# An open, analysis-ready, and quality controlled resource for pediatric brain white-matter research

**Adam Richie-Halford[1][†][*], Matthew Cieslak[2,3][†][*], Lei Ai[4], Sendy Caffarra[1,5,6], Sydney Covitz[2,3], Alexandre R. Franco[4,7], Iliana I. Karipidis[5,8,9], John Kruper[10], Michael Milham[4,7], Bárbara Avelar-Pereira[8], Ethan Roy[5], Valerie J. Sydnor[2,3], Jason Yeatman[1,5], The Fibr Community Science Consortium[12], Theodore D. Satterthwaite[2,3][‡], Ariel Rokem[10,11][‡]**

**\*For correspondence:**
adamrh@stanford.edu (ARH);
matthew.cieslak@pennmedicine.upenn.edu (MC)

[†]These authors contributed equally to this work
[‡]These authors also contributed equally to this work

[1]Stanford University, Division of Developmental and Behavioral Pediatrics, Stanford, California, 94305, USA; [2]University of Pennsylvania, Department of Psychiatry, Philadelphia, Pennsylvania, 19104, USA; [3]University of Pennsylvania, Lifespan Informatics and Neuroimaging Center, Philadelphia, Pennsylvania, 19104, USA; [4]Child Mind Institute, Center for the Developing Brain, New York City, New York, 10022, USA; [5]Stanford University, Graduate School of Education, Stanford, California, 94305, USA; [6]University of Modena and Reggio Emilia, Department of Biomedical, Metabolic and Neural Sciences, 41125 Modena, Italy; [7]Nathan Kline Institute for Psychiatric Research, Center for Biomedical Imaging and Neuromodulation, Orangeburg, New York, 10962, USA; [8]Stanford University, Department of Psychiatry and Behavioral Sciences, School of Medicine, Stanford, California, 94305, USA; [9]University of Zurich, Department of Child and Adolescent Psychiatry and Psychotherapy, University Hospital of Psychiatry Zurich, Zurich, 8032, Switzerland; [10]University of Washington, Department of Psychology, Seattle, Washington, 98195, USA; [11]University of Washington, eScience Institute, Seattle, Washington, 98195, USA; [12]The Fibr Community Science Consortium

## Abstract

We created resources to facilitate research on the role of human brain microstructure in the development of mental health disorders, based on openly-available diffusion MRI (dMRI) data from the Healthy Brain Network (HBN) study. First, we curated the HBN dMRI data (N=2747) into the Brain Imaging Data Structure and preprocessed it according to best-practices, including denoising and correcting for motion effects, susceptibility-related distortions, and eddy currents. Preprocessed, analysis-ready data was made openly available. Data quality plays a key role in the analysis of dMRI, and we provide automated quality control (QC) scores for every scan, as part of the data release. To scale QC to this large dataset, we trained a neural network through the combination of a small data subset scored by experts and a larger set scored by community scientists. The network performs QC highly concordant with that of experts on a held out set (ROC-AUC = 0.947). A further analysis of the neural network demonstrates that it relies on image features with relevance to QC. Altogether, this work both delivers a resource for transdiagnostic research in brain connectivity and pediatric mental health and serves as a novel tool for automated QC of large datasets.

## Introduction

Childhood and adolescence are characterized by rapid dynamic change to human brain structure and function (*Lebel and Deoni, 2018*). This period of development is also a time during which the symptoms of many mental health disorders emerge (*Paus et al., 2008*). Understanding how individual differences in brain development relate to the onset and progression of psychopathology inevitably requires large datasets (*Paus, 2010*; *Fair et al., 2021*). The Healthy Brain Network (HBN) is a landmark pediatric mental health study that is designed to eventually include MRI images along with detailed clinical and cognitive phentoyping from over 5000 New York City area children and adolescents (*Alexander et al., 2017*). The HBN dataset takes a trans-diagnostic approach and provides a broad range of phenotypic and brain imaging data for each individual. One of the brain imaging measurements acquired is diffusion MRI (dMRI), which is the dominant technology for inferring the physical properties of white matter (*Wandell, 2016*). The dMRI data is openly available in its raw form through the Functional Connectomes Project and the International Neuroimaging Data-Sharing Initiative (FCP-INDI), spurring collaboration on open and reproducible science (*Mennes et al., 2013*). However, this raw, publicly available data requires extensive processing and quality assurance before it can be fruitfully analyzed.

The analysis of a large, multi-site dMRI dataset must take into account the inevitable variability in scanning parameters across scanning sessions. Critical preprocessing steps, such as susceptibility distortion correction (*Jones and Cercignani, 2010*) require additional MRI acquisitions besides dMRI and accurate metadata accompanying each image. A session missing an acquisition or important metadata can either be processed to the extent its available data allows or excluded entirely. In addition, the quality of preprocessed data is heavily affected by differences in acquisition parameters (*Yeh et al., 2019*) and by differences in preprocessing steps. Here we address these problems by meticulously curating the HBN data according to the Brain Imaging Data Specification (BIDS) (*Gorgolewski et al., 2016*) and processing the data using the *QSIPrep* (*Cieslak et al., 2021*) BIDS App (*Gorgolewski et al., 2017*). *QSIPrep* automatically builds and executes benchmarked workflows that adhere to best practices in the field given the available BIDS data. The results include automated data quality metrics, visual reports and a description of the processing steps automatically chosen to process each session.

This preprocessing requires a costly compute infrastructure and is both time-consuming and error-prone. Requiring researchers to process dMRI data on their own introduces both a practical barrier to access and an extra source of heterogeneity into the data, devaluing its scientific utility. We provide the preprocessed data as a transparent and open resource, thereby reducing barriers to data access and allowing researchers to spend more of their time answering questions in brain development and psychopathology rather than recapitulating preprocessing.

In addition to requiring extensive preprocessing, dMRI data must be thoroughly checked for quality. dMRI measurements are susceptible to a variety of artifacts that affect the quality of the signals and the ability to make accurate inferences from them. In small studies, with few participants, it is common to thoroughly examine the data from every participant as part of a quality control (QC) process. However, expert examination is time consuming and is prohibitive in large datasets such as HBN. This difficulty could be ameliorated through the automation of QC. Given their success in other visual recognition tasks, machine learning and computer vision methods, such as convolutional deep artificial neural networks or "deep learning" (*LeCun et al., 2015*), are promising avenues for automation of QC. However, one of the challenges of these new methods is that they require a large training dataset to attain accurate performance. In previous work, we demonstrated that deep learning can accurately emulate expert QC of T1-weighted (T1w) anatomical brain images (*Keshavan et al., 2019*). To obtain a large enough training dataset of T1w images in our prior study, we deployed a community science tool [1] that collected quality control scores of

---

[1] While the term "citizen science" evokes a sense of civic duty in scientific engagement, it can also imply a barrier for community members who want to contribute to science but may not be citizens of a particular country. In this manuscript we use the more modern term "community science."

parts of the dataset from volunteers through a web application. The scores were then calibrated using a gold standard expert-scored subset of these images. A deep learning neural network was trained on the calibrated and aggregated score, resulting in very high concordance with expert ratings on a separate test dataset. We termed this approach "hybrid QC", because it combined information from experts with information from community scientists to create a scalable machine learning algorithm that can be applied to future data collection.

However, the hybrid QC proof-of-concept left lingering questions about its applicability to other datasets because it was trained on a single-site, single-modality dataset. Here, we expand the hybrid-QC approach to a large multi-site dMRI dataset. Moreover, one of the common critiques of deep learning is that it can learn irrelevant features of the data and does not provide information that is transparent enough to interpret (*Lipton, 2017*; *Salahuddin et al., 2022*; *Zech et al., 2018*). To confirm that the hybrid-QC deep learning algorithm uses meaningful features of the diffusion-weighted images (DWI) to perform accurate QC, we used machine learning interpretation methods that pry open the "black box" of the neural network, thereby highlighting the features that lead to a specific QC score (*Sundararajan et al., 2017*; *Murdoch et al., 2019*).

Taken together, the combination of curated BIDS data, preprocessed images, and quality control scores generated by the deep learning algorithm provides researchers with a rich and accessible data resource. We anticipate that these HBN Preprocessed Open Diffusion Derivatives (HBN-POD2) will accelerate translational research on both normal and abnormal brain development.

## Results

The aims of this study were fourfold: (i) curate the HBN MRI data into a fully BIDS-compliant MRI dataset, (ii) perform state-of-the-art diffusion MRI (dMRI) preprocessing using *QSIPrep*, (iii) assign QC scores to each participant, and (iv) provide unrestricted public release to the outputs from each of these steps. We started with MRI data from 2747 HBN participants available through FCP-INDI, curating these data for compliance with the Brain Imaging Data Structure (BIDS) specification (*Gorgolewski et al., 2016*). We preprocessed the structural MRI (sMRI) and diffusion MRI (dMRI) data using *QSIPrep*. Participants that could not be curated to comply with the BIDS standard or that did not have dMRI data were excluded, resulting in 2134 participants with preprocessed, BIDS-compliant dMRI data (Figure 1). HBN neuroimaging data was collected at four sites: Staten Island (SI, $N = 300$), Rutgers University Brain Imaging Center (RU, $N = 873$), the CitiGroup Cornell Brain Imaging Center (CBIC, $N = 887$), and the City University of New York Advanced Science Research Center (CUNY, $N = 74$), where numbers in parentheses represent participant counts in HBN-POD2. Figure 2 depicts the age distribution of study participants by sex for each of these scan site as well as pairwise distributions for the automated quality metrics that are described in the next section.

## Healthy Brain Network Preprocessed Open Diffusion Derivatives

Curated BIDS data and their corresponding *QSIPrep* outputs are provided in the FCP-INDI Amazon Web Services (AWS) S3 bucket [2]. This public resource can be accessed by anyone using standard S3 access tools.

The curation process accounts for the acquisition variability inherent in large multi-site datasets by identifying unique *variants* in the HBN dMRI and fieldmap acquisitions. Each session was grouped according to metadata parameters that affect the dMRI signal (PhaseEncodingDirection, EchoTime, VoxelSize, FlipAngle, PhasePartialFourier, NumberOfVolumes, Fieldmap availability). Using the "Curation of BIDS" (CuBIDS) package (*Covitz et al., 2022*), we identified a total of 20 unique DWI acquisitions across HBN-POD2, where about 5% of acquisitions were different from the most common DWI acquisition at their site. The specific variant of each session is provided as a column in the participant.tsv file and a summary of variants with participant counts is provided in Appendix 1.

---

[2] Curated BIDS data is available at s3://fcp-indi/data/Projects/HBN/BIDS_curated/ and *QSIPrep* outputs are available at s3://fcp-indi/data/Projects/HBN/BIDS_curated/derivatives/qsiprep/.
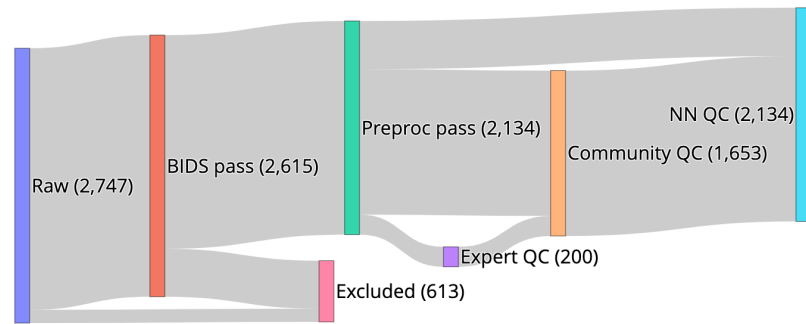
**Figure 1. HBN-POD2 data provenance**: Imaging data for 2747 participants, aged 5-21 years and collected at four sites in the New York City area, was made available through the Functional Connectomes Project and the International Neuroimaging Data-Sharing Initiative (FCP-INDI). These data were curated for compliance to the BIDS specification (*Gorgolewski et al., 2016*) and availability of imaging metadata in json format. 2615 participants met this specification. Imaging data was preprocessed using *QSIPrep* (*Cieslak et al., 2021*) to group, distortion correct, motion correct, denoise, coregister and resample MRI scans. Of the BIDS curated participants, 2134 passed this step, with the majority of failures coming from participants with missing dMRI scans. Expert raters assigned QC scores to 200 of these participants, creating a "gold standard" QC subset. Community raters then assigned binary QC ratings to a superset of the gold standard containing 1653 participants. An image classification algorithm was trained on a combination of automated quality metrics from *QSIPrep* and community scientist reviews to "extend" the expert ratings to the community science subset. Finally, a deep learning QC model was trained on the community science subset to assign QC scores to the entire dataset and to future releases from HBN. The HBN-POD2 dataset, including QC ratings, is openly available through FCP-INDI.

---

¹³⁵      The processed diffusion derivatives are standard *QSIPrep* outputs, which contain preprocessed
¹³⁶ imaging data along with the corresponding QC metrics:

¹³⁷     • *Anatomical Data* Preprocessed images, segmentations and transforms for spatial normaliza-
¹³⁸       tion are located in the `anat/` directory of each session. The gray matter, white matter and
¹³⁹       cerebrospinal fluid (`GM`, `WM`, `CSF`) probabilistic segmentations are provided in nifti format with
¹⁴⁰       the `_probtissue` suffix. The deterministic segmentation is in `_dseg.nii.gz`. All images are
¹⁴¹       in alignment with AC-PC-aligned `sub-X_desc-preproc_T1w.nii.gz` image unless they have
¹⁴²       `space-MNI152NLin2009cAsym` in their file name, in which case they are aligned to the MNI
¹⁴³       Nonlinear T1-weighted asymmetric brain template (version 2009c; (*Fonov et al., 2009a*)). The
¹⁴⁴       spatial transform between the AC-PC T1w image and the MNI space brain is in the ITK/ANTs
¹⁴⁵       format file named `sub-X_from-MNI152NLin2009cAsym_to-T1w_mode-image_xfm.h5`. The brain
¹⁴⁶       mask from `ANTsBrainExtraction.sh` is included in the file with the `_desc-brain_mask.nii.gz`
¹⁴⁷       suffix.

¹⁴⁸     • *Diffusion Data* The preprocessed dMRI scan and accompanying metadata are located in the
¹⁴⁹       `dwi/` directory of each session. The fully-preprocessed dMRI data is named according to the
¹⁵⁰       file pattern `sub-X_space-T1w_desc-preproc_dwi.nii.gz`. These images all have an isotropic
¹⁵¹       voxel size of 1.7 mm and are aligned in world coordinates with the anatomical image located
¹⁵²       at `anat/sub-X_desc-preproc_T1w.nii.gz`. Gradient information is provided in `bval/bvec` for-
¹⁵³       mat compatible with DIPY and DSI Studio and the `.b` format compatible with MRtrix3. Volume-
¹⁵⁴       wise QC metrics including head motion parameters are included in the `confounds.tsv` file.
¹⁵⁵       Automatically computed quality measures for the entire image series are provided in the
¹⁵⁶       `ImageQC.csv` file, which includes the neighboring DWI Correlation, number of bad slices and
¹⁵⁷       head motion summary statistics. Figure 2 depicts pairwise distributions for the three of these
¹⁵⁸       automated data quality metrics that were most informative in QC models described later (see
¹⁵⁹       Appendix 3 for further details). The `desc-brain_mask` file is a dMRI-based brain mask that
¹⁶⁰       should only be used when the T1w-based brain mask is inappropriate (i.e. when no suscep-
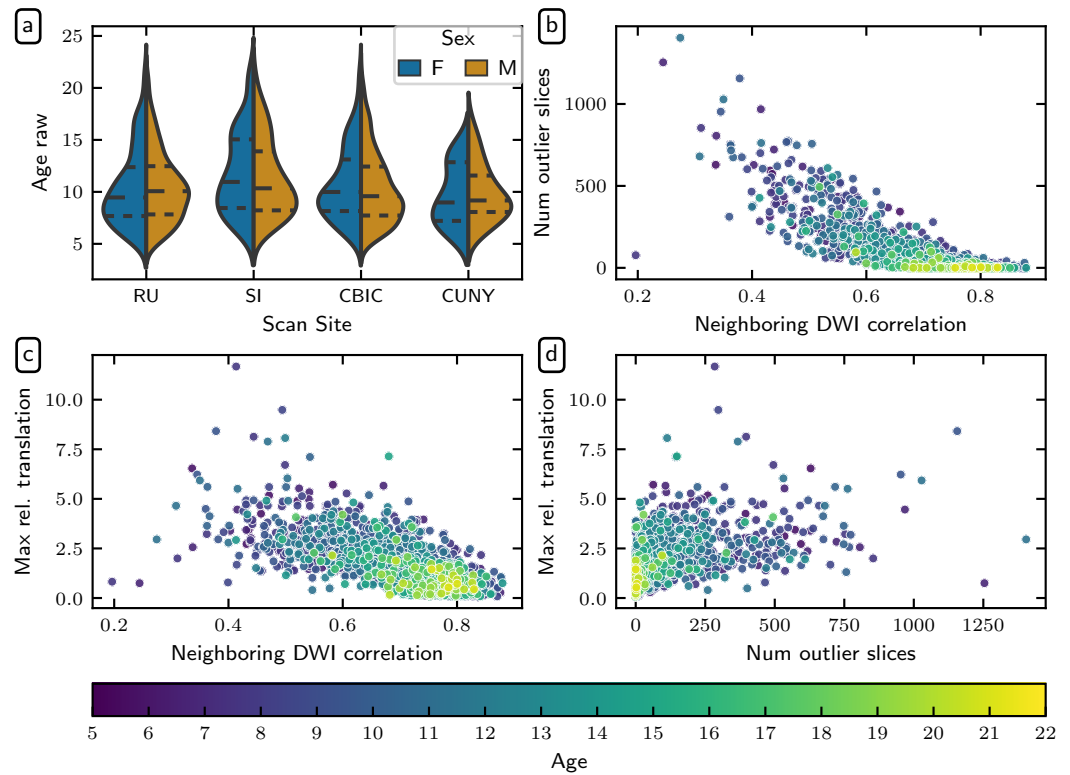¹⁶¹       tibility distortion correction has been applied).

**Figure 2. Demographic and *QSIPrep* quality metric distributions**: **(a)** HBN age distributions by sex for each scanning site. Dashed lines indicate age quartiles. The remaining plots show associations between **(b)** neighboring diffusion-weighted imaging (DWI) correlation (*Yeh et al., 2019*) and the number of outlier slices, **(c)** neighboring DWI correlation and maximum relative translation, and **(d)** the number of outlier slices and maximum relative translation. The number of outlier slices is positively associated with the maximum relative translation, while neighboring DWI correlation is negatively associated with the other two metrics. These plots are colored by age, and reveal that older participants generally have higher quality data.

## Quality Control

To QC all available HBN dMRI data, we adopted a hybrid QC approach that combines expert rating, community science, and deep learning, drawing on the success of a previous application in assessing the quality of HBN's structural T1w MRI data (*Keshavan et al., 2019*). This method (i) starts with dMRI expert raters labelling a small subset of participants, the "gold standard" dataset; (ii) amplifies these labels using a community science web application to extend expert ratings to a much larger subset of the data, the community science subset and (iii) trains a deep learning model on the community science subset to predict expert decisions on the entire dataset.

