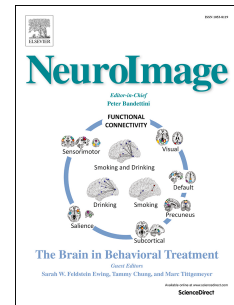


# Accepted Manuscript

Replication and generalization in applied neuroimaging

Garikoitz Lerma-Usabiaga, Pratik Mukherjee, Zhimei Ren, Michael L. Perry, Brian A. Wandell



PII: S1053-8119(19)30629-9

DOI: <https://doi.org/10.1016/j.neuroimage.2019.116048>

Article Number: 116048

Reference: YNIMG 116048

To appear in: *NeuroImage*

Received Date: 18 December 2018

Revised Date: 29 April 2019

Accepted Date: 22 July 2019

Please cite this article as: Lerma-Usabiaga, G., Mukherjee, P., Ren, Z., Perry, M.L., Wandell, B.A., Replication and generalization in applied neuroimaging, *NeuroImage* (2019), doi: <https://doi.org/10.1016/j.neuroimage.2019.116048>.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# 1 Replication and generalization in applied neuroimaging

2

3 **Short Title:** Applied neuroimaging

4

5 **Authors:** Garikoitz Lerma-Usabiaga<sup>1,2</sup>, Pratik Mukherjee<sup>3,4</sup>, Zhimei Ren<sup>5</sup>, Michael L. Perry<sup>1</sup>, Brian A.  
6 Wandell<sup>1</sup>

7

8

9 **Affiliation:**

10 <sup>1</sup> Department of Psychology, Stanford University, 450 Serra Mall, Jordan Hall Building, 94305 Stanford,  
11 California, USA

12 <sup>2</sup> BCBL. Basque Center on Cognition, Brain and Language. Mikeletegi Pasealekua 69, Donostia - San  
13 Sebastián, 20009 Gipuzkoa, Spain

14 <sup>3</sup> Radiology and Biomedical Imaging, and <sup>4</sup> Bioengineering and Therapeutic Sciences, University of  
15 California, San Francisco, California, USA

16 <sup>5</sup> Department of Statistics, Stanford University, 390 Serra Mall, Sequoia Hall Building, 94305 Stanford,  
17 California, USA

18

19 **Corresponding Author:**

20 Garikoitz Lerma-Usabiaga

21 Email: garikoitz@gmail.com

22 450 Serra Mall

23 Room 488, Jordan Hall, Building 420, Main Quad

24 94305 Stanford CA, United States

## 1 Abstract

2 There is much interest in translating neuroimaging findings into meaningful clinical diagnostics. The goal  
3 of scientific discoveries differs from clinical diagnostics. Scientific discoveries must replicate under a  
4 specific set of conditions; to translate to the clinic we must show that findings using purpose-built  
5 scientific instruments will be observable in clinical populations and instruments. Here we describe and  
6 evaluate data and computational methods designed to translate a scientific observation to a clinical  
7 setting. Using diffusion weighted imaging (DWI), Wahl et al., (2010) observed that across subjects the  
8 mean fractional anisotropy (FA) of homologous pairs of tracts is highly correlated. We hypothesize that  
9 this is a fundamental biological trait that should be present in most healthy participants, and deviations  
10 from this assessment may be a useful diagnostic metric. Using this metric as an illustration of our  
11 methods, we analyzed six pairs of homologous white matter tracts in nine different DWI datasets with 44  
12 subjects each. Considering the original FA measurement as a baseline, we show that the new metric is  
13 between 2 and 4 times more precise when used in a clinical context. Our framework to translate research  
14 findings into clinical practice can be applied, in principle, to other neuroimaging results.

15

16

## 17 Keywords

18 Replication; Generalization; Generalizability; Computational Reproducibility; Structural MRI; DWI;  
19 White Matter Tracts; Biomarker

20

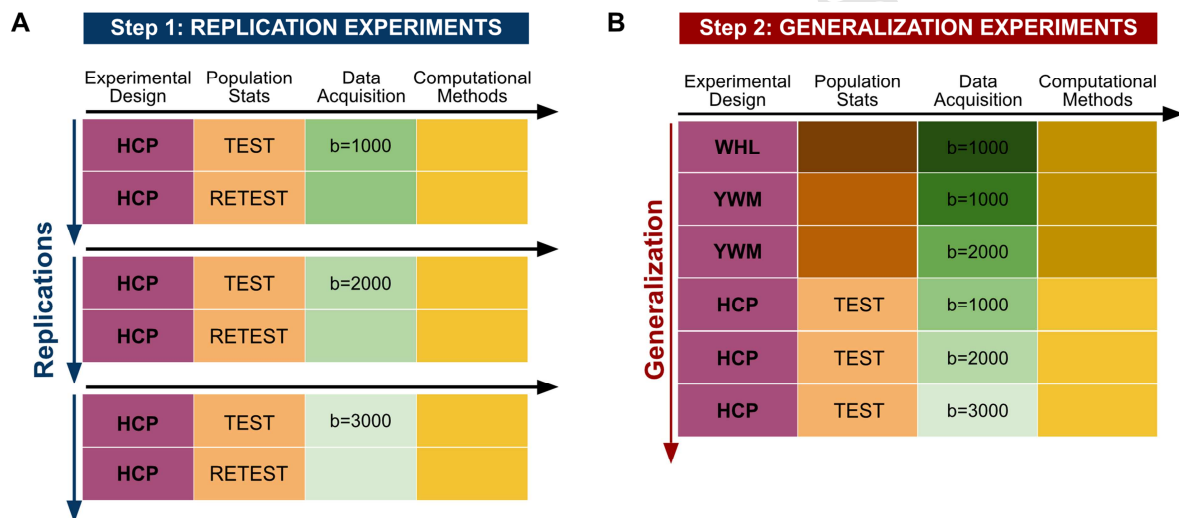
## 1 1.- Introduction

2 We describe methods to translate magnetic resonance imaging (MRI) scientific findings into clinical  
3 practice. The goal of scientific discoveries differs from clinical diagnostics. Clinical applications should  
4 be based on quantitative measurements that replicate in controlled laboratory conditions. These  
5 applications must also be applicable to a clinical environment where data acquisition methods, subject  
6 populations, and computational methods can vary substantially.

7  
8 We base our methods on the ideas of replication and generalization. Because these terms, along with  
9 reproducibility, re-execution, and robustness are used in various ways in the literature (Goodman et al.,  
10 2016; Kennedy et al., 2019; McNaught and Wilkinson, 1997; Patil et al., 2016; Plessner, 2017), we begin  
11 by explaining our usage. Scientific experimentalists typically set out to make a measurement that can be  
12 replicated. For example, a team makes a measurement using a specific rig and experimental conditions.  
13 Other scientists check the work by following the published instructions that define how to construct the  
14 rig and implement the experimental conditions. Scientific *replication* means repeating the experiment as  
15 precisely as possible. This approach is appropriate for investigations that test theories or quantify  
16 important phenomena, but replication is not a realistic possibility for extending discoveries into clinical  
17 applications. These applications do not have access to the carefully calibrated instruments that have been  
18 purpose-built for scientific measurements (for example, the Human Connectome Project scanners). For a  
19 scientific discovery to become clinically relevant, the finding must *generalize* across variations in the  
20 population and instruments.

21  
22 Replication and generalization are contrasted in [Figure 1](#). Panel A emphasizes scientific discovery and  
23 replication. An experimental design is chosen and measurements are made with a selected population,  
24 data acquisition instruments and methods, and a computational method. We measure the precision of the  
25 measurement when the experiment is repeated (test-retest). In this case three replication experiments are

1 illustrated using different data acquisition parameters. If the scientific measurements replicate with  
 2 sufficient precision, we might carry out generalization measurements (Panel B), to test the extent of  
 3 applicability of said measurements. The panel illustrates generalization experiments that share the same  
 4 experimental design, but use different populations (e.g., geographic locales, age and gender), different  
 5 data acquisition methods, (e.g., pulse sequences and vendors), and different computational methods (e.g.,  
 6 pre-processing software). Translating a scientific measurement into a clinical application is a two step  
 7 process: beginning with an experiment that replicates, we test how well the experiment generalizes.  
 8



**Figure 1. Summary of the individual experiments, organized as replication or generalization experiments.**

The columns correspond to the experimental pipeline steps; every row corresponds to an experiment. Different colors represent different steps in the experimental pipeline; different shades represent implementation differences within the step. A) Three replication experiments, based on the Human Connectome Project (HCP) test-retest datasets. The difference among the experiments is the b-value used in the acquisition. In a replication experiment, the intention is to repeat the original methods as far as possible, hence the same shades; the test-retest case goes uses the same population and instrumentation at different times. B) The generalization experiment reflects the transition to the clinical environment. The goal is to evaluate whether the measurements are robust to expected variations in the measurement conditions. The generalization is undertaken after validating the results in the replication experiment. The datasets are from Wahl *et al.* (2010) (WHL), Yeatman *et al.* (2014) (YWM), and (Glasser *et al.*, 2013) (HCP).

1  
2 This paper applies replication and generalization to a neuroimaging measurement that has the potential to  
3 become clinically relevant: identifying lateralized white matter disease in individual subjects. Wahl *et al.*  
4 (2010) used diffusion-weighted imaging (DWI) to measure white matter tracts in healthy adults; they  
5 observed that across subjects the mean fractional anisotropy (FA) of homologous pairs of tracts is highly  
6 correlated. We hypothesized that this finding might be a fundamental biological trait in healthy  
7 participants, that can be measured in research labs and clinical settings. We investigated if the relation  
8 between homologous tract pairs is a more useful clinical measure than assessing the measurements from  
9 each tract separately. We find that using the relation between homologous left and right tracts does  
10 provide a potential clinical measure.

11

## 2.- Materials and methods

To evaluate the replication and generalization of the DWI finding, we obtained data from multiple sources. We use nine datasets that we group into three categories.

- **WHL:** Original 44 subject dataset used in (Wahl et al., 2010). The authors shared the original DICOM files for this work, and we performed the analysis using our computational methods. We obtain 1 dataset, called WHL1000.
- **YWM:** We selected a 44 subject subset of the data reported by (Yeatman et al., 2014). The subjects matched the mean age (but not the age range) of the WHL dataset. We obtained two datasets: YWM1000 and YWM2000 that differ in data acquisition parameters (b-values, number of directions).
- **HCP:** We selected a 44 subjects whose test-retest data are available from the 1200 Human Connectome Project (HCP) release (Glasser et al., 2013). HCP1000, HCP2000 and HCP3000 differ only in data acquisition parameters (b-values). HCP1000RETEST, HCP2000RETEST and HCP3000RETEST are the corresponding retest data.

