

Evaluating the reliability of human brain white matter tractometry

John Kruper^{a,b}, Jason D. Yeatman^{c,d}, Adam Richie-Halford^b, David Bloom^{a,b}, Mareike Grotheer^{e,f}, Sedy Caffarra^{c,d,g}, Gregory Kiar^h, Iliana I. Karipidisⁱ, Ethan Roy^c, and Ariel Rokem^{1 a,b}

^aDepartment of Psychology, University of Washington, Seattle, WA, 98195, United States of America

^bScience Institute, University of Washington, Seattle, WA, 98195, United States of America

^cGraduate School of Education, Stanford University, Stanford, CA, 94305, United States of America

^dDivision of Developmental-Behavioral Pediatrics, Stanford University School of Medicine, Stanford, CA, 94305, United States of America

^eCenter for Mind, Brain and Behavior - CMBB, Hans-Meerwein-Straße 6, Marburg 35032, Germany

^fDepartment of Psychology, University of Marburg, Marburg 35039, Germany

^gBasque Center on Cognition, Brain and Language, BCBL, 20009, Spain

^hDepartment of Biomedical Engineering, McGill University, Montreal, H3A 0E9, Canada

ⁱCenter for Interdisciplinary Brain Sciences Research, Department of Psychiatry and Behavioral Sciences, Stanford University School of Medicine, Stanford, CA, 94305, United States of America

The validity of research results depends on the reliability of analysis methods. In recent years, there have been concerns about the validity of research that uses diffusion-weighted MRI (dMRI) to understand human brain white matter connections *in vivo*, in part based on reliability of the analysis methods used in this field. We defined and assessed three dimensions of reliability in dMRI-based tractometry, an analysis technique that assesses the physical properties of white matter pathways: (1) reproducibility, (2) test-retest reliability and (3) robustness. To facilitate reproducibility, we provide software that automates tractometry (<https://yeatmanlab.github.io/pyAFQ>). In measurements from the Human Connectome Project, as well as clinical-grade measurements, we find that tractometry has high test-retest reliability that is comparable to most standardized clinical assessment tools. We find that tractometry is also robust: showing high reliability with different choices of analysis algorithms. Taken together, our results suggest that tractometry is a reliable approach to analysis of white matter connections. The overall approach taken here both demonstrates the specific trustworthiness of tractometry analysis and outlines what researchers can do to demonstrate the reliability of computational analysis pipelines in neuroimaging.

Diffusion MRI | Brain Connectivity | Tractography | Reproducibility | Robustness

Correspondence: arokem@uw.edu

Introduction

The white matter of the brain contains the long-range connections between distant cortical regions. The integration and coordination of brain activity through the fascicles containing these connections is important for information processing and for brain health (1, 2). Using voxel-specific directional diffusion information from diffusion-weighted MRI (dMRI), computational tractography produces three-dimensional trajectories through the white matter within the MRI volume that are called “streamlines” (3, 4). Collections of streamlines that match the location and direction of major white matter pathways within an individual can be generated with different strategies: using probabilistic (5, 6) or streamline-based (7, 8) atlases, or known anatomical landmarks (9–12). Because these are models of the anatomy, we refer to these estimates as “bundles” to distinguish them from the anatomical path-

ways themselves. The delineation of well-known anatomical pathways overcomes many of the concerns about confounds in dMRI-based tractography (13, 14), because “brain connections derived from diffusion MRI tractography can be highly anatomically accurate – if we know where white matter pathways start, where they end, and where they do not go” (15).

The physical properties of the tissue affect the diffusion of water within the brain and the microstructure of tissue within the white matter along the length of computationally-generated bundles can be assessed using a variety of models (16, 17). Taken together, computational tractography, bundle segmentation and diffusion modeling provide so-called “tract profiles”: estimates of microstructural properties of tissue along the length of major pathways. This is the basis of tractometry: statistical analysis that compares different groups, or assesses individual variability in brain connection structure (9, 18? –20). For the inferences made from tractometry to be valid and useful, tract profiles need to be reliable.

In the present work, we provide an assessment of three different ways in which scientific results can be reliable: reproducibility, test-retest reliability, and robustness. These terms are often debated and conflicting definitions for these terms have been proposed (21, 22). Here, we use the definitions proposed in (23). *Reproducibility* is defined as the case in which data and methods are fully accessible and usable: running the same code with the same data should produce an identical result. Use of different data (e.g., in a test-retest experiment) resulting in quantitatively comparable results would denote *test-retest reliability* (TRR). In clinical science and psychology in general, TRR (e.g., in the form of inter-rater reliability) is considered a key metric of the reliability of a measurement. Use of a different analysis approach or different analysis system (e.g., different software implementation of the same ideas) could result in similar conclusions, denoting their *robustness* against implementation details. The recent findings of Botvinik-Nezer *et al* (24) show that even when full computational reproducibility is achieved, the results of analysing a single fMRI dataset can vary significantly between teams and analysis pipelines, demonstrating issues

of robustness.