Expert quality control

To create a gold standard QC dataset, we first developed *dmriprep-viewer*, a dMRI data viewer and QC rating web application to display *QSIPrep* outputs and collect expert ratings (*Richie-Halford et al., 2022*). Six of the co-authors, who are all dMRI experts, rated a 200-participant subset of the HBN-POD2 data using extensive visual examination of each participant's dMRI data, including the preprocessed diffusion weighting imaging (DWI) time series, a plot of motion parameters throughout the DWI scan, and full 3D volumes depicting (i) the brain mask and $b = 0$ to T1w registration and (ii) a directionally encoded color fractional anisotropy (DEC-FA) image laid over the $b = 0$ volume. See Appendix 2 for an example of the *dmriprep-viewer* interface. The experts rated participants using a five-point scale with ratings of "definitely fail," "probably fail," "not sure," "probably pass," and "definitely pass." The distribution of scores given by the experts demonstrates that the gold
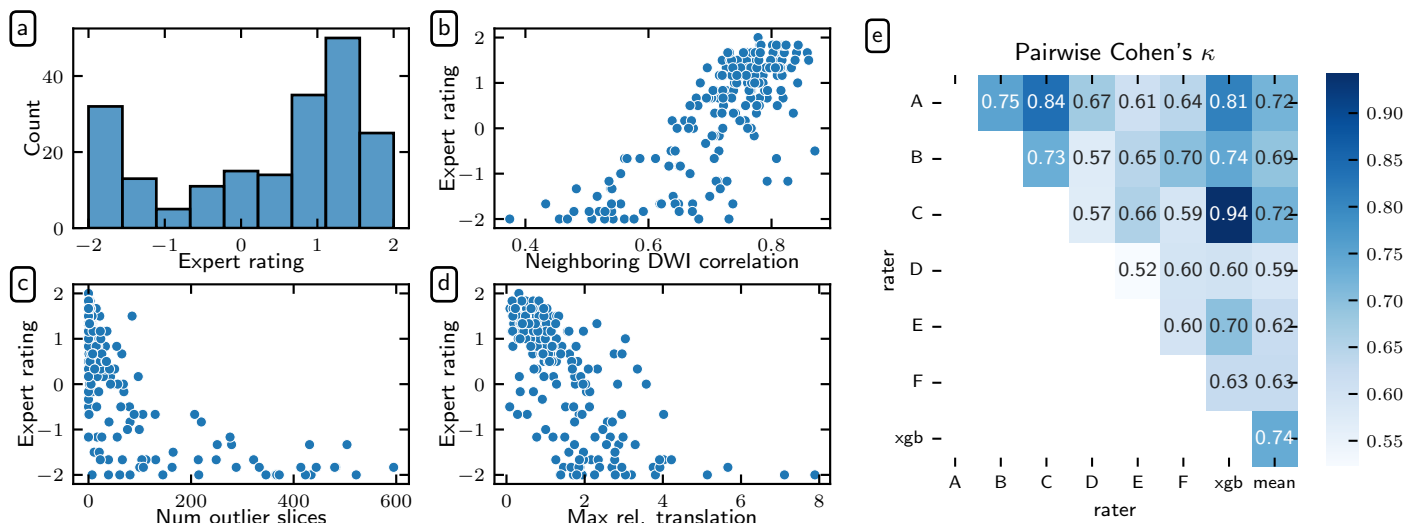
**Figure 3. Expert QC results**: Six dMRI experts rated a subset of 200 participants. Experts agreed with *QSIPrep*'s automated QC metrics. Here we show the distribution of mean expert QC ratings **(a)** and associations between the mean expert QC rating and the *QSIPrep* metrics **(b)** neighboring diffusion-weighted imaging (DWI) correlation (*Yeh et al., 2019*), **(c)** maximum relative translation, and **(d)** number of outlier slices. As expected, neighboring DWI correlation is directly correlated with expert rating while the other two metrics are inversely correlated with expert rating. **(e)** Experts agreed with each other. Here we show the pairwise Cohen's $\kappa$ measure of inter-rater reliability (see text for ICC calculations). The XGB model has an inter-rater reliability (quantified here as Cohen's $\kappa$) that is indistinguishable from the other raters

181　standard dataset included a range of data quality (Figure 3a). Mean expert ratings correlated with
182　the three *QSIPrep* automated QC metrics that were most informative for the XGB model described
183　in the next section: neighboring DWI correlation (*Yeh et al., 2019*) (Figure 3b), maximum relative
184　translation (Figure 3c), and number of outlier slices (Figure 3d). The neighboring DWI correlation
185　characterizes the pairwise spatial correlation between pairs of DWI volumes that sample neigh-
186　boring points in $q$-space. Since lower values indicate reduced data quality, it is reassuring that the
187　neighboring DWI correlation correlated directly with expert ratings (Pearson CC: $0.797$). Conversely,
188　high relative translation and a high number of motion outlier slices reflect poor data quality and
189　these metrics were inversely related to mean expert rating (Pearson CC: $-0.692$ and Pearson CC:
190　$-0.695$, respectively).

191　　In addition to agreeing qualitatively with *QSIPrep*'s automated QC metrics on average, the expert
192　raters also tended to agree with each other (Figure 3e). We assessed inter-rater reliability (IRR)
193　using the pairwise Cohen's $\kappa$ (*Di Eugenio and Glass, 2004*), which exceeded 0.52 in all cases, and
194　with a mean value of 0.648. In addition to the pairwise Cohen's $\kappa$, we also computed the intra-class
195　correlation (ICC) (*Hallgren, 2012*) as a measure of IRR. ICC3k is the appropriate variant of the ICC to
196　use when a fixed set of $k$ raters each code an identical set of participants, as is the case here. ICC3k
197　for inter-rater reliability among the experts was 0.930 (95% CI: [0.91, 0.94]), which is qualitatively
198　considered an "excellent" level of IRR (*Cicchetti, 1994*). The high IRR provides confidence that the
199　average of the expert ratings for each image in the gold standard is an appropriate target to use
200　for training a machine learning model that predicts the expert scores.

201　Community science quality control
202　Although the expert raters achieved high IRR and yielded intuitive associations with *QSIPrep*'s au-
203　tomated QC metrics, generating expert QC labels for the entire HBN-POD2 dataset would be pro-
204　hibitively time consuming. To assess the image quality of the remaining participants, we deployed
205　*Fibr* (https://fibr.dev), a community science web application in which users assigned binary pass/fail
206　labels assessing the quality of horizontal slice DEC-FA images overlaid on the $b = 0$ image (see Ap-
207　pendix 2 for an example). Specifically, *Fibr* users saw individual slices or an animated sequence of
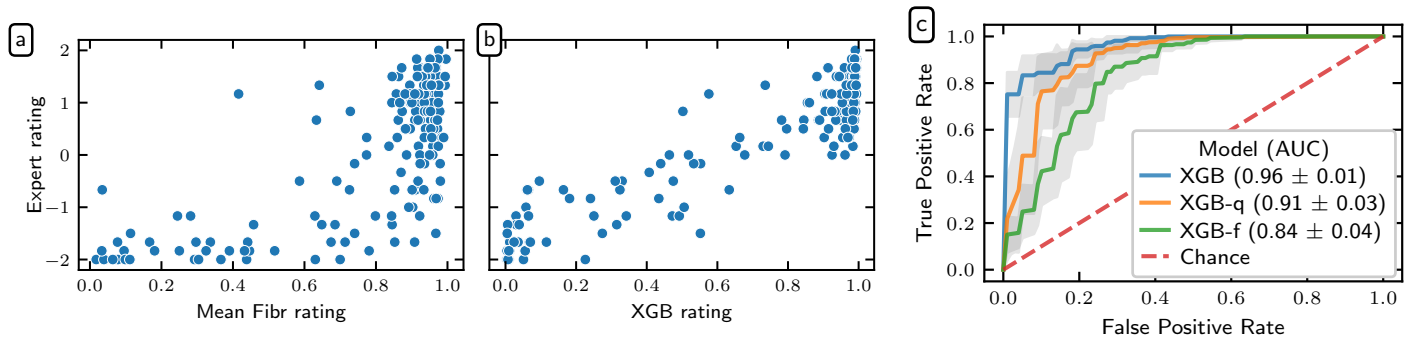
**Figure 4. Community science predictions of the expert ratings**: Scatter plots showing the relationship between mean expert rating and both mean *Fibr* rating **(a)** and XGB prediction **(b)**. *Fibr* raters overestimated the quality of images compared to expert raters. But the XGB prediction compensated for this by incorporating automated QC metrics and weighting more valuable *Fibr* raters. **(c)** ROC curves for the XGB, XGB-q, and XGB-f models. Translucent bands represent one standard deviation from the mean of the cross-validation splits.

208 ten slices taken from the entire DEC-FA volume that the expert raters saw. The *Fibr* users, there-
209 fore, saw only a subset of the imaging data that the dMRI experts had access to for a given partici-
210 pant, but they saw data from many more participants. In total, 374 community scientists provided
211 587,778 ratings for a mean of $> 50$ ratings per slice (or $> 200$ ratings per participant) from 1653 par-
212 ticipants. Of the community scientists, 145 raters provided $> 3,000$ ratings each and are included
213 in the *Fibr* Community Science Consortium as co-authors on this paper (***Ward-Fear et al., 2020***)
214 (see Appendix 4 for a list of consortium co-authors).

215  There are three issues to account for when comparing *Fibr* and expert QC ratings. First, the
216 unadjusted *Fibr* ratings were overly optimistic; i.e., on average, community scientists were not as
217 conservative as the expert raters (Figure 4a). Second, different community scientists provide data
218 of differing accuracy. That is, they were less consistent across different views of the same image,
219 and/or were less consistent with expert ratings for the same data). This means that data from some
220 *Fibr* raters was more informative than others. Third, important information about data quality was
221 provided in the *QSIPrep* data quality metrics, which were not available to *Fibr* raters. To account for
222 rater variability and take advantage of the information provided by *QSIPrep*, we trained gradient
223 boosted decision trees (***Chen and Guestrin, 2016a***) to predict expert scores, scaled to the range
224 $[0, 1]$ and binarized with a $0.5$ threshold, based on a combination of community science ratings and
225 automated *QSIPrep* QC metrics. One can think of the gradient boosting model as assigning more
226 weight to *Fibr* raters who reliably agree with the expert raters, thereby resolving the aforesaid
227 issues with community rater accuracy. We refer to this gradient boosting model as XGB.

228  To clarify the contributions of the automated QC metrics and the community science raters, we
229 trained two additional gradient boosting models: (i) one trained only on the automated *QSIPrep*
230 data quality metrics, which we call XGB-q and (ii) one trained on only the *Fibr* ratings, which we
231 call XGB-f. XGB-f may be viewed as a data-driven weighting of community scientists' ratings, while
232 XGB-q may be viewed as a generalization of data quality metric exclusion criteria. XGB, combining
233 information from both *Fibr* ratings and *QSIPrep* data quality metrics attained a cross-validated area
234 under the receiver operating curve (ROC-AUC) of $0.96 \pm 0.01$ on the "gold standard," where the $\pm$
235 indicates the standard deviation of scores from repeated $k$-fold cross-validation (Figure 4b). In
236 contrast, XGB-q attained an ROC-AUC of $0.91 \pm 0.03$ and XGB-f achieved an ROC-AUC of $0.84 \pm 0.04$.
237 The enhanced performance of XGB-q over XGB-f shows that community scientists alone are not as
238 accurate as automated data quality metrics are at predicting expert ratings. And yet, the increased
239 performance of XGB over XGB-q demonstrates that there is additional image quality information
240 to be gained by incorporating community scientist input.

241  As a way of evaluating the quality of the XGB predictions, consider the fact that the average
242 Cohen's $\kappa$ between XGB and the expert raters was 0.74, which is higher than the average Cohen's

243 $\kappa$ between any of the other raters and their human peers (Figure 3). This is not surprising, given
244 that the XGB model was fit to optimize this match, but further demonstrates the goodness of fit of
245 this model.

246 Nevertheless, this provides confidence in using the XGB scores in the next step of analysis,
247 where we treat the XGB model as an additional coder and extend XGB ratings to participants with-
248 out *Fibr* ratings. In this case, when a subset of participants is coded by multiple raters and the
249 reliability of their ratings is meant to generalize to other participants rated by only one coder, the
250 single-measure ICC3, as opposed to ICC3k, should be used. When adding XGB to the existing expert
251 raters as a seventh expert, we achieved **ICC3** $= 0.709(95\%CI : [0.66, 0.75])$. The high ICC3 value after
252 inclusion of the XGB model justifies using the XGB scores as the target for training an image-based
253 deep learning network.

## Automated quality control labelling through deep learning

255 While the XGB "rater" does a good job of extending QC ratings to the entire community science
256 subset, this approach requires *Fibr* scores; without community science *Fibr* scores, only the less
257 accurate XGB-q prediction can be employed. Consequently, a new, fully automated QC approach
258 is needed that can be readily applied to new data releases from HBN.

259 We therefore trained a deep convolutional neural network to predict the XGB ratings directly
260 from *QSIPrep* outputs. We modified an existing 3D convolutional neural network (CNN) architec-
261 ture (*Zunair et al., 2020*)—previously applied to the ImageCLEF Tuberculosis Severity Assessment
262 2019 benchmark (*Dicente Cid et al., 2019*)—to accept multichannel input generated from the pre-
263 processed dMRI: the $b = 0$ reference diffusion image, each of the three cardinal axis components
264 of the DEC-FA image, and, optionally, automated QC metrics from *QSIPrep*. We trained this net-
265 work on XGB scores and validated it against the gold standard expert-scored dataset. We refer
266 to the convolutional neural network model trained only on imaging data as CNN-i and the model
267 that incorporates automated QC metrics as CNN-i+q. The two models performed nearly identically
268 and achieved an ROC-AUC of $0.947 \pm 0.004$ (Figure 5a). The near-identical performance suggests
269 that *QSIPrep*'s automated data quality metrics provided information that was redundant with in-
270 formation available in the imaging data. Both CNN-i and CNN-i+q outperformed XGB-q, which was
271 trained only on automated QC metrics, but both modestly underperformed relative to the full XGB
272 model, that uses *Fibr* scores in addition to the *QSIPrep* data quality metrics.

273 The openly available HBN-POD2 data released with this paper provides four QC ratings: the
274 mean expert QC ratings, XGB-q and XGB predicted scores, as well as the CNN-i predicted score.
275 However, we treat the CNN-i score as the definitive QC score because it is available for all partici-
276 pants, can be easily calculated for new participants in future HBN releases, and is more accurate
277 than XGB-q in predicting expert ratings in the "gold standard" report set. When we refer to a par-
278 ticipant's QC score without specifying a generating model, the CNN-i score is assumed. Figure 5
279 depicts the distribution of these QC scores by age (Figure 5b), sex (Figure 5c), and scanning site
280 (Figure 5d). QC distributions are similar for each scan site and for male and female participants [3].

## Attribution masks for the deep learning classifier

282 We generated post-hoc attribution maps that highlight regions of the input volume that are rele-
283 vant for the QC score. The integrated gradient method (*Sundararajan et al., 2017*) is a gradient-
284 based attribution method (*Ancona et al., 2019*) that aggregates gradients for synthetic images in-
285 terpolating between a baseline image and the input image. It has been used to interpret deep
286 learning models applied to retinal imaging in diabetic retinopathy (*Sayres et al., 2019*) and glau-
287 coma (*Mehta et al., 2021*) prediction, as well as in multiple sclerosis prediction from brain MRI
288 (*Wargnier-Dauchelle et al., 2021*). Our goal is to confirm that the CNN-i model was driven by the
289 same features that would drive the expert rating, thereby bolstering the decision to apply it to new
290 data.

---

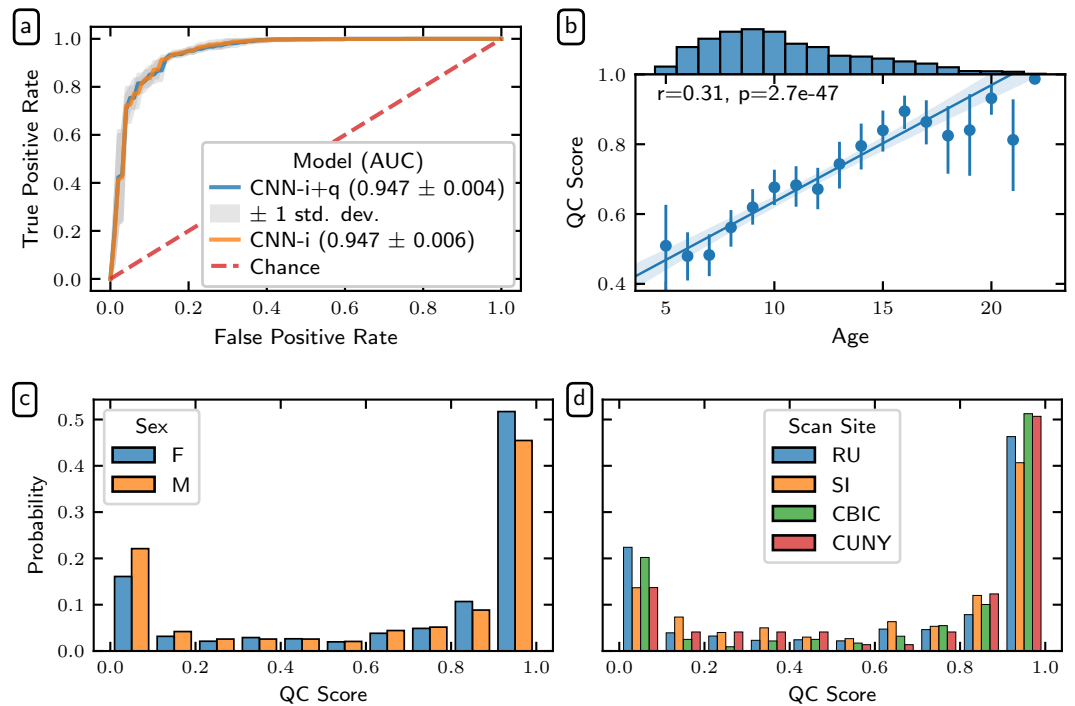[3]Responses for the sex variable in HBN phenotypic data are limited to "male" and "female."

**Figure 5. Deep learning QC scores**: **(a)** ROC curves for two deep learning models trained on imaging data: one trained with additional automated data quality metrics from *QSIPrep* (blue) and one trained without (orange). The models performed roughly identically, reflecting that the data quality metrics are derived from the imaging data and are therefore redundant. Both outperformed the XGB-q predictions, indicating the added value of the diffusion weighted images. However, both models underperformed the XGB predictions, which also incorporate information from *Fibr* ratings for each scan. The error bands represent one standard deviation from the mean of the cross-validation splits. **(b)** Joint distributions showing a strong direct association between age and QC score (Pearson CC: $0.31$). This likely reflects the well-known negative association between age and head motion in pediatric neuroimaging. The dots encode the mean QC score for each year of age with error bands representing a bootstrapped 95% confidence interval. The line depicts a linear regression relating age and QC score with translucent bands encoding a bootstrapped 95% confidence interval. Histograms showing the relationship between participants QC scores and their sex **(c)** and scan site **(d)**. QC distributions are independent of sex and scanning site.

Figure 6 shows attribution maps for example participants from each confusion class: true positive, true negative, false positive, and false negative. The columns correspond to the different channels of the deep learning input volume: the $b = 0$ reference image and the DEC-FA in the $x$, $y$, and $z$ directions. The blue voxels indicate positive attribution, that is, data that supports a passing QC classification. Conversely, the red voxels indicate negative attribution, data that supports a failing QC classification. The true positive map indicates that the network was looking at the entire brain rather than focusing on any one anatomical region (Figure 6a). Moreover, the model identified white matter fascicles that travel along the direction of the input channel: lateral for $x$, anterior-posterior for $y$, and superior-inferior for $z$. The true negative attribution map (Figure 6b) reveals that when the reference $b = 0$ volume contains motion artifacts, such as banding, the network ignored the otherwise positive attributions for the clearly identifiable white matter tracts in the DEC-FA channels. The false positive map (Figure 6c) and the false negative map (Figure 6d) should be interpreted differently since they come from low confidence predictions; the probability of passing hovered on either side of the pass/fail threshold. For example, in the false positive case, the network was confused enough that it treated voxels that are outside of the brain to be as informative as voxels in the major white matter bundles.
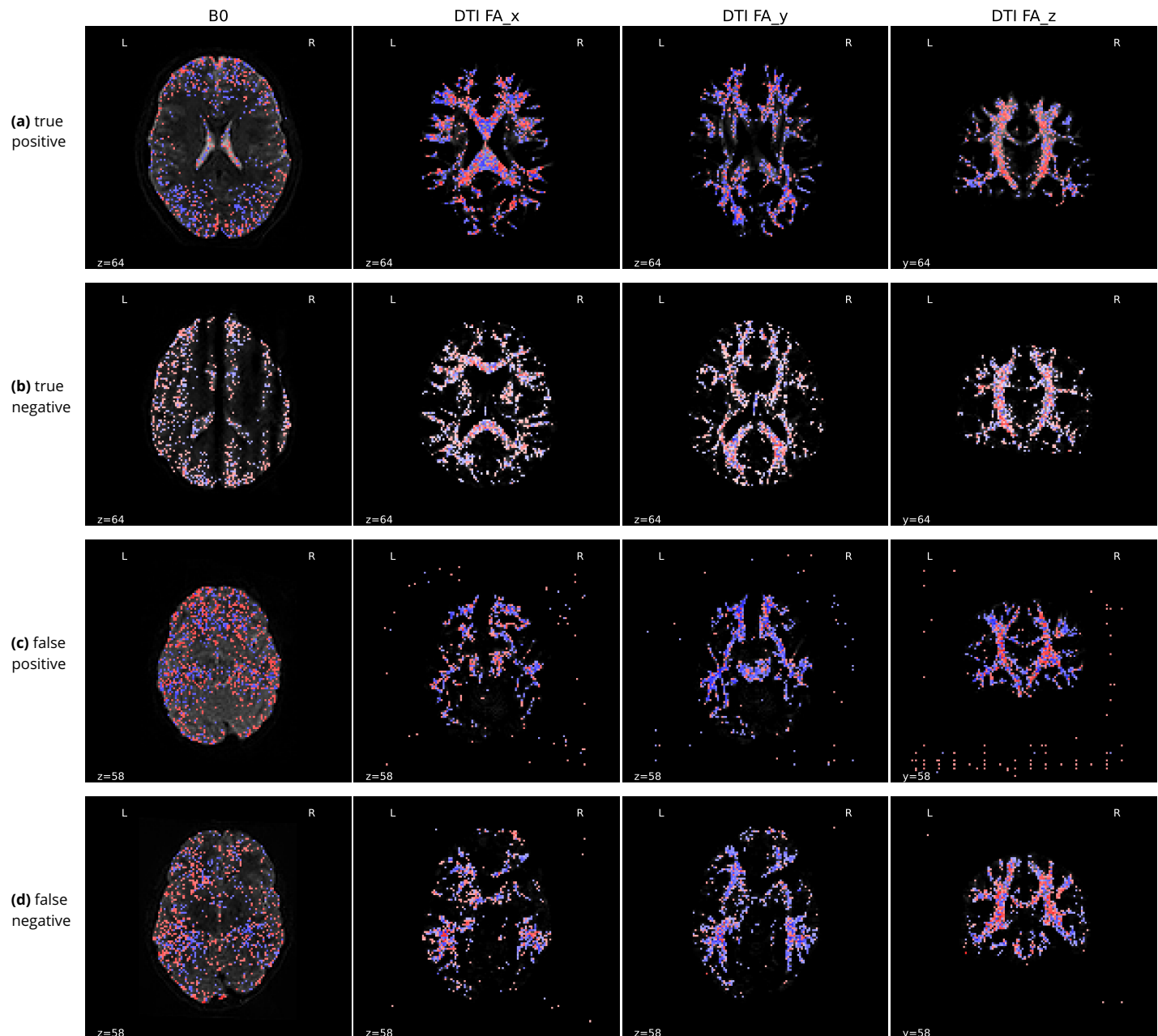
**Figure 6. Integrated gradients attribution maps for the deep learning classifier**: Each column depicts a different channel of the input tensor: the $b = 0$ DWI volume and the DEC-FA images in the $x$, $y$, and $z$ directions. The first three columns show an axial slice while the last column shows a coronal slice. Blue voxels indicate positive attribution (i.e. evidence for passing the participant), while red voxels indicate negative attribution (i.e. evidence for QC failure). The underlying grayscale depicts the input channel. Each row depicts a representative participant from each confusion class: **(a)** Attribution maps for a true positive prediction. The model looked at the entire brain and focused on known white matter bundles in the DEC-FA channels. In particular, it focused on lateral bundles in the $x$ direction, anterior-posterior bundles in the $y$ direction, and superior-inferior bundles in the $z$ direction. **(b)** Attribution maps for a true negative prediction. The model focused primarily on the $b = 0$ channel, suggesting that it ignores DEC-FA when motion artifacts like banding are present. **(c)** Attribution maps for a false positive prediction. Both the false positive and negative predictions were low confidence predictions. This is reinforced by the fact that the model viewed some voxels that are outside of the brain as just as informative as those in major white matter tracts. **(d)** Attribution maps for a false negative prediction. The model failed to find long-range white matter tracts in the anterior-posterior and lateral directions. We also speculate that the model expected left-right symmetry in the DEC-FA channels and assigned negative attribution to asymmetrical features.

