0.033	0.744	0.95
MID	0.202	0.013	0.738	0.944
rs_fMRI	0.233	0.042	0.749	0.955
sMRI	0.208	0.04	0.763	0.969
DTI	0.19	0.033	0.757	0.972
Out-of-sample predictive ability of multimodal MRI: follow-up samples				
Models	$r$	$R^2$	MAE	RMSE
Stacked_All	0.414	0.166	0.719	0.913
Stacked_Complete	0.427	0.168	0.651	0.829
Stacked_Best	0.438	0.175	0.666	0.846
Stacked_NoBest	0.317	0.1	0.794	1
N-Back	0.383	0.118	0.687	0.875
SST	0.145	-0.004	0.76	0.961
MID	0.148	-0.003	0.757	0.955
rs_fMRI	0.251	0.055	0.754	0.954
sMRI	0.226	0.049	0.764	0.965
DTI	0.212	0.045	0.771	0.978
Mean (SD) of out-of-site predictive ability of multimodal MRI				
Models	$r$	$R^2$	MAE	RMSE
Stacked	0.46 (0.057)	0.21 (0.052)	0.698 (0.023)	0.888 (0.029)
N-Back	0.408 (0.069)	0.167 (0.055)	0.718 (0.028)	0.91 (0.031)
MID	0.227 (0.096)	0.05 (0.05)	0.772 (0.021)	0.973 (0.025)
SST	0.139 (0.071)	0.019 (0.024)	0.783 (0.014)	0.988 (0.012)
rs_fMRI	0.255 (0.061)	0.064 (0.03)	0.765 (0.015)	0.966 (0.016)
sMRI	0.248 (0.092)	0.061 (0.046)	0.763 (0.024)	0.967 (0.024)
DTI	0.223 (0.076)	0.049 (0.037)	0.766 (0.016)	0.974 (0.019)

Abbreviations: MAE, mean absolute error;  $R$ , Pearson's correlation;  $R^2$ , coefficient of determination; RMSE, root mean squared error.

multimodal MRI for both baseline and follow-up samples. Performance of Stacked All, Stacked Complete and Stacked Best was among the top with Pearson's  $r$  over 0.4 and  $R^2$  around 0.19. Importantly, the superior performance of stacked models was found across baseline and follow-up test sets at a similar magnitude, suggesting their longitudinal stability. Note that given that the N-back task-based fMRI had the highest performance among modality-specific models, we set the missing values of the Stacked Best to be the same as those of the N-back task-based fMRI. Moreover, the opportunistic stacking (Engemann et al., 2020) algorithm that led to the stacked model with at least one modality present, Stacked All, was robust against missing values as the performance

of Stacked All was similar to that of the stacked model with all modalities present, Stacked Complete.

Figure 3 shows the proportion of missing data in the two test sets. sMRI had the lowest missing observations, while the three task-based fMRI data had the highest. Missing observations in Stacked All were around 3%–6%, while those in Stacked Complete were up to 78.79%. Figure 3 also shows the differences in the  $g$ -factor between participants with versus without missing values for each model in the two test sets. Participants with missing values had a significantly lower  $g$ -factor than those without missing values, as indicated by Welch's  $t$ -test, for all models, except for Stacked No Best, which showed the opposite direction. Yet, numerically these differences in



**FIGURE 2** Out-of-sample predictive ability of multimodal MRI as a function of modalities in the test sets for baseline (a) and follow-up (b) samples. Stacked all required the test data with at least one modality present. Stacked complete required the test data with all modalities present. Stacked best had the same missing values with the modality with the best prediction (N-back task-based fMRI). Stacked no best did not have any test data from the modality with the best prediction and had at least one modality present

the  $g$ -factor were weaker in magnitude in the Stacked All than in other models with high predictive performance (such as the N-back task-based fMRI and Stacked Complete) as indicated by Cohen's  $d$ . Accordingly, by imputing the data via the opportunistic stacking (Engemann et al., 2020), we were able to include more participants, and thus, less likely to exclude participants with a lower  $g$ -factor.

### 3.2.2 | Comparing out-of-sample predictive ability of multimodal MRI between the stacked model and N-back task-based fMRI

N-back task-based fMRI provided the best out-of-sample predictive ability for both baseline and follow-up test sets, relative to other modality-specific models. Figure 4 compared the predictive ability between the Stacked Best and N-back task-based fMRI using bootstrapped differences. The Stacked Best had significantly higher performance in both baseline and follow-up test sets, reflected by higher Pearson's  $r$  and  $R^2$  and lower MAE and RMSE. This indicates the boost in predictive performance when multiple modalities were integrated, at around 12% for the baseline data and 6% for the follow-up data. Accordingly, the stacked model performed better than the best single modality.

### 3.2.3 | Out-of-site predictive ability of multimodal MRI

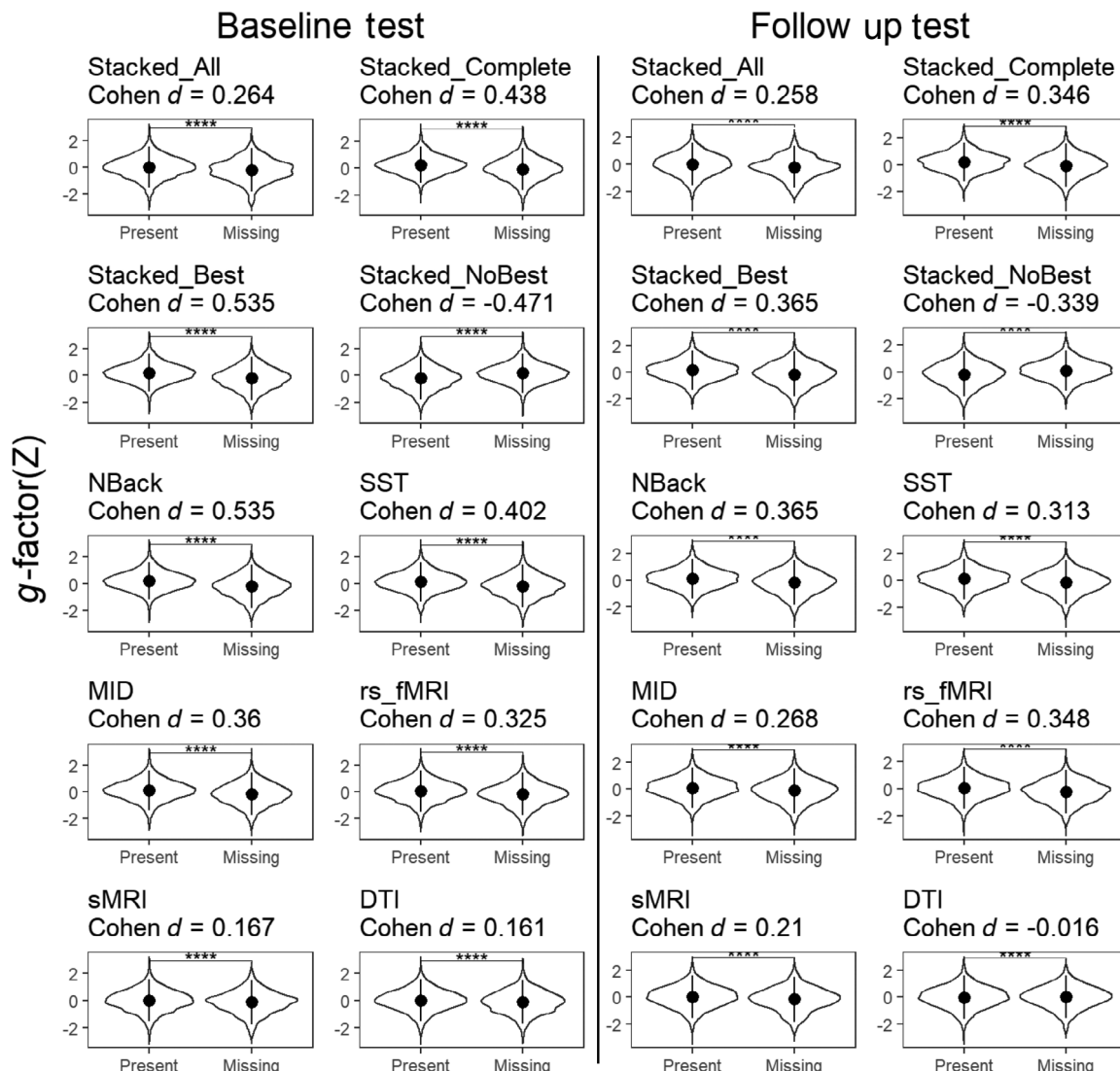
Based on leave-one-site-out cross-validation, the out-of-site predictive ability of the stacked model was highest, explaining on-average

21% (SD = 5.2) of the variance in the  $g$ -factor across 21 sites (Table 1 and Figure 5). This confirmed the generalisability of the stacked model and ensured its use for subsequent mediation analyses.