Figure 1 represents three replication experiments (panel A) and a generalization experiment (panel B). The replication experiments compare test-retest values of the mean tract FA at three different b-values; they were collected using the same subjects, instruments and computational methods at the HCP. This replication analysis bounds the precision of the estimated mean tract FA: the generalization precision shouldn't be better than the replication precision.

The generalization experiment compares the mean tract FA across different subjects, instruments and computational methods (Figure 1B). The precision derived from these six experiments assesses generalization. The HCP RETEST experiments are omitted from the generalization experiment to avoid a HCP bias. In addition to subject and instrument differences, the WHL and YWM differ in computational

1 processing. Some unmeasured variability is introduced by non-deterministic aspects of these  
2 computations.

3

4 In the following sections, we describe three different aspects of the experimental pipeline. The Population  
5 statistics section shows that cohorts are similar, but not identical. The Data acquisition section includes  
6 MRI pulse sequences and parameter choices that are different between sites and vendors, as is often the  
7 case in clinical settings. The Computational methods section describes the infrastructure we used to  
8 implement computational reproducibility, as well as a detailed description of the data analysis pipeline  
9 and numerical calculations.

## 10 **2.1.- Population statistics**

11 The population statistics for the three datasets are similar, but not exactly the same (see [Table 1](#)). All the  
12 groups include 44 subjects of a similar mean age, ranging from 30.7 to 31.8. The age range of the YWM  
13 dataset is the largest, with a standard deviation of 14.4. The HCP dataset age standard deviations is 3.2,  
14 which is an approximation: the HCP ages are binned to protect participant privacy. The YMW and WHL  
15 datasets are matched in male-female ratio, but the HCP dataset has more females than males. The original  
16 publications include more information about the populations (Glasser et al., 2013; Wahl et al., 2010;  
17 Yeatman et al., 2014).

18

Dataset	Count	Age	Gender	Age
WHL	44	30.8±7.8	20 female	29.5±7.5
			24 male	31.9±7.9
YWM	44	31.8±14.4	24 female	29.5±2.1
			20 male	34.7±3.6
HCP	44	30.7±3.2	31 female	31.9±3.2
			13 male	27.8±3.2



Table 1. Descriptive statistics of the three different populations used across the datasets.

## 2.2.- Data acquisition

Table 2 shows the main characteristics of the DWI data acquisition, emphasizing the differences between sites and experiments. As a practical matter, measurements made across multiple sites are very likely to have different MRI scanner models that are calibrated using different tools. The MRI vendors compete on intellectual property concerning the pulse sequences, making a perfect replication either extremely inconvenient or impossible. For example, the scanner used by the HCP site was specially designed and this type of instrument is unlikely to become available to the thousands of clinical sites around the world (Glasser et al., 2016, 2013). The datasets differ with respect to the number of acquisition channels, gradient strength, diffusion directions, b-value and voxel size. Such differences are unavoidable because not all sites can implement the same acquisition parameters. In addition to vendor differences, data are acquired over time, technology evolves, and people make choices.

Dataset	Scanner Vendor; Model; Location	Magnetic Field; Head Coil Receivers; Max. Gradient Strength	Main Sequence Characteristics	Experiment Codename
WHL	GE Signa EXCITE UCSF	3T 8 channels 40 mT/m	55 dirs., 1.8 mm <sup>3</sup> voxels b = 1000 s/mm <sup>2</sup>	WHL1000
	YWM	GE Discovery 750 Stanford CNI	3T 32 channels	30 dirs., 2 mm <sup>3</sup> voxels b = 1000 s/mm <sup>2</sup>
40 mT/m			96 dirs., 2 mm <sup>3</sup> voxels b = 2000 s/mm <sup>2</sup>	YMN2000
HCP	Siemens	3T	90 dirs., 1.25 mm <sup>3</sup> vox	HCP1000 &
	Connectom CMRR/WASH	32 channels	b = 1000 s/mm <sup>2</sup>	HCP1000RETEST

	WashU	100 mT/m	90 dirs., 1.25 mm <sup>3</sup> vox b = 2000 s/mm <sup>2</sup>	<b>HCP2000 &amp; HCP2000RETEST</b>
			90 dirs., 1.25 mm <sup>3</sup> vox b = 3000 s/mm <sup>2</sup>	<b>HCP3000 &amp; HCP3000RETEST</b>

Table 2. Main characteristics of the data acquisition parameters across datasets.

## 1 2.3.- Computational methods

2 The computational methods are divided into two parts: (1) *the infrastructure*: required for a  
3 computationally reproducible system, sometimes called the neuroinformatics platform (Marcus et al.,  
4 2011); and, (2) *the data analysis pipeline*: comprises all the steps starting with the DICOM images  
5 generated in the MRI scanner (the acquisition device) to the final published results.

### 6 2.3.1.- Infrastructure for computational reproducibility

7 The data management and computational infrastructure use a technology (Flywheel.io) that (a)  
8 implements reproducible computational methods, (b) tracks provenance of the data, and (c) facilitates  
9 data sharing. For reproducibility, all computational methods were performed using containerized  
10 methods. These are small virtual machines that include all dependencies and runs the same computation  
11 across platforms. The analytical methods implemented in the containers are open-source, and we provide  
12 links to the containers in the following sections. To track the provenance, the computational system  
13 stores: (a) the input data, (b) the container version that was executed, (c) the container input parameters,  
14 and (c) the output files. The analyses are fully reproducible by anyone with IRB authorization to access  
15 the system. More details about the infrastructure and implementation can be found at Lerma-Usabiaga et  
16 al. (2019).

## 1 2.3.2.- Data analysis pipeline

2 The diffusion-weighted imaging analysis methods consisted of two main steps, implemented in two  
3 containers: preprocessing and tractography. Both were applied to WHL and YWM datasets. The HCP  
4 dataset was preprocessed by that consortium (Andersson et al., 2003; Andersson and Sotiropoulos, 2016,  
5 2015) and only the tractography container was applied.

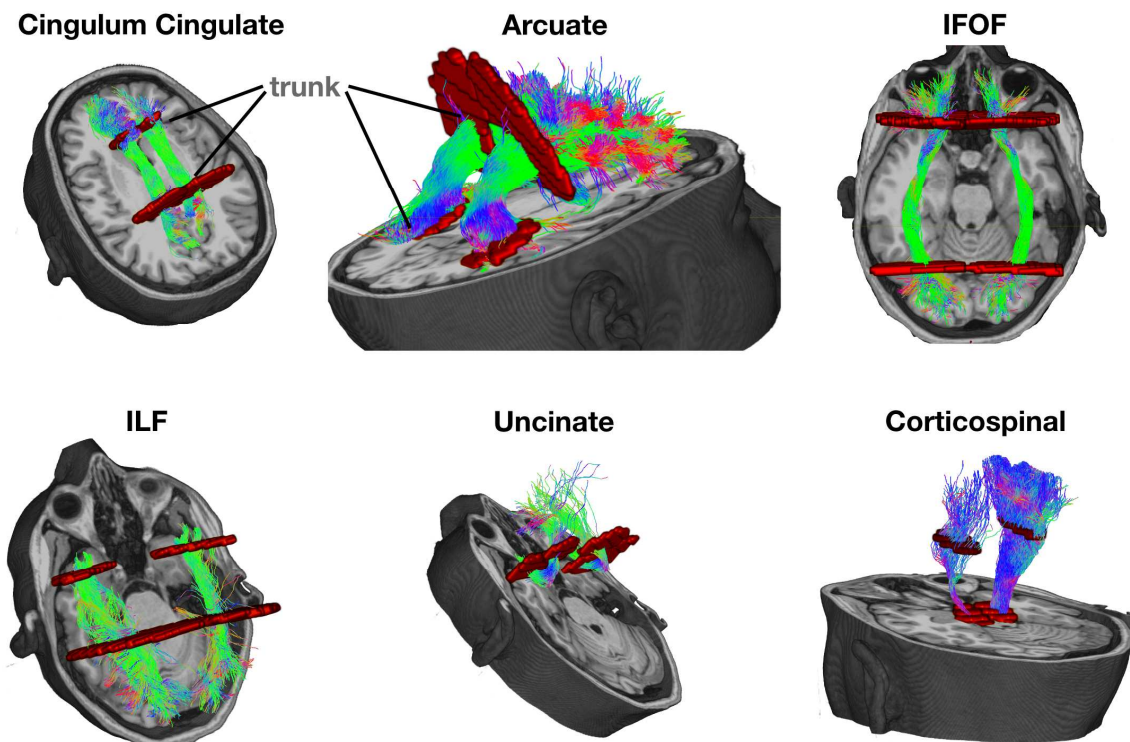
### 6 2.3.2.1.- Preprocessing

7 The preprocessing consists of the data preparation required to do the tractography and fractional  
8 anisotropy (FA) analyses. The preprocessing container comprises the following steps: first, using the tools  
9 provided by MRtrix ([github.com/MRtrix3/mrtrix3](https://github.com/MRtrix3/mrtrix3)), we perform a principal component analysis (PCA)  
10 based denoising of the data; second, additional Rician based denoising and Gibbs ringing corrections  
11 were applied (Kellner et al., 2016; Veraart et al., 2016a, 2016b); third, FSL's eddy current correction was  
12 applied (Andersson and Sotiropoulos, 2016); fourth, we performed bias correction using the ANTs  
13 package (Tustison et al., 2010); fifth, we applied a Rician background noise removal using MRtrix tools  
14 again. The code and parameters are available through GitHub ([github.com/vistalab/RTP-preproc](https://github.com/vistalab/RTP-preproc)) and  
15 Docker Hub ([hub.docker.com/r/vistalab/RTP-preproc/](https://hub.docker.com/r/vistalab/RTP-preproc/)).