307 QC prediction models can generalize to unseen sites

308 Site harmonization is a major issue for any multisite neuroimaging study and developing auto-
309 mated QC tools that generalize between sites has been a perennial issue (*Esteban et al., 2017*).
310 Furthermore, the ability to generalize between sites in a single multisite study would signal the
311 promise of generalizing to other datasets altogether. To better understand the ability of our QC
312 models to generalize across scanning sites, we trained several variants of the XGB-q and CNN-i
313 models on partitions of the data with different sites held out (Figure 7). ROC-AUC for generalization
314 is uniformly high for both the XGB-q and the CNN-i models (Table 1). However, more importantly,
315 accuracy and balanced accuracy vary substantially: depending on the site that was used for train-
316 ing, balanced accuracy could be as low as guess rate, particularly for the CNN-i model. Notably,
317 it seems that including the RU site in the training data led to relatively high balanced accuracy in
318 both models. The XGB-q model balanced accuracy was less dependent on the specific sites used
319 for training, but also displayed some variability across permutations of this experiment. In partic-
320 ular, the benefit from including the "right site" in the training data, namely RU, eclipsed the slight
321 benefit conferred by including more than one site in the training data.

## Quality control improves inference

323 To demonstrate the effect that quality control has on inference, we analyzed tract profile data
324 derived from HBN-POD2 data. Tract profiling (*Yeatman et al., 2012*; *Jones et al., 2005*; *Colby et al.,*
325 *2012*; *O'Donnell et al., 2009*; *Kruper et al., 2021*) is a subset of tractometry (*Jones et al., 2005*; *Bells*
326 *et al., 2011*), which uses the results of dMRI tractography to quantify properties of the white matter
327 along major pathways. Tract-profiling retains the values of diffusion metrics along the trajectory of
328 each bundle of tractography streamlines, rather than computing summary statistics summarized
329 at the level of each bundle. In Figure 8, we plot mean diffusivity tract profiles grouped into four QC
330 bins along the length of twenty-four bundles: While some bundles, such as the cingulum cingulate
331 (CGC) and the inferior longitudinal fasciculus (ILF), appear insensitive to QC score, others, such
332 as the uncinate (UNC) and the orbital portion of the corpus callosum, exhibit strong differences
333 between QC bins. In most bundles, low QC scores tend to flatten the profile, indicating that mean
334 diffusivity appears artifactually homogeneous across the bundle.

335 The effect of QC score on white matter
336 bundle profiles indicates that researchers
337 using HBN-POD2 should incorporate QC in
338 their analyses, either by applying a QC cut-
339 off when selecting participants or by explic-
340 itly adding QC score to their inferential mod-
341 els. Failure to do so may cause spurious
342 associations or degrade predictive perfor-
343 mance. To demonstrate this, we selected
344 participant age as a representative pheno-
345 typic benchmark because (i) it operates on
346 a natural scale with meaningful units and
347 (ii) despite the unique methodological chal-
348 lenges it presents for biomarker identifica-
349 tion (*Nelson et al., 2020*), brain age pre-
350 diction may be diagnostic of overall brain
351 health (*Franke et al., 2010*; *Cole et al., 2019*;
352 *Richie-Halford et al., 2021*). We observed
353 the effect of varying QC cutoff on the predic-
354 tive performance of an age prediction model.
355 Cross-validated $R^2$ scores for an age predic-
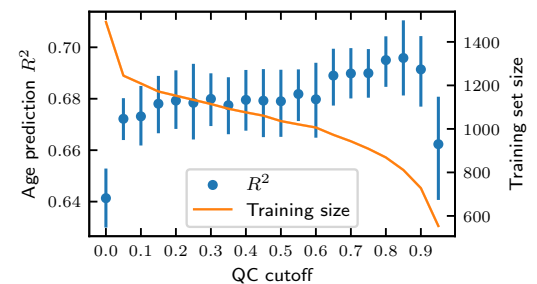356 tion model varied depending on the QC cut-



**Figure 9. Imposing a QC cutoff improves age prediction**: Cross validated $R^2$ scores (left axis, blue dots) from an age prediction model increase after screening participants by QC score. We see the most dramatic increase in $R^2$ after imposing even the lowest cutoff of $0.05$. Thereafter, the $R^2$ scores trend upward until a cutoff of $\sim 0.95$, where the training set size (right axis, orange line) becomes too small to sustain model performance. The error bands represent a bootstrapped 95% confidence interval.
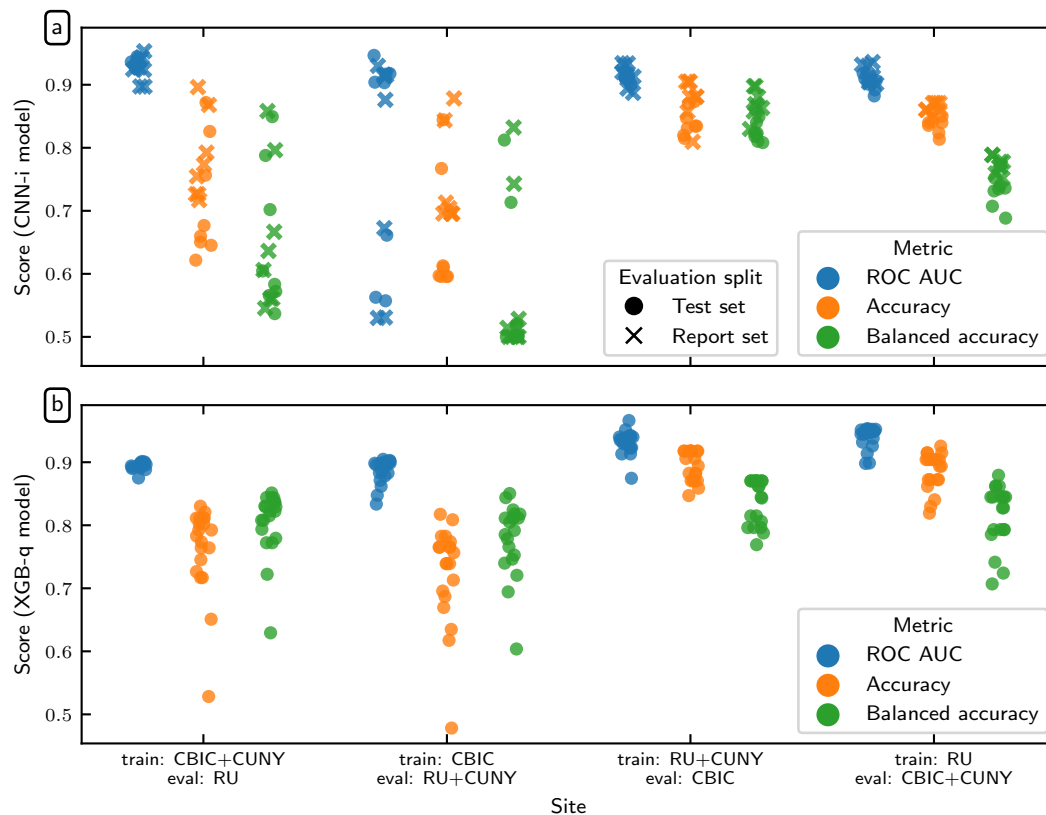
**Figure 7. Generalization of QC scores to unseen sites**: In each experiment, CNN-i (**a**) and XGB-q (**b**) models were trained with some sites held out and evaluated only on data from these held out sites. Model performance is quantified as ROC-AUC (blue), accuracy (orange) and balanced accuracy (green). For XGB-q, the targets are the expert ratings on data from the held out site. For CNN-i, performance is scored against XGB scores (as used before; test set in filled circles), or expert ratings on the data from the held out site (report set in crosses). Summary statistics for this plot are listed in Table 1.

**Table 1. Site generalization summary statistics**: Below we list the mean ± standard deviation of the site generalization evaluation metrics displayed in Figure 7. For each of the CNN-i and XGB-q model families and each of the site generalization splits, we report the accuracy, balanced accuracy, and ROC-AUC.

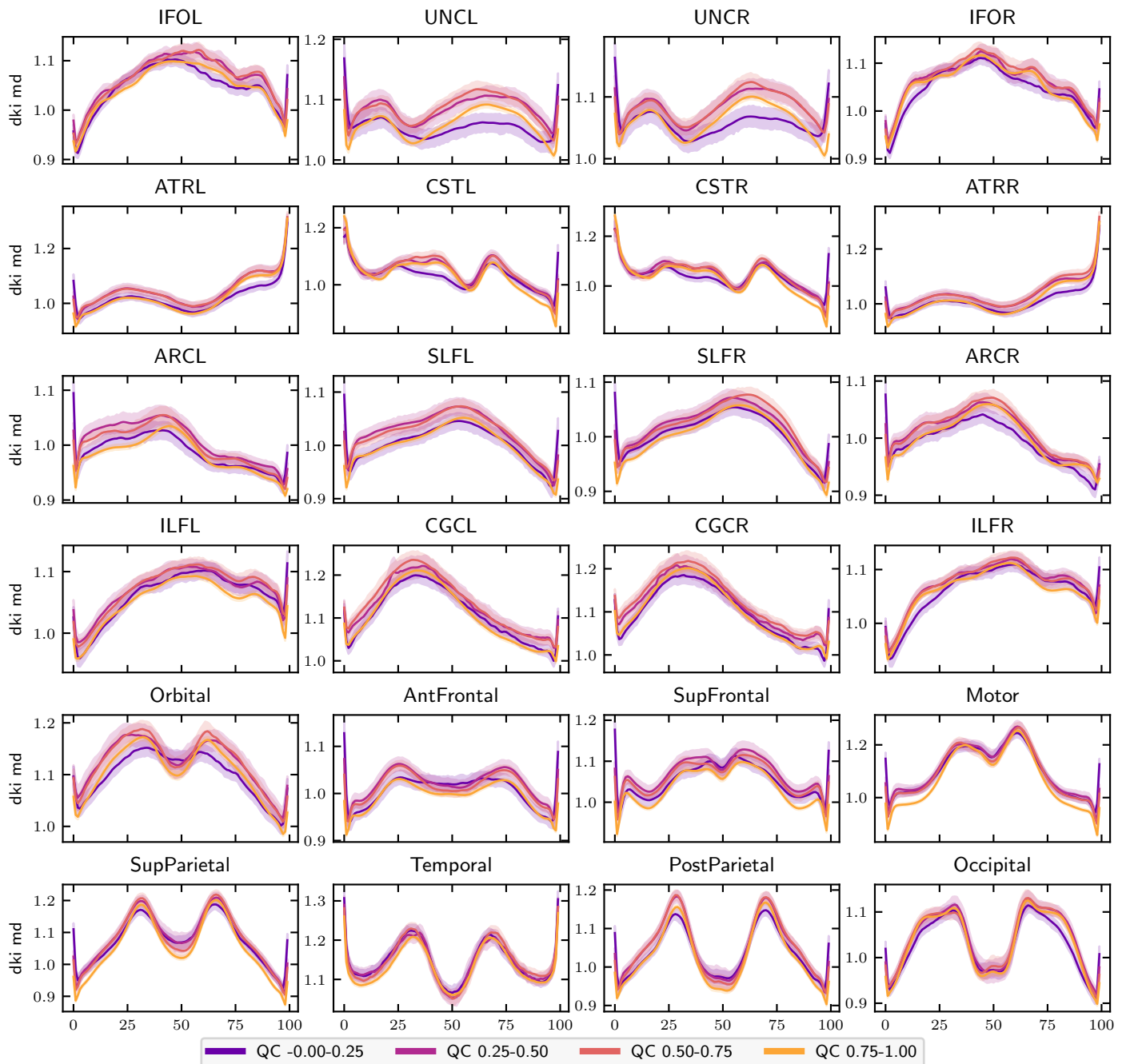| Model | Site | Accuracy | Balanced accuracy | ROC-AUC |
|---|---|---|---|---|
| CNN-i | train: CBIC + CUNY, test: RU | $0.748 \pm 0.086$ | $0.652 \pm 0.112$ | $0.930 \pm 0.015$ |
| | train: CBIC, test: RU + CUNY | $0.696 \pm 0.095$ | $0.574 \pm 0.123$ | $0.791 \pm 0.169$ |
| | train: RU + CUNY, test: CBIC | $0.859 \pm 0.033$ | $0.847 \pm 0.030$ | $0.912 \pm 0.013$ |
| | train: RU, test: CBIC + CUNY | $0.851 \pm 0.018$ | $0.753 \pm 0.029$ | $0.910 \pm 0.014$ |
| XGB-q | train: CBIC+CUNY, test: RU | $0.763 \pm 0.071$ | $0.805 \pm 0.052$ | $0.895 \pm 0.006$ |
| | train: CBIC, test: RU+CUNY | $0.725 \pm 0.079$ | $0.779 \pm 0.058$ | $0.886 \pm 0.019$ |
| | train: RU+CUNY, test: CBIC | $0.894 \pm 0.024$ | $0.838 \pm 0.036$ | $0.931 \pm 0.018$ |
| | train: RU, test: CBIC+CUNY | $0.886 \pm 0.030$ | $0.816 \pm 0.048$ | $0.940 \pm 0.017$ |

**Figure 8. MD bundle profiles show large QC group differences**: MD profiles binned by QC score in twenty-four major while matter bundles. The *x*-axis represents distance along the length of the fiber bundle. The left and right uncinate bundles were the most sensitive to QC score. Generally, QC score tended to flatten bundle profiles. Error bands represent bootstrapped 95% confidence intervals. Bundle abbreviations for lateralized bundles contain a trailing "L" or "R" indicating the hemisphere. Bundle abbreviations: inferior fronto-occipital fasciculus (IFO), uncinate (UNC), thalamic radiation (ATR), corticospinal (CST), arcuate (ARC), superior longitudinal fasciculus (SLF). inferior longitudinal fasciculus (ILF), cingulum cingulate (CGC), orbital corpus callosum (Orbital), anterior frontal corpus callosum (AntFrontal), superior frontal corpus callosum (SupFrontal), motor corpus callosum (Motor), superior parietal corpus callosum (SupParietal), temporal corpus callosum (Temporal), post-parietal corpus callosum (PostParietal), and occipital corpus callosum (Occipital).

**Figure 8–Figure supplement 1.** FA bundle profiles

357  off (Figure 9). An initial large improvement was achieved by excluding the 200 participants with the

358  lowest QC scores, followed by a gradual increase in performance. Finally, when a large number of

359  participants is excluded, performance deteriorated again.

## Discussion

361  We present HBN-POD2, one of the largest child and adolescent diffusion imaging datasets with

362  preprocessed derivatives that is currently openly available. The dataset was designed to comply

363  with the best practices of the field. For example, it complies with the current draft of the BIDS

364  diffusion derivative specification (*Pestilli et al., 2021*). It will grow continuously as the HBN study

365  acquires more data, eventually reaching its 10,000 participant goal.

### Preprocessing and quality control increase the impact of openly-available data

367  The most immediate contribution of this work is a large analysis-ready dMRI data resource, openly

368  accessible to the public. In the past decade, projects such as the Human Connectome Project (HCP)

369  (*Van Essen et al., 2013*), UK Biobank (*Miller et al., 2016*), ABCD (*Jernigan and Brown, 2018*), and Cam-

370  CAN (*Taylor et al., 2017*; *Shafto et al., 2014*) and of course FCP-INDI (which includes HBN) (*Mennes*

371  *et al., 2013*) have ushered a culture of data sharing in open big-data neuroscience. The adoption

372  and reuse of these datasets reduces or eliminates the data collection burden on downstream re-

373  searchers. Some projects, such as the HCP (*Glasser et al., 2013*), also provide preprocessed deriva-

374  tives, further reducing researchers' burden and extending the benefits of data-sharing from data

375  collection to preprocessing and secondary analysis. Following the example of the HCP, HBN-POD2

376  provides analysis-ready dMRI derivatives. This avoids duplication of and heterogeneity across the

377  preprocessing effort while also ensuring a minimum standard of data quality for HBN researchers.

378  We also provide the CuBIDS variant annotation in the participants.tsv file, allowing researchers to

379  account for the imaging heterogeneity inherent in a dataset of this size. Making MRI derivatives

380  accessible not only reduces the burden of processing large datasets for research groups with lim-

381  ited resources (*Laird, 2021*), but also aids research performed by clinicians who are interested in

382  brain-behavior relationships but may be lacking the technical training to process large-scale dMRI

383  data.

384      The data is amenable to many different analyses, including tractometry (*Yeatman et al., 2012*,

385  *2018*; *Kruper et al., 2021*), graph theoretical analysis (*Yeh et al., 2020*), and combinations with func-

386  tional MRI data and other data types for the same participants. The availability of standardized

387  preprocessed diffusion data will allow researchers to create and test hypotheses on the white mat-

388  ter properties underlying behavior and disease, from reading and math acquisition to childhood

389  adversity and mental health. As such, this dataset will accelerate discovery at the nexus of white

390  matter microstructure and neurodevelopmental and learning disorders.

391      In large developmental datasets, it is critically important to perform accurate and reliable QC

392  of the data. QC is associated not just with age, but with many phenotypic variables of interest in

393  cognition and psychopathology (*Siegel et al., 2017*). HBN-POD2 provides four separate QC scores

394  alongside its large dataset of pediatric neuroimaging diffusion derivatives, paving the way for users

395  of the data to incorporate considerations of data quality into their analysis of the processed data.

396  Unsurprisingly, QC scores are strongly correlated with age (Figure 5). This accords with the negative

397  association between head motion and age in developmental studies, which is well established both

398  in general (*Power et al., 2012*; *Satterthwaite et al., 2012*; *Fair et al., 2012*; *Yendiki et al., 2014*)

399  and specifically for resting-state fMRI in the HBN dataset (*Alexander et al., 2017*). Moreover, it

400  is important that QC has bundle-specific and spatially localized effects (Figure 8). Analysis of this

401  data that does not incorporate QC is likely to find replicable but invalid effects. For example, in

402  patient-control studies, patients are likely to have lower quality data. And analysis of such patient

403  data that does not control for QC will find spatially-localized and replicable group differences that

404  are due to data quality, not necessarily underlying neuroanatomical differences.