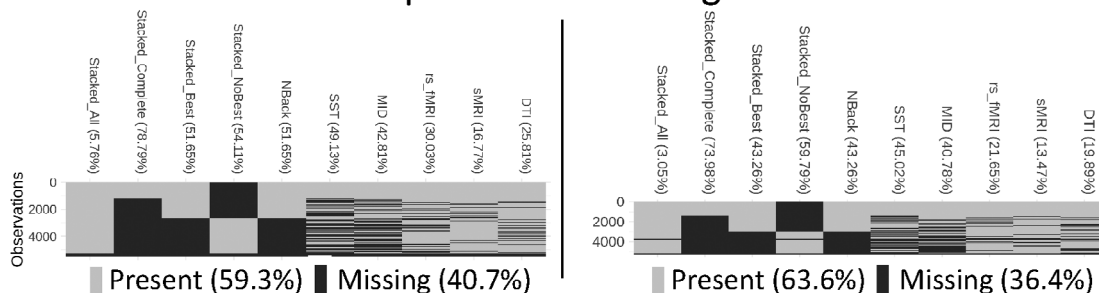
### 3.3 | Feature importance of multimodal MRI models

Figure 6 shows the feature importance of both the modality-specific and stacked models. For the modality-specific models, we applied eNetXplorer (Candia & Tsang, 2019a) to show brain features that significantly (empirical  $p < .05$ ) contributed to the prediction. For N-back task-based fMRI, the  $g$ -factor prediction was driven by activity in areas, such as the precuneus, sulcus intermedius primus, superior frontal sulci and dorsal cingulate. For MID task-based fMRI, the prediction was driven by activity in several areas in the parietal, frontal and temporal regions. For SST, the prediction was contributed by activity in areas such as the supramarginal gyrus and inferior precentral sulcus. For rs-fMRI, the prediction was driven by connectivity within cinguloparietal and sensory-motor-hand as well as between networks that were connected with frontoparietal, default-mode and sensory-motor-hand networks. For sMRI, the prediction was driven by the volume/thickness at several areas, such as the insula, middle frontal gyrus and lingual sulcus. For DTI, the prediction was driven by FA at several white matter tracts, such as the superior longitudinal fasciculus, forceps minor, uncinata and parahippocampal cingulum. For the stacked model, we applied the CPI (Strobl et al., 2008) and SHAP (Lundberg & Lee, 2017) to examine which of the modalities contributed strongly to the prediction. CPI and SHAP provided similar results. N-back task-based fMRI by far had the highest importance score.

(a) Differences in  $g$ -Factor( $Z$ ) as a function of missing data

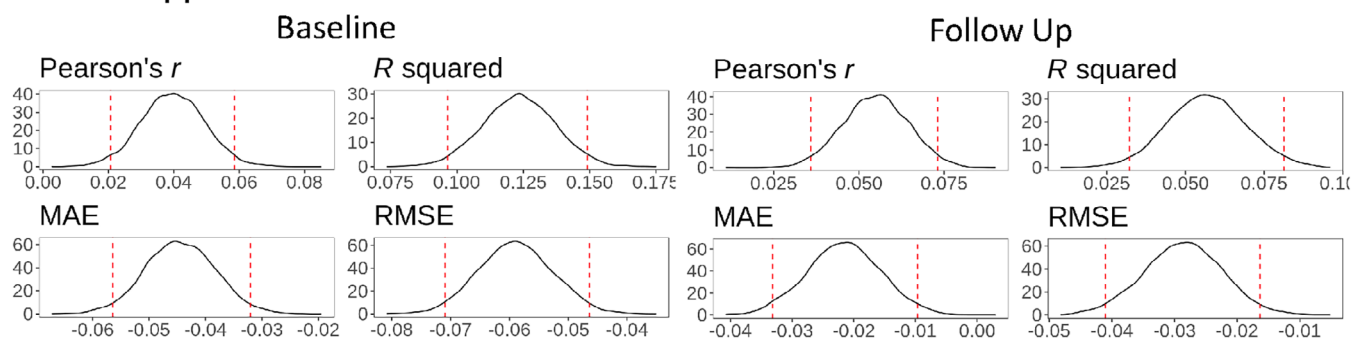


(b) Proportion of missing data



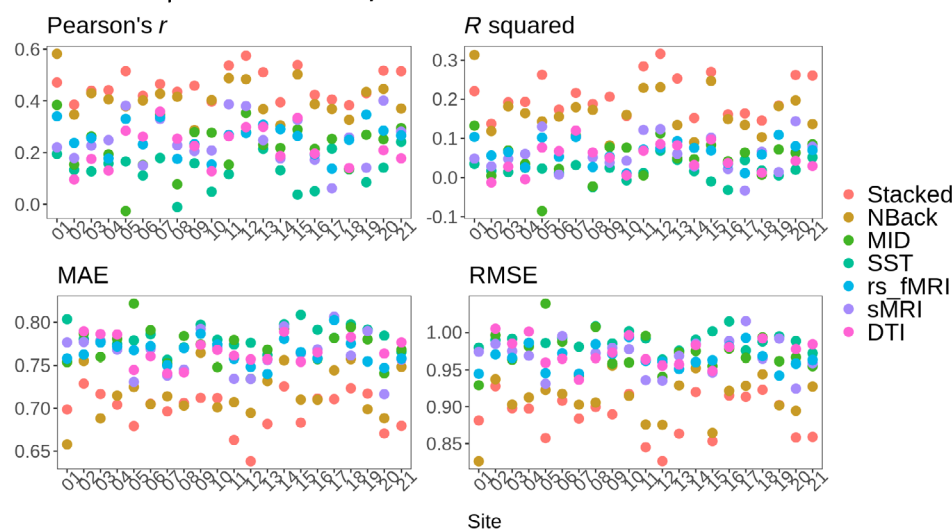
**FIGURE 3** Missing values in each predictive model in the baseline and follow-up test sets. (a) Shows the differences in the  $g$ -factor between participants with versus without missing values for each predictive model in the two test sets. \*\*\*\* indicates  $p$ -value  $< .001$  based on Welch's  $t$ -test. Positive Cohen's  $d$  indicates that participants without missing values had a higher  $g$ -factor than participants without missing values. Dot and line are the mean and standard deviation  $\times 2$  of the  $g$ -factor, respectively. (b) Shows the proportion of missing data for each predictive model in the two test sets

## Comparing Out-of-sample Predictive Ability of Multimodal MRI: Bootstrapped Distribution of Stacked Best > N-Back



**FIGURE 4** Comparing out-of-sample predictive ability of multimodal MRI between stacked best and the model with the best modality (N-back task-based fMRI). Here, we separately applied bootstrapping on the baseline and follow-up test sets. At each of 5000 iterations, we computed performance indices (including  $r$ ,  $R^2$ , MAE and RMSE) of stacked best and N-back task-based fMRI models and subtracted performance indices of N-back task-based fMRI from that of stacked best. Dotted lines indicate 95% confidence intervals. MAE, mean absolute error;  $R^2$ , coefficient of determination; RMSE, root mean squared error

### Out-of-site predictive ability of multimodal MRI



**FIGURE 5** Out-of-site predictive ability of multimodal MRI via leave-one-site-out cross-validation. We evaluated out-of-site predictive ability between predicted versus observed  $g$ -factor in the hold-out site. Note that DTI data were not available from three sites (sites 1, 17 and 19). MAE, mean absolute error;  $R^2$ , coefficient of determination; RMSE, root mean squared error

### 3.4 | Did the brain-based predictive models mediate the relationships of the behaviourally derived $g$ -factor with socio-demographic, psychological and genetic factors?

#### 3.4.1 | Key socio-demographic and psychological factors

Based on leave-one-site-out cross-validation, socio-demographic and psychological factors explained on-average 29.7% (SD = 8.1) of the variance in the behaviourally derived  $g$ -factor across sites (see Figure 7). The top features in the Elastic-Net models that had the magnitude of their standardised coefficients over 0.1 included parents' education and income along with child's attention and social problems as well as extracurricular activities.

#### 3.4.2 | Polygenic scores

Figure 8a,b shows the relationship between the behaviourally derived  $g$ -factor and the PGS of cognitive abilities at different thresholds. While the behaviourally derived  $g$ -factor was significantly related to the PGS of cognitive abilities across thresholds, the relationship at the  $p < .01$  PGS threshold was the numerically strongest ( $r = 0.21$ ,  $p < .001$  [CI95% = 0.18–0.24]). Accordingly, we used PGS at the  $p < .01$  PGS threshold as the gene-based  $g$ -factor for the mediation analyses.

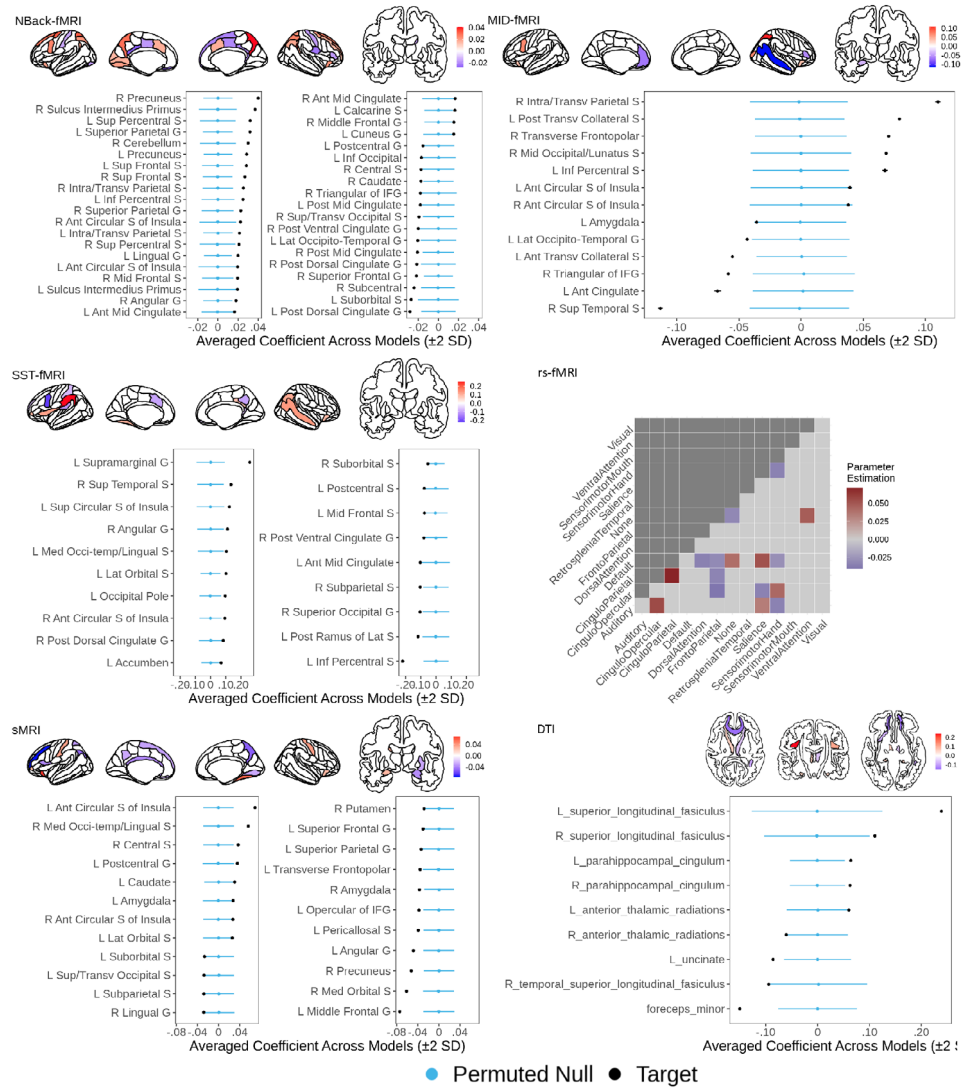
#### 3.4.3 | Mediation analyses

We tested whether brain-based  $g$ -factor mediated the relationships between the behaviourally derived  $g$ -factor and socio-demography-and-