### 16 2.3.2.2.- DWI processing and tractography

17 The tractography container takes the preprocessed DWI data and an un-preprocessed anatomical T1-  
18 weighted file as input. It outputs the FA of the selected 6 homologous tract pairs. The algorithms in the  
19 container perform the following steps: first, the diffusion data are aligned and resliced to the anatomical  
20 image ([https://github.com/vistalab/vistasoft\\_dtiInit](https://github.com/vistalab/vistasoft_dtiInit)); second, the whole brain white matter streamlines are  
21 estimated using the Ensemble Tractography (ET) method (Takemura et al., 2016). ET invokes MRtrix's  
22 constrained spherical deconvolution (CSD) implementation once and the tractography tool 5 times,  
23 constructing whole brain tractograms with a range of minimum angle parameters (values 47.2, 23.1, 11.5,  
24 5.7, 2.9). The LiFE (Linear Fascicle Evaluation) method evaluates the tractogram streamlines and retains

1 those that meaningfully contribute to predicting variance in the DWI data (Pestilli et al., 2014). Finally,  
2 the Automated Fiber Quantification (AFQ) method (Yeatman et al., 2012) segments streamlines into  
3 tracts (Figure 2). The code and parameters are available through GitHub ([github.com/vistalab/RTP-](https://github.com/vistalab/RTP-pipeline)  
4 [pipeline](https://github.com/vistalab/RTP-pipeline)) and the container through Docker Hub ([hub.docker.com/r/vistalab/RTP-pipeline](https://hub.docker.com/r/vistalab/RTP-pipeline)).  
5



**Figure 2. Six pairs of homologous tracts and their defining ROIs.**

The streamlines serve as a model of white matter tracts; they are selected by fitting to the diffusion weighted imaging (DWI) measurements. The tracts are defined by regions of interest (ROIs, red) that select specific streamlines from the whole brain tractogram. The region between the two ROIs is relatively stable and called the trunk. We estimate a core fiber from the collection of streamlines and sample 100 equally spaced segments. The FA of the core fiber is calculated by combining FA transverse to the core fiber at every sample point, using a Gaussian weighting scheme over distance. The set of sample points is the tract profile; the average of the FA values of the core fiber is the mean tract FA.

## 1 2.3.2.3.- Mean tract FA values

2 We analyzed the six homologous-tract pairs analyzed in (Wahl et al., 2010) (Figure 2). The ROIs used to  
3 identify the streamlines that form the tracts are shown in red. The mean tract FA is calculated in several  
4 steps. A core fiber, representing the central tendency of all the streamlines in the tract, is identified.  
5 Equally spaced positions along the fiber between the two defining ROIs are sampled (N=100). The FA  
6 values of streamlines at locations transverse to each sample position are measured and combined. The  
7 value is a Gaussian-weighted sum where the weight depends on the distance from the sample point  
8 (Yeatman et al., 2012). The sampling and transverse averaging generates a tract profile of 100 FA values.  
9 The mean tract FA is the average of these values.

10

## 11 2.3.3.- Data preparation and statistical analysis

12 The data preparation, statistical analysis and plotting scripts read the input data directly from the Flywheel  
13 neuroinformatics platform using a software development kit (SDK). To maintain reproducibility and data  
14 provenance, these scripts are stored and versioned in a GitHub repository, and the input data are stored  
15 and the specific version that was executed is stored in the neuroinformatic platform. The scripts read the  
16 files containing the FA values for each subject and each tract, categorize it for the different experiments,  
17 create the descriptive plots and calculate the metrics. The scripts to replicate the figures and calculations  
18 can be found at <https://github.com/garikoitz/paper-reproducibility>.

ACCEPTED MANUSCRIPT

## 1 3.- Results

2 We first illustrate replication and generalization analyses for the FA measurement of individual tracts and  
3 evaluate the usefulness of this measure as a clinical application. Next, we evaluate a metric based on the  
4 homologous tract correlation reported by Wahl et al. (2010). The main figures describe one illustrative  
5 tract, the inferior fronto-occipital fasciculus (IFOF), and in total we report findings for six pairs of  
6 homologous tracts. We selected the mean FA of a tractogram as an example because it is useful to explain  
7 our methods, but the analysis can be applied to many other measures. For example, Wahl et al. report four  
8 DWI measures (FA, MD, AD, RD).

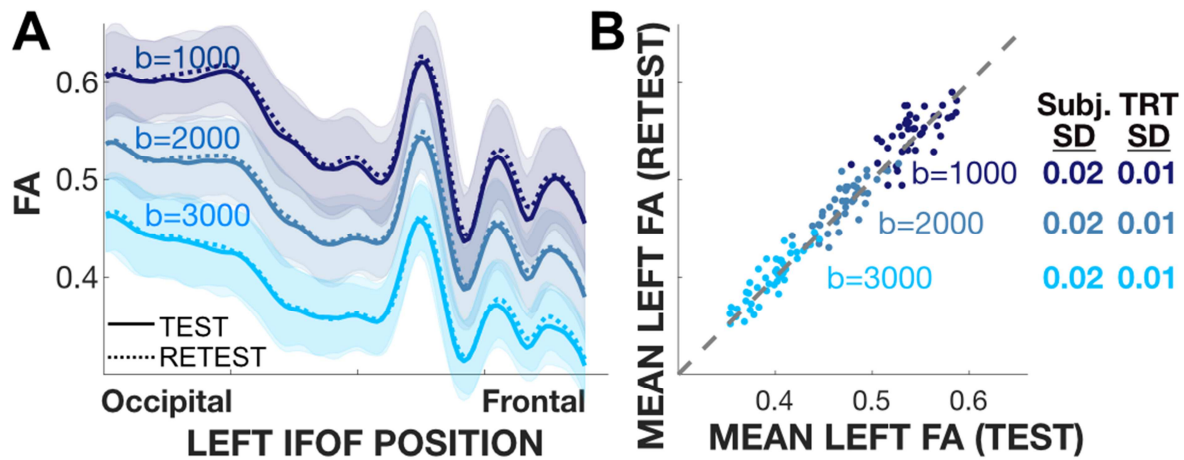
### 9 3.1.- FA measurement

#### 10 3.1.1.- Replication experiment

11 Figure 3A shows the mean tract FA profiles at three b-values for the streamlines that model the IFOF.  
12 The solid and dashed lines show the mean tract profile across subjects for the test (solid) and retest  
13 (dashed) acquisitions. The profiles are similar at each b-value; consistent with prior measurements the FA  
14 values decrease as b-value increases (Farrell et al., 2007a; Jones and Basser, 2004; Landman et al., 2007;  
15 Mukherjee et al., 2008a, 2008b). The shaded regions indicate the range ( $\pm 1$  SD) across the population  
16 of participants.

17

18



**Figure 3. Replication analyses of the tract profile and mean tract FA.**

Analyses are shown for a representative tract (left IFOF), and based on the HCP test-retest data. **A)** Tract profiles of the subject average FA in the test (solid) and retest (dashed) experiments. The mean profile (thin line) and  $\pm 1$  SD (shaded band) are shown. The profiles at each b-value match very closely; across b-values the profiles have a similar shape but different absolute values. **B)** Test-retest scatter plot. For all b-values, the SD of the difference between the test-retest pairs of FA values is 0.01 (TRT SD), and the SD of the distribution of FA values is 0.02 (Subj. SD). (Tract profiles and scatter plots for 11 other tracts are similar and reported in Figures S1a, S2).

1

2 The test-retest analyses for the mean tract FA of the IFOF are shown in Figure 3B. Each point is a subject,  
 3 and the three types of symbols show test-retest at three b-values. The test-retest mean tract FA values are  
 4 distributed near the identity line. For each b-value the mean tract FA varies between subjects (standard  
 5 deviation, 0.025). The scatter about the identity line is smaller, (standard deviation, 0.01-0.02). The  
 6 scatter around the identity line is similar for measurements at the three b-values, suggesting that the noise  
 7 level is similar (Rokem et al., 2015).

8

9 The replication analyses for an additional 11 tracts follow the same trends as the IFOF (see Supplemental  
 10 material, Figures S1a-S2). The FA values decrease with increasing b-value, and the between-subject  
 11 standard deviation is larger than the within-subject test-retest standard deviation. Considering all tracts,

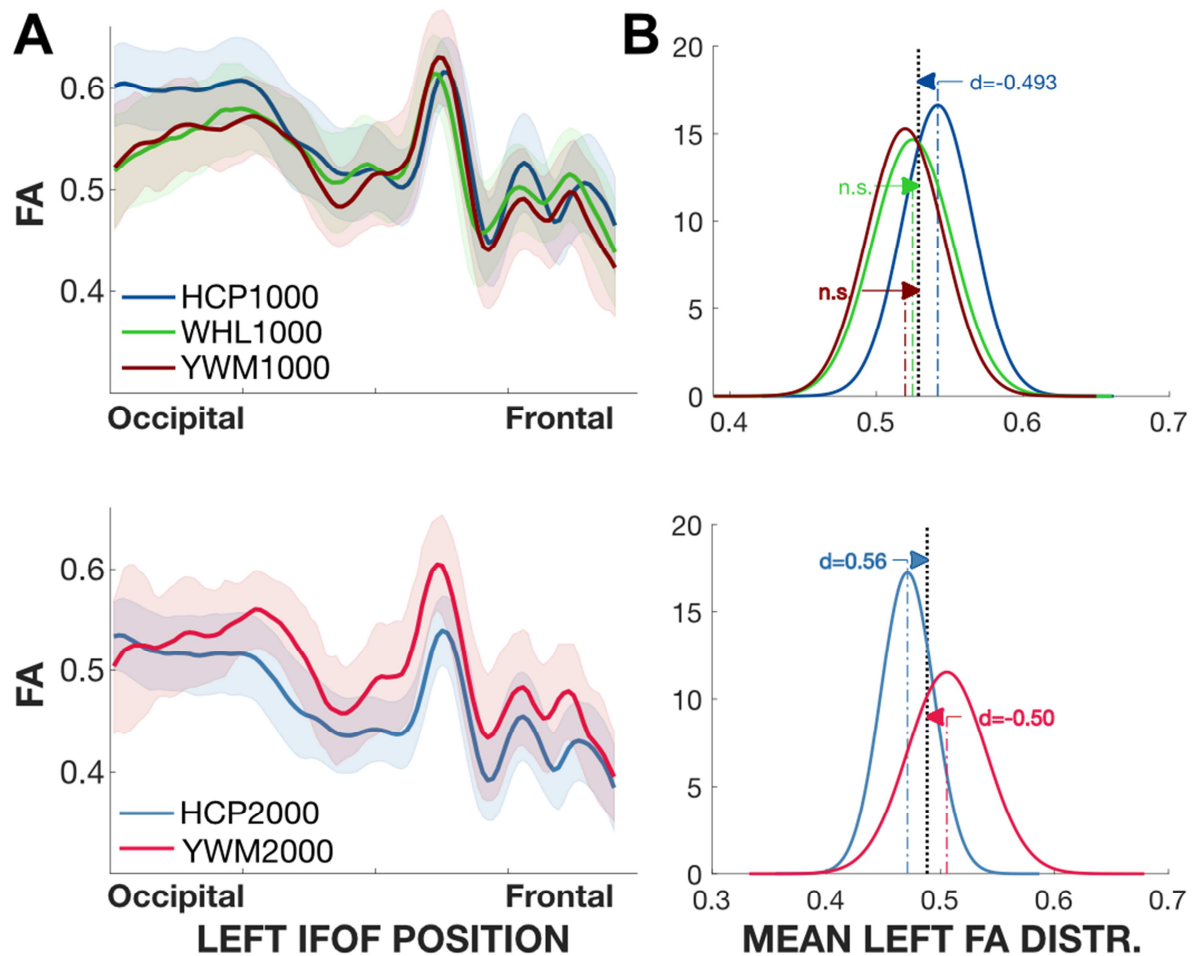


1 the largest between-subject standard deviation is for the arcuate fasciculus, and the smallest is for the  
2 corticospinal tract. In all cases, the shape of the tract profiles remain similar across b-values. This  
3 supports the idea that tract profiles are a useful target for further investigation (St-Jean et al., 2019;  
4 Yeatman et al., 2014).