405 　We further demonstrated the impact of QC in a benchmark age prediction task (Figure 9). In
406 this case, the increase in model performance from imposing a QC cutoff is intuitive: we know
407 from Figure 8 that participants with low QC scores have reduced MD, but MD also decreases as
408 participants mature (*Yeatman et al., 2014*; *Richie-Halford et al., 2021*). Eliminating participants with
409 low QC therefore removes the ones who may look artificially older from the analysis, improving
410 overall performance. The most noticeable improvement in performance comes after imposing
411 the most modest cutoff of 0.05, suggesting that inferences may benefit from *any* QC screening. On
412 the other hand, QC screening inherently introduces a tradeoff between the desire for high quality
413 data and the desire for a large sample size. In this case, after a QC cutoff of around 0.9, the training
414 set size is reduced such that it degrades predictive performance. Importantly, we do not expect
415 the sensitivity analysis of an age prediction model to generalize to other analyses and therefore
416 recommend that researchers using HBN-POD2 choose the most appropriate QC cutoff for their
417 research question and consider including QC score as a model covariate in their analyses.

418 **Automated quality control: scalability, interpretability, and generalization**
419 The predictive performance of the CNN-i model (Figure 5a) gives us confidence that it could ac-
420 curately classify unseen data from the same sites, justifying its extension to the entire HBN-POD2
421 dataset and to future releases of HBN. However, one limitation of this model is that it does not satis-
422 factorily explain its decisions. As deep learning models have been increasingly applied to medical
423 image analysis, there is an evolving interest in the interpretability of these models (*Salahuddin
424 et al., 2022*; *Lipton, 2017*; *Zech et al., 2018*; *Ghassemi et al., 2021*). While an exhaustive interpre-
425 tation of deep learning QC models is beyond the scope of this work, we provided a preliminary
426 qualitative interpretation of the CNN-i model (Figure 6) that demonstrates the intuitive nature of
427 its decisions.

428 　The accuracy in generalizing to unseen data from HBN also suggested the tantalizing possibility
429 that the QC models would be able to generalize to similar data from other datasets. To assess this,
430 we trained the models with unseen sites held out (Figure 7). Both the CNN-i model and the XGB-q
431 model do sometimes generalize to data from unseen sites, suggesting that they would be able to
432 generalize to some other datasets as well. However, they do not reliably generalize, implying that
433 they should not currently be used in this way. Future work could build upon the work that we have
434 done here to establish a procedure whereby the models that we fit in HBN would be applied to data
435 from other studies, but comprehensive calibration and validation would have to be undertaken as
436 part of this procedure.

437 　We recognize that decisions about QC inclusion must balance accuracy, interpretability, gener-
438 alization to new data, and scalability to ever larger datasets. We therefore provide three additional
439 scores: (i) the mean expert QC score for the 200 participants in the gold standard dataset, (ii) the
440 scores predicted by the XGB model, which outperformed all other models when evaluated against
441 the gold standard ratings, but which are only available for participants that have community sci-
442 ence scores, and (iii) the scores predicted by the XGB-q model, which underperformed the deep
443 learning generated scores, but which rely only on the automated QC metrics output by *QSIPrep*.
444 We view the XGB-q scores, which are available for all participants, as a more interpretable and scal-
445 able fallback because the XGB-q model ingests *QSIPrep* output without any further postprocess-
446 ing. XGB-q also provides slightly more uniform performance in generalization to unseen HBN sites
447 (Figure 7). Because the XGB-q model most readily generalizes to other *QSIPrep* outputs, we pack-
448 age it as an independent QC service in the QSIQC software package (*Richie-Halford and Rokem,
449 2022b*), available both as a docker image at `ghcr.io/richford/qsiqc` and as a Streamlit app at
450 https://share.streamlit.io/richford/qsiqc/main/app.py. The decision to use a more interpretable but
451 slightly less performant method of generating QC scores was also advocated by (*Tobe et al., 2021*),
452 who noted that the Euler number of T1-weighed images (*Rosen et al., 2018*) in the NKI-Rockland
453 dataset can reliably predict scores generated with *Braindr*, the community science application de-
454 veloped in our previous work (*Keshavan et al., 2019*).

455 We also note that the issue of algorithmic impact in choosing a QC method is not exclusive to
456 the deep learning model. We have chosen models that most reliably reproduce the gold standard
457 ratings, but a reliable algorithm might still negatively influence researcher's decisions. For example,
458 excluding participants by QC score could spur them to exclude populations deserving of study,
459 as when QC score is highly correlated with age or socio-economic status. We therefore caution
460 researchers to examine interactions between the QC scores we provide and their phenotype of
461 interest.

462 More generally, QC in the dataset that we have produced is fundamentally anchored to the
463 decisions made by the expert observers. While Cohen's $\kappa$ between some pairs of experts can be
464 as low as 0.52, IRR quantified across all of the experts with ICC3k is excellent. Nevertheless, it is
465 possible that improvements to the final QC scores could be obtained through improvements to
466 IRR, or by designing a more extensive expert QC protocol. The tradeoff between more extensive
467 QC for each participant and more superficial QC on more participants was not explored in this
468 study, but could also be the target for future research.

## Transparent pipelines provide an extensible baseline for future methods

470 While the primary audience of HBN-POD2 is researchers in neurodevelopment who will use the
471 dMRI derivatives in their studies, other researchers may use HBN-POD2 to develop new prepro-
472 cessing algorithms or quality control methods. In this respect, HBN-POD2 follows *Avesani et al.*
473 (*2019*), who recognized the diverse interests that different scientific communities have in reusing
474 neuroimaging data and coined the term *data upcycling* to promote multiple-use data sharing for
475 purposes secondary to those of the original project. Complementing the approach taken in Avesani
476 et al.'s work, which provided dMRI from a small number of participants preprocessed with many
477 pipelines, HBN-POD2 contains many participants, all processed with a single state of the art pipeline,
478 *QSIPrep*. For researchers developing new preprocessing algorithms, HBN-POD2 provides a large,
479 openly available baseline to which they can compare their results.

480 Similarly, neuroimaging QC methods developers will benefit from a large benchmark dataset
481 of expert, community science, and automated QC ratings, with which to test new methods. Im-
482 portantly, the architecture and parameters of the deep learning network used for QC are also
483 provided as part of this work, allowing application of this network to future releases of HBN data,
484 and allowing other researchers to build upon our efforts. Indeed, in this work, we have extended
485 our previous work on what we now call "hybrid QC". This approach, which we originally applied
486 to the first two releases of the HBN T1-weighted data (*Keshavan et al., 2019*) (using the *Braindr*
487 web app: https://braindr.us) was extended here in several respects. First, the *Braindr* study used a
488 smaller dataset of approximately 700 participants, while we extended this approach to well over
489 2000 participants. Second, *Braindr* relied on approximately 80,000 ratings from 261 users. Here,
490 we received more than 500,000 ratings from 374 community scientists. As our understanding of
491 the role of community scientist contributions has evolved, we decided that we would include as col-
492 lective co-authors community scientists who contributed more than 3000 ratings (*Ward-Fear et al.,*
493 *2020*). Third, *Braindr* used data from only a single site. Here, multi-site data was used. This opens
494 up multiple possibilities for deeper exploration of between-site quality differences, and also for har-
495 monization of QC across sites, as we have attempted here. Last, the most challenging extension of
496 hybrid QC from *Braindr* to this study entailed developing an approach that would encompass multi-
497 volume dMRI data. On the one hand, this meant that the task performed by the expert observers
498 was more challenging, because it required examination of the full dMRI time-series for every scan.
499 To wit, expert inter-rater reliability was considerably higher for the T1-weighted only data in *Ke-*
500 *shavan et al.* (*2019*) than for the dMRI data used (Figure 3e). On the other hand, it also meant that
501 the 4D data had to be summarized into 2D data to be displayed in the *Fibr* web application. This
502 was achieved by summarizing the entire time-series as a DEC-FA + $b = 0$ image and presenting
503 community scientists with animated sections of these images that showed how the data extended
504 over several horizontal slices. In addition, the extension to 4D data required developing new deep

505 learning architectures for analysis of 4D images, including upstream contributions to *Nobrainer*, a
506 community-developed software library for deep learning in neuroimaging data (*Kaczmarzyk et al.,*
507 *2021*). These extensions demonstrate that the hybrid QC approach generalizes very well to a vari-
508 ety of different circumstances. Future applications of this approach could generalize to functional
509 MRI data, as well as other large datasets from other kinds of measurements and other research
510 domains.

### Future work and open problems

512 The HBN study plans to acquire imaging data for over 5000 participants, necessitating future data
513 releases. Since future releases of HBN will also require future releases of HBN-POD2, a plan for
514 these is essential. This is a general issue affecting multi-year neuroimaging projects for which
515 derivative data is being released before study completion. The use of *QSIPrep*, *cloudknot* and the
516 containerization of the QC score assignment process facilitate running the exact pipeline described
517 in this paper on newly released participants. However, this approach is somewhat unsatisfactory
518 because it fails to anticipate improvements in preprocessing methodology. That is, what should we
519 do when *QSIPrep* is inevitably updated between HBN releases? Enforce standardization by using an
520 outdated pipeline or use state-of-the-art preprocessing at the expense of standardized processing
521 between releases? Because the use of *cloudknot* and AWS Spot Instances renders preprocessing
522 fast and relatively inexpensive, we propose a third way: if improvements to the preprocessing
523 pipeline are available with a new HBN release, we plan to execute the improved pipeline on the
524 entire HBN dataset, while preserving the previous baseline release in an archived BIDS derivative
525 dataset.

526 Undertaking the processing and QC effort to generate HBN-POD2 required construction and
527 deployment of substantial informatics infrastructure, including tools for cloud computing, web
528 applications for expert annotation and for community science rating and analysis software. All of
529 these tools are provided openly, so that this approach can be generalized even more widely in
530 other projects and in other scientific fields.

### Methods and Materials

532 To facilitate replicability, Jupyter notebooks (*Kluyver et al., 2016*) and Dockerfiles (*Merkel, 2014*) nec-
533 essary to reproduce the methods described herein are provided in the HBN-POD2 GitHub reposi-
534 tory at https://github.com/richford/hbn-pod2-qc. The specific version of the repository used in this
535 study is documented in *Richie-Halford and Rokem* (*2022a*). The `make` or `make help` commands will
536 list the available commands and `make build` will build the requisite Docker images to analyze HBN-
537 POD2 QC data. In order to separate data from analysis code (*Wilson et al., 2017*), we provide inter-
538 mediate data necessary to analyze the QC results in an OSF (*Foster and Deardorff, 2017*) project
539 (*Richie-Halford and Rokem, 2021*), which can be downloaded using the `make data` command in the
540 root of the HBN-POD2 GitHub repository. Most of the code in this repository uses Pandas (*McKin-*
541 *ney, 2010*; *pandas development team, 2020*), Numpy (*Harris et al., 2020*), Matplotlib (*Hunter, 2007*),
542 and Seaborn (*Waskom, 2021*).

### Inputs

544 Inputs for this study consisted of MRI data from the Healthy Brain Network pediatric mental health
545 study (*Alexander et al., 2017*), containing dMRI data from 2747 participants aged 5-21 years. These
546 data were measured using a 1.5 T Siemens mobile scanner on Staten Island (SI) and three fixed 3 T
547 Siemens MRI scanners at sites in the New York area: Rutgers University Brain Imaging Center (RU),
548 the CitiGroup Cornell Brain Imaging Center (CBIC), and the City University of New York Advanced
549 Science Research Center (CUNY). Informed consent was obtained from each participant aged 18 or
550 older. For participants younger than 18, written consent was obtained from their legal guardians
551 and written assent was obtained from the participant. Voxel resolution was $1.8\,\text{mm} \times 1.8\,\text{mm} \times$
552 $1.8\,\text{mm}$ with 64 non-colinear directions measured for each of $b = 1000$ s/mm$^2$ and $b = 2000$ s/mm$^2$.

### BIDS curation

We curated the imaging metadata for 2615 of the 2747 currently available HBN participants. Using dcm2bids and custom scripts, we conformed the data to the Brain Imaging Data Structure (BIDS; (*Gorgolewski et al., 2016*)) specification. The BIDS-curated dataset is available on FCP-INDI and can be accessed via AWS S3 at s3://fcp-indi/data/Projects/HBN/BIDS_curated/.

After conforming the data to BIDS, we used the "Curation of BIDS" (CuBIDS) package (*Covitz et al., 2022*) to identify unique combinations, or "variants" of imaging parameters in the curated dataset. CuBIDS is a Python-based software package that provides a sanity-preserving workflow to help users reproducibly parse, validate, curate, and understand heterogeneous BIDS imaging datasets. CuBIDS includes a robust implementation of the BIDS Validator that scales to large samples and incorporates DataLad (*Halchenko et al., 2021*), a distributed data management system, to ensure reproducibility and provenance tracking throughout the curation process. CuBIDS tools also employ agglomerative clustering to identify the aforementioned variants of imaging parameters. Users may then test BIDS-Apps on a subset of participants that represent the full range of acquisition parameters that are present. These variants are listed in the participants.tsv file in the BIDS-curated dataset.

### Preprocessing

We performed dMRI preprocessing on 2615 participants, using *QSIPrep* (*Cieslak et al., 2021*) 0.12.1, which is based on *Nipype* 1.5.1 (*Gorgolewski et al., 2011, 2018*), RRID:SCR_002502. *QSIPrep* a robust and scalable pipeline to group, distortion correct, motion correct, denoise, coregister and resample MRI scans. In total, 417 participants failed this preprocessing step, largely due to missing dMRI files. In keeping with the BIDS specification, the preprocessed dataset is available as a derivative dataset within the BIDS-curated dataset and can be access on AWS S3 at s3://fcp-indi/data/Projects/HBN/BIDS_curated/derivatives/qsiprep/. *QSIPrep* fosters reproducibility by automatically generating thorough methods boilerplate for later use in scientific publications, which we use for the remainder of this subsection to document each preprocessing step.

- *Anatomical data preprocessing* The T1-weighted (T1w) image was corrected for intensity non-uniformity (INU) using `N4BiasFieldCorrection` (*Tustison et al., 2010*) (ANTs 2.3.1), and used as T1w-reference throughout the workflow. The T1w-reference was then skull-stripped using `antsBrainExtraction.sh` (ANTs 2.3.1), using OASIS as target template. Spatial normalization to the ICBM 152 Nonlinear Asymmetrical template version 2009c (*Fonov et al.* (*2009b*), RRID:SCR_008796) was performed through nonlinear registration with `antsRegistration` (*Avants et al.* (*2008*), ANTs 2.3.1, RRID:SCR_004757), using brain-extracted versions of both T1w volume and template. Brain tissue segmentation of cerebrospinal fluid (CSF), white-matter (WM) and gray-matter (GM) was performed on the brain-extracted T1w using `FAST` (*Zhang et al., 2001*), FSL 6.0.3:b862cdd5, RRID:SCR_002823.

- *Diffusion data preprocessing*
  Any images with a $b$-value less than 100 s/mm$^2$ were treated as a $b = 0$ image. MP-PCA denoising as implemented in MRtrix3's `dwidenoise` (*Veraart et al., 2016*) was applied with a 5-voxel window. After MP-PCA, B1 field inhomogeneity was corrected using `dwibiascorrect` from MRtrix3 with the N4 algorithm (*Tustison et al., 2010*). After B1 bias correction, the mean intensity of the DWI series was adjusted so all the mean intensity of the $b = 0$ images matched across each separate DWI scanning sequence.
  FSL (version 6.0.3:b862cdd5)'s eddy was used for head motion correction and Eddy current correction (*Andersson and Sotiropoulos, 2016*). Eddy was configured with a $q$-space smoothing factor of 10, a total of 5 iterations, and 1000 voxels used to estimate hyperparameters. A linear first level model and a linear second level model were used to characterize Eddy current-related spatial distortion. $q$-space coordinates were forcefully assigned to shells. Field offset was attempted to be separated from participant movement. Shells were aligned post-

eddy. Eddy's outlier replacement was run (*Andersson et al., 2016*). Data were grouped by slice, only including values from slices determined to contain at least 250 intracerebral voxels. Groups deviating by more than four standard deviations from the prediction had their data replaced with imputed values. Data was collected with reversed phase-encode blips, resulting in pairs of images with distortions going in opposite directions. Here, $b = 0$ reference images with reversed phase encoding directions were used along with an equal number of $b = 0$ images extracted from the DWI scans. From these pairs the susceptibility-induced off-resonance field was estimated using a method similar to that described in (*Andersson et al., 2003*). The fieldmaps were ultimately incorporated into the Eddy current and head motion correction interpolation. Final interpolation was performed using the `jac` method.

Several confounding time-series were calculated based on the *preprocessed DWI*: framewise displacement (FD) using the implementation in *Nipype* following the definitions by (*Power et al., 2014*). The DWI time-series were resampled to ACPC, and their corresponding gradient directions were rotated accordingly to generate a *preprocessed DWI run in ACPC space*.

Many internal operations of *QSIPrep* use *Nilearn* 0.6.2 (*Abraham et al., 2014*), RRID:SCR_001362 and *DIPY* (*Garyfallidis et al., 2014*). For more details of the pipeline, see the section corresponding to workflows in *QSIPrep*'s documentation.

### Cloud-based distributed preprocessing

The containerization of *QSIPrep* provided a consistent preprocessing pipeline for each participant but the number of participants made serial processing of each participant prohibitive on a single machine. We used *cloudknot*, a previously developed cloud-computing library (*Richie-Halford and Rokem, 2018*) to parallelize the preprocessing over individual participants on spot instances in the Amazon Web Services Batch service. *Cloudknot* takes as input a user-defined Python function and creates the necessary AWS infrastructure to map that function onto a range of inputs, in this case, the participant IDs. The Python preprocessing function was a thin wrapper around *QSIPrep*'s command line interface and is provided in a Jupyter notebook in the HBN-POD2 GitHub repository in the "notebooks" directory. Using *cloudknot* and AWS Batch Spot Instances, the preprocessing cost less than $1.00 per participant.

### Expert quality control

The expert QC "gold standard" subset was created by randomly selecting 200 participants from the preprocessed dataset, sampled such that the proportional site distribution in the gold standard subset matched that of the preprocessed dataset.

We created a web application for expert quality control of preprocessed dMRI, called *dmriprep-viewer* (*Richie-Halford et al., 2022*). The viewer ingests *QSIPrep* outputs and generates a browser-based interface for expert QC. It provides a study overview displaying the distributions of *QSIPrep*'s automated data quality metrics (described at https://qsiprep.readthedocs.io/en/latest/preprocessing.html#quality-control-data). Each datum on the study overview page is interactively linked to a participant-level QC page that provides an interactive version of *QSIPrep*'s visual reports (described at https://qsiprep.readthedocs.io/en/latest/preprocessing.html#visual-reports). The viewer allows users to assign a rating of $-2$ (definitely fail), $-1$ (probably fail), $0$ (not sure), $1$ (probably pass), or $2$ (definitely pass) to a participant. To standardize rater expectations before rating, expert raters watched a tutorial video (available on YouTube at https://youtu.be/SQ0v-O-e5b8 and in the OSF project). They then rated each participant and saved their scores and sent them to the lead author for compilation.

To compute the pairwise Cohen's $\kappa$ scores in Figure 3e, we used the *scikit-learn* (*Pedregosa et al., 2011*) `cohen_kappa_score` function with quadratic weights. To compute intra-class correlation, we used the *pingouin* statistical package (*Vallat, 2018*) `intraclass_corr` function. The expert rating analysis can be replicated using the `make expert-qc` command in the HBN-POD2 GitHub repository.