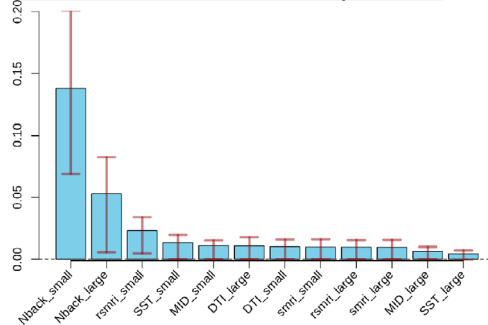
**FIGURE 6** Feature importance of the modality-specific and stacked models. For the modality-specific models, we applied eNetXplorer (Candia & Tsang, 2019a) permutation and only plotted brain features with empirical  $p < .05$ . For the stacked model, we applied conditional permutation importance (CPI) (Debeer & Strobl, 2020) and SHapley additive exPlanations (SHAP) (Lundberg & Lee, 2017). Both CPI and SHAP were computed based on the second-layer training set. Error bars in the CPI plot show an interval between 0.25 and 0.75 quantiles of the CPI for each tree in the random forests. The ‘\_large’ and ‘\_small’ suffixes indicate whether the missing values were coded as a large (1000) or small (−1000) number, respectively. For SHAP, we combined Shapley values across the two coded features of the same modality. We then ranked the modalities according to the absolute value of SHAP; the highest one was N-back task-based fMRI. Note the grey colour indicates observations with a missing value (coded as 1000 or −1000). ant, anterior; G, gyrus; IFG, inferior frontal gyrus; L, left; Lat, lateral; med, medial; R, right; S, sulcus; Sup, superior

**Feature Importance of the Modality-Specific Models: eNetXplorer**

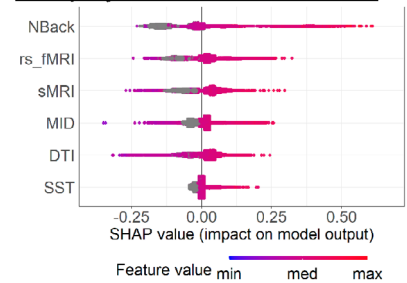


**Feature Importance of the Stacked Model:**

**Conditional Permutation Importance**



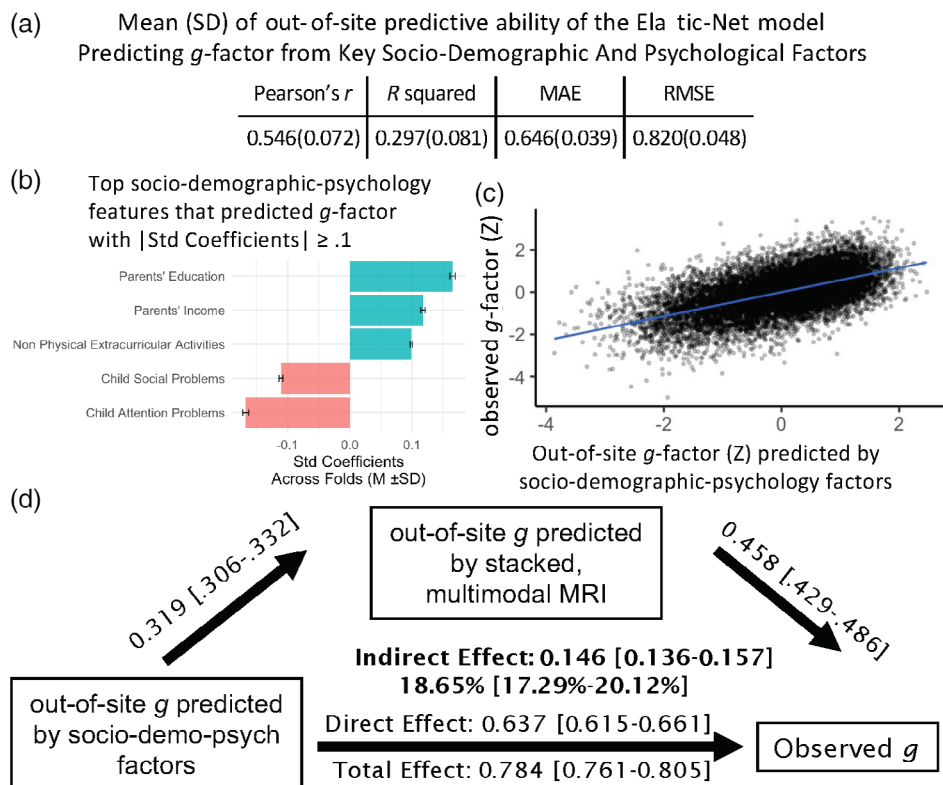
**SHapley Additive exPlanations**



psychology-based and gene-based  $g$ -factors. We found significant indirect effects (1) when the socio-demography-and-psychology-based  $g$ -factor was the sole independent variable (Figure 7d proportion mediated = 19.1%), (2) when the gene-based  $g$ -factor was the sole independent variable (Figure 8c, proportion mediated = 15.6%) and (3) when both socio-demography-and-psychology-based  $g$ -factor (Figure 9, proportion mediated = 15%) and gene-based  $g$ -factor (Figure 9, proportion mediated = 10.75%) were the covered independent variables.

**4 | DISCUSSION**

Following the RDoC's integrative approach for cognitive abilities (Morris & Cuthbert, 2012), we aimed to develop brain-based predictive models that can (a) improve our current ability to predict children's cognitive abilities and (b) account for the relationships between cognitive abilities and socio-demographic, psychological and genetic factors. Here, we showed that incorporating data from different MRI modalities



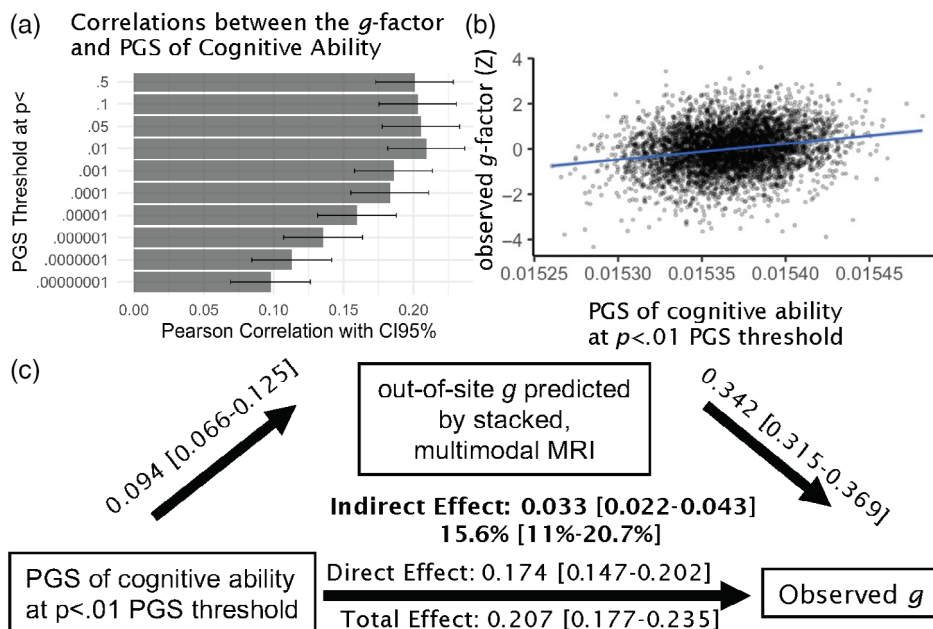
**FIGURE 7** Key socio-demographic and psychological factors. (a) Shows the out-of-site predictive ability of the elastic-net model predicting the  $g$ -factor from key socio-demographic and psychological factors, based on leave-one-site-out cross-validation. (b) Shows the top socio-demographic and psychological features with the magnitude of standardised coefficient over 0.1 based on the elastic-net model. Blue indicates a positive relationship while red indicates a negative relationship. (c) Shows a scatter plot between out-of-site predicted values of the  $g$ -factor based on key socio-demographic and psychological factors and the observed (i.e., real) values of the  $g$ -factor. (d) Shows a mediation analysis where (1) the socio-demography-and-psychology-based  $g$ -factor (the out-of-site predicted values of the  $g$ -factor based on the key socio-demographic and psychological factors at all hold-out sites) is the independent variable, (2) the brain-based  $g$ -factor (the out-of-site predicted values of the  $g$ -factor of the stacked model based on multimodal MRI data at all hold-out sites) is the mediator and (3) the behaviourally derived  $g$ -factor (the observed  $g$ -factor) is the dependent variable. % under the indirect effect indicates proportion mediated. [] indicates a 95% confidence interval based on bootstrapping. MAE, mean absolute error;  $R^2$ , coefficient of determination; RMSE, root mean squared error

into stacked models substantially improved our ability to predict cognitive abilities, operationalised as the behaviourally derived  $g$ -factor. Our brain-based, stacked predictive models were stable across years and generalisable to different sites while being able to handle missing values. Moreover, we showed that the brain-based, stacked models significantly, albeit partially, mediated the relationships of the behaviourally derived  $g$ -factor with socio-demographic, psychological and genetic factors. Thus, our brain-based predictive models for children's  $g$ -factor demonstrated construct validity according to the RDoC framework (Insel et al., 2010; Morris & Cuthbert, 2012; NIMH, n.d.-a; n.d.-c).