### 5 3.1.2.- Generalization experiments

6 We assess the generalization of the FA measure by comparing the HCP data with those from YWM and  
7 WHL. Because of the large differences in FA, we separate the analysis by b-value (Figure 4).

8



**Figure 4. FA Analyses for the generalization experiment and selected tracts.**

Top: Left Corticospinal. Bottom: Left IFOF. **A)** The curves show the average FA tract profiles for different experiments. The shaded region is  $\pm 1$  SD. **B)** Normal distribution summary of the mean FA values in each experiment. The mean is the average of the FA values of each participant's profile. The grey plot shows the distribution for all experiments; the curves are scaled so that the sum of the areas of the experiments equals the grey area. The arrows show the difference between each of the means and the group mean, and the numbers express effect size (Cohen's  $d$ ). The distributions were estimated using 10,000 bootstrap samples. **C)** Mean FA values and 90% experimental confidence intervals. *n.s.*: non-significant. Plots for additional tracts are in the Supplementary Materials (Figure S1b-S1c-S3a-S3b).

- 1
- 2 The HCP, YWM and WHL data obtained at  $b=1000$  are compared in the top two panels. We use the IFOF
- 3 tract, but the conclusions are the same for other tracts (see Supplementary material Figures S1b,S1c-S3a-

1 S3b). The tract profile from the HCP dataset is the same as that shown in [Figure 2](#), and the green and red  
2 curves are from the WHL and YWM datasets, respectively. Over much of the tract the three data sets  
3 agree in the sense that they are closer than the between-subject variance. The HCP tractogram profile  
4 diverges from the YWM and WHL on the left side of the graph (occipital end), and this appears to be the  
5 largest source of the difference between the three datasets.

6  
7 The distribution of HCP mean tract FA values are about 1 standard deviation larger than the values in the  
8 YWM and WHL data set, and this causes the precision of the generalization to be substantially lower than  
9 the precision of the replication (Figure 4B, top). It is notable that at  $b=1000$  the mean FA tract values for  
10 the IFOF in the WHL and YWM data sets contain values that are never observed in the HCP data set ( $FA$   
11  $< 0.47$ ). The expansion of the range of FA values provides an indication of what one would observe in a  
12 clinical application compared to measurements obtained at a single site.

13  
14 The HCP and YWM data obtained at  $b=2000$  are compared in the two bottom panels. In this  
15 measurement the HCP FA values are generally lower than the YWM FA values. This difference is seen  
16 in the mean FA distributions, which are again separated by about 1 standard deviation. It is notable that  
17 at  $b=2000$  the mean FA values for the IFOF include values in the YWM data that are never observed in  
18 the HCP data (e.g.,  $FA > 0.55$ ). Again, the generalization analysis shows that combining data from  
19 multiple sites extends the range of FA values one would observe from healthy participants.

### 20 3.1.3.- Evaluation

21 The analyses of replication and generalization do not force a conclusion about whether the technique may  
22 have value in practice. Rather, the analyses define the range of values one might observe using a  
23 restricted set of instruments and methods (replication), compared to the range of values observed as we  
24 measure in clinical applications (generalization). For most tracts, the range of the mean tract FA value

- 1 increases by about a factor of two as we include data from different, but typical, instruments and sites.  
 2 Adding more sites, or expanding the population, can only increase this factor.

### 3 **3.2.- Homologous tract FA values**

4 The evaluation of mean tract FA motivated us to search for a dependent measure with better  
 5 generalization. The high positive correlation in FA between pairs of homologous tracts (Wahl et al.,  
 6 2010), measured across subjects, suggests an alternative measure. The correlation implies that a  
 7 participant with a relatively high FA value in the left tract will have a relatively high FA in the  
 8 homologous right tract. Using this type of measure has the potential to improve generalization because  
 9 measurements of the two tracts depend on common experimental factors. Qualitatively, the measurements  
 10 of the left tract serve as calibration data to predict the FA measurement of the right tract. This is  
 11 analogous to the use of image contrast rather than image value.

#### 12 **3.2.1.- Homologous tracts linear model**

13 The next question we address is how to convert the observed correlations, obtained from multiple  
 14 participants, into a measurement that can be applied to individual participants. The initial approach is to  
 15 use the linear model implicit in the correlation. Specifically, the correlation between homologous tracts  
 16 means that there is an affine transform that predicts the mean tract FA in the right from knowledge of the  
 17 left.

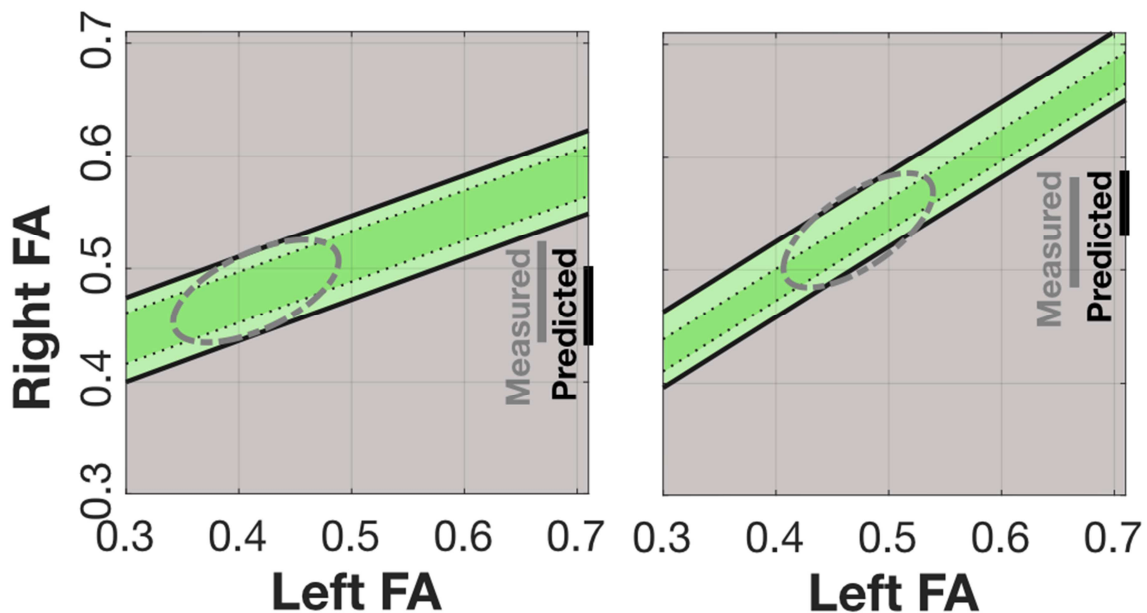
$$19 \quad \text{PredictedRight}_{FA} = \alpha \text{MeasuredLeft}_{FA} + \beta$$

20  
 21 The prediction error (residuals) are the difference between the measured and predicted FA,

$$22 \quad \text{Residuals} = \text{MeasuredRight}_{FA} - \text{PredictedRight}_{FA},$$

23

1  
 2 and bootstrapping with replacement from the residuals we estimate the FA range where we expect to find  
 3 some percentage, say 95%, of the data (Figure 5: the vertical black line represents this range). If we  
 4 calculate the range of possible  $PredictedRight_{FA}$  values for all  $MeasuredLeft_{FA}$  values, we obtain a  
 5 band of likely  $PredictedRight_{FA}$  values (green bands). The center of the band is the linear prediction  
 6 and the dashed (solid) lines represents the 68% (95%) limits. Given a measurement of the left FA, the  
 7 band defines the range of expected values for the MeasuredRightFA in a healthy participant.



**Figure 5. Representation of the relation between homologous tracts**

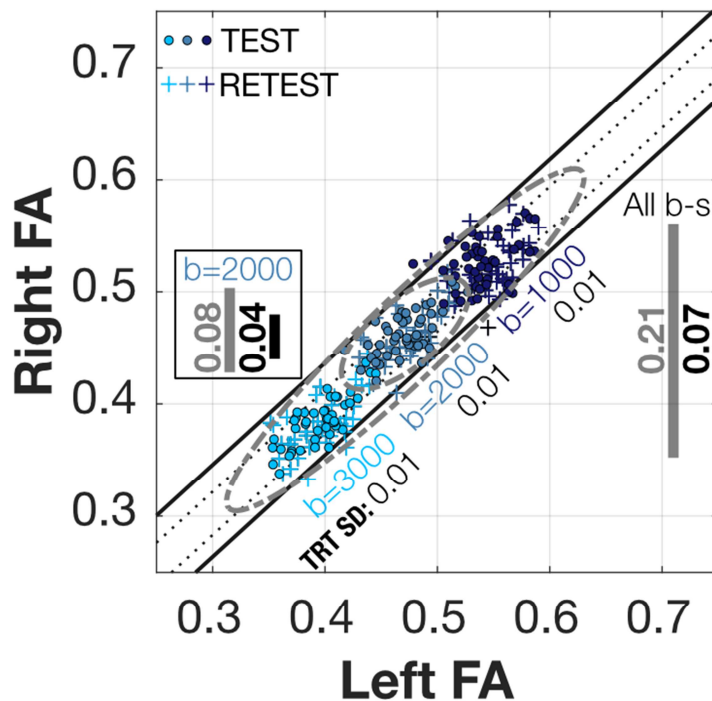
The linear correlation between mean FA in homologous white matter tracts defines a band of predicted right FA values given a left FA value. Measurements across clinically relevant cases, including variations in population, data acquisition, and computational methods, define the correlation and the size and shape of this region. For each tract, a participant's data may fall inside or outside the green region, and this serves as a diagnostic of their white matter health. *Measured*: the range of Right FA values. *Predicted*: the size of the range of predicted Right FA values given a Left FA value (vertical height of the green band).