The mean expert ratings were scaled to the range 0 to 1, so that a mean rating from 0 to 0.2 corresponds to an expert rating of "definitely fail", a mean rating from 0.2 to 0.4 corresponds to "probably fail", from 0.4 to 0.6 corresponds to "not sure", from 0.6 to 0.8 corresponds to "probably pass", and 0.8 to 1.0 corresponds to "definitely pass." These expert scores are available in the "expert_qc_score" column of the `participants.tsv` file on FCP-INDI.
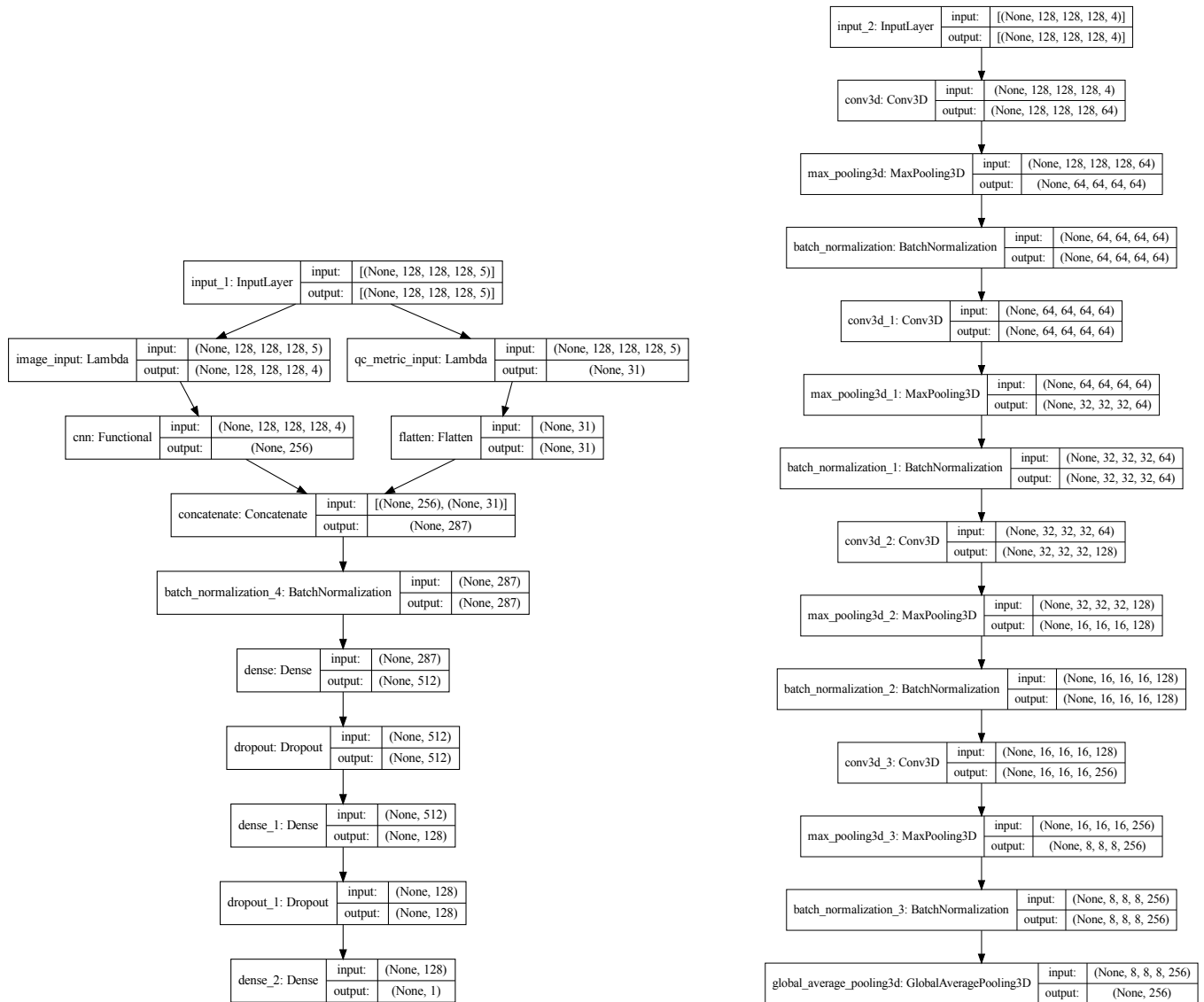
### Community scientist quality control

The community science web application is based on the SwipesForScience framework https://swipesforscience.org/, which generates a web application for community science given an open repository of images to be labelled and a configuration file. The source code for the *Fibr* web application is available at https://github.com/richford/fibr. After a brief tutorial, community scientists provided binary pass/fail ratings based on the DEC-FA from a fit of a DTI model to each participant's preprocessed dMRI data. These images were generated using a *DIPY* (*Garyfallidis et al., 2014*) `TensorModel` in a *cloudknot*-enabled Jupyter notebook that is available in the "notebooks" directory of the *Fibr* GitHub repository. *Fibr* saves each community rating to its Google Firebase backend, the contents of which have been archived to the HBN-POD2 OSF project.

The *Fibr* ratings were then combined with 31 automated *QSIPrep* data quality metrics to train the gradient boosted trees models XGB, XGB-f, and XGB-q. See Appendix 3 for a list of these automated QC metrics and a measure of their global feature importance in the XGB and XGB-q models. These models were implemented as binary classifiers using the XGBoost library (*Chen and Guestrin, 2016b*). The targets for these classifiers were the mean expert ratings in the gold standard dataset, rescaled to the range $[0, 1]$ and binarized with a threshold of $0.5$. Using repeated stratified K-fold cross-validation, with three splits and two repeats, we evaluated the models' performance in predicting the gold standard ratings. In each fold, the best model hyperparameters were chosen using the scikit-optimize (*Head et al., 2021*) `BayesSearchCV` class. Saved model checkpoints for each cross-validation split are available in the HBN-POD2 OSF project. Since each split resulted in a different XGB model and we required a single QC score to train the deep learning model, we combined the models from each cross-validation split using a voting classifier, computing a weighted averaged of the predicted probability of passing from each model, weighted by its out-of-sample ROC-AUC. This was implemented using scikit-learn's `VotingClassifier` class. Treating the voting classifier as another "expert" rater, we reassessed the pairwise Cohen's $\kappa$ and ICC scores as in the expert QC subsection. The community ratings analysis can be replicated using the `make community-qc` command in the HBN-POD2 GitHub repository. The XGB model's positive class probabilities are available in the "xgb_qc_score" column of the `participants.tsv` file on FCP-INDI, while the XGB-q model's positive class probabilities are available in the "xgb_qsiprep_qc_score" column.

### Deep learning to predict quality control

The binarized voting classifier's predictions were then used as targets to train a deep learning binary classifier to predict QC scores based on each participant's preprocessed dMRI data. We trained two different model architectures: (i) CNN-i, which took only preprocessed dMRI images as input and (ii) CNN-i+q, whose input also included *QSIPrep*'s automated data quality metrics. Both models were implemented in Tensorflow 2 (*Abadi et al., 2015*) using the Keras module (*Chollet et al., 2015*). The image processing part of the model architecture was identical for both models: a modification of an existing 3D CNN (*Zunair et al., 2020*) previously applied to assess tuberculosis severity (*Dicente Cid et al., 2019*). It accepts a 3D volume as input with four channels: (i) the $b = 0$ reference volume, (ii) DEC-FA in the $x$-direction, (iii) DEC-FA in the $y$-direction and (iv) DEC-FA in the $z$-direction. The *QSIPrep*'s automated QC metrics were included as an additional fifth channel. The CNN-i+q model architecture is summarized in Figure 10. Upon input, the CNN-i+q model extracts the imaging channels and passes them through the CNN architecture. The remaining data quality metrics channel is flattened and passed "around" the CNN architecture and concatenated with the output of the convolutional layers. This concatenated output is then passed through a

**(a)** Slicing and combining the input channels

**(b)** CNN architecture

**Figure 10. Deep learning model architecture**: **(a)** The CNN-i+q model accepts multichannel input that combined four imaging channels with a fifth channel containing 31 *QSIPrep* automated data quality metrics. The imaging channels are separated from the data quality channel using `Lambda` layers. The imaging channels are passed through a CNN **(b)**, the output of which is concatenated with the data quality metrics, batch normalized and passed through two fully-connected (FC) layers, with rectified linear unit (ReLu) activation functions and with 512 and 128 units respectively. Each FC layer is followed by a dropout layer which drops 40% of the input units. The final layer contains a single unit with a sigmoid activation function and outputs the probability of passing QC. **(b)** The CNN portion of the model passes the imaging input through four convolutional blocks. Each block consists of a 3D convolutional layer with a kernel size of 3 and a ReLu activation, a 3D max pooling layer with a pool size of 2, and a batch normalization layer with Tensorflow's default parameters. The number of filters in the convolutional layers in each block are 64, 64, 128, and 256 respectively. The output of the final block is passed through a 3D global average pooling layer with Tensorflow's default parameters.

**Figure 10–Figure supplement 1.** Deep learning model loss curves

699  fully-connected layer to produce a single output, the probability of passing QC. This architecture
700  has 1,438,783 trainable parameters.
701      We used *DIPY* (*Garyfallidis et al., 2014*) and *cloudknot* (*Richie-Halford and Rokem, 2018*) to gen-
702  erate these multichannel volumes for each participant and save them as NIfTI-1 files (*Cox et al.,*
703  *2004*). These NIfTI files were then converted to the Tensorflow TFRecord format using the *Nobrainer*
704  deep learning framework (*Kaczmarzyk et al., 2021*). The Jupyter notebooks used to create these
705  NIfTI and TFRecord files are available in the "notebooks" directory of the *Fibr* GitHub repository.
706      We trained each model using the Google Cloud AI Platform Training service; the HBN-POD2
707  GitHub repository contains Docker services to launch training (with `make dl-train`) and prediction
708  (with `make dl-predict`) jobs on Google Cloud, if the user has provided the appropriate credentials
709  in an environment file and placed the TFRecord files on Google Cloud Storage. To estimate the
710  variability in model training, we trained ten separate models using different training and valida-
711  tion splits of the data. The gold standard dataset was not included in any of these splits and was
712  reserved for reporting final model performance. Models were optimized for binary crossentropy
713  loss using the Adam optimizer (*Kingma and Ba, 2017*) with an initial learning rate of 0.0001. We
714  reduced the learning rate by a factor of 0.5 when the validation loss plateaued for more than two
715  epochs. We also stopped training when the validation loss failed to improve by more than 0.001
716  for twenty consecutive epochs. These two adjustments were made using the `ReduceLROnPlateau`
717  and `EarlyStopping` callbacks in Tensorflow 2 (*Abadi et al., 2015*) respectively. The training and
718  validation loss curves for both the CNN-i and CNN-i+q models are depicted in *Figure 10–Figure*
719  *Supplement 1*. While the CNN-i+q model achieved better validation loss, it did not outperform the
720  CNN-i model on the held out gold standard dataset. The CNN-i+q model's positive class probabili-
721  ties are available in the "dl_qc_score" column of the `participants.tsv` file on FCP-INDI.
722      To generate the attribution maps, we followed Tensorflow's integrated gradients tutorial (*Ten-*
723  *sorFlow Authors, 2021*) with a black baseline image and 128 steps in the Riemann sum approxima-
724  tion of the integral (i.e. `m_steps = 128`). In the HBN-POD2 GitHub repository, we provide a Docker
725  service to compute integrated gradient attribution maps on Google Cloud, which can be invoked
726  using the `make dl-integrated-gradients` command.

## Site generalization experiments

728  To simulate the generalization of the XGB-q and CNN-i models to new scanning sites, we trained
729  multiple versions of XGB-q and CNN-i with different scanning sites held out and then evaluated
730  those models on the held out sites. These models were therefore evaluated on data from "un-
731  seen" sites. We constructed these train/evaluate splits from combinations of the HBN sites with
732  3T scanners (RU, CBIC, and CUNY), and excluded CUNY as a standalone training or test site be-
733  cause of its low number of participants ($N = 74$). This left four combinations of site-generated
734  training splits: CBIC + CUNY (eval: RU), CBIC (eval: RU + CUNY), RU + CUNY (eval: CBIC), and RU
735  (eval: CBIC + CUNY).
736      We trained eight models (with distinct random seeds) from the CNN-i family of models using
737  the global XGB scores as targets, just as with the full CNN-i model. Similarly, we trained twenty
738  models (with distinct random seeds) from the XGB-q family of models using the expert scores as
739  targets, just as with the full XGB-q model. For each model, we reported three evaluation metrics:
740  ROC-AUC, accuracy, and balanced accuracy. Because the distribution of QC scores was imbalanced
741  (Figures 3a and 5d), we included balanced accuracy as an evaluation metric. Balanced accuracy
742  avoids inflated accuracy estimates on imbalanced data (*Velez et al., 2007*), and in the binary clas-
743  sification case, it is the mean of the sensitivity and specificity. For the CNN-i family, we further
744  decomposed the evaluation split into a report set, for which expert scores were available, and a
745  test set, with participants who were not in the "gold standard" dataset. For the report set, we eval-
746  uated the model using the expert scores as the ground truth. For the test set, we evaluated each
747  model using the XGB scores as ground truth.
748      Aside from the specification of train and evaluation splits, model training followed exactly the

749 same procedure as for the full dataset. For example, we use the same cross validation and hyperpa-
750 rameter optimization procedure for the XGB-q family as for the original XGB-q model and the same
751 architecture, input format, and early stopping criteria for the CNN-i family as for the CNN-i model.
752 In the HBN-POD2 GitHub repository, we provide a Docker service to conduct the CNN-i site general-
753 ization experiments Google Cloud, which can be invoked using the `make dl-site-generalization`
754 command. The XGB-q site generalization experiments can be replicated locally using the `make`
755 `site-generalization` command, which will also plot the results of the CNN-i experiments.

## QC bundle profiles

757 To generate bundle profiles, reconstruction was performed using the *QSIprep* 0.12.1 preconfigured
758 reconstruction workflow `mrtrix_multishell_msmt`, modified to generate two million streamlines
759 rather than the default ten million. Multi-tissue fiber response functions were estimated using the
760 dhollander algorithm. Fiber orientation distributions (FODs) were estimated via constrained spher-
761 ical deconvolution (CSD, (*Tournier et al., 2004, 2008*)) using an unsupervised multi-tissue method
762 (*Dhollander et al., 2019, 2016*). Reconstruction was done using MRtrix3 (*J-Donald et al., 2019*). FODs
763 were intensity-normalized using mtnormalize (*Raffelt et al., 2017*).

764 These tractograms were then used as input to the Python Automated Fiber Quantification tool-
765 box (pyAFQ) (*Kruper et al., 2021*). Twenty-four major tracts, which are enumerated in Figure 8,
766 were identified using multiple criteria: streamlines are selected as candidates for inclusion in a
767 bundle of streamlines that represents a tract if they pass through known inclusion ROIs and do
768 not pass through exclusion ROIs (*Wakana et al., 2007*). In addition, a probabilistic atlas is used
769 to exclude streamlines which are unlikely to be part of a tract and to adjudicate in cases where a
770 streamline could belong to more than one tract (*Hua et al., 2008*). Each streamline is resampled
771 to 100 nodes and the robust mean at each location is calculated by estimating the 3D covariance
772 of the location of each node and excluding streamlines that are more than 5 standard deviations
773 from the mean location in any node. Finally, a bundle profile of tissue properties in each bundle
774 was created by interpolating the value of MRI maps of these tissue properties to the location of
775 the nodes of the resampled streamlines designated to each bundle. In each of 100 nodes, the val-
776 ues are summed across streamlines, weighting the contribution of each streamline by the inverse
777 of the mahalanobis distance of the node from the average of that node across streamlines. This
778 means that streamlines that are more representative of the tract contribute more to the bundle
779 profile, relative to streamlines that are on the edge of the tract.

780 These processes create bundle profiles, in which diffusion measures are quantified and av-
781 eraged along twenty-four major fiber tracts. We retain only the mean diffusivity (MD) and the
782 fractional anisotropy (FA) from a diffusion kurtosis imaging (DKI) model (*Jensen et al., 2005*), im-
783 plemented in DIPY (*Henriques et al., 2021*), and impute missing bundles using median imputation
784 as implemented by *scikit-learn*'s `SimpleImputer` class. Because the HBN-POD2 bundle profiles ex-
785 hibit strong site effects (*Richie-Halford et al., 2021*), we used the ComBat harmonization method
786 to robustly adjust for site effects in the tract profiles. Initially designed to correct for site effects
787 in gene expression studies (*Johnson et al., 2007*), ComBat employs a parametric empirical Bayes
788 approach to adjust for batch effects and has since been applied to multi-site cortical thickness
789 measurements (*Fortin et al., 2018*), multi-site DTI studies (*Fortin et al., 2017*), and functional MRI
790 data in the Adolescent Brain Cognitive Development Study (ABCD) (*Nielson et al., 2018*). We rely
791 on the *neurocombat_sklearn* library (*Pinaya, 2020*), to apply ComBat in before plotting bundle pro-
792 files in Figure 8 using plotting functions from the AFQ-Insight package (*Richie-Halford et al., 2019*).
793 The bundle profile analysis can be replicated using the `make bundle-profiles` command in the
794 HBN-POD2 GitHub repository.

## Brain age prediction

796 We evaluated the effect of varying the QC cutoff on model performance by observing cross-validated
797 $R^2$ values of gradient boosted trees models implemented using XGBoost. The input feature space

798 for each model consisted of 4800 features per participant, comprising 100 nodes for each of MD
799 and FA in the twenty-four major tracts. We imputed missing bundles and harmonized the different
800 scanning sites as above. The XGBoost models' hyperparameters were hand-tuned to values that
801 have been performant in the authors' previous experience. Within the limited age range of the HBN
802 study, MD and FA follow logarithmic maturation trajectories (*Yeatman et al., 2014*). We therefore
803 log-transformed each participant's age before prediction using the `TransformedTargetRegressor`
804 class from *scikit-learn* . For each value of the QC cutoff between 0 and 0.95, in steps of 0.05, we
805 computed the cross-validated $R^2$ values using *scikit-learn*'s `cross_val_score` function with repeated
806 K-fold cross-validation using five folds and five repeats.

## Author contributions statement

808 The last two authors named share senior authorship. The first two authors named share lead
809 authorship. The remaining authors are listed in alphabetical order, with the exception of the *Fibr*
810 Community Science Consortium, whose members provided community science QC ratings and are
811 listed in Appendix 4. We describe contributions to the paper using the CRediT taxonomy (*Brand*
812 *et al., 2015*; *Allen et al., 2014*): Conceptualization: A.R-H., A.R., T.S., and M.C.; Methodology: A.R-H.
813 and A.R.; Software: A.R-H., M.C., and S.C.; Validation: A.R-H., M.C., and S.C.; Formal Analysis: A.R-H.
814 and M.C.; Investigation: A.R-H. and M.C.; Resources: A.R., T.S., and M.M.; Data Curation: S.C., M.C.,
815 V.J.S., I.I.K., B.A-P. and L.A.; Writing – Original Draft: A.R-H. and A.R.; Writing – Review & Editing: A.R-
816 H., A.R., M.C., A.F., T.S., V.J.S., I.I.K, B.A-P., and S.C.; Visualization: A.R-H.; Supervision: A.R. and T.S.;
817 Project Administration: A.R-H. and A.R.; Funding Acquisition: A.R. and T.S.

## Acknowledgments

## References

829 **Abadi M**, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin M, Ghemawat
830 S, Goodfellow I, Harp A, Irving G, Isard M, Jia Y, Jozefowicz R, Kaiser L, Kudlur M, Levenberg J, et al., Tensor-
831 Flow: Large-Scale Machine Learning on Heterogeneous Systems; 2015. https://www.tensorflow.org/, software
832 available from tensorflow.org.

833 **Abraham A**, Pedregosa F, Eickenberg M, Gervais P, Mueller A, Kossaifi J, Gramfort A, Thirion B, Varoquaux G.
834 Machine learning for neuroimaging with scikit-learn. Frontiers in Neuroinformatics. 2014; 8. https://www.
835 frontiersin.org/articles/10.3389/fninf.2014.00014/full, doi: 10.3389/fninf.2014.00014.

836 **Alexander LM**, Escalera J, Ai L, Andreotti C, Febre K, Mangone A, Vega-Potler N, Langer N, Alexander A, Kovacs
837 M, Litke S, O'Hagan B, Andersen J, Bronstein B, Bui A, Bushey M, Butler H, Castagna V, Camacho N, Chan
838 E, et al. An open resource for transdiagnostic research in pediatric mental health and learning disorders.
839 Scientific Data. 2017 Dec; 4:170181. https://doi.org/10.1038/sdata.2017.181, doi: 10.1038/sdata.2017.181.

840 **Allen L**, Scott J, Brand A, Hlava M, Altman M. Publishing: Credit where credit is due. Nature. 2014 Apr;
841 508(7496):312–313. http://dx.doi.org/10.1038/508312a, doi: 10.1038/508312a.

842 **Ancona M**, Ceolini E, Öztireli C, Gross M. Gradient-Based Attribution Methods. In: Samek W, Montavon G,
843 Vedaldi A, Hansen LK, Müller KR, editors. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*
844 Cham: Springer International Publishing; 2019.p. 169–191. https://doi.org/10.1007/978-3-030-28954-6_9,
845 doi: 10.1007/978-3-030-28954-6\_9.

846 **Andersson JL**, Graham MS, Zsoldos E, Sotiropoulos SN. Incorporating outlier detection and replacement into
847 a non-parametric framework for movement and distortion correction of diffusion MR images. Neuroimage.
848 2016; 141:556–572.

849 **Andersson JL**, Skare S, Ashburner J. How to correct susceptibility distortions in spin-echo echo-planar images:
850 application to diffusion tensor imaging. Neuroimage. 2003; 20(2):870–888.

851 **Andersson JL**, Sotiropoulos SN. An integrated approach to correction for off-resonance effects and subject
852 movement in diffusion MR imaging. Neuroimage. 2016; 125:1063–1078.

853 **Avants BB**, Epstein CL, Grossman M, Gee JC. Symmetric diffeomorphic image registration with cross-
854 correlation: Evaluating automated labeling of elderly and neurodegenerative brain. Medical Image
855 Analysis. 2008; 12(1):26–41. http://www.sciencedirect.com/science/article/pii/S1361841507000606, doi:
856 10.1016/j.media.2007.06.004.