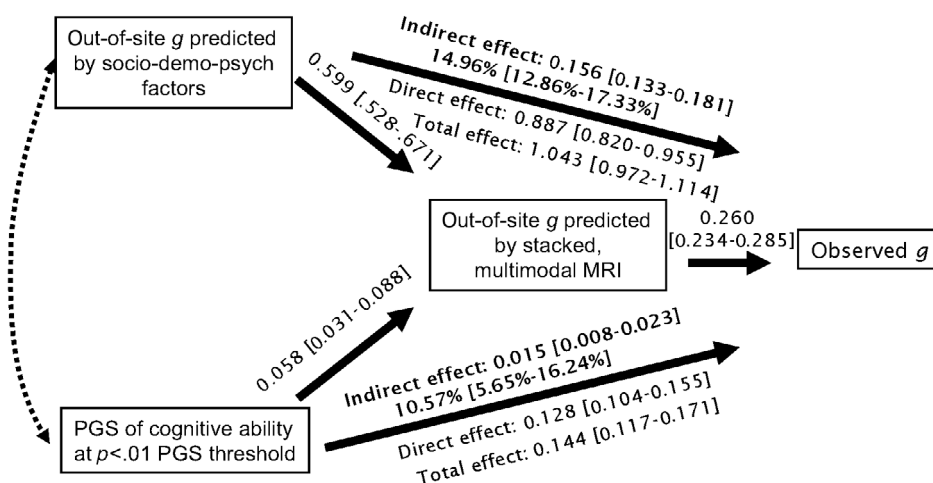
#### 4.1 | The brain-based, stacked predictive models for the $g$ -factor were (1) predictive, (2) longitudinal stable, (3) robust against missing values and (4) explainable

We developed longitudinal predictive models for children's  $g$ -factor from MRI data of different modalities. We built models from the

baseline MRI data and tested them on unseen children at the same age and 2 years older. We found similar predictive abilities across these two test sets for all modality-specific and stacked models. That is, the models that had high out-of-sample prediction on same-age children also had high out-of-sample prediction on older children, suggesting the longitudinal stability of MRI for many modalities. The best model across all performance indices (Pearson's  $r$ ,  $R^2$ , MAE and RMSE) was the stacked model that incorporated all six modalities, which was followed closely by the N-back task-related fMRI model. Apart from the SST task-related fMRI model, other models (including the MID task-related, rs-fMRI, sMRI and DTI) performed moderately well. We also found a similar magnitude for out-of-site predictive ability based on leave-one-site-out cross-validation, suggesting the generalisability of MRI not only across ages but also across data collection sites. Overall, the stacked model partially predicted the children's  $g$ -factor at around 20% of the variance. This made the stacked model the most generalisable model to out-of-sample, out-of-site children as well as the most longitudinally stable model.



**FIGURE 8** Polygenic scores (PGSs) of cognitive abilities. (a) Shows Pearson's correlations between the *g*-factor and PGS of cognitive abilities at different PGS thresholds. (b) Shows a scatter plot between the PGS of cognitive abilities at the *p* < .01 PGS threshold and the observed (i.e., real) values of the *g*-factor. (c) Shows a mediation analysis where (1) gene-based *g*-factor (the PGS of cognitive abilities at the *p* < .01 PGS threshold) is the independent variable, (2) the brain-based *g*-factor (the predicted values of the *g*-factor of the stacked model based on multimodal MRI data at all hold-out sites) is the mediator and (3) the behaviourally derived *g*-factor (the observed *g*-factor) is the dependent variable. Not shown in the figure are four PCs included as the control variables to adjust for population stratification. % under the indirect effect indicates proportion mediated. [] indicates a 95% confidence interval based on bootstrapping



**FIGURE 9** Mediation analysis with both key socio-demographic and psychological factors as well as genetic factors as independent variables. Specifically, this model treated (1) the socio-demography-and-psychology-based *g*-factor (i.e., the out-of-site predicted values of the *g*-factor based on the key socio-demographic and psychological factors at all hold-out sites) and (2) the gene-based *g*-factor (i.e., the PGS of cognitive abilities at the *p* < .01 PGS threshold) as two separate independent variables. It treated the brain-based *g*-factor (i.e., the predicted values of the out-of-site predicted values of the *g*-factor of the stacked model based on multimodal MRI data at all hold-out sites) as the mediator and the behaviourally derived *g*-factor (i.e., the observed *g*-factor) as the dependent variable. Not shown in the figure are four PCs included as the control variables to adjust for population stratification. % under the indirect effect indicates proportion mediated. [] indicates a 95% confidence interval based on bootstrapping. The dotted, double arrowed line indicates covariation between the two independent variables. PGS, polygenic score

Beyond generalisability across ages and sites, the stacked model based on opportunistic stacking (Engemann et al., 2020) also allowed us to handle missingness in the data. This is especially important for children's MRI data given high levels of noise in certain modalities

(Fassbender et al., 2017). If we were to use data only from children with all modalities present (i.e., the Stacked Complete), the model would not apply to around 80% of the children. The opportunistic stacking allowed us to use the data as long as one modality was

present (i.e., the Stacked All), leaving the exclusion to just around 5%. Importantly, the predictive performance of Stacked Complete and Stacked All were both relatively high, ensuring the ability of opportunistic stacking to deal with the missing data. Furthermore, handling missingness in the data via opportunistic stacking also heightened the chance of including participants with a wider range of the  $g$ -factor, including those with a lower  $g$ -factor who usually had missingness in the MRI data (perhaps due to high movement artefacts [Fassbender et al., 2017]). Moreover, in the case when the best modality was not available, using the stacked model (i.e., the Stacked No Best) could be helpful. While the predictive ability of the Stacked No Best was not as strong as the Stacked Complete, Stacked All and Stacked Best, its performance measures of variance (Pearson's  $r$  and  $R^2$ ) appeared stronger in magnitude than any other non-optimal modalities by themselves. Accordingly, in settings where not all of the modalities are available, researchers/practitioners can still take advantage of the boosted predictive ability of the stacked models over unimodal models.

The stacked model improved predictive ability over and above the best modality, which was the N-back task-based fMRI. This is based on bootstrapping distributions of the differences in performance indices between the N-back task-based fMRI and the stacked model with the same participants (i.e., the Stacked Best). Our finding is consistent with previous studies showing the enhanced predictive power of the stacked model (Engemann et al., 2020; Rasero et al., 2021). Yet, it is important to note that, while the improvement in performance was statistically significant, the magnitude of this improvement was somewhat modest. For instance, in the case of the baseline samples, the Stacked Best led to  $r = 0.442$  and  $R^2 = 0.195$ , which was improved from the N-back task-based fMRI at  $r = 0.402$  and  $R^2 = 0.072$ , rendering the improvement at around  $r \sim 0.04$  and  $R^2 \sim 0.123$ . Accordingly for researchers who have access to all MRI modalities and several fMRI tasks, including the N-back task, using the stacked model should provide the best possible performance for predicting the  $g$ -factor. However, if resources are constrained, the next best option would be using the N-back task-based fMRI along with other modalities that are available.

In addition to predictability, our machine learning framework allowed for easy-to-explain models, highlighting the neurobiological bases of children's  $g$ -factor. Explainability is used in a specific machine-learning sense (Molnar, 2019), referring to the extent to which a technique applied allows us to explain the contribution of each brain feature to the prediction. Here, CPI (Debeer & Strobl, 2020) and SHAP (Lundberg & Lee, 2017) allowed us to infer that prediction from the stacked model was driven primarily by N-back task-related fMRI. This indicates the important role of working memory. eNetXplorer permutation (Candia & Tsang, 2019a) further showed us that contribution from fMRI activity in the parietal and frontal areas during the N-back task drove the prediction. These areas were similar to the areas previously found in a recent study in adults (Sripada, Angstadt, et al., 2020). Similarly, we also found brain indices from other modalities, from activity during other tasks to the cortical thickness and white matter density, that contributed to the prediction of the  $g$ -factor, albeit with lower predictive performance.