8  
 9 A more general formulation, beyond the linear relation, assesses the distribution of left-right FA values in  
 10 the plane. These distributions form a cloud of points in the plane that can be reasonably approximated by

1 a bivariate Gaussian. Consequently, the likely locations of the points are circumscribed by an ellipse. The  
 2 distance of any single point from the center of the ellipse, say measured by the Mahalanobis distance, can  
 3 serve as a measure of the participant's health in a clinical application. This formulation has the added  
 4 benefit incorporating additional information: the absolute value of tract mean FA.

### 5 3.2.2.- Replication of the linear model

6 A scatterplot of the mean tract FA of the left and right IFOF for six HCP data sets (three b-values, test-  
 7 retest) is in [Figure 6](#). The different blue colors represent measurements at different b-values, and the  
 8 different shapes represent test (circles) and retest (crosses) measurements. The slope of the linear relation  
 9 between the mean tract FA of the left-right IFOF tracts is slightly less than one. Each pair of tracts has its  
 10 own best-fitting line (see Figure S4a).



**Figure 6. Left-Right IFOF FA scatterplots and iso-residual contour lines for HCP**

Scatterplot of the Left-Right IFOF HCP mean FA values. Inside the square, the grey line (0.08) shows the 95% range of all Right FA values for  $b=2000$ , and the black line (0.04) shows the range of possible values for any given Left FA value. Although not pictured, the values when using the  $b=2000$  test-retest data points increases to 0.09 and 0.05. Outside the square to the right, the grey line (0.21) shows the 95% range of all right FA values for the combined six HCP Test-Retest values. The diagonal bands are the contour lines holding the 68% and 95% of the residuals from the linear model fitted to all the six datasets. See Figure S4a for the rest of the tracts.

1  
 2 The test-retest data points thoroughly intermingle, which is a replication of the left-right linear relation.  
 3 The mean tract FA of a single tract replicates with a precision of 0.01 s.d. (Figure 3), and the separation in  
 4 the FA plane for mean tract FA of left-right homologous tracts (corresponding circles and crosses)  
 5 replicates with the same precision (0.01 s.d.).  
 6  
 7 The data obtained at the three different b-values fall along roughly the same line. Consequently, this left-  
 8 right measurement generalizes well across b-values, despite the fact that the mean tract FA values do not  
 9 (Figure 3). Considering the data from the three b-values, 95% of the FA measurements fall within 0.21  
 10 FA (gray line at right). Correspondingly, for the left-right difference 95% of the measurements fall within  
 11 0.07 FA (black line at right).  
 12  
 13 In certain cases, different sites may adopt measurement protocols at a single b-value. In that case, the  
 14 range of the left-right difference is reduced. For example, in the  $b=2000$  data set the FA range would be  
 15 reduced to 0.04 FA, which is smaller than the FA range across subjects (0.08 FA).  
 16  
 17 There are different causes for the range of FA values between subjects. Some of the differences are likely  
 18 to be the natural variation between subjects. Additional variation may be due to uncontrolled instrumental

1 factors. The co-linear relation between data obtained at the different b-values suggests that some  
2 differences arise because the nominal and true gradient (b-value) differs between subjects.

### 3 3.2.3.- Generalization of the linear model

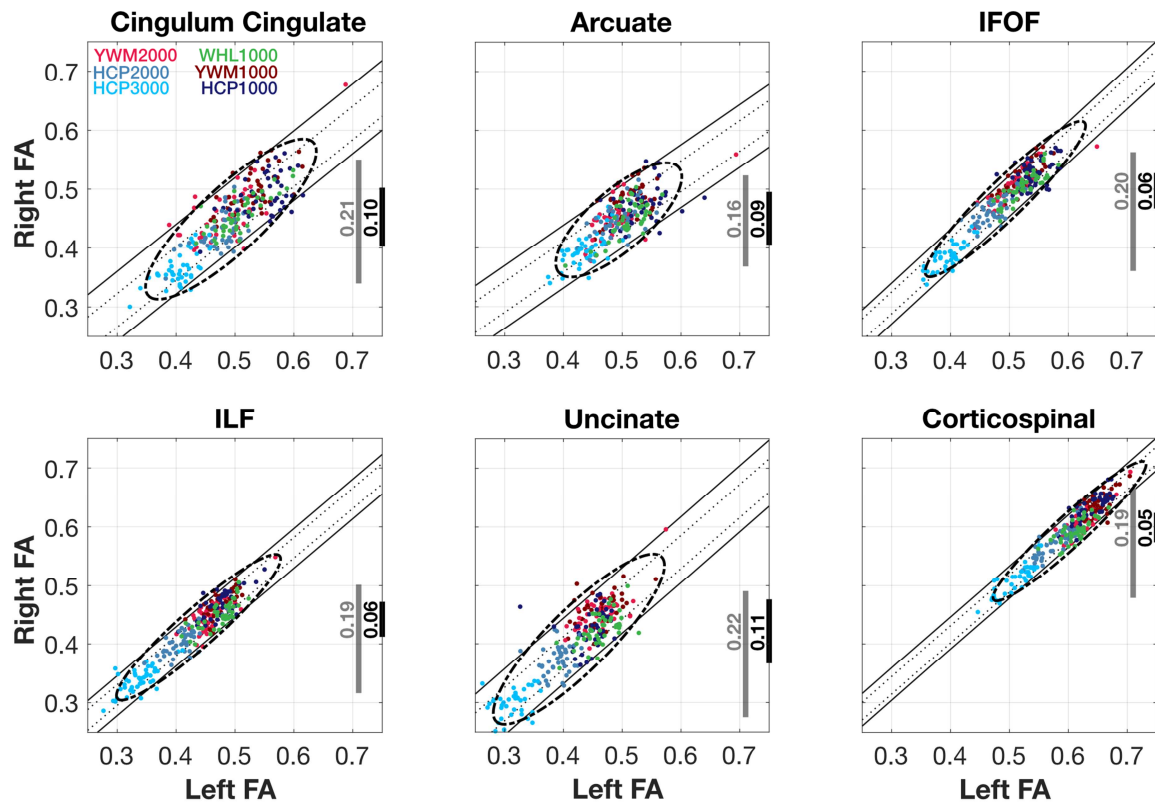
4 To assess generalization we combined the six datasets at three b-values (b=1000, 2000, 3000) and three  
5 sites (WHL,YWM,HCP). The left-right scatterplots, one for each of the six pairs of tracts, are shown in  
6 (Figure 7). There are qualitative similarities between data from different tracts, but each has its own  
7 parameters and precision.

8  
9 The left-right scatterplots of the IFOF, ILF and CST are the most compact. Given a measurement of the  
10 left mean tract FA, the right mean tract FA falls within about 0.05 FA. For the Cingulum, Arcuate and  
11 Uncinate the left mean tract FA predicts the right mean tract FA within about 0.10 FA. In all cases the  
12 slopes of the linear regions (orientation of the principal axis of the ellipse) are near one.

13

14





**Figure 7. Homologous tract Left-Right FA scatterplots and iso-residual contour lines**

Scatterplot of the Left-Right mean FA for all tracts and all projects. The grey vertical lines shows the 95% range of all Right FA values, and the black line the range of possible values for any given Left FA value. The diagonal bands are the iso-residual contour lines holding the 68% and 95% of the residuals from the linear model fitted to all the six datasets.

1

2 The left-right relation generalizes across the different sites and b-values. For each of the tracts, there is no  
 3 substantial loss of FA precision when calculating the left-right difference using the data at a single  
 4 nominal b-value or data from all b-values at all sites.

5

6 The scatter plots reveal outliers in the cohort, and one particular sample point stands out. This point arises  
 7 from a single subject at b=2000 who is an outlier in all of the tracts (YWM2000 data, red dot). In a  
 8 clinical setting, this subject would be subject to more scrutiny. We can compare this subject's data to the

1 acquisition at  $b=1000$  (YWM1000). The subject's FA values in the YWM1000 acquisition are normal, so  
2 we assume that something went wrong in the YWM2000 acquisition and/or analyses. Such outliers occur,  
3 and it is not unexpected that one of 264 data points might be problematic

4  
5 Some of the variation in the mean tract FA arises from the tractography algorithms. For example, the  
6 Arcuate and Uncinate are more curved than the other tracts, and previously several groups observed that  
7 the right Arcuate is not well-recovered from DWI data (Catani et al., 2007; Lebel and Beaulieu, 2009;  
8 Wahl et al., 2010; Yeatman et al., 2011). Other differences may arise because of differences in the length  
9 of the trunks used to estimate the mean FA of each tract (see Figure 2).

10  
11 Similar variability was observed in five of the homologous tract pair correlations in the original Wahl et  
12 al (2010) experiment (Cingulum Cingulate: 0.57, Arcuate: 0.5, IFOF: 0.88, ILF: 0.73, Uncinate: 0.7),  
13 with the one exception of the corticospinal tract (0.62). The original Wahl et al. result for corticospinal  
14 may be due to their method of identifying the corticospinal tract; because using our tractography methods  
15 on the original data (WHL1000) the value is higher (0.71).

16  
17 The correlation values of the data combined across b-values are very high (Cingulum Cingulate: 0.85,  
18 Arcuate: 0.76, IFOF: 0.94, ILF: 0.94, Uncinate: 0.87, Corticospinal: 0.95). This suggests that as hoped the  
19 same left-right relation is revealed at different b-values and that using the relation rather than absolute FA  
20 levels compensates for variations in the data acquisition.

21

22

## 1 4.- Discussion

2 We write in support of the idea that modern neuroimaging is sufficiently mature to develop useful  
3 quantitative applications for structural neuroimaging. As an example, we showed how the test-retest MRI  
4 scans produce highly reliable diffusion measures, even when accounting for instrumental noise, system  
5 calibration between scans, and repeating the probabilistic numerical processing in the computational  
6 methods. On the other hand, the experiments confirm prior reports that the compliance range of the data  
7 acquisition parameters for FA does not extend to changes in the diffusion gradient b-value (Chou et al.,  
8 2013; Farrell et al., 2007b; Hutchinson et al., 2017; Landman et al., 2007). For this reason, we proposed:  
9 (i) a two-step assessment system (measure replication, measure generalization) to translate MRI metrics  
10 with potential to be useful in the clinic; and, (ii) a simple method for improving the precision of our  
11 metrics by using the relationship between two measurements that compensates for the acquisition  
12 differences.