857 **Avesani P**, McPherson B, Hayashi S, Caiafa CF, Henschel R, Garyfallidis E, Kitchell L, Bullock D, Patterson A,
858 Olivetti E, Sporns O, Saykin AJ, Wang L, Dinov I, Hancock D, Caron B, Qian Y, Pestilli F. The open diffusion data
859 derivatives, brain data upcycling via integrated publishing of derivatives and reproducible open cloud ser-
860 vices. Scientific data. 2019 May; 6(1):69. http://dx.doi.org/10.1038/s41597-019-0073-y, doi: 10.1038/s41597-
861 019-0073-y.

862 **Bells S**, Cercignani M, Deoni S, Assaf Y, Pasternak O, Evans C, Leemans A, Jones D. Tractometry–comprehensive
863 multi-modal quantitative assessment of white matter along specific tracts. In: *Proc. ISMRM*, vol. 678; 2011.
864 p. 1.

865 **Brand A**, Allen L, Altman M, Hlava M, Scott J. Beyond authorship: attribution, contribution, collaboration, and
866 credit. Learned publishing: journal of the Association of Learned and Professional Society Publishers. 2015
867 Apr; 28(2):151–155. http://doi.wiley.com/10.1087/20150211, doi: 10.1087/20150211.

868 **Chen T**, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: *Proceedings of the 22nd ACM SIGKDD Interna-
869 tional Conference on Knowledge Discovery and Data Mining* KDD '16, New York, NY, USA: Association for Comput-
870 ing Machinery; 2016. p. 785–794. https://doi.org/10.1145/2939672.2939785, doi: 10.1145/2939672.2939785.

871 **Chen T**, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: *Proceedings of the 22nd ACM SIGKDD Interna-
872 tional Conference on Knowledge Discovery and Data Mining* KDD '16, New York, NY, USA: ACM; 2016. p. 785–794.
873 http://doi.acm.org/10.1145/2939672.2939785, doi: 10.1145/2939672.2939785.

874 **Chollet F**, et al., Keras; 2015. https://keras.io.

875 **Cicchetti DV**. Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment
876 instruments in psychology. Psychol Assess. 1994 Dec; 6(4):284–290.

877 **Cieslak M**, Cook PA, He X, Yeh FC, Dhollander T, Adebimpe A, Aguirre GK, Bassett DS, Betzel RF, Bourque J,
878 Cabral LM, Davatzikos C, Detre JA, Earl E, Elliott MA, Fadnavis S, Fair DA, Foran W, Fotiadis P, Garyfallidis E,
879 et al. QSIPrep: an integrative platform for preprocessing and reconstructing diffusion MRI data. Nature
880 methods. 2021 Jul; 18(7):775–778. http://dx.doi.org/10.1038/s41592-021-01185-5, doi: 10.1038/s41592-021-
881 01185-5.

882 **Colby JB**, Soderberg L, Lebel C, Dinov ID, Thompson PM, Sowell ER. Along-tract statistics allow for enhanced
883 tractography analysis. NeuroImage. 2012 Feb; 59(4):3227–3242. http://www.sciencedirect.com/science/article/
884 pii/S1053811911012833, doi: 10.1016/j.neuroimage.2011.11.004.

885 **Cole JH**, Marioni RE, Harris SE, Deary IJ. Brain age and other bodily 'ages': implications for neuropsychi-
886 atry. Molecular psychiatry. 2019 Feb; 24(2):266–281. http://dx.doi.org/10.1038/s41380-018-0098-1, doi:
887 10.1038/s41380-018-0098-1.

888 **Covitz S**, Tapera T, Adebimpe A, Ai L, Alexander-Bloch A, Bertolero M, Fair D, Feczko E, Franco A, Gur R, Gur
889 R, Hendrickson T, Houghton A, Mehta K, Murtha K, Perrone A, Robert-Fitzgerald T, Schabdach J, Shinohara
890 R, Vogel J, et al. Curation of BIDS (CuBIDS): a sanity-preserving workflow and software package for curating
891 large BIDS datasets; 2022, in preparation.

892 **Cox RW**, Ashburner J, Breman H, Fissell K, Haselgrove C, Holmes CJ, Lancaster JL, Rex DE, Smith SM, Woodward
893 JB, Strother SC. A (sort of) new image data format standard: NiFTI-1. In: *10th Annual Meeting of the Organization
894 for Human Brain Mapping*; 2004. https://nifti.nimh.nih.gov/nifti-1/documentation/hbm_nifti_2004.pdf.

895 **Dhollander T**, Mito R, Raffelt D, Connelly A. Improved white matter response function estimation for 3-tissue
896 constrained spherical deconvolution. In: *Proc. Intl. Soc. Mag. Reson. Med*; 2019. p. 555.

**Dhollander T**, Raffelt D, Connelly A. Unsupervised 3-tissue response function estimation from single-shell or multi-shell diffusion MR data without a co-registered T1 image. In: *ISMRM Workshop on Breaking the Barriers of Diffusion MRI*, vol. 5; 2016. p. 5.

**Di Eugenio B**, Glass M. The kappa statistic: a second look. Computational Linguistics. 2004 Mar; 30(1):95–101. https://doi.org/10.1162/089120104773633402, doi: 10.1162/089120104773633402.

**Dicente Cid Y**, Liauchuk V, Klimuk D, Tarasau A, Kovalev V, Müller H. Overview of ImageCLEFtuberculosis 2019 - Automatic CT-based Report Generation and Tuberculosis Severity Assessment. In: *CLEF*; 2019. .

**Esteban O**, Birman D, Schaer M, Koyejo OO, Poldrack RA, Gorgolewski KJ. MRIQC: Advancing the automatic prediction of image quality in MRI from unseen sites. PloS one. 2017 Sep; 12(9):e0184661. http://dx.doi.org/10.1371/journal.pone.0184661, doi: 10.1371/journal.pone.0184661.

**Fair DA**, Dosenbach NU, Moore AH, Satterthwaite T, Milham MP. Developmental Cognitive Neuroscience in the Era of Networks and Big Data: Strengths, Weaknesses, Opportunities, and Threats. Annual Review of Developmental Psychology. 2021; 3.

**Fair DA**, Nigg JT, Iyer S, Bathula D, Mills KL, Dosenbach NUF, Schlaggar BL, Mennes M, Gutman D, Bangaru S, Buitelaar JK, Dickstein DP, Di Martino A, Kennedy DN, Kelly C, Luna B, Schweitzer JB, Velanova K, Wang YF, Mostofsky S, et al. Distinct neural signatures detected for ADHD subtypes after controlling for micro-movements in resting state functional connectivity MRI data. Frontiers in systems neuroscience. 2012; 6:80. http://dx.doi.org/10.3389/fnsys.2012.00080, doi: 10.3389/fnsys.2012.00080.

**Fonov VS**, Evans AC, McKinstry RC, Almli CR, Collins DL. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. Neuroimage. 2009 Jul; 47:S102.

**Fonov V**, Evans A, McKinstry R, Almli C, Collins D. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. NeuroImage. 2009; 47, Supplement 1:S102. http://www.sciencedirect.com/science/article/pii/S1053811909708845, doi: 10.1016/S1053-8119(09)70884-5.

**Fortin JP**, Cullen N, Sheline YI, Taylor WD, Aselcioglu I, Cook PA, Adams P, Cooper C, Fava M, McGrath PJ, McInnis M, Phillips ML, Trivedi MH, Weissman MM, Shinohara RT. Harmonization of cortical thickness measurements across scanners and sites. NeuroImage. 2018 Feb; 167:104–120. http://dx.doi.org/10.1016/j.neuroimage.2017.11.024, doi: 10.1016/j.neuroimage.2017.11.024.

**Fortin JP**, Parker D, Tunç B, Watanabe T, Elliott MA, Ruparel K, Roalf DR, Satterthwaite TD, Gur RC, Gur RE, Schultz RT, Verma R, Shinohara RT. Harmonization of multi-site diffusion tensor imaging data. NeuroImage. 2017 Nov; 161:149–170. http://dx.doi.org/10.1016/j.neuroimage.2017.08.047, doi: 10.1016/j.neuroimage.2017.08.047.

**Foster ED MSLS**, Deardorff A MLIS. Open Science Framework (OSF). Journal of the Medical Library Association: JMLA. 2017 Apr; 105(2). http://jmla.pitt.edu/ojs/jmla/article/view/88, doi: 10.5195/jmla.2017.88.

**Franke K**, Ziegler G, Klöppel S, Gaser C, Alzheimer's Disease Neuroimaging Initiative. Estimating the age of healthy subjects from T1-weighted MRI scans using kernel methods: exploring the influence of various parameters. NeuroImage. 2010 Apr; 50(3):883–892. http://dx.doi.org/10.1016/j.neuroimage.2010.01.005, doi: 10.1016/j.neuroimage.2010.01.005.

**Garyfallidis E**, Brett M, Amirbekian B, Rokem A, Van Der Walt S, Descoteaux M, Nimmo-Smith I. DIPY, a library for the analysis of diffusion MRI data. Frontiers in neuroinformatics. 2014; 8:8.

**Ghassemi M**, Oakden-Rayner L, Beam AL. The false hope of current approaches to explainable artificial intelligence in health care. The Lancet Digital health. 2021 Nov; 3(11):e745–e750. http://dx.doi.org/10.1016/S2589-7500(21)00208-9, doi: 10.1016/S2589-7500(21)00208-9.

**Glasser MF**, Sotiropoulos SN, Wilson JA, Coalson TS, Fischl B, Andersson JL, Xu J, Jbabdi S, Webster M, Polimeni JR, Van Essen DC, Jenkinson M, WU-Minn HCP Consortium. The minimal preprocessing pipelines for the Human Connectome Project. NeuroImage. 2013 Oct; 80:105–124. http://dx.doi.org/10.1016/j.neuroimage.2013.04.127, doi: 10.1016/j.neuroimage.2013.04.127.

**Gorgolewski K**, Burns CD, Madison C, Clark D, Halchenko YO, Waskom ML, Ghosh S. Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in Python. Frontiers in Neuroinformatics. 2011; 5:13. doi: 10.3389/fninf.2011.00013.

946 **Gorgolewski KJ**, Alfaro-Almagro F, Auer T, Bellec P, Capotă M, Chakravarty MM, Churchill NW, Cohen AL, Crad-
947 dock RC, Devenyi GA, Eklund A, Esteban O, Flandin G, Ghosh SS, Guntupalli JS, Jenkinson M, Keshavan A,
948 Kiar G, Liem F, Raamana PR, et al. BIDS apps: Improving ease of use, accessibility, and reproducibility of
949 neuroimaging data analysis methods. PLoS Comput Biol. 2017 Mar; 13(3):e1005209.

950 **Gorgolewski KJ**, Auer T, Calhoun VD, Craddock RC, Das S, Duff EP, Flandin G, Ghosh SS, Glatard T, Halchenko
951 YO, Handwerker DA, Hanke M, Keator D, Li X, Michael Z, Maumet C, Nichols BN, Nichols TE, Pellman J, Po-
952 line JB, et al. The brain imaging data structure, a format for organizing and describing outputs of neu-
953 roimaging experiments. Scientific data. 2016 Jun; 3:160044. http://dx.doi.org/10.1038/sdata.2016.44, doi:
954 10.1038/sdata.2016.44.

955 **Gorgolewski KJ**, Esteban O, Markiewicz CJ, Ziegler E, Ellis DG, Notter MP, Jarecka D, Johnson H, Burns C,
956 Manhães-Savio A, Hamalainen C, Yvernault B, Salo T, Jordan K, Goncalves M, Waskom M, Clark D, Wong J,
957 Loney F, Modat M, et al. Nipype. Software. 2018; doi: 10.5281/zenodo.596855.

958 **Halchenko YO**, Meyer K, Poldrack B, Solanky DS, Wagner AS, Gors J, MacFarlane D, Pustina D, Sochat V,
959 Ghosh SS, Mönch C, Markiewicz CJ, Waite L, Shlyakhter I, de la Vega A, Hayashi S, Häusler CO, Poline JB,
960 Kadelka T, Skytén K, et al. DataLad: distributed system for joint management of code, data, and their re-
961 lationship. Journal of Open Source Software. 2021; 6(63):3262. https://doi.org/10.21105/joss.03262, doi:
962 10.21105/joss.03262.

963 **Hallgren KA**. Computing Inter-Rater Reliability for Observational Data: An Overview and Tutorial. Tutorials
964 in quantitative methods for psychology. 2012; 8(1):23–34. http://dx.doi.org/10.20982/tqmp.08.1.p023, doi:
965 10.20982/tqmp.08.1.p023.

966 **Harris CR**, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith
967 NJ, Kern R, Picus M, Hoyer S, van Kerkwijk MH, Brett M, Haldane A, del Río JF, Wiebe M, Peterson P, Gérard-
968 Marchant P, et al. Array programming with NumPy. Nature. 2020 Sep; 585(7825):357–362. https://doi.org/
969 10.1038/s41586-020-2649-2, doi: 10.1038/s41586-020-2649-2.

970 **Head T**, Kumar M, Nahrstaedt H, Louppe G, Shcherbatyi I, scikit-optimize/scikit-optimize. Zenodo; 2021. https:
971 //doi.org/10.5281/zenodo.5565057, doi: 10.5281/zenodo.5565057.

972 **Henriques RN**, Correia MM, Marrale M, Huber E, Kruper J, Koudoro S, Yeatman JD, Garyfallidis E, Rokem A.
973 Diffusional Kurtosis Imaging in the Diffusion Imaging in Python Project. Front Hum Neurosci. 2021; 15:390.

974 **Hua K**, Zhang J, Wakana S, Jiang H, Li X, Reich DS, Calabresi PA, Pekar JJ, van Zijl PCM, Mori S. Tract
975 probability maps in stereotaxic spaces: analyses of white matter anatomy and tract-specific quantifi-
976 cation. NeuroImage. 2008 Jan; 39(1):336–347. http://dx.doi.org/10.1016/j.neuroimage.2007.07.053, doi:
977 10.1016/j.neuroimage.2007.07.053.

978 **Hunter JD**. Matplotlib: A 2D graphics environment. Computing in Science & Engineering. 2007; 9(3):90–95. doi:
979 10.1109/MCSE.2007.55.

980 **J-Donald**, Smith R, Raffelt D, Tabbara R, Dhollander T, Pietsch M, Christiaens D, Jeurissen B, Yeh CH, Connelly
981 A. MRtrix3: A fast, flexible and open software framework for medical image processing and visualisation.
982 NeuroImage. 2019; 202:116137.

983 **Jensen JH**, Helpern JA, Ramani A, Lu H, Kaczynski K. Diffusional kurtosis imaging: the quantification of non-
984 gaussian water diffusion by means of magnetic resonance imaging. Magnetic resonance in medicine: official
985 journal of the Society of Magnetic Resonance in Medicine / Society of Magnetic Resonance in Medicine. 2005
986 Jun; 53(6):1432–1440. http://dx.doi.org/10.1002/mrm.20508, doi: 10.1002/mrm.20508.

987 **Jernigan TL**, Brown SA. Introduction. Developmental cognitive neuroscience. 2018; 32:1–3. http://www.
988 sciencedirect.com/science/article/pii/S1878929317301883, doi: 10.1016/j.dcn.2018.02.002.

989 **Johnson WE**, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical
990 Bayes methods. Biostatistics. 2007 Jan; 8(1):118–127. http://dx.doi.org/10.1093/biostatistics/kxj037, doi:
991 10.1093/biostatistics/kxj037.

992 **Jones DK**, Cercignani M. Twenty-five pitfalls in the analysis of diffusion MRI data. NMR in biomedicine. 2010
993 Aug; 23(7):803–820. http://dx.doi.org/10.1002/nbm.1543, doi: 10.1002/nbm.1543.

994 **Jones DK**, Travis AR, Eden G, Pierpaoli C, Basser PJ. PASTA: pointwise assessment of streamline tractography
995 attributes. Magnetic resonance in medicine: official journal of the Society of Magnetic Resonance in Medicine
996 / Society of Magnetic Resonance in Medicine. 2005 Jun; 53(6):1462–1467. http://dx.doi.org/10.1002/mrm.
997 20484, doi: 10.1002/mrm.20484.

**998** **Kaczmarzyk J**, McClure P, Zulfikar W, Rana A, Rajaei H, Richie-Halford A, Bansal S, Jarecka D, Lee J, Ghosh
**999** S, neuronets/nobrainer: 0.2.0. Zenodo; 2021. https://doi.org/10.5281/zenodo.5803350, doi: 10.5281/zen-
**1000** odo.5803350.

**1001** **Keshavan A**, Yeatman JD, Rokem A. Combining Citizen Science and Deep Learning to Amplify Expertise in
**1002** Neuroimaging. Frontiers in neuroinformatics. 2019 May; 13:29. http://dx.doi.org/10.3389/fninf.2019.00029,
**1003** doi: 10.3389/fninf.2019.00029.

**1004** **Kingma DP**, Ba J, Adam: A Method for Stochastic Optimization; 2017.

**1005** **Kluyver T**, Ragan-Kelley B, Pérez F, Granger B, Bussonnier M, Frederic J, Kelley K, Hamrick J, Grout J, Corlay S,
**1006** Ivanov P, Avila D, Abdalla S, Willing C, Jupyter Development Team. Jupyter Notebooks – a publishing format
**1007** for reproducible computational workflows. In: *Positioning and Power in Academic Publishing: Players, Agents*
**1008** *and Agendas* Amsterdam, NY: IOS Press; 2016.p. 87–90. https://ebooks.iospress.nl/publication/42900, doi:
**1009** 10.3233/978-1-61499-649-1-87.

**1010** **Kruper J**, Yeatman JD, Richie-Halford A, Bloom D, Grotheer M, Caffarra S, Kiar G, Karipidis II, Roy E, Chan-
**1011** dio BQ, Garyfallidis E, Rokem A. Evaluating the reliability of human brain white matter tractome-
**1012** try. Aperture Neuro. 2021 Oct; https://www.biorxiv.org/content/early/2021/02/24/2021.02.24.432740, doi:
**1013** 10.1101/2021.02.24.432740.

**1014** **Laird AR**. Large, open datasets for human connectomics research: Considerations for reproducible and re-
**1015** sponsible data use. NeuroImage. 2021 Dec; 244:118579. http://dx.doi.org/10.1016/j.neuroimage.2021.118579,
**1016** doi: 10.1016/j.neuroimage.2021.118579.

**1017** **Lebel C**, Deoni S. The development of brain white matter microstructure. Neuroimage. 2018 Nov; 182:207–218.

**1018** **LeCun Y**, Bengio Y, Hinton G. Deep learning. Nature. 2015 May; 521(7553):436–444. http://dx.doi.org/10.1038/
**1019** nature14539, doi: 10.1038/nature14539.

**1020** **Lipton ZC**, The Doctor Just Won't Accept That!; 2017. http://arxiv.org/abs/1711.08037.

**1021** **Lundberg SM**, Erion G, Chen H, DeGrave A, Prutkin JM, Nair B, Katz R, Himmelfarb J, Bansal N, Lee SI. From
**1022** Local Explanations to Global Understanding with Explainable AI for Trees. Nature machine intelligence. 2020
**1023** Jan; 2(1):56–67. http://dx.doi.org/10.1038/s42256-019-0138-9, doi: 10.1038/s42256-019-0138-9.

**1024** **Lundberg SM**, Lee SI. A Unified Approach to Interpreting Model Predictions. In: Guyon I, Luxburg
**1025** UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R, editors. *Advances in Neural Informa-*
**1026** *tion Processing Systems 30* Curran Associates, Inc.; 2017.p. 4765–4774. http://papers.nips.cc/paper/
**1027** 7062-a-unified-approach-to-interpreting-model-predictions.pdf.

**1028** **McKinney W**. Data Structures for Statistical Computing in Python. In: Stéfan van der Walt, Jarrod Millman,
**1029** editors. *Proceedings of the 9th Python in Science Conference*; 2010. p. 56–61. doi: 10.25080/Majora-92bf1922-
**1030** 00a.

**1031** **Mehta P**, Petersen CA, Wen JC, Banitt MR, Chen PP, Bojikian KD, Egan C, Lee SI, Balazinska M, Lee AY, Rokem A.
**1032** Automated detection of glaucoma with interpretable machine learning using clinical data and multi-modal
**1033** retinal images. Am J Ophthalmol. 2021 May; .

**1034** **Mennes M**, Biswal BB, Castellanos FX, Milham MP. Making data sharing work: the FCP/INDI experi-
**1035** ence. NeuroImage. 2013 Nov; 82:683–691. http://dx.doi.org/10.1016/j.neuroimage.2012.10.064, doi:
**1036** 10.1016/j.neuroimage.2012.10.064.

**1037** **Merkel D**. Docker: lightweight linux containers for consistent development and deployment. Linux journal.
**1038** 2014; 2014(239):2.