Unlike previous unimodal (Dubois et al., 2018; Genç et al., 2018; Góngora et al., 2020; Gray et al., 2003; Narr et al., 2007; Pamplona et al., 2015; Sripada, Rutherford, et al., 2020; Waiter et al., 2009) and multimodal studies (Jiang et al., 2020; Rasero et al., 2021), we were able to compare the ability of task-based fMRI with other modalities in predicting the  $g$ -factor. We found that one of the three task-based fMRI models, the N-back, performed exceptionally well. Based on the CPI (Debeer & Strobl, 2020) and SHAP (Lundberg & Lee, 2017), the N-back task-related fMRI appeared to drive the prediction of the stacked model. This finding is consistent with a recent study using adults' data from the Human Connectome Project, showing superior performance of the N-back task in predicting the  $g$ -factor, compared to rs-fMRI (Sripada, Angstadt, et al., 2020) and other tasks. Showing that task-based fMRI from a certain task could capture cognitive ability across a 2-year gap provided a promising outlook for the use of task-based fMRI as a predictive tool. Our finding is contradictory to a more common practice in cognitive neuroscience that usually relies on sMRI (McDaniel, 2005; Mihalik et al., 2019; Pietschnig et al., 2015) or rs-fMRI (Dubois et al., 2018; Rasero et al., 2021; Sripada, Angstadt, et al., 2020) when predicting cognitive abilities. These sMRI and rs-fMRI studies often result in poorer predictive performance (at  $r < 0.4$ ) than what was found here. Accordingly, we are in agreement with a recent movement (Finn, 2021) for studies on individual differences to move from rs-fMRI and embrace other MRI modalities, including task-based fMRI.

It is important to note that not all fMRI tasks were suitable for predicting certain targets. The N-back task and SST, for instance, were designed to capture working memory (Barch et al., 2013; Casey et al., 2018) and inhibitory control (Casey et al., 2018; Whelan et al., 2012), respectively. Accordingly, both should be related to the  $g$ -factor, especially on memory recall and mental flexibility portions of the  $g$ -factor. Yet, only the N-back task showed good predictive ability. This may be due to different cognitive processes in each task (i.e., working memory vs. inhibitory control) or to different task configurations. It is entirely possible, for instance, that the block design used in the N-back, as opposed to the event-related design used in the SST, allowed the N-back to have higher predictive power. Accordingly, while task-based fMRI can have high predictive power, systematic comparisons are required in future research to better understand the characteristics of some tasks that make them more suitable for predicting the  $g$ -factor and other individual differences.

#### 4.2 | The brain-based, stacked predictive models for the $g$ -factor demonstrated construct validity, according to the RDoC framework (Insel et al., 2010)

Here we tested the construct validity of the brain-based, stacked predictive models for the  $g$ -factor according to the RDoC framework (Insel et al., 2010). The RDoC proposes that cognitive abilities are affected by socio-demographic and psychological factors (Morris & Cuthbert, 2012; NIMH, n.d.-b). The RDoC also proposes that cognitive abilities as measured by brain differences belong to the same domain as cognitive abilities as measured by gene differences (Insel

et al., 2010; Morris & Cuthbert, 2012; NIMH, n.d.-a, n.d.-c). Accordingly, to satisfy these presuppositions, our brain-based, stacked predictive models for the *g*-factor should be able to capture the relationship between the behaviourally derived *g*-factor and socio-demographic, psychological and genetic factors.

We first built a predictive model of the *g*-factor using 70 socio-demographic and psychological features (Kirlic et al., 2021), resulting in the socio-demography-and-psychology-based *g*-factor. This model had relatively high performance, accounting for around 30% of the *g*-factor. Moreover, the top contributing features are consistent with previous studies, including socio-demographics (Farah et al., 2006) (e.g., parents' education and income) along with children's mental health (Biederman et al., 2004; Goodall et al., 2018) (e.g., attention and social problems) and children's extracurricular activities (Kirlic et al., 2021). More importantly, our mediation analysis showed that the brain-based *g*-factor captured approximately 19% of the relationship between the behaviourally derived *g*-factor and the socio-demography-and-psychology-based *g*-factor.

As for the genetic factor, we first showed that the PGS based on adults' cognitive abilities (Lee et al., 2018) was related to children's *g*-factor, consistent with a recent study (Allegrini et al., 2019). This enabled us to use the PGS of cognitive abilities as the gene-based *g*-factor. Similar to the socio-demography-and-psychology-based *g*-factor, our mediation analysis showed that the brain-based *g*-factor accounted for approximately 16% of the relationship between the behaviourally derived *g*-factor and the gene-based *g*-factor. In fact, mediation from the brain-based *g*-factor was still significant when having both socio-demography-and-psychology-based and gene-based *g*-factors together as independent variables in the model. Altogether, our brain-based, stacked predictive models for the *g*-factor demonstrated the construct validity of cognitive abilities that is in line with the RDoC framework (Insel et al., 2010).

### 4.3 | Applications, limitations and disclaimers

For applications, our brain-based predictive models for the *g*-factor facilitate the development of a robust, transdiagnostic research tool for cognition at the neural level in keeping with the RDoC (Morris & Cuthbert, 2012). Cognitive abilities are one of RDoC's six major transdiagnostic domains (Morris & Cuthbert, 2012), relating to a number of psychiatric disorders (Sheffield et al., 2018; Shilyansky et al., 2016; Thaler et al., 2013). Based on RDoC (Morris & Cuthbert, 2012), to improve our understanding of cognitive abilities, we need research tools that allow us to integrate different units of analysis, from behavioural down to neural and genetic levels, and that reflect the influences of socio-demographical and psychological factors across the lifespan (Insel et al., 2010; NIMH, n.d.-a). Our brain-based predictive models satisfied many presuppositions of RDoC (Morris & Cuthbert, 2012). Our brain-based predictive models for the *g*-factor were not only longitudinal stable (Insel et al., 2010; NIMH, n.d.-a), but they also captured the influences of socio-demographical, psychological and genetic factors on cognitive abilities (Insel et al., 2010;

Morris & Cuthbert, 2012; NIMH, n.d.-a, n.d.-c). In fact, the predictive ability of our brain-based predictive models in capturing the behavioural performance of cognitive tasks was considerably higher than that of PGS (multimodal MRI's  $r \sim 0.4$  and  $R^2 \sim 0.2$  vs. PGS's  $r \sim 0.21$  in our study and  $R^2 < 0.1$  in another study [Allegrini et al., 2019]), suggesting the potential use of brain-based predictive models for a robust, transdiagnostic, brain-based marker for cognitive abilities.

With opportunistic stacking, those who wish to adapt our brain-based predictive models to compute a transdiagnostic brain-based marker for cognition in their own data, but do not have as many modalities as the ABCD, can still use our models. That is, they can still use the model built from the ABCD and impute missing values of certain modalities to fits with their study. Accordingly, our use of opportunistic stacking provides a scalable and flexible approach for future researchers following the RDoC framework (Morris & Cuthbert, 2012).

Our study is not without limitations. We relied on the ABCD study's curated, preprocessed data (Casey et al., 2018; Hagler et al., 2019; Yang & Jernigan, n.d.). This provided certain advantages. For instance, given that the curated data provided by the ABCD have already been preprocessed, other studies that wish to apply our model of the *g*-factor to the ABCD data can readily do so without concerns about differences in preprocessing steps. Preprocessed data also enabled us to apply the manual quality control done by the study, a process that required time and well-trained labour (Casey et al., 2018; Hagler et al., 2019; Yang & Jernigan, n.d.). Preprocessing large-scale multi-modal data ourselves would not only demand significant computer power and time but is prone to error. However, using the preprocessed data only allowed us to follow the choices of processing done by the study. For example, ABCD Release 3 only provided Freesurfer's parcellation (Destrieux et al., 2010; Fischl et al., 2002) for task-based fMRI. While this popular method allowed us to explain task-based activity on subject-specific anatomical landmarks, the regions are relatively large compared to other parcellations. Future studies will need to examine if smaller and/or different parcellations would improve predictive performance. Next, our predictive modelling framework was designed to predict the out-of-sample *g*-factor, but not the developmental changes in the *g*-factor, from multimodal MRI. More specifically, we standardised MRI and cognitive data within each age group to satisfy the assumption of our machine-learning algorithms (Zou & Hastie, 2005) and to force behavioural performance from different cognitive tasks onto the same scale. This unfortunately made our predictive models inappropriate for predicting the developmental changes in cognition over years (Moeller, 2015). Future research that aims to capture the developmental changes in cognition would need to employ a different strategy for standardisation (Moeller, 2015).

In terms of important disclaimers, research reporting on cognitive abilities can be misunderstood or misquoted for alien purposes (Suzuki & Aronson, 2005). It is therefore important to clarify the following. First, the fact that measurements taken from the brain were related to cognitive abilities should not be equated with assertions that variability in cognitive abilities is 'purely biological'. Here, we showed that the predictive model for the *g*-factor based on socio-demographic and psychological variables that were available in the



ABCD (Kirlic et al., 2021) already accounted for a larger variance of the  $g$ -factor ( $\sim 30\%$ ) than the predictive models based on the brain ( $\sim 20\%$ ) or genes ( $< 10\%$  [Allegrini et al., 2019]). Moreover, our mediation analysis showed that the brain-based predictive models could only account for approximately 19% of the relationship between cognitive abilities and socio-demographic and psychological factors. Accordingly, it is very plausible that social-demographic and psychological circumstances, broadly construed, have at least partial aetiological primacy. Second, it should be clear that social-demographic, psychological and genetic circumstances may not be independent of one another, as suggested by studies on the complex interplay of genes and environments on cognitive abilities over the course of cognitive development (Tucker-Drob et al., 2013; Tucker-Drob & Briley, 2014). This is shown in our mediation analyses. Here, the brain-based  $g$ -factor showed less proportion mediated for the influences of social-demographic, psychological factors and genes when they were included together in the model, compared to when they were included in separate models. This suggests the interdependency among the brain, genes, social-demographic and psychological factors as proposed by the RDoC (Insel et al., 2010; NIMH, n.d.-a). Third, under no circumstances should the results of this article be interpreted as entailing a value judgement about how people vary in measurements of cognitive abilities. Indeed, it is important to reflect on the fact that the way we measured cognitive abilities, for example, through the  $g$ -factor here, reflects norms that are entrenched in cultures and societies of a certain time in history, rather than reflecting some universal truth or a supra-historical marker of cognitive abilities (Flynn, 2009). The value of the  $g$ -factor here is as a marker (present in early life) of a series of other important life outcomes in current societal circumstances.