### 13 4.1.- Replication-generalization tradeoff

14 There is a tradeoff between replication and generalization in neuroimaging. Over the past decade, the two  
15 extremes have been represented by: (1) the HCP for high-quality highly replicable anatomic, diffusion  
16 and functional imaging using custom-designed hardware (the Connectome scanner) and software (e.g.,  
17 multiband echo planar sequences) that were not generalizable to other platforms (Glasser et al., 2016);  
18 and (2) the Enhancing NeuroImaging Genetics through Meta-Analysis (ENIGMA) consortium that began  
19 with low quality and low precision imaging metrics that were primarily limited to gross macroscopic  
20 features such as total intracranial volume, but were platform-independent and did not require standardized  
21 sequences and therefore generalized for worldwide data aggregation.

22

1 The tension between these two goals is being addressed by specifying standardized pulse sequences  
2 across a broad range of scanners for multicenter studies. This approach is exemplified by the HCP  
3 Lifespan protocol (Bookheimer et al., 2019; Harms et al., 2018; Somerville et al., 2018), the ENIGMA  
4 protocol (Acheson et al., 2017; Adhikari et al., 2018; Kochunov et al., 2017), and protocols for Precision  
5 Medicine studies such as ADNI3 for Alzheimer disease (Reid et al., 2017; Zavaliangos-Petropulu et al.,  
6 2019) and TRACK-TBI for traumatic brain injury (Yuh et al., 2013). This approach is applicable to  
7 coordinated multi-center studies.

8  
9 The harmonization of measurements puts a strong emphasis on replication, hoping to limit the problem of  
10 generalization. There are economic and technology trademark issues that will prevent the widespread  
11 distribution of the most advanced instruments. Because there will be variations in clinical instrumentation  
12 and methods, we advocate for investigators to design tools and experiments that directly address  
13 generalization. The approach in this paper emphasizes collecting multiple datasets and then evaluating  
14 different dependent measures to select the ones that generalize. In this approach, it becomes important to  
15 specify the precision and the compliance range when reporting results for potential application, as  
16 different pathologies will have different requirements.

#### 17 **4.2.- Explicit measures of generalization and context of use**

18 Clinical applications should be based on measurements that replicate with confidence intervals that are  
19 compact enough to support a meaningful diagnostic. This attribute is crucial for the validation of  
20 “biomarkers” that can be widely used for biomedical science and clinical translation. A biomarker is  
21 defined by the US National Institutes of Health (NIH) and the US Food & Drug Administration (FDA) as  
22 “a defined characteristic that is measured as an indicator of normal biological processes, pathogenic  
23 processes, or responses to an exposure or intervention, including therapeutic interventions” (Naylor,  
24 2003). This definition encompasses brain imaging (Mayeux, 2004). Precision Medicine is “an emerging

1 approach for disease treatment and prevention that takes into account individual variability in genes,  
2 environment, and lifestyle for each person” (Collins and Varmus, 2015). Objectively quantifying these  
3 individual differences in order to tailor treatment and prevention strategies for specific patients requires  
4 validated biomarkers. Ensuring the reliability of these biomarkers over time in individual subjects is  
5 crucial for adequately testing the efficacy of precision medicine therapies (Senn, 2018). Throughout this  
6 work, we provided a valid range of normal values that gives the precision at which departures from  
7 normality can be measured. This range, like all the measurements we used in this work, is given in FA  
8 units, which is directly interpretable by any researcher or clinical practitioner.

9  
10 In addition, but less appreciated, is that neuroimaging applications deployed in the field will use a range  
11 of instruments, participant populations, and measurement protocols (Goodman et al., 2016); the range of  
12 conditions in the field will be wider than that encountered in scientific studies. It is important, therefore,  
13 to assess how effectively an applied measurement generalizes across the clinical conditions. For an  
14 applied measurement to scale from the lab to the clinic, the result must generalize across these  
15 measurement conditions. This range of conditions where the measurement is valid for a proposed  
16 application should be specified contained in the “context of use” that the FDA requires as part of the  
17 biomarker qualification process (Goodsaid and Mendrick, 2010). After our experiment, we could claim  
18 that the context of use of our metric is circumscribed to 3 Tesla MRI magnets and b-values between 1000  
19 and 3000. We think that the symmetry in homologous mean tract FA is a fundamental human biological  
20 trait, but we should extend our generalization experiments to extend its context of use.

### 21 **4.3 - A continuous aggregation platform**

22 The relatively recent increase in complexity of neuroimaging is a major complicating factor that impacts  
23 reproducible research. New MR instruments, analysis algorithms and the use of special participants have  
24 increased the size and complexity of datasets. In many neuroimaging publications, there is no realistic

1 chance that a reader can repeat the experimental data acquisition or even the computational analyses  
2 (Buckheit and Donoho, 1995; Sandve et al., 2013; Wilson et al., 2017). The best we can hope for is to be  
3 able to repeat, check, and explore portions of the computational analysis of the published data (Peng,  
4 2011; Stodden et al., 2014).

5  
6 The increase in computational power has also led to an increase in algorithm complexity and the number  
7 of user-defined parameters. Several authors have analyzed the effect of pipeline parameters and reported  
8 large impacts on fMRI data; the variations in the result as a function of the parameters can be quite  
9 significant. For example, the position of the peak activation may range over a cortical area of 25 cm<sup>2</sup>  
10 (Carp, 2012). (Yarkoni and Westfall, 2017) observe that we are often uncertain about critical parameters  
11 that must be in computational models. We can confirm that the general point also applies to DWI  
12 methods. It is our experience, too, that scientists find it very difficult to keep track of the specific  
13 parameters used in any particular analysis, and even fewer scientists record the combinations of  
14 parameters they used during data exploration (Baker, 2016).

15  
16 To overcome most of these problems, the system we used in this paper encapsulates the software and its  
17 dependencies in a container; it also stores the history of which analyses (and with what configuration)  
18 were run in the database. This approach overlaps with many of the proposals for scientific reproducibility.  
19 For example, (Poldrack et al., 2017) describe desiderata for reproducible research tools that closely align  
20 with those we have implemented.

21  
22 ... The entire analysis workflow (including both successful and failed analyses) would be  
23 completely automated in a workflow engine and packaged in a software container or virtual  
24 machine to ensure computational reproducibility. All data sets and results would be assigned  
25 version numbers to enable explicit tracking of provenance ... (page 124).

26

1 Furthermore, our system is extensible. We can add datasets to our neuroinformatics platform, analyze  
2 them with identical computational methods, and check how the compliance range of our measurement  
3 changes. Analogously, we can containerize a computational tool from another group, process our data  
4 again, and do the same checks. Therefore, new results sets can be continuously aggregated. In the long  
5 term, this continuous aggregation will continue to inform the compliance range, and it will naturally work  
6 towards the harmonization of measurement protocols: settings that worsen the compliance range will be  
7 abandoned. We think that this continuous aggregating and improvement process will provide a useful  
8 approach for translating scientific research to the clinic.

#### 9 **4.4.- Related research**

10 A particularly related recent investigation of DWI generalization considered data from 13 different 3T  
11 MRI scanners throughout the USA, representing all three major vendors (GE, Philips and Siemens),  
12 found a coefficient of variation (CoV) of 4.2% for the FA of whole-brain white matter, with the FA CoV  
13 varying from 2% to 6% for individual major white matter tracts (Palacios et al., 2017). That study was  
14 limited to a single subject, to scanners with similar hardware capabilities, and to a harmonized DTI  
15 protocol in which all major acquisition parameters are as similar as possible.

16  
17 This study extends that work by probing generalization across a wider range of acquisition parameters  
18 (e.g., spatial resolution, b-value, and the number of diffusion directions) using scanners with different  
19 hardware capabilities (e.g., 8 receiver channels vs 32 and 40 mT/m maximum gradient amplitude vs 100  
20 mT/m), and in different participant populations. The scope of our tests is for a very modest set of  
21 instruments, data acquisition parameters, and population statistics; but the generalization could have  
22 proved much worse. A fundamental difference in our work is the intention to vary the experimental  
23 conditions instead of harmonizing them, assessing how the instrumental variations affect the precision  
24 range.

1  
2 Furthermore, the generalization issues in neuroimaging applications are similar to those in other human  
3 research fields (He et al., 2015; Shavelson et al., 1989; Shavelson and Webb, 1991; Tipton, 2014); the  
4 issues are also closely linked to meta-analysis, which aggregate the outcomes of multiple studies  
5 (Evangelou and Ioannidis, 2013; Simpson and Pearson, 1904). The unique features of neuroimaging  
6 applications we discuss are that they are motivated by the observation that these applications are likely to  
7 arise from experimental measures that are not precisely controlled.

#### 8 **4.5.- Limitations and opportunities**

9 All the datasets were obtained from research environments. We obtained data from different sites to  
10 illustrate our point, but for a real experiment, more datasets with more variability should be included.  
11 Further generalization could come from scanners (models, mean field strength), acquisition sequences  
12 (e.g. dual-spin echo), population (e.g. age range) or computational methods (e.g. Tracula (Yendiki et al.,  
13 2011)). The database system we use is extensible; we can add data and re-evaluate the generalization  
14 should new dataset become available.

15  
16 This work assesses one type of structural data which eliminated the need to analyze the impact of  
17 experimental design. Developing a deeper understanding of such factors is important for clinical  
18 assessments using task-based functional MRI, say for psychiatric disorders. Such analyses introduce  
19 many new parameters including factors ranging from stimulus selection and delivery and subject  
20 instructions and compliance.

21  
22 Some functional experiments quantify characteristics of individual participants (e.g. defining V1). A  
23 much larger set of the scientific literature uses group comparisons. In many cases, it will not be clear how  
24 to convert a group comparison experiment into a clinical assessment of individual participants.



ACCEPTED MANUSCRIPT

1

## 2 5.- Conclusion

3 This paper illustrates an approach for translating neuroimaging findings from the lab to the clinic. We  
4 describe software tools designed for large data sets and computational reproducibility that are helpful  
5 calculating the impact of increasing the number of sites, experiments, different subjects, and/or the impact  
6 of higher quality instrumentation. We consider a full approach, from the definition of the data set for  
7 replication and generalization experiments, to the neuroinformatics platform and computational methods  
8 required to define an evaluate metrics with diagnostic value.

9

10

## 1 Acknowledgements

2 This work was supported by a Marie Skłodowska-Curie (H2020-MSCA-IF-2017-795807-ReCiModel)  
3 grant to G.L.-U. We thank the Simons Foundation Autism Research Initiative and Weston Havens  
4 foundation for support. We acknowledge research grant support from the James S. McDonnell  
5 Foundation, the Charles A. Dana Foundation, the American Society of Neuroradiology, the U.S. National  
6 Institutes of Health (R01 NS060776), and the Academic Senate of the University of California, San  
7 Francisco for the Wahl 2010 et al. dataset.