**1039** **Miller KL**, Alfaro-Almagro F, Bangerter NK, Thomas DL, Yacoub E, Xu J, Bartsch AJ, Jbabdi S, Sotiropoulos SN,
**1040** Andersson JLR, Griffanti L, Douaud G, Okell TW, Weale P, Dragonu I, Garratt S, Hudson S, Collins R, Jenkin-
**1041** son M, Matthews PM, et al. Multimodal population brain imaging in the UK Biobank prospective epidemi-
**1042** ological study. Nature neuroscience. 2016 Nov; 19(11):1523–1536. http://dx.doi.org/10.1038/nn.4393, doi:
**1043** 10.1038/nn.4393.

**1044** **Murdoch WJ**, Singh C, Kumbier K, Abbasi-Asl R, Yu B. Definitions, methods, and applications in interpretable
**1045** machine learning. Proceedings of the National Academy of Sciences. 2019; 116(44):22071–22080. https:
**1046** //www.pnas.org/content/116/44/22071, doi: 10.1073/pnas.1900654116.

**1047** Nelson PG, Promislow DEL, Masel J. Biomarkers for Aging Identified in Cross-sectional Studies Tend
**1048** to Be Non-causative. The journals of gerontology Series A, Biological sciences and medical sciences.
**1049** 2020 Feb; 75(3):466–472. http://academic.oup.com/biomedgerontology/article/75/3/466/5540066, doi:
**1050** 10.1093/gerona/glz174.

**1051** Nielson DM, Pereira F, Zheng CY, Migineishvili N, Lee JA, Thomas AG, Bandettini PA. Detecting and harmonizing
**1052** scanner differences in the ABCD study - annual release 1.0. bioRxiv. 2018; https://www.biorxiv.org/content/
**1053** early/2018/05/02/309260, doi: 10.1101/309260.

**1054** O'Donnell LJ, Westin CF, Golby AJ. Tract-based morphometry for white matter group analysis.
**1055** NeuroImage. 2009 Apr; 45(3):832–844. http://dx.doi.org/10.1016/j.neuroimage.2008.12.023, doi:
**1056** 10.1016/j.neuroimage.2008.12.023.

**1057** Paus T. Population neuroscience: why and how. Hum Brain Mapp. 2010 Jun; 31(6):891–903.

**1058** Paus T, Keshavan M, Giedd JN. Why do many psychiatric disorders emerge during adolescence? Nat Rev
**1059** Neurosci. 2008 Dec; 9(12):947–957.

**1060** Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R,
**1061** Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine
**1062** Learning in Python. Journal of Machine Learning Research. 2011; 12:2825–2830.

**1063** Pestilli F, Poldrack R, Rokem A, Satterthwaite T, Feingold F, Duff E, Pernet C, Smith R, Esteban O, Cieslak M, A
**1064** community-driven development of the Brain Imaging Data Standard (BIDS) to describe macroscopic brain
**1065** connections. OSF; 2021. osf.io/u4g5p, doi: 10.17605/OSF.IO/U4G5P.

**1066** Pinaya WHL, NeuroCombat-sklearn; 2020. https://github.com/Warvito/neurocombat_sklearn.

**1067** Power JD, Barnes KA, Snyder AZ, Schlaggar BL, Petersen SE. Spurious but systematic correlations in functional
**1068** connectivity MRI networks arise from subject motion. NeuroImage. 2012 Feb; 59(3):2142–2154. http://dx.
**1069** doi.org/10.1016/j.neuroimage.2011.10.018, doi: 10.1016/j.neuroimage.2011.10.018.

**1070** Power JD, Mitra A, Laumann TO, Snyder AZ, Schlaggar BL, Petersen SE. Methods to detect, characterize, and
**1071** remove motion artifact in resting state fMRI. NeuroImage. 2014; 84(Supplement C):320–341. http://www.
**1072** sciencedirect.com/science/article/pii/S1053811913009117, doi: 10.1016/j.neuroimage.2013.08.048.

**1073** Raffelt D, Dhollander T, Tournier JD, Tabbara R, Smith RE, Pierre E, Connelly A. Bias field correction and intensity
**1074** normalisation for quantitative analysis of apparent fibre density. In: *Proc. Intl. Soc. Mag. Reson. Med*, vol. 25;
**1075** 2017. p. 3541.

**1076** Richie-Halford A, Keshavan A, Cieslak M, Esteban O, , Yeatman J, Rokem A. dmriprep-viewer: a web application
**1077** for quality control of large neuroimaging datasets; 2022, in preparation.

**1078** Richie-Halford A, Rokem A. Cloudknot: A Python Library to Run your Existing Code on AWS Batch. Proceedings
**1079** of the 17th Python in Science Conference. 2018; p. 8–14. http://conference.scipy.org/proceedings/scipy2018/
**1080** adam_richie-halford.html, doi: 10.25080/Majora-4af1f417-001.

**1081** Richie-Halford A, Rokem A, HBN-POD2 QC. OSF; 2021. osf.io/8cy32, doi: 10.17605/OSF.IO/8CY32.

**1082** Richie-Halford A, Rokem A, HBN-POD2-QC: Code accompanying the HBN-POD2 manuscript. Zenodo; 2022.
**1083** https://doi.org/10.5281/zenodo.5949280, doi: 10.5281/zenodo.5949280.

**1084** Richie-Halford A, Rokem A, QSIQC: Predict diffusion MRI quality ratings. Zenodo; 2022. https://doi.org/10.
**1085** 5281/zenodo.5949269, doi: 10.5281/zenodo.5949269.

**1086** Richie-Halford A, Rokem A, Simon N, Yeatman J, richford/AFQ-Insight: AFQ-Insight. Zenodo; 2019. https://doi.
**1087** org/10.5281/zenodo.3585942, doi: 10.5281/zenodo.3585942.

**1088** Richie-Halford A, Yeatman JD, Simon N, Rokem A. Multidimensional analysis and detection of informative
**1089** features in human brain white matter. PLoS computational biology. 2021 Jun; 17(6):e1009136. http://dx.doi.
**1090** org/10.1371/journal.pcbi.1009136, doi: 10.1371/journal.pcbi.1009136.

**1091** Rosen AFG, Roalf DR, Ruparel K, Blake J, Seelaus K, Villa LP, Ciric R, Cook PA, Davatzikos C, Elliott MA, Garcia de
**1092** La Garza A, Gennatas ED, Quarmley M, Schmitt JE, Shinohara RT, Tisdall MD, Craddock RC, Gur RE, Gur RC,
**1093** Satterthwaite TD. Quantitative assessment of structural image quality. NeuroImage. 2018 Apr; 169:407–418.
**1094** http://dx.doi.org/10.1016/j.neuroimage.2017.12.059, doi: 10.1016/j.neuroimage.2017.12.059.

1095 **Salahuddin Z**, Woodruff HC, Chatterjee A, Lambin P. Transparency of deep neural networks
1096 for medical image analysis: A review of interpretability methods. Computers in biology and
1097 medicine. 2022 Jan; 140:105111. https://www.sciencedirect.com/science/article/pii/S0010482521009057, doi:
1098 10.1016/j.compbiomed.2021.105111.

1099 **Satterthwaite TD**, Wolf DH, Loughead J, Ruparel K, Elliott MA, Hakonarson H, Gur RC, Gur RE. Impact of in-
1100 scanner head motion on multiple measures of functional connectivity: relevance for studies of neurodevel-
1101 opment in youth. NeuroImage. 2012 Mar; 60(1):623–632. http://dx.doi.org/10.1016/j.neuroimage.2011.12.063,
1102 doi: 10.1016/j.neuroimage.2011.12.063.

1103 **Sayres R**, Taly A, Rahimy E, Blumer K, Coz D, Hammel N, Krause J, Narayanaswamy A, Rastegar Z, Wu D, Xu S,
1104 Barb S, Joseph A, Shumski M, Smith J, Sood AB, Corrado GS, Peng L, Webster DR. Using a Deep Learning Algo-
1105 rithm and Integrated Gradients Explanation to Assist Grading for Diabetic Retinopathy. Ophthalmology. 2019
1106 Apr; 126(4):552–564. http://dx.doi.org/10.1016/j.ophtha.2018.11.016, doi: 10.1016/j.ophtha.2018.11.016.

1107 **Shafto MA**, Tyler LK, Dixon M, Taylor JR, Rowe JB, Cusack R, Calder AJ, Marslen-Wilson WD, Duncan J, Dalgleish
1108 T, Henson RN, Brayne C, Matthews FE, Cam-CAN. The Cambridge Centre for Ageing and Neuroscience (Cam-
1109 CAN) study protocol: a cross-sectional, lifespan, multidisciplinary examination of healthy cognitive ageing.
1110 BMC neurology. 2014 Oct; 14:204. http://dx.doi.org/10.1186/s12883-014-0204-1, doi: 10.1186/s12883-014-
1111 0204-1.

1112 **Siegel JS**, Mitra A, Laumann TO, Seitzman BA, Raichle M, Corbetta M, Snyder AZ. Data Quality Influences Ob-
1113 served Links Between Functional Connectivity and Behavior. Cerebral cortex. 2017 Sep; 27(9):4492–4502.
1114 http://dx.doi.org/10.1093/cercor/bhw253, doi: 10.1093/cercor/bhw253.

1115 **Sundararajan M**, Taly A, Yan Q. Axiomatic Attribution for Deep Networks. In: Precup D, Teh YW, editors. *Pro-
1116 ceedings of the 34th International Conference on Machine Learning*, vol. 70 of Proceedings of Machine Learning
1117 Research PMLR; 2017. p. 3319–3328. https://proceedings.mlr.press/v70/sundararajan17a.html.

1118 **Taylor JR**, Williams N, Cusack R, Auer T, Shafto MA, Dixon M, Tyler LK, Cam-Can, Henson RN. The Cambridge
1119 Centre for Ageing and Neuroscience (Cam-CAN) data repository: Structural and functional MRI, MEG, and
1120 cognitive data from a cross-sectional adult lifespan sample. NeuroImage. 2017 Jan; 144(Pt B):262–269. http:
1121 //dx.doi.org/10.1016/j.neuroimage.2015.09.018, doi: 10.1016/j.neuroimage.2015.09.018.

1122 **pandas development team T**, pandas-dev/pandas: Pandas. Zenodo; 2020. https://doi.org/10.5281/zenodo.
1123 3509134, doi: 10.5281/zenodo.3509134.

1124 **TensorFlow Authors T**, Integrated gradients tutorial; 2021. Accessed: 2021-11-15. https://www.tensorflow.org/
1125 tutorials/interpretability/integrated_gradients.

1126 **Tobe RH**, MacKay-Brandt A, Lim R, Kramer M, Breland MM, Trautman KD, Hu C, Sangoi R, Tu L, Alexan-
1127 der L, Gabbay V, Castellanos FX, Leventhal BL, Craddock RC, Colcombe SJ, Franco AR, Milham MP. A
1128 Longitudinal Resource for Studying Connectome Development and its Psychiatric Associations During
1129 Childhood. medRxiv. 2021; https://www.medrxiv.org/content/early/2021/03/12/2021.03.09.21253168, doi:
1130 10.1101/2021.03.09.21253168.

1131 **Tournier JD**, Calamante F, Gadian DG, Connelly A. Direct estimation of the fiber orientation density function
1132 from diffusion-weighted MRI data using spherical deconvolution. NeuroImage. 2004; 23(3):1176–1185.

1133 **Tournier JD**, Yeh CH, Calamante F, Cho KH, Connelly A, Lin CP. Resolving crossing fibres using constrained
1134 spherical deconvolution: validation using diffusion-weighted imaging phantom data. Neuroimage. 2008;
1135 42(2):617–625.

1136 **Tustison NJ**, Avants BB, Cook PA, Zheng Y, Egan A, Yushkevich PA, Gee JC. N4ITK: Improved N3 Bias Correction.
1137 IEEE Transactions on Medical Imaging. 2010; 29(6):1310–1320. doi: 10.1109/TMI.2010.2046908.

1138 **Vallat R**. Pingouin: statistics in Python. Journal of Open Source Software. 2018; 3(31):1026. https://doi.org/10.
1139 21105/joss.01026, doi: 10.21105/joss.01026.

1140 **Van Essen DC**, Smith SM, Barch DM, Behrens TEJ, Yacoub E, Ugurbil K, WU-Minn HCP Consortium. The WU-
1141 Minn Human Connectome Project: an overview. NeuroImage. 2013 Oct; 80:62–79. http://dx.doi.org/10.
1142 1016/j.neuroimage.2013.05.041, doi: 10.1016/j.neuroimage.2013.05.041.

1143 **Velez DR**, White BC, Motsinger AA, Bush WS, Ritchie MD, Williams SM, Moore JH. A balanced accuracy function
1144 for epistasis modeling in imbalanced datasets using multifactor dimensionality reduction. Genetic epidemi-
1145 ology. 2007 May; 31(4):306–315. http://dx.doi.org/10.1002/gepi.20211, doi: 10.1002/gepi.20211.

**Veraart J**, Novikov DS, Christiaens D, Ades-Aron B, Sijbers J, Fieremans E. Denoising of diffusion MRI using random matrix theory. NeuroImage. 2016; 142:394–406.

**Wakana S**, Caprihan A, Panzenboeck MM, Fallon JH, Perry M, Gollub RL, Hua K, Zhang J, Jiang H, Dubey P, Blitz A, van Zijl P, Mori S. Reproducibility of quantitative tractography methods applied to cerebral white matter. NeuroImage. 2007 Jul; 36(3):630–644. http://dx.doi.org/10.1016/j.neuroimage.2007.02.049, doi: 10.1016/j.neuroimage.2007.02.049.

**Wandell BA**. Clarifying Human White Matter. Annual review of neuroscience. 2016; 39(1):103–128. https://doi.org/10.1146/annurev-neuro-070815-013815, doi: 10.1146/annurev-neuro-070815-013815.

**Ward-Fear G**, Pauly GB, Vendetti JE, Shine R. Authorship Protocols Must Change to Credit Citizen Scientists. Trends Ecol Evol. 2020 Mar; 35(3):187–190.

**Wargnier-Dauchelle V**, Grenier T, Durand-Dubief F, Cotton F, Sdika M. A More Interpretable Classifier For Multiple Sclerosis. In: *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*; 2021. p. 1062–1066. http://dx.doi.org/10.1109/ISBI48211.2021.9434074, doi: 10.1109/ISBI48211.2021.9434074.

**Waskom ML**. seaborn: statistical data visualization. Journal of Open Source Software. 2021; 6(60):3021. https://doi.org/10.21105/joss.03021, doi: 10.21105/joss.03021.

**Wilson G**, Bryan J, Cranston K, Kitzes J, Nederbragt L, Teal TK. Good enough practices in scientific computing. PLoS computational biology. 2017 Jun; 13(6):e1005510. http://dx.doi.org/10.1371/journal.pcbi.1005510, doi: 10.1371/journal.pcbi.1005510.

**Yeatman JD**, Dougherty RF, Myall NJ, Wandell BA, Feldman HM. Tract Profiles of White Matter Properties: Automating Fiber-Tract Quantification. PloS one. 2012 Nov; 7(11):e49790. https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0049790, doi: 10.1371/journal.pone.0049790.

**Yeatman JD**, Richie-Halford A, Smith JK, Keshavan A, Rokem A. A browser-based tool for visualization and analysis of diffusion MRI data. Nature communications. 2018 Mar; 9(1):940. https://www.nature.com/articles/s41467-018-03297-7, doi: 10.1038/s41467-018-03297-7.

**Yeatman JD**, Wandell BA, Mezer AA. Lifespan maturation and degeneration of human brain white matter. Nature communications. 2014 Sep; 5:4932. https://www.nature.com/articles/ncomms5932, doi: 10.1038/ncomms5932.

**Yeh CH**, Jones DK, Liang X, Descoteaux M, Connelly A. Mapping Structural Connectivity Using Diffusion MRI: Challenges and Opportunities. Journal of magnetic resonance imaging: JMRI. 2020 Jun; http://dx.doi.org/10.1002/jmri.27188, doi: 10.1002/jmri.27188.

**Yeh FC**, Zaydan IM, Suski VR, Lacomis D, Richardson RM, Maroon JC, Barrios-Martinez J. Differential tractography as a track-based biomarker for neuronal injury. NeuroImage. 2019 Nov; 202:116131. http://dx.doi.org/10.1016/j.neuroimage.2019.116131, doi: 10.1016/j.neuroimage.2019.116131.

**Yendiki A**, Koldewyn K, Kakunoori S, Kanwisher N, Fischl B. Spurious group differences due to head motion in a diffusion MRI study. NeuroImage. 2014 Mar; 88:79–90. http://dx.doi.org/10.1016/j.neuroimage.2013.11.027, doi: 10.1016/j.neuroimage.2013.11.027.

**Zech JR**, Badgeley MA, Liu M, Costa AB, Titano JJ, Oermann EK. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study. PLoS Med. 2018 Nov; 15(11):e1002683.

**Zhang Y**, Brady M, Smith S. Segmentation of brain MR images through a hidden Markov random field model and the expectation-maximization algorithm. IEEE Transactions on Medical Imaging. 2001; 20(1):45–57. doi: 10.1109/42.906424.

**Zunair H**, Rahman A, Mohammed N, Cohen JP. Uniformizing Techniques to Process CT Scans with 3D CNNs for Tuberculosis Prediction. In: *Predictive Intelligence in Medicine* Springer International Publishing; 2020. p. 156–168. http://dx.doi.org/10.1007/978-3-030-59354-4_15, doi: 10.1007/978-3-030-59354-4\_15.

## Appendix 1

### CuBIDS variant annotation

We identified 20 unique dMRI acquisitions across HBN-POD2, which are summarized in Table 1. Site CBIC has two acquisition types: "64dir," which shares it's pulse sequence with sites RU and CUNY, and "ABCD64dir," with acquisition parameters that better match the ABCD study (TE=0.089 s and TR=4.1 s). The "Most_Common" variant identifies the most common combination of acquisition parameters for a given site and acquisition. The "Low_Volume" variant identifies participants from all sites with less that 129 DWI volumes, which is the number of volumes in the most common variants. All remaining variants names identify the acquisition parameter(s) that differ from those of the most common variant. For example, the "MultibandAccelerationFactor" variant has a different multiband acceleration factor than that of the the most common variant but all participants within that variant share the same multiband acceleration factor. Variants that differ by multiple acquisition parameters have names that are composed of concatenated parameters. For example, the variant "Dim3SizeVoxelSizeDim3" varies both in the number of voxels in dimension 3 ("Dim3Size") and in the voxel size in dimension 3 ("VoxelSizeDim3").