## 5 | CONCLUSION

In conclusion, we developed brain-based stacked, predictive models for children's cognitive abilities that were longitudinally stable, generalisable and robust against missingness. More importantly, our brain-based models were able to partially mediate the relationships of childhood cognitive abilities with the socio-demographic, psychological and genetic factors. Accordingly, our approach should pave the way for future researchers to employ multimodal MRI as a useful research tool for integrative, RDoC-inspired research in cognition and mental health.

## ACKNOWLEDGEMENTS

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multi-site, 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 is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041022, U01DA041028, U01DA041048, U01DA041089, U01DA041106,

U01DA041117, U01DA041120, U01DA041134, U01DA041148, U01DA041156, U01DA041174, U24DA041123, U24DA041147, U01DA041093 and U01DA041025. 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/scientists/workgroups/>. 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. We thank the developers of several R libraries, including semTools (Sunthud Pornprasertmanit), eNetXplorer (Julián Candia) and ggseg (Athanasia M. Mowinckel), for their technical advice. Narun Pat and Yue Wang were supported by Health Research Council Funding (21/618) and by University of Otago.

## CONFLICT OF INTEREST

The authors declare no conflict of interests.

## DATA AVAILABILITY STATEMENT

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multi-site, 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 is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041022, U01DA041028, U01DA041048, U01DA041089, U01DA041106, U01DA041117, U01DA041120, U01DA041134, U01DA041148, U01DA041156, U01DA041174, U24DA041123, U24DA041147, U01DA041093, and U01DA041025. 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/scientists/workgroups/>. 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.

## ORCID

Narun Pat  <https://orcid.org/0000-0003-1459-5255>

Richard Anney  <https://orcid.org/0000-0002-6083-407X>

Lucy Riglin  <https://orcid.org/0000-0002-5124-5230>

Anita Thapar  <https://orcid.org/0000-0002-3689-737X>

Argyris Stringaris  <https://orcid.org/0000-0002-6264-8377>

## REFERENCES

- 1000 Genomes Project Consortium, Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P., Kang, H. M., Korbel, J. O., Marchini, J. L., McCarthy, S., McVean, G., & Abecasis, G. R. (2015). A global reference for human genetic variation. *Nature*, 526, 68–74.
- Achenbach, T. M., Ivanova, M. Y., & Rescorla, L. A. (2017). Empirically based assessment and taxonomy of psychopathology for ages 1½–90+ years: Developmental, multi-informant, and multicultural findings. *Comprehensive Psychiatry*, 79, 4–18.

- Adolescent and School Health | CDC. (2020). Youth risk behavior surveillance system (YRBSS) and data. <https://www.cdc.gov/healthyouth/data/yrbs/index.htm>
- Alexander, A. L., Lee, J. E., Lazar, M., & Field, A. S. (2007). Diffusion tensor imaging of the brain. *Neurotherapeutics*, 4, 316–329.
- Allegrini, A. G., Selzam, S., Rimfeld, K., von Stumm, S., Pingault, J. B., & Plomin, R. (2019). Genomic prediction of cognitive traits in childhood and adolescence. *Molecular Psychiatry*, 24, 819–827.
- Ang, Y.-S., Frontero, N., Belleau, E., & Pizzagalli, D. A. (2020). Disentangling vulnerability, state and trait features of neurocognitive impairments in depression. *Brain*, 143, 3865–3877. <https://doi.org/10.1093/brain/awaa314>
- Bagot, K. S., Matthews, S. A., Mason, M., Squeglia, L. M., Fowler, J., Gray, K., Herting, M., May, A., Colrain, I., Godino, J., Tapert, S., Brown, S., & Patrick, K. (2018). Current, future and potential use of mobile and wearable technologies and social media data in the ABCD study to increase understanding of contributors to child health. *Developmental Cognitive Neuroscience*, 32, 121–129.
- Barch, D. M., Burgess, G. C., Harms, M. P., Petersen, S. E., Schlaggar, B. L., Corbetta, M., Glasser, M. F., Curtiss, S., Dixit, S., Feldt, C., Nolan, D., Bryant, E., Hartley, T., Footer, O., Bjork, J. M., Poldrack, R., Smith, S., Johansen-Berg, H., Snyder, A. Z., ... WU-Minn HCP Consortium. (2013). Function in the human connectome: Task-fMRI and individual differences in behavior. *NeuroImage*, 80, 169–189.
- Bauer, P. J., Dikmen, S. S., Heaton, R. K., Mungas, D., Slotkin, J., & Beaumont, J. L. (2013). III. NIH toolbox cognition battery (CB): Measuring episodic memory. *Monographs of the Society for Research in Child Development*, 78, 34–48.
- Biederman, J., Monuteaux, M. C., Doyle, A. E., Seidman, L. J., Wilens, T. E., Ferrero, F., Morgan, C. L., & Faraone, S. V. (2004). Impact of executive function deficits and attention-deficit/hyperactivity disorder (ADHD) on academic outcomes in children. *Journal of Consulting and Clinical Psychology*, 72, 757–766.
- Bissett, P. G., Hagen, M. P., & Poldrack, R. A. (2020). Design issues and solutions for stop-signal data from the Adolescent Brain Cognitive Development (ABCD) study. *eLife*, 10, e60185.
- Bleck, T. P., Nowinski, C. J., Gershon, R., & Koroshetz, W. J. (2013). What is the NIH toolbox, and what will it mean to neurology? *Neurology*, 80, 874–875.
- Bogdan, R., Baranger, D. A. A., & Agrawal, A. (2018). Polygenic risk scores in clinical psychology: Bridging genomic risk to individual differences. *Annual Review of Clinical Psychology*, 14, 119–157.
- Bolt, T., Nomi, J. S., Yeo, B. T. T., & Uddin, L. Q. (2017). Data-driven extraction of a nested model of human brain function. *The Journal of Neuroscience*, 37, 7263–7277.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32.
- Bruni, O., Ottaviano, S., Guidetti, V., Romoli, M., Innocenzi, M., Cortesi, F., & Giannotti, F. (1996). The sleep disturbance scale for children (SDSC) construction and validation of an instrument to evaluate sleep disturbances in childhood and adolescence. *Journal of Sleep Research*, 5, 251–261.
- Candia, J., & Tsang, J. S. (2019a). eNetXplorer: An R package for the quantitative exploration of elastic net families for generalized linear models. *BMC Bioinformatics*, 20, 189.
- Candia, J., & Tsang, J. S. (2019b). eNetXplorer: An R package for the quantitative exploration of elastic net families for generalized linear models. *BMC Bioinformatics*, 20, 1–11.
- Carlozzi, N. E., Tulskey, D. S., Kail, R. V., & Beaumont, J. L. (2013). NIH toolbox cognition battery (CB): Measuring processing speed. *Monographs of the Society for Research in Child Development*, 78, 88–102.
- Carver, C. S., & White, T. L. (1994). Behavioral inhibition, behavioral activation, and affective responses to impending reward and punishment: The BIS/BAS scales. *Journal of Personality and Social Psychology*, 67, 319–333.
- Casey, B. J., Cannonier, T., Conley, M. I., 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., ... ABCD Imaging Acquisition Workgroup. (2018). The adolescent brain cognitive development (ABCD) study: Imaging acquisition across 21 sites. *Developmental Cognitive Neuroscience*, 32, 43–54.
- Chilcoat, H. D., & Anthony, J. C. (1996). Impact of parent monitoring on initiation of drug use through late childhood. *Journal of the American Academy of Child & Adolescent Psychiatry*, 35, 91–100.
- 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. *Developmental Cognitive Neuroscience*, 32, 143–154.
- 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, 179–194.
- Daniel, M. H. & Wahlstrom, D. Equivalence of Q-interactive™ and paper administrations of cognitive tasks: WISC®-V. 13 (2014).
- Davies, G., Tenesa, A., Payton, A., Yang, J., Harris, S. E., Liewald, D., Ke, X., le Hellard, S., Christoforou, A., Luciano, M., McGhee, K., Lopez, L., Gow, A. J., Corley, J., Redmond, P., Fox, H. C., Haggarty, P., Whalley, L. J., McNeill, G., ... Deary, I. J. (2011). Genome-wide association studies establish that human intelligence is highly heritable and polygenic. *Molecular Psychiatry*, 16, 996–1005.
- Debeer, D., & Strobl, C. (2020). Conditional permutation importance revisited. *BMC Bioinformatics*, 21, 307.
- Destrieux, C., Fischl, B., Dale, A., & Halgren, E. (2010). Automatic parcellation of human cortical gyri and sulci using standard anatomical nomenclature. *NeuroImage*, 53, 1–15.
- Dick, A. S., Garcia, N. L., Pruden, S. M., Thompson, W. K., Hawes, S. W., Sutherland, M. T., Riedel, M. C., Laird, A. R., & Gonzalez, R. (2019). No evidence for a bilingual executive function advantage in the ABCD study. *Nature Human Behaviour*, 3, 692–701.
- Dubois, J., Galdi, P., Paul, L. K., & Adolphs, R. (2018). A distributed brain network predicts general intelligence from resting-state human neuroimaging data. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 373, 20170284.
- Duncan, L., Shen, H., Gelaye, B., Meijssen, J., Ressler, K., Feldman, M., Peterson, R., & Domingue, B. (2019). Analysis of polygenic risk score usage and performance in diverse human populations. *Nature Communications*, 10, 3328.
- Echeverria, S. E., Diez-Roux, A. V., & Link, B. G. (2004). Reliability of self-reported neighborhood characteristics. *Journal of Urban Health*, 81, 682–701.
- Engemann, D. A., Kozynets, O., Sabbagh, D., Lemaître, G., Varoquaux, G., Liem, F., & Gramfort, A. (2020). Combining magnetoencephalography with magnetic resonance imaging enhances learning of surrogate-biomarkers. *eLife*, 9, e54055.
- Epskamp, S. (2015). semPlot: Unified visualizations of structural equation models. *Structural Equation Modeling*, 22, 474–483.
- Eriksen, B. A., & Eriksen, C. W. (1974). Effects of noise letters upon the identification of a target letter in a nonsearch task. *Perception & Psychophysics*, 16, 143–149.
- Farah, M. J., Shera, D. M., Savage, J. H., Betancourt, L., Giannetta, J. M., Brodsky, N. L., Malmud, E. K., & Hurt, H. (2006). Childhood poverty: Specific associations with neurocognitive development. *Brain Research*, 1110, 166–174.