8

## 9 Competing financial interests

10 The authors declare that the research was conducted in the absence of any commercial or financial  
11 relationships that could be construed as a potential conflict of interest. Brian Wandell is a co-founder of  
12 Flywheel.io.

13

14

## 1 References

- 2 Acheson, A., Wijtenburg, S.A., Rowland, L.M., Winkler, A., Mathias, C.W., Hong, L.E., Jahanshad, N.,  
3 Patel, B., Thompson, P.M., McGuire, S.A., Sherman, P.M., Kochunov, P., Dougherty, D.M., 2017.  
4 Reproducibility of tract-based white matter microstructural measures using the ENIGMA-DTI  
5 protocol. *Brain Behav.* 7, e00615.
- 6 Adhikari, B.M., Jahanshad, N., Shukla, D., Turner, J., Grotegerd, D., Dannlowski, U., Kugel, H.,  
7 Engelen, J., Dietsche, B., Krug, A., Kircher, T., Fieremans, E., Veraart, J., Novikov, D.S., Boedhoe,  
8 P.S.W., van der Werf, Y.D., van den Heuvel, O.A., Ipser, J., Uhlmann, A., Stein, D.J., Dickie, E.,  
9 Voineskos, A.N., Malhotra, A.K., Pizzagalli, F., Calhoun, V.D., Waller, L., Veer, I.M., Walter, H.,  
10 Buchanan, R.W., Glahn, D.C., Hong, L.E., Thompson, P.M., Kochunov, P., 2018. A resting state  
11 fMRI analysis pipeline for pooling inference across diverse cohorts: an ENIGMA rs-fMRI protocol.  
12 *Brain Imaging Behav.* <https://doi.org/10.1007/s11682-018-9941-x>
- 13 Andersson, J.L.R., Skare, S., Ashburner, J., 2003. How to correct susceptibility distortions in spin-echo  
14 echo-planar images: application to diffusion tensor imaging. *Neuroimage* 20, 870–888.
- 15 Andersson, J.L.R., Sotiropoulos, S.N., 2016. An integrated approach to correction for off-resonance  
16 effects and subject movement in diffusion MR imaging. *Neuroimage* 125, 1063–1078.
- 17 Andersson, J.L.R., Sotiropoulos, S.N., 2015. Non-parametric representation and prediction of single- and  
18 multi-shell diffusion-weighted MRI data using Gaussian processes. *Neuroimage* 122, 166–176.
- 19 Baker, M., 2016. 1,500 scientists lift the lid on reproducibility. *Nature* 533, 452–454.
- 20 Bookheimer, S.Y., Salat, D.H., Terpstra, M., Ances, B.M., Barch, D.M., Buckner, R.L., Burgess, G.C.,  
21 Curtiss, S.W., Diaz-Santos, M., Elam, J.S., Fischl, B., Greve, D.N., Hagy, H.A., Harms, M.P., Hatch,  
22 O.M., Hedden, T., Hodge, C., Japardi, K.C., Kuhn, T.P., Ly, T.K., Smith, S.M., Somerville, L.H.,  
23 Uğurbil, K., van der Kouwe, A., Van Essen, D., Woods, R.P., Yacoub, E., 2019. The Lifespan  
24 Human Connectome Project in Aging: An overview. *Neuroimage* 185, 335–348.

- 1 Buckheit, J.B., Donoho, D.L., 1995. WaveLab and Reproducible Research. *Wavelets and Statistics*.  
2 [https://doi.org/10.1007/978-1-4612-2544-7\\_5](https://doi.org/10.1007/978-1-4612-2544-7_5)
- 3 Carp, J., 2012. On the plurality of (methodological) worlds: estimating the analytic flexibility of fMRI  
4 experiments. *Front. Neurosci.* 6, 149.
- 5 Catani, M., Allin, M.P.G., Husain, M., Pugliese, L., Mesulam, M.M., Murray, R.M., Jones, D.K., 2007.  
6 Symmetries in human brain language pathways correlate with verbal recall. *Proc. Natl. Acad. Sci. U.*  
7 *S. A.* 104, 17163–17168.
- 8 Chou, M.C., Kao, E.F., Mori, S., 2013. Effects of b-value and echo time on magnetic resonance diffusion  
9 tensor imaging-derived parameters at 1.5 t: A voxel-wise study. *J. Med. Biol. Eng.* 33, 45–50.
- 10 Collins, F.S., Varmus, H., 2015. A new initiative on precision medicine. *N. Engl. J. Med.* 372, 793–795.
- 11 Evangelou, E., Ioannidis, J.P.A., 2013. Meta-analysis methods for genome-wide association studies and  
12 beyond. *Nat. Rev. Genet.* 14, 379–389.
- 13 Farrell, J.A.D., Landman, B.A., Jones, C.K., Smith, S.A., Prince, J.L., van Zijl, P.C.M., Mori, S., 2007a.  
14 Effects of signal-to-noise ratio on the accuracy and reproducibility of diffusion tensor imaging--  
15 derived fractional anisotropy, mean diffusivity, and principal eigenvector measurements at 1.5 T.  
16 *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for*  
17 *Magnetic Resonance in Medicine* 26, 756–767.
- 18 Farrell, J.A.D., Landman, B.A., Jones, C.K., Smith, S.A., Prince, J.L., Van Zijl, P.C.M., Mori, S., 2007b.  
19 Effects of signal-to-noise ratio on the accuracy and reproducibility of diffusion tensor imaging--  
20 derived fractional anisotropy, mean diffusivity, and principal eigenvector measurements at 1.5T. *J.*  
21 *Magn. Reson. Imaging* 26, 756–767.
- 22 Glasser, M.F., Smith, S.M., Marcus, D.S., Andersson, J., Auerbach, E.J., Behrens, T.E.J., Coalson, T.S.,  
23 Harms, M.P., Jenkinson, M., Moeller, S., Robinson, E.C., Sotiropoulos, S.N., Xu, J., Yacoub, E.,  
24 Ugurbil, K., Van Essen, D.C., 2016. The Human Connectome Project's neuroimaging approach. *Nat.*  
25 *Neurosci.* In press, 1175–1187.
- 26 Glasser, M.F., Sotiropoulos, S.N., Wilson, J.A., Coalson, T.S., Fischl, B., Andersson, J.L., Xu, J., Jbabdi,

- 1 S., Webster, M., Polimeni, J.R., Van Essen, D.C., Jenkinson, M., WU-Minn HCP Consortium, 2013.  
2 The minimal preprocessing pipelines for the Human Connectome Project. *Neuroimage* 80, 105–124.
- 3 Goodman, S.N., Fanelli, D., Ioannidis, J.P.A., 2016. What does research reproducibility mean? *Sci.*  
4 *Transl. Med.* 8, 341ps12.
- 5 Goodsaid, F.M., Mendrick, D.L., 2010. Translational medicine and the value of biomarker qualification.  
6 *Sci. Transl. Med.* 2, 47ps44.
- 7 Harms, M.P., Somerville, L.H., Ances, B.M., Andersson, J., Barch, D.M., Bastiani, M., Bookheimer,  
8 S.Y., Brown, T.B., Buckner, R.L., Burgess, G.C., Coalson, T.S., Chappell, M.A., Dapretto, M.,  
9 Douaud, G., Fischl, B., Glasser, M.F., Greve, D.N., Hodge, C., Jamison, K.W., Jbabdi, S., Kandala,  
10 S., Li, X., Mair, R.W., Mangia, S., Marcus, D., Mascali, D., Moeller, S., Nichols, T.E., Robinson,  
11 E.C., Salat, D.H., Smith, S.M., Sotiropoulos, S.N., Terpstra, M., Thomas, K.M., Tisdall, M.D.,  
12 Ugurbil, K., van der Kouwe, A., Woods, R.P., Zöllei, L., Van Essen, D.C., Yacoub, E., 2018.  
13 Extending the Human Connectome Project across ages: Imaging protocols for the Lifespan  
14 Development and Aging projects. *Neuroimage* 183, 972–984.
- 15 He, Z., Chandar, P., Ryan, P., Weng, C., 2015. Simulation-based Evaluation of the Generalizability Index  
16 for Study Traits. *AMIA Annu. Symp. Proc.* 2015, 594–603.
- 17 Hutchinson, E.B., Avram, A.V., Irfanoglu, M.O., Koay, C.G., Barnett, A.S., Komlosh, M.E., Özarlan,  
18 E., Schwerin, S.C., Juliano, S.L., Pierpaoli, C., 2017. Analysis of the effects of noise, DWI sampling,  
19 and value of assumed parameters in diffusion MRI models. *Magn. Reson. Med.* 78, 1767–1780.
- 20 Jones, D.K., Basser, P.J., 2004. “Squashing peanuts and smashing pumpkins”: how noise distorts  
21 diffusion-weighted MR data. *Magnetic Resonance in Medicine: An Official Journal of the*  
22 *International Society for Magnetic Resonance in Medicine* 52, 979–993.
- 23 Kellner, E., Dhital, B., Kiselev, V.G., Reisert, M., 2016. Gibbs-ringing artifact removal based on local  
24 subvoxel-shifts. *Magn. Reson. Med.* 76, 1574–1581.
- 25 Kennedy, D.N., Abraham, S.A., Bates, J.F., Crowley, A., Ghosh, S., Gillespie, T., Goncalves, M., Grethe,  
26 J.S., Halchenko, Y.O., Hanke, M., Haselgrove, C., Hodge, S.M., Jarecka, D., Kaczmarzyk, J.,