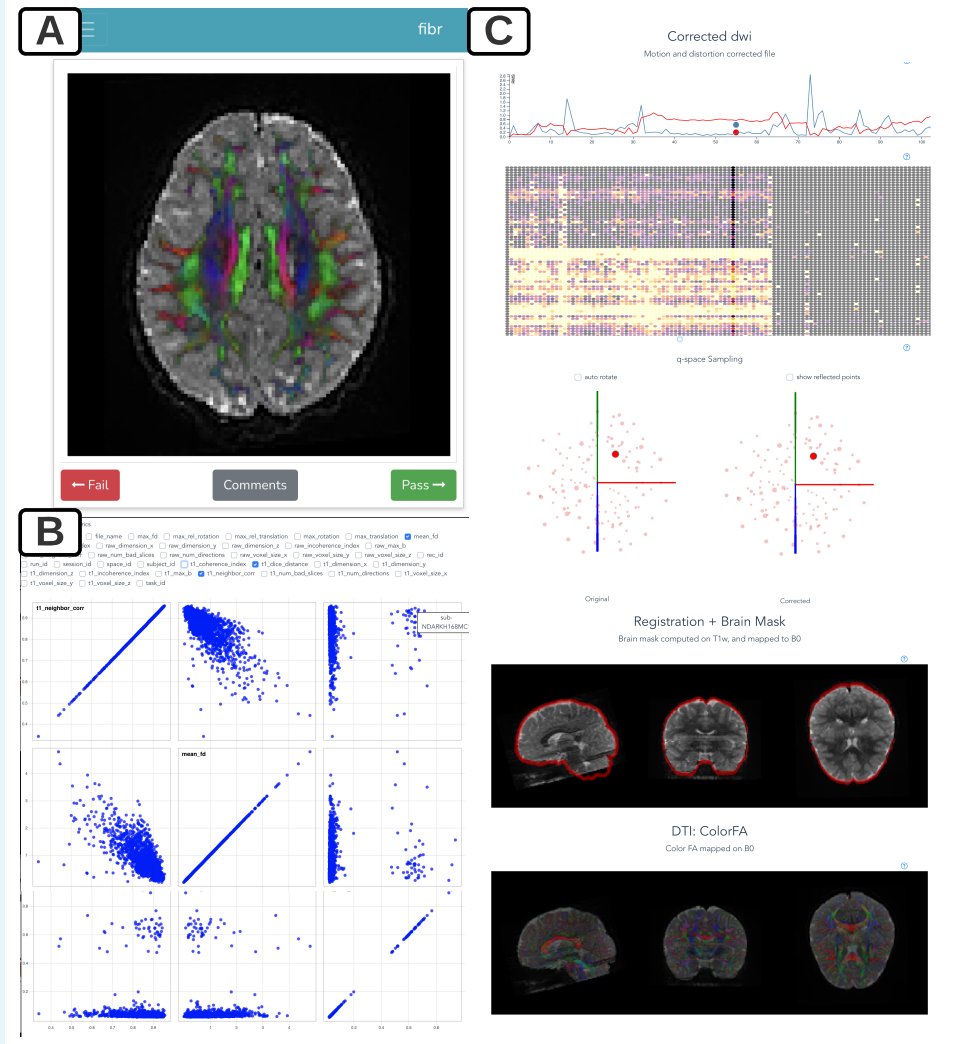| Site | Acquisition | Variant | Count |
|------|-------------|---------|-------|
| CBIC | 64dir | Most_Common | 828 |
| CBIC | 64dir | Obliquity | 32 |
| CBIC | 64dir | VoxelSizeDim1VoxelSizeDim2 | 1 |
| CBIC | ABCD64dir | Most_Common | 15 |
| CBIC | ABCD64dir | HasFmap | 2 |
| CBIC | ABCD64dir | MultibandAccelerationFactor | 1 |
| CBIC | ABCD64dir | Obliquity | 1 |
| CUNY | 64dir | Most_Common | 68 |
| CUNY | 64dir | Dim3SizeVoxelSizeDim3 | 4 |
| CUNY | 64dir | Obliquity | 2 |
| RU | 64dir | Most_Common | 859 |
| RU | 64dir | NoFmap | 5 |
| RU | 64dir | Obliquity | 8 |
| RU | 64dir | PhaseEncodingDirection | 1 |
| SI | 64dir | EchoTime | 1 |
| SI | 64dir | EchoTimePhaseEncodingDirection | 9 |
| SI | 64dir | Most_Common | 269 |
| SI | 64dir | NoFmap | 2 |
| SI | 64dir | Obliquity | 12 |
| All Sites | All Acquisitions | Low_Volume_Count | 14 |

**Appendix 1 Table 1.** Participant counts for HBN-POD2 variants.

## Appendix 2

1211
1212
1213
1214
1215
1216

### HBN-POD2 quality control instruments

We created quality control web applications for both community raters and expert raters. These apps are publicly accessible at https://fibr.dev, for the community science instrument and at http://www.nipreps.org/dmriprep-viewer/ for the expert rating instrument. We encourage readers to try these web applications on their own but have included screenshots and a summary of the interfaces in Figure 1.



**Appendix 2 Figure 1. HBN-POD2 quality control instruments**: **(A)** The user interface for community science QC app *Fibr*. After a tutorial, users are asked to give binary pass/fail ratings to each subject's DEC-FA image. The intuitive swipe or click interface allows community scientists to review more images than is practical for expert reviewers. Expert reviewers use the more advanced *dmriprep-viewer* interface, where they can **(B)** view the distribution of data quality metrics for the entire study using interactive scatterplots and violin plots, and **(C)** inspect individual participants' preprocessing results, including corrected dMRI images, frame displacement, q-space sampling distributions, registration information, and a DTI model.

## Appendix 3

### XGB feature importance

SHAP is a method to explain individual predictions based on game theoretically optimal Shapley values (*Lundberg and Lee, 2017*). To estimate global feature importance for the XGB and XGB-q models, we use the `shap` library's `TreeExplainer` (*Lundberg et al., 2020*) and average the absolute Shapley value per feature across each individual prediction. Tables 1 and 2 list the *QSIPrep* automated QC metric features in order of decreasing mean absolute shap value for the XGB and XGB-q models, respectively. We chose the top three metrics from Table 1 to plot metric distributions in Figure 2 and correlations with the expert QC results in Figure 3.

| feature | mean abs shap | feature | mean abs shap |
|---|---|---|---|
| raw_neighbor_corr | 0.666429 | raw_neighbor_corr | 0.767536 |
| max_rel_translation | 0.348662 | raw_incoherence_index | 0.453897 |
| raw_num_bad_slices | 0.288937 | raw_num_bad_slices | 0.430422 |
| t1_neighbor_corr | 0.282198 | t1_coherence_index | 0.382218 |
| raw_incoherence_index | 0.229733 | max_rel_translation | 0.363052 |
| raw_coherence_index | 0.162103 | raw_coherence_index | 0.320438 |
| max_rel_rotation | 0.118963 | t1_neighbor_corr | 0.250948 |
| mean_fd | 0.116457 | t1_dice_distance | 0.248104 |
| max_fd | 0.099359 | t1_incoherence_index | 0.242348 |
| max_rotation | 0.078774 | max_rel_rotation | 0.135590 |
| t1_coherence_index | 0.035553 | mean_fd | 0.128642 |
| t1_dice_distance | 0.034510 | max_translation | 0.120815 |
| max_translation | 0.032323 | max_fd | 0.119739 |
| t1_incoherence_index | 0.030225 | max_rotation | 0.101209 |
| raw_voxel_size_x | 0.000000 | t1_num_bad_slices | 0.007075 |
| raw_voxel_size_y | 0.000000 | raw_dimension_y | 0.000000 |
| raw_voxel_size_z | 0.000000 | raw_dimension_z | 0.000000 |
| raw_num_directions | 0.000000 | raw_voxel_size_x | 0.000000 |
| raw_max_b | 0.000000 | raw_voxel_size_y | 0.000000 |
| raw_dimension_y | 0.000000 | raw_voxel_size_z | 0.000000 |
| raw_dimension_z | 0.000000 | raw_max_b | 0.000000 |
| t1_voxel_size_x | 0.000000 | t1_voxel_size_x | 0.000000 |
| t1_dimension_x | 0.000000 | raw_num_directions | 0.000000 |
| t1_dimension_y | 0.000000 | t1_dimension_x | 0.000000 |
| t1_dimension_z | 0.000000 | t1_dimension_y | 0.000000 |
| t1_voxel_size_y | 0.000000 | t1_dimension_z | 0.000000 |
| t1_voxel_size_z | 0.000000 | t1_voxel_size_y | 0.000000 |
| t1_max_b | 0.000000 | t1_voxel_size_z | 0.000000 |
| t1_num_bad_slices | 0.000000 | t1_max_b | 0.000000 |
| t1_num_directions | 0.000000 | t1_num_directions | 0.000000 |
| raw_dimension_x | 0.000000 | raw_dimension_x | 0.000000 |

**Appendix 3 Table 1.** XGB mean absolute shap values

**Appendix 3 Table 2.** XGB-q mean absolute shap values

### The *Fibr* Community Science Consortium

1242

1243 The following community raters provided $> 3,000$ ratings each and elected to be included

1244 in the *Fibr* Community Science Consortium as co-authors on this paper.

| Name | ORCID iD |
| --- | --- |
| Nicholas J. Abbott | 0000-0003-1466-0352 |
| John A. E. Anderson | 0000-0001-6511-1957 |
| Gagana B. | |
| MaryLena Bleile | 0000-0002-0762-2596 |
| Peter S. Bloomfield | 0000-0002-8356-7701 |
| Vince Bottom | |
| Josiane Bourque | |
| Rory Boyle | 0000-0003-0787-6892 |
| Julia K. Brynildsen | 0000-0002-1627-6576 |
| Navona Calarco | 0000-0002-4391-0472 |
| Jaime J. Castrellon | 0000-0001-5834-7101 |
| Natasha Chaku | 0000-0003-0944-6159 |
| Bosi Chen | 0000-0002-0117-9757 |
| Sidhant Chopra | 0000-0003-0866-3477 |
| Emily B. J. Coffey | 0000-0001-8249-7396 |
| Nigel Colenbier | 0000-0003-0928-2668 |
| Daniel J. Cox | |
| James Elliott Crippen | |
| Jacob J. Crouse | 0000-0002-3805-2936 |
| Szabolcs David | 0000-0003-0316-3895 |
| Benjamin De Leener | 0000-0002-1378-2756 |
| Gwyneth Delap | |
| Zhi-De Deng | 0000-0001-8925-0871 |
| Jules Roger Dugre | 0000-0003-4946-0350 |
| Anders Eklund | 0000-0001-7061-7995 |
| Kirsten Ellis | 0000-0002-7570-0939 |
| Arielle Ered | 0000-0002-8386-4423 |
| Harry Farmer | 0000-0002-3684-0605 |
| Joshua Faskowitz | 0000-0003-1814-7206 |
| Jody E. Finch | 0000-0003-2457-1345 |
| Guillaume Flandin | 0000-0003-0077-7859 |
| Matthew W. Flounders | 0000-0001-7014-4665 |
| Leon Fonville | 0000-0001-8874-7843 |
| Dea Garic | 0000-0003-3595-4210 |
| Patricia Garrido-Vásquez | 0000-0002-9561-8983 |
| Gabriel Gonzalez-Escamilla | 0000-0002-7209-1736 |
| Shannon E. Grogans | 0000-0003-0383-4601 |
| Mareike Grotheer | 0000-0002-8653-1157 |
| David C. Gruskin | 0000-0001-6504-191X |
| Guido I. Guberman | |
| Edda Briana Haggerty | 0000-0003-0597-7956 |

| Name | ORCID |
|---|---|
| Younghee Hahn | |
| Elizabeth H. Hall | |
| Jamie L. Hanson | 0000-0002-0469-8886 |
| Yann Harel | 0000-0002-8970-1983 |
| Bruno Hebling Vieira | 0000-0002-8770-7396 |
| Meike D. Hettwer | 0000-0002-7973-6752 |
| Corey Horien | 0000-0001-6738-1029 |
| Fan Huang | |
| Zeeshan M. Huque | |
| Anthony R. James | 0000-0002-5297-2229 |
| Isabella Kahhale | 0000-0002-0963-9738 |
| Sarah L. H. Kamhout | |
| Arielle S. Keller | 0000-0003-4708-1672 |
| Harmandeep Singh Khera | 0000-0001-6840-4616 |
| Gregory Kiar | 0000-0001-8915-496X |
| Peter Alexander Kirk | 0000-0003-0786-3039 |
| Simon H. Kohl | 0000-0003-0949-6754 |
| Stephanie A. Korenic | |
| Cole Korponay | 0000-0003-2562-9617 |
| Alyssa K. Kozlowski | |
| Nevena Kraljevic | 0000-0003-0869-648X |
| Alberto Lazari | 0000-0002-8688-581X |
| Mackenzie J. Leavitt | 0000-0002-6100-3235 |
| Zhaolong Li | 0000-0003-2246-4116 |
| Giulia Liberati | 0000-0002-5684-4443 |
| Elizabeth S. Lorenc | 0000-0003-1311-726X |
| Annabelle Julina Lossin | 0000-0001-5921-1353 |
| Leon D. Lotter | 0000-0002-2337-6073 |
| David M. Lydon-Staley | 0000-0001-8702-3923 |
| Christopher R. Madan | 0000-0003-3228-6501 |
| Neville Magielse | 0000-0002-6777-4225 |
| Hilary A. Marusak | 0000-0002-0771-6795 |
| Julien Mayor | 0000-0001-9827-542 |
| Amanda L. McGowan | 0000-0003-3422-0135 |
| Kahini P. Mehta | |
| Steven Lee Meisler | 0000-0002-8888-1572 |
| Cleanthis Michael | 0000-0002-5300-473X |
| Mackenzie E. Mitchell | 0000-0002-0225-6320 |
| Simon Morand-Beaulieu | 0000-0002-5880-3688 |
| Benjamin T. Newman | 0000-0002-0668-2853 |
| Jared A. Nielsen | 0000-0002-2717-193X |
| Shane M. O'Mara | |
| Amar Ojha | 0000-0002-1038-0225 |
| Adam Omary | |
| Evren Özarslan | 0000-0003-0859-1311 |
| Linden Parkes | 0000-0002-9329-7207 |
| Madeline Peterson | |
| Adam Robert Pines | |
| Claudia Pisanu | 0000-0002-9151-4319 |

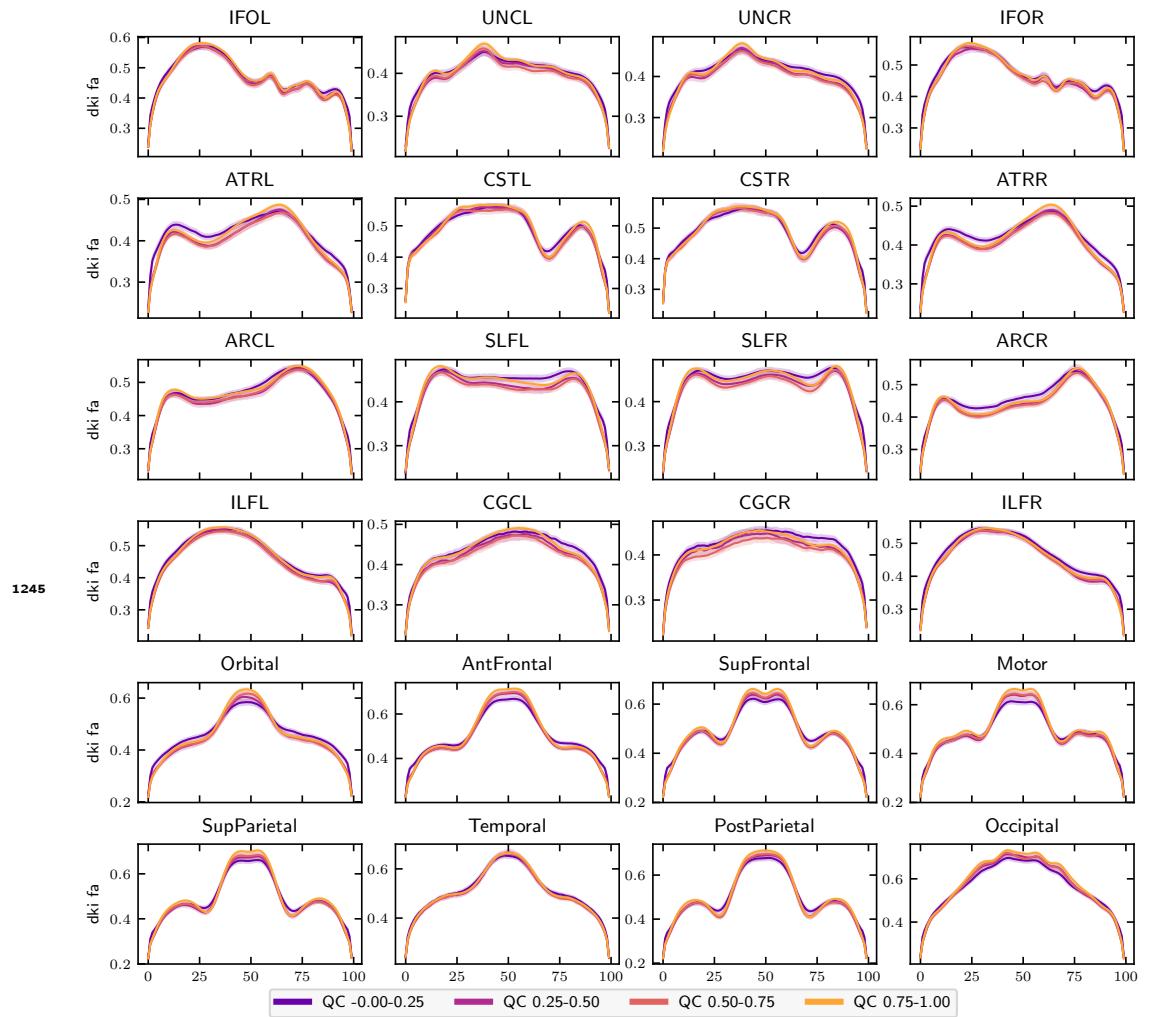| Ryan R. Rich | 0000-0001-9495-3184 |
| Ashish K. Sahoo | 0000-0003-1815-6655 |
| Amjad Samara | 0000-0002-6001-7395 |
| Farah Sayed | |
| Jonathan Thore Schneider | 0000-0002-1925-6669 |
| Lindsay S. Shaffer | 0000-0002-0642-1717 |
| Ekaterina Shatalina | 0000-0001-8900-0792 |
| Sara A. Sims | 0000-0001-7107-1891 |
| Skyler Sinclair | 0000-0003-3010-6431 |
| Jae W. Song | 0000-0002-3127-6427 |
| Griffin Stockton Hogrogian | 0000-0003-2877-078X |
| Ursula A. Tooley | 0000-0001-6377-3885 |
| Vaibhav Tripathi | |
| Hamid B. Turker | 0000-0002-2670-4036 |
| Sofie Louise Valk | 0000-0003-2998-6849 |
| Matthew B. Wall | 0000-0002-0493-6274 |
| Cheryl K. Walther | |
| Yuchao Wang | 0000-0001-9871-3006 |
| Bertil Wegmann | 0000-0003-2193-6003 |
| Thomas Welton | 0000-0002-9503-2093 |
| Alex I. Wiesman | 0000-0003-0917-1570 |
| Andrew G. Wiesman | |
| Mark Wiesman | |
| Drew E. Winters | 0000-0002-0701-9658 |
| Ruiyi Yuan | |
| Sadie J. Zacharek | 0000-0001-8770-4614 |
| Chris Zajner | 0000-0002-0204-6497 |
| Ilya Zakharov | 0000-0001-7207-9641 |
| Gianpaolo Zammarchi | 0000-0002-9733-380X |
| Dale Zhou | 0000-0001-9240-1327 |
| Benjamin Zimmerman | 0000-0003-2570-8198 |
| Kurt Zoner | |

**Figure 8–Figure supplement 1. FA bundle profiles binned by QC score**: FA profiles binned by QC score in twenty-four major while matter bundles. The *x*-axis represents distance along the length of the fiber bundle. Error bands represent bootstrapped 95% confidence intervals. Bundle abbreviations are as in Figure 8
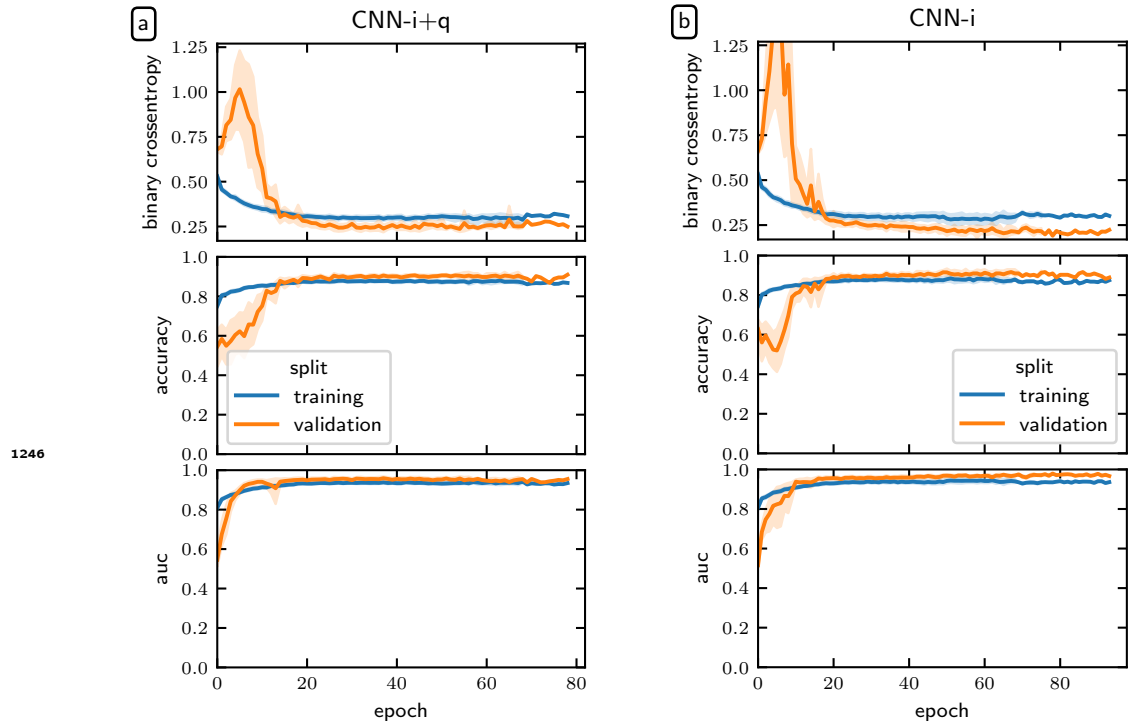
**Figure 10–Figure supplement 1. Deep learning model loss curves**: The binary cross-entropy loss (top), accuracy (middle), and ROC-AUC (bottom) for **(a)** the CNN-i+q model and **(a)** the CNN-i model. Model performance typically plateaued after twenty epochs but was allowed continue until meeting the early stopping criterion. The error bands represent a bootstrapped 95% confidence interval.