- Fassbender, C., Mukherjee, P., & Schweitzer, J. B. (2017). Minimizing noise in pediatric task-based functional MRI: Adolescents with developmental disabilities and typical development. *NeuroImage*, *149*, 338–347.
- Finn, E. S. (2021). Is it time to put rest to rest? *Trends in Cognitive Sciences*, *25*, 1021–1032.
- 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. *Neuron*, *33*, 341–355.
- Flynn, J. R. (2009). *What is intelligence? Beyond the Flynn effect*. Cambridge University Press.
- Fortin, J.-P., Parker, D., Tunç, B., Watanabe, T., Elliott, M. A., Ruparel, K., Roalf, D. R., Satterthwaite, T. D., Gur, R. C., Gur, R. E., Schultz, R. T., Verma, R., & Shinohara, R. T. (2017). Harmonization of multi-site diffusion tensor imaging data. *NeuroImage*, *161*, 149–170.
- Friedman, J. H., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, *33*, 1–22.
- Frostenson, S., & Kliff, S. on April 6, 2016, 8:50 a.m. ET. Where is the lead exposure risk in your community? *Vox.com* <http://www.vox.com/a/lead-exposure-risk-map>
- Garavan, H., Bartsch, H., Conway, K., Decastro, A., Goldstein, R. Z., Heeringa, S., Jernigan, T., Potter, A., Thompson, W., & Zabs, D. (2018). Recruiting the ABCD sample: Design considerations and procedures. *Developmental Cognitive Neuroscience*, *32*, 16–22.
- Garavan, H., Hahn, S., Charani, B., Juliano, A., Allgaier, N., Yuan, D. K., Weigard, A., Orr, C., Watts, R., Wager, T. D., de Leon, O. R., Hagler, D. J., Jr., & Potter, A. (2020). The ABCD stop signal data: Response to Bissett et al. *bioRxiv*, 2020.07.27.223057. <https://doi.org/10.1101/2020.07.27.223057>
- Genç, E., Fraenz, C., Schlüter, C., Friedrich, P., Hossiep, R., Voelkle, M. C., Ling, J. M., Güntürkün, O., & Jung, R. E. (2018). Diffusion markers of dendritic density and arborization in gray matter predict differences in intelligence. *Nature Communications*, *9*, 1905.
- Gershon, R. C., Cook, K. F., Mungas, D., Manly, J. J., Slotkin, J., Beaumont, J. L., & Weintraub, S. (2014). Language measures of the NIH toolbox cognition battery. *Journal of the International Neuropsychological Society*, *20*, 642–651.
- Góngora, D., Vega-Hernández, M., Jahansahi, M., Valdés-Sosa, P. A., & Bringas-Vega, M. L. (2020). Crystallized and fluid intelligence are predicted by microstructure of specific white-matter tracts. *Human Brain Mapping*, *41*, 906–916.
- Goodall, J., Fisher, C., Hetrick, S., Phillips, L., Parrish, E. M., & Allott, K. (2018). Neurocognitive functioning in depressed young people: A systematic review and meta-analysis. *Neuropsychology Review*, *28*, 216–231.
- Goodman, R., Meltzer, H., & Bailey, V. (2003). The strengths and difficulties questionnaire: A pilot study on the validity of the self-report version. *International Review of Psychiatry*, *15*, 173–177.
- 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. *Cerebral Cortex*, *26*, 288–303.
- Gray, J. R., Chabris, C. F., & Braver, T. S. (2003). Neural mechanisms of general fluid intelligence. *Nature Neuroscience*, *6*, 316–322.
- Hagler, D. J., Hatton, S. N., 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., Sicut, C. S., Kuperman, J., Bartsch, H., Xue, F., ... Dale, A. M. (2019). Image processing and analysis methods for the adolescent brain cognitive development study. *NeuroImage*, *202*, 116091.
- Hagler, D. J., Jr., Ahmadi, M. E., Kuperman, J., Holland, D., McDonald, C., Halgren, E., & Dale, A. M. (2009). Automated white-matter tractography using a probabilistic diffusion tensor atlas: Application to temporal lobe epilepsy. *Human Brain Mapping*, *30*, 1535–1547.
- Insel, T., Cuthbert, B., Garvey, M., Heinssen, R., Pine, D. S., Quinn, K., Sanislow, C., & Wang, P. (2010). Research domain criteria (RDoC): Toward a new classification framework for research on mental disorders. *Am J Psychiatry*, *167*, 748–751.
- Jiang, R., Calhoun, V. D., Cui, Y., Qi, S., Zhuo, C., Li, J., Jung, R., Yang, J., du, Y., Jiang, T., & Sui, J. (2020). Multimodal data revealed different neurobiological correlates of intelligence between males and females. *Brain Imaging and Behavior*, *14*, 1979–1993.
- Jorgensen, T. D., S Pornprasertmanit, AM. Schoemann, Y Rosseel, P Miller, C Quick, M Garnier-Villarreal, J Selig, A Boulton, K Preacher, D Coffman, M Rhemtulla, A Robitzsch, C Enders, R Arslan, B Clinton, P Panko, E Merkle, S Chesnut, J Byrnes, ... AR. Johnson (2018). semTools: Useful tools for structural equation modeling. *R package version 0.5-1*.
- Josse, J., Prost, N., Scornet, E., & Varoquaux, G. (2020). On the consistency of supervised learning with missing values. *arXiv:1902.06931 [cs, math, stat]*.
- Kessler, R. C., Avenevoli, S., Costello, E. J., Green, J. G., Gruber, M. J., Heeringa, S., Merikangas, K. R., Pennell, B. E., Sampson, N. A., & Zaslavsky, A. M. (2009). Design and field procedures in the US National Comorbidity Survey Replication Adolescent Supplement (NCS-A). *International Journal of Methods in Psychiatric Research*, *18*, 69–83.
- Kind, A. J. H., Jencks, S., Brock, J., Yu, M., Bartels, C., Ehlenbach, W., Greenberg, C., & Smith, M. (2014). Neighborhood socioeconomic disadvantage and 30-day rehospitalization. *Annals of Internal Medicine*, *161*, 765–774.
- Kirlic, N., Colaizzi, J. M., Cosgrove, K. T., Cohen, Z. P., Yeh, H. W., Breslin, F., Morris, A. S., Uppercle, R. L., Singh, M. K., & Paulus, M. P. (2021). Extracurricular activities, screen media activity, and sleep May be modifiable factors related to Children's cognitive functioning: Evidence from the ABCD study®. *Child Development*, *92*, 2035–2052.
- Knutson, B., Westdorp, A., Kaiser, E., & Hommer, D. (2000). fMRI visualization of brain activity during a monetary incentive delay task. *NeuroImage*, *12*, 20–27.
- Lee, J. J., Wedow, R., Okbay, A., Kong, E., Maghziyan, O., Zacher, M., Nguyen-Viet, T. A., Bowers, P., Sidorenko, J., Karlsson Linnér, R., Fontana, M. A., Kundu, T., Lee, C., Li, H., Li, R., Royer, R., Timshel, P. N., Walters, R. K., Willoughby, E. A., ... Cesarini, D. (2018). Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nature Genetics*, *50*, 1112–1121.
- Luciana, M., Bjork, J. M., Nagel, B. J., Barch, D. M., Gonzalez, R., Nixon, S. J., & Banich, M. T. (2018). Adolescent neurocognitive development and impacts of substance use: Overview of the adolescent brain cognitive development (ABCD) baseline neurocognition battery. *Developmental Cognitive Neuroscience*, *32*, 67–79.
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc..
- MacKinnon, D. P., Fairchild, A. J., & Fritz, M. S. (2007). Mediation analysis. *Annual Review of Psychology*, *58*, 593–614.
- McDaniel, M. A. (2005). Big-brained people are smarter: A meta-analysis of the relationship between in vivo brain volume and intelligence. *Intelligence*, *33*, 337–346.
- Merikangas, K. R., Avenevoli, S., Costello, E. J., Koretz, D., & Kessler, R. C. (2009). National Comorbidity Survey Replication Adolescent Supplement (NCS-A): I. Background and measures. *Journal of the American Academy of Child & Adolescent Psychiatry*, *48*, 367–379.
- Mihalik, A., Brudfors, M., Robu, M., Ferreira, F. S., Lin, H., Rau, A., Wu, T., Blumberg, S. B., Kanber, B., Tariq, M., Garcia, M. E., Zor, C., Nikitichev, D. I., Mourão-Miranda, J., & Oxtoby, N. P. (2019). ABCD neurocognitive prediction challenge 2019: Predicting individual fluid intelligence scores from structural MRI using probabilistic