- 1 Keator, D.B., Meyer, K., Martone, M.E., Padhy, S., Poline, J.-B., Preuss, N., Sincomb, T., Travers,  
2 M., 2019. Everything Matters: The ReprONim Perspective on Reproducible Neuroimaging. *Front.*  
3 *Neuroinform.* 13, 1.
- 4 Kochunov, P., Dickie, E.W., Viviano, J.D., Turner, J., Kingsley, P.B., Jahanshad, N., Thompson, P.M.,  
5 Ryan, M.C., Fieremans, E., Novikov, D., Veraart, J., Hong, E.L., Malhotra, A.K., Buchanan, R.W.,  
6 Chavez, S., Voineskos, A.N., 2017. Integration of routine QA data into mega-analysis may improve  
7 quality and sensitivity of multisite diffusion tensor imaging studies. *Hum. Brain Mapp.* 1–9.
- 8 Landman, B.A., Farrell, J.A.D., Jones, C.K., Smith, S.A., Prince, J.L., Mori, S., 2007. Effects of diffusion  
9 weighting schemes on the reproducibility of DTI-derived fractional anisotropy, mean diffusivity, and  
10 principal eigenvector measurements at 1.5T. *Neuroimage* 36, 1123–1138.
- 11 Lebel, C., Beaulieu, C., 2009. Lateralization of the arcuate fasciculus from childhood to adulthood and its  
12 relation to cognitive abilities in children. *Human Brain Mapping*. <https://doi.org/10.1002/hbm.20779>
- 13 Lerma-Usabiaga, G., Perry, M., Wandell, B.A., 2019. Reproducible Tract Profiles (RTP): from diffusion  
14 MRI acquisition to publication. <https://doi.org/10.1101/680173>
- 15 Marcus, D., Harwell, J., Olsen, T., Hodge, M., Glasser, M., Prior, F., Jenkinson, M., Laumann, T.,  
16 Curtiss, S., Van Essen, D., 2011. Informatics and data mining tools and strategies for the human  
17 connectome project. *Front. Neuroinform.* 5, 4.
- 18 Mayeux, R., 2004. Biomarkers: potential uses and limitations. *NeuroRx* 1, 182–188.
- 19 McNaught, A.D., Wilkinson, A. (Eds.), 1997. IUPAC. Compendium of Chemical Terminology (the  
20 “Gold Book”) - Reproducibility, Second. ed. Blackwell Scientific Publications, Oxford.
- 21 Mukherjee, P., Berman, J.I., Chung, S.W., Hess, C.P., Henry, R.G., 2008a. Diffusion tensor MR imaging  
22 and fiber tractography: theoretic underpinnings. *AJNR Am. J. Neuroradiol.* 29, 632–641.
- 23 Mukherjee, P., Chung, S.W., Berman, J.I., Hess, C.P., Henry, R.G., 2008b. Diffusion tensor MR imaging  
24 and fiber tractography: technical considerations. *AJNR Am. J. Neuroradiol.* 29, 843–852.
- 25 Naylor, S., 2003. Biomarkers: current perspectives and future prospects. *Expert Rev. Mol. Diagn.* 3, 525–  
26 529.

- 1 Palacios, E.M., Martin, A.J., Boss, M.A., Ezekiel, F., Chang, Y.S., Yuh, E.L., Vassar, M.J., Schnyer,  
2 D.M., MacDonald, C.L., Crawford, K.L., Irimia, A., Toga, A.W., Mukherjee, P., TRACK-TBI  
3 Investigators, 2017. Toward Precision and Reproducibility of Diffusion Tensor Imaging: A  
4 Multicenter Diffusion Phantom and Traveling Volunteer Study. *AJNR Am. J. Neuroradiol.* 38, 537–  
5 545.
- 6 Patil, P., Peng, R.D., Leek, J., 2016. A statistical definition for reproducibility and replicability. *bioRxiv.*  
7 <https://doi.org/10.1101/066803>
- 8 Peng, R.D., 2011. Reproducible research in computational science. *Science* 334, 1226–1227.
- 9 Pestilli, F., Yeatman, J.D., Rokem, A., Kay, K.N., Wandell, B.A., 2014. Evaluation and statistical  
10 inference for human connectomes. *Nat. Methods* 11, 1058–1063.
- 11 Plesser, H.E., 2017. Reproducibility vs. Replicability: A Brief History of a Confused Terminology. *Front.*  
12 *Neuroinform.* 11, 76.
- 13 Poldrack, R.A., Baker, C.I., Durnez, J., Gorgolewski, K.J., Matthews, P.M., Munafò, M.R., Nichols, T.E.,  
14 Poline, J.B., Vul, E., Yarkoni, T., 2017. Scanning the horizon: Towards transparent and reproducible  
15 neuroimaging research. *Nat. Rev. Neurosci.* 18, 115–126.
- 16 Reid, R.I., Borowski, B.J., Thostenson, K., Arani, A., Thomas, D.L., Cash, D.M., Zhang, H., Gunter, J.L.,  
17 Bernstein, M.A., DeCarli, C.S., Fox, N.C., Thompson, P.M., Tosun, D., Weiner, M., Jack, C.R.,  
18 2017. THE ADNI3 DIFFUSION MRI PROTOCOL: BASIC ADVANCED. *Alzheimer's &*  
19 *Dementia.* <https://doi.org/10.1016/j.jalz.2017.06.1542>
- 20 Rokem, A., Yeatman, J.D., Pestilli, F., Kay, K.N., Mezer, A., Van Der Walt, S., Wandell, B.A., 2015.  
21 Evaluating the accuracy of diffusion MRI models in white matter. *PLoS One* 10, 1–26.
- 22 Sandve, G.K., Nekrutenko, A., Taylor, J., Hovig, E., 2013. Ten simple rules for reproducible  
23 computational research. *PLoS Comput. Biol.* 9, e1003285.
- 24 Senn, S., 2018. Statistical pitfalls of personalized medicine. *Nature* 563, 619–621.
- 25 Shavelson, R.J., Webb, N.M., 1991. *Generalizability Theory: A Primer.* SAGE.
- 26 Shavelson, R.J., Webb, N.M., Rowley, G.L., 1989. Generalizability theory. *Am. Psychol.* 44, 922.



- 1 Simpson, R.J.S., Pearson, K., 1904. Report On Certain Enteric Fever Inoculation Statistics. *Br. Med. J.* 2,  
2 1243–1246.
- 3 Somerville, L.H., Bookheimer, S.Y., Buckner, R.L., Burgess, G.C., Curtiss, S.W., Dapretto, M., Elam,  
4 J.S., Gaffrey, M.S., Harms, M.P., Hodge, C., Kandala, S., Kastman, E.K., Nichols, T.E., Schlaggar,  
5 B.L., Smith, S.M., Thomas, K.M., Yacoub, E., Van Essen, D.C., Barch, D.M., 2018. The Lifespan  
6 Human Connectome Project in Development: A large-scale study of brain connectivity development  
7 in 5-21 year olds. *Neuroimage* 183, 456–468.
- 8 St-Jean, S., Chamberland, M., Viergever, M.A., Leemans, A., 2019. Reducing variability in along-tract  
9 analysis with diffusion profile realignment. *arXiv [q-bio.QM]*.
- 10 Stodden, V., Leisch, F., Peng, R.D., 2014. *Implementing Reproducible Research*. CRC Press.
- 11 Takemura, H., Caiafa, C.F., Wandell, B.A., Pestilli, F., 2016. Ensemble Tractography. *PLoS Comput.*  
12 *Biol.* 12, 1–22.
- 13 Tipton, E., 2014. How Generalizable Is Your Experiment? An Index for Comparing Experimental  
14 Samples and Populations. *J. Educ. Behav. Stat.* 39, 478–501.
- 15 Tustison, N.J., Avants, B.B., Cook, P.A., Zheng, Y., Egan, A., Yushkevich, P.A., Gee, J.C., 2010.  
16 N4ITK: improved N3 bias correction. *IEEE Trans. Med. Imaging* 29, 1310–1320.
- 17 Veraart, J., Fieremans, E., Novikov, D.S., 2016a. Diffusion MRI noise mapping using random matrix  
18 theory. *Magn. Reson. Med.* 76, 1582–1593.
- 19 Veraart, J., Novikov, D.S., Christiaens, D., Ades-Aron, B., Sijbers, J., Fieremans, E., 2016b. Denoising of  
20 diffusion MRI using random matrix theory. *Neuroimage* 142, 394–406.
- 21 Wahl, M., Li, Y.-O., Ng, J., Lahue, S.C., Cooper, S.R., Sherr, E.H., Mukherjee, P., 2010. Microstructural  
22 correlations of white matter tracts in the human brain. *Neuroimage* 51, 531–541.
- 23 Wilson, G., Bryan, J., Cranston, K., Kitzes, J., Nederbragt, L., Teal, T.K., 2017. Good enough practices in  
24 scientific computing. *PLoS Comput. Biol.* 13, e1005510.
- 25 Yarkoni, T., Westfall, J., 2017. Choosing Prediction Over Explanation in Psychology: Lessons From  
26 Machine Learning. *Perspect. Psychol. Sci.* 12, 1100–1122.

- 1 Yeatman, J.D., Dougherty, R.F., Myall, N.J., Wandell, B.A., Feldman, H.M., 2012. Tract profiles of  
2 white matter properties: automating fiber-tract quantification. *PLoS One* 7, e49790.
- 3 Yeatman, J.D., Dougherty, R.F., Rykhlevskaia, E., Sherbondy, A.J., Deutsch, G.K., Wandell, B.A., Ben-  
4 Shachar, M., 2011. Anatomical properties of the arcuate fasciculus predict phonological and reading  
5 skills in children. *J. Cogn. Neurosci.* 23, 3304–3317.
- 6 Yeatman, J.D., Wandell, B.A., Mezer, A., 2014. Maturation and degeneration of human white matter.  
7 *Nat. Commun.* 5, 1–12.
- 8 Yendiki, A., Panneck, P., Srinivasan, P., Stevens, A., Zöllei, L., Augustinack, J., Wang, R., Salat, D.,  
9 Ehrlich, S., Behrens, T.E.J., Jbabdi, S., Gollub, R., Fischl, B., 2011. Automated probabilistic  
10 reconstruction of white-matter pathways in health and disease using an atlas of the underlying  
11 anatomy. *Front. Neuroinform.* 5, 23.
- 12 Yuh, E.L., Mukherjee, P., Lingsma, H.F., Yue, J.K., Ferguson, A.R., Gordon, W.A., Valadka, A.B.,  
13 Schnyer, D.M., Okonkwo, D.O., Maas, A.I.R., Manley, G.T., TRACK-TBI Investigators, 2013.  
14 Magnetic resonance imaging improves 3-month outcome prediction in mild traumatic brain injury.  
15 *Ann. Neurol.* 73, 224–235.
- 16 Zavaliangos-Petropulu, A., Nir, T.M., Thomopoulos, S.I., Reid, R.I., Bernstein, M.A., Borowski, B., Jack,  
17 C.R., Jr, Weiner, M.W., Jahanshad, N., Thompson, P.M., 2019. Diffusion MRI Indices and Their  
18 Relation to Cognitive Impairment in Brain Aging: The Updated Multi-protocol Approach in ADNI3.  
19 *Front. Neuroinform.* 13, 2.

20

- Reproducible white matter diagnostics, generalizable to clinical conditions
- Data collected from multiple scanners into a searchable and computable database
- Data analysis software implemented as platform-independent, reproducible containers
- Strategy to minimize measurement variance across a likely span of clinical scanners

ACCEPTED MANUSCRIPT