- segmentation and kernel ridge regression. In K. M. Pohl, W. K. Thompson, E. Adeli, & M. G. Linguraru (Eds.), *Adolescent brain cognitive development neurocognitive prediction* (pp. 133–142). Springer. [https://doi.org/10.1007/978-3-030-31901-4\\_16](https://doi.org/10.1007/978-3-030-31901-4_16)
- Moeller, J. (2015). A word on standardization in longitudinal studies: Don't. *Frontiers in Psychology*, *6*, 1389.
- Molnar, C. (2019). *Interpretable machine learning: A guide for making black box models explainable*. Leanpub.
- Moos, R. H., & Humphrey, B. (1974). *Preliminary manual for family environment scale, work environment scale, group environment scale*. Consulting Psychologists Press.
- Morris, S. E., & Cuthbert, B. N. (2012). Research domain criteria: Cognitive systems, neural circuits, and dimensions of behavior. *Dialogues in Clinical Neuroscience*, *14*, 29–37.
- Mowinckel, A. M., & Vidal-Piñeiro, D. (2020). Visualization of brain statistics with R packages *ggseg* and *ggseg3d*. *Advances in Methods and Practices in Psychological Science*, *3*, 466–483.
- Narr, K. L., Woods, R. P., Thompson, P. M., Szeszko, P., Robinson, D., Dimtcheva, T., Gurbani, M., Toga, A. W., & Bilder, R. M. (2007). Relationships between IQ and regional cortical Gray matter thickness in healthy adults. *Cerebral Cortex*, *17*, 2163–2171.
- National Institute of Mental Health (NIMH). (n.d.-a). About RDoC <https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/about-rdoc>
- Nielson, D. M., Pereira, F., Zheng, C. Y., Migineishvili, N., Lee, J. A., Thomas, A. G., & Bandettini, P. A. (2018). Detecting and harmonizing scanner differences in the ABCD study—Annual release 1.0. *bioRxiv*, 309260. <https://doi.org/10.1101/309260>
- NIMH. (n.d.-b). Developmental and Environmental Aspects <https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/developmental-and-environmental-aspects>
- NIMH. (n.d.-c). Units of Analysis <https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/units/units-of-analysis>.
- Niu, X., Zhang, F., Kounios, J., & Liang, H. (2020). Improved prediction of brain age using multimodal neuroimaging data. *Human Brain Mapping*, *41*, 1626–1643.
- Pamplona, G. S. P., Santos Neto, G. S., Rosset, S. R. E., Rogers, B. P., & Salmon, C. E. G. (2015). Analyzing the association between functional connectivity of the brain and intellectual performance. *Frontiers in Human Neuroscience*, *9*, 61.
- Pat, N., Riglin, L., Anney, R., Wang, Y., Barch, D. M., Thapar, A., & Stringaris, A. (2021) Epub September 7, 2021). Motivation and cognitive abilities as mediators between polygenic scores and psychopathology in children. *Journal of the American Academy of Child & Adolescent Psychiatry*, *61*, 782–795.e3. <https://doi.org/10.1016/j.jaac.2021.08.019>
- Pietschnig, J., Penke, L., Wicherts, J. M., Zeiler, M., & Voracek, M. (2015). Meta-analysis of associations between human brain volume and intelligence differences: How strong are they and what do they mean? *Neuroscience & Biobehavioral Reviews*, *57*, 411–432.
- Plomin, R., & Deary, I. J. (2015). Genetics and intelligence differences: Five special findings. *Molecular Psychiatry*, *20*, 98–108.
- Rasero, J., Sentis, A. I., Yeh, F.-C., & Verstynen, T. (2021). Integrating across neuroimaging modalities boosts prediction accuracy of cognitive ability. *PLoS Computational Biology*, *17*, e1008347.
- Rosseel, Y. (2012). Lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, *48*, 1–36.
- Roth, A. E. (Ed.). (1988). *The Shapley value: Essays in honor of Lloyd S. Shapley*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511528446>
- Schaefer, E. S. (1965). A configurational analysis of children's reports of parent behavior. *Journal of Consulting Psychology*, *29*, 552–557.
- Sheffield, J. M., Karcher, N. R., & Barch, D. M. (2018). Cognitive deficits in psychotic disorders: A lifespan perspective. *Neuropsychology Review*, *28*, 509–533.
- Shilyansky, C., Williams, L. M., Gyurak, A., Harris, A., Usherwood, T., & Etkin, A. (2016). Effect of antidepressant treatment on cognitive impairments associated with depression: A randomised longitudinal study. *The Lancet Psychiatry*, *3*, 425–435.
- Snellen, H. (1862). *Letterproeven Tot Bepaling Der Gezichtscherpte* (Dutch edition). Utrecht, the Netherlands: Weyers. Also published in many languages as *Optotypi ad Visum Determinandum*.
- Sripada, C., Angstadt, M., Rutherford, S., Taxali, A., & Shedden, K. (2020). Toward a “treadmill test” for cognition: Improved prediction of general cognitive ability from the task activated brain. *Human Brain Mapping*, *41*, 3186–3197.
- 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. *Molecular Psychiatry*, *25*, 3413–3421.
- Stover, P. J., Harlan, W. R., Hammond, J. A., Hendershot, T., & Hamilton, C. M. (2010). PhenX: A toolkit for interdisciplinary genetics research. *Current Opinion in Lipidology*, *21*, 136–140.
- Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable importance for random forests. *BMC Bioinformatics*, *9*, 307.
- Strobl, C., Boulesteix, A.-L., Zeileis, A., & Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*, *8*, 25.
- Sui, J., Jiang, R., Bustillo, J., & Calhoun, V. (2020). Neuroimaging-based individualized prediction of cognition and behavior for mental disorders and health: Methods and promises. *Biological Psychiatry*, *88*, 818–828.
- Suzuki, L., & Aronson, J. (2005). The cultural malleability of intelligence and its impact on the racial/ethnic hierarchy. *Psychology, Public Policy, and Law*, *11*, 320–327.
- Thaler, N. S., Bello, D. T., & Etkoff, L. M. (2013). WISC-IV profiles are associated with differences in symptomatology and outcome in children with ADHD. *Journal of Attention Disorders*, *17*, 291–301.
- Thompson, W. K., Barch, D. M., Bjork, J. M., Gonzalez, R., Nagel, B. J., Nixon, S. J., & Luciana, M. (2019). The structure of cognition in 9 and 10 year-old children and associations with problem behaviors: Findings from the ABCD study's baseline neurocognitive battery. *Developmental Cognitive Neuroscience*, *36*, 100606.
- Tucker-Drob, E. M. (2009). Differentiation of cognitive abilities across the life span. *Developmental Psychology*, *45*, 1097–1118.
- Tucker-Drob, E. M., & Briley, D. A. (2014). Continuity of genetic and environmental influences on cognition across the life span: A meta-analysis of longitudinal twin and adoption studies. *Psychological Bulletin*, *140*, 949–979.
- Tucker-Drob, E. M., Briley, D. A., & Harden, K. P. (2013). Genetic and environmental influences on cognition across development and context. *Current Directions in Psychological Science*, *22*, 349–355.
- Uban, K. A., Horton, M. K., Jacobus, J., Heyser, C., Thompson, W. K., Tapert, S. F., Madden, P. A. F., Sowell, E. R., & Adolescent Brain Cognitive Development Study. (2018). Biospecimens and the ABCD study: Rationale, methods of collection, measurement and early data. *Developmental Cognitive Neuroscience*, *32*, 97–106.
- United States Department of Justice. Office of Justice Programs. Federal Bureau of Investigation. (2012). Uniform crime reporting program data: County-level detailed arrest and offense data, 2010: Version 2. <https://doi.org/10.3886/ICPSR33523.V2>
- Waiter, G. D., Deary, I. J., Staff, R. T., Murray, A. D., Fox, H. C., Starr, J. M., & Whalley, L. J. (2009). Exploring possible neural mechanisms of intelligence differences using processing speed and working memory tasks: An fMRI study. *Intelligence*, *37*, 199–206.
- Whelan, R., Conrod, P. J., Poline, J. B., Lourdasamy, A., Banaschewski, T., Barker, G. J., Bellgrove, M. A., Büchel, C., Byrne, M., Cummins, T. D., Fauth-Bühler, M., Flor, H., Gallinat, J., Heinz, A., Ittermann, B.,



- Mann, K., Martino, J. L., Lalor, E. C., Lathrop, M., ... IMAGEN Consortium. (2012). Adolescent impulsivity phenotypes characterized by distinct brain networks. *Nature Neuroscience*, *15*, 920–925.
- Wolpert, D. H. (1992). Stacked generalization. *Neural Networks*, *5*, 241–259.
- Wright, M. N., & Ziegler, A. (2017). Ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, *77*, 1–17.
- Yang, R. & Jernigan, T. L. (n.d.) Adolescent Brain Cognitive Development Study (ABCD)—Annual Release 3.0. <https://doi.org/10.15154/1519007>
- Youngstrom, E. A., Frazier, T. W., Demeter, C., Calabrese, J. R., & Findling, R. L. (2008). Developing a 10-item mania scale from the parent general behavior inventory for children and adolescents. *The Journal of Clinical Psychiatry*, *69*, 831–839.
- Zapolski, T. C. B., Stairs, A. M., Settles, R. F., Combs, J. L., & Smith, G. T. (2010). The measurement of dispositions to rash action in children. *Assessment*, *17*, 116–125.
- Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, *67*, 301–320.
- Zucker, R. A., Gonzalez, R., Feldstein Ewing, S. W., Paulus, M. P., Arroyo, J., Fuligni, A., Morris, A. S., Sanchez, M., & Wills, T. (2018). Assessment of culture and environment in the adolescent brain and cognitive development study: Rationale, description of measures, and early data. *Developmental Cognitive Neuroscience*, *32*, 107–120.

#### SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Pat, N., Wang, Y., Anney, R., Riglin, L., Thapar, A., & Stringaris, A. (2022). Longitudinally stable, brain-based predictive models mediate the relationships between childhood cognition and socio-demographic, psychological and genetic factors. *Human Brain Mapping*, *43*(18), 5520–5542. <https://doi.org/10.1002/hbm.26027>