The contribution of the present work is three-fold: To support reproducible research using tractometry, we developed an open-source software library called Automated Fiber Quantification in Python (pyAFQ; <https://yeatmanlab.github.io/pyAFQ>). Given dMRI data that has undergone standard preprocessing (e.g., using QSIprep (25)), pyAFQ automatically performs tractography, classifies streamlines into bundles representing the major tracts, and extracts tract profiles of diffusion properties along those bundles, producing “tidy” CSV output files (26) that are amenable to further statistical analysis (Fig. S1). The library implements the major functionality provided by a previous MATLAB implementation of tractometry analysis (9), and offers a menu of configurable algorithms allowing researchers to tune the pipeline to their specific scientific questions (Fig. S2). Second, we use pyAFQ to assess test-retest reliability of tractometry results. Third, we assess robustness of tractometry results to variations across different models of the diffusion in individual voxels, across different bundle recognition approaches, and across different implementations.

Materials and Methods

pyAFQ. We developed an open-source tractometry software library to support computational reproducibility: Python Automated Fiber Quantification (pyAFQ; <https://github.com/yeatmanlab/pyAFQ>). The software relies heavily on methods implemented in DIPY (27) and is also based on a previous Matlab implementation of tractometry (9). More details are available in the ‘Automated Fiber Quantification in Python (pyAFQ)’ section of Supplementary Methods.

Tractometry. The pyAFQ software is configurable, allowing users to specify methods and parameters for different stages of the analysis (Fig. S2). Here, we will describe the default setting. In the first step, computational tractography methods, implemented in DIPY (27), are used to generate streamlines throughout the brain white matter (Fig. S1A). Next, the T1-weighted MNI template (28, 29) is registered to the anisotropic power map (APM) (30) computed from the diffusion data, that has a T1-like contrast (Fig. S1B) using the symmetric image normalization method (31) implemented in DIPY (27). The next step is to perform bundle recognition, where each tractography streamline is classified as either belonging to a particular bundle, or discarded. We use the transform found during registration to bring canonical anatomical landmarks, such as waypoint regions of interest (ROIs) and probability maps, from template space to the individual subject’s native space. Waypoint ROIs are used to delineate the trajectory of the bundles (32). See Table S1 for the bundle abbreviations we use in this paper. Streamlines that pass through inclusion waypoint ROIs for a particular bundle, and do not pass through exclusion ROI, are selected as candidates to include in the bundle. In addition, a probabilistic atlas (33) is used as a tie-breaker to determine whether a streamline is more likely

to belong to one bundle or another (in cases where the streamline matches the criteria for inclusion in either). For example, the corticospinal tract is identified by finding streamlines that do pass through an axial waypoint ROI in the brainstem and another ROI axially oriented in the white matter of the corona radiata, but that do not pass through the midline (Fig. S1C). The final step is to extract the tract profile: each streamline is resampled to a fixed number of points and the mean value of a diffusion-derived scalar (e.g., FA, MD) is found for each one of these nodes. The values are summarized by weighting the contribution of each streamline, based on how concordant the trajectory of this streamline is with respect to the other streamlines in the bundle (Fig. S1D).

Data. We used two datasets with test-retest measurements. We used Human Connectome Project test-retest measurements of dMRI for 44 neurologically healthy subjects aged 22-35 (HCP-TR) (34). The other is an experimental dataset, with dMRI from 48 children, 5 years old in age, collected at the University of Washington (UW-PREK). More details about the measurement are available in the ‘Data’ section of Supplementary Methods.

HCP-TR Configurations. We processed HCP-TR with three different pyAFQ configurations. In the first configuration, we used DKI as the ODF model. In the second configuration, we used CSD as the ODF model. For the final configuration, we used Recobundles (8) for bundle recognition instead of the default waypoint ROI approach, and DKI as the ODF model. More details are available in the ‘Configurations’ section of Supplementary Methods.

Measures of Reliability. Tract segmentation of each bundle was compared across measurements and methods using the Dice coefficient, weighted by streamline count (wDSC) (35). Tract profiles were compared with three measures: (1) Profile reliability: mean Pearson’s correlation between each point in the tract profile for different data; (2) Subject reliability: Pearson’s correlation between the mean of the tract profiles across individuals, which quantifies reliable differences between individuals per bundle for mean tract profiles; (3) an adjusted contrast index profile to directly compare the values of individual nodes in the tract profiles in different measurements. To estimate test-retest reliability (TRR), the above measures were calculated for each individual across different measurements. To estimate robustness, these were calculated for each individual across different analysis methods. More details are available in the ‘Measures of Reliability’ section of Supplementary Methods.

Results

Tractometry using pyAFQ classifies streamlines into bundles that represent major anatomical pathways. The streamlines are used to sample dMRI-derived scalars into bundle profiles that are calculated for every individual and can be summarized for a group of subjects. An example of the process and

result of the tract profile extraction process is shown in Supplementary Fig. S3, together with the results of this process across the 18 major white matter pathways for all subjects in the HCP-TR dataset.

Assessing test-retest reliability of tractometry. In datasets with scan-rescan data we can assess test-retest reliability (TRR) at several different levels of tractometry. For example, the correlation between two profiles provides a measure of the reliability of the overall tract profile in that subject. Analyzing the Human Connectome Project’s test-retest dataset (HCP-TR), we find that for fractional anisotropy (FA) calculated using the diffusion kurtosis model (DKI), the values of *profile reliability* vary across subjects (Figure 1A), but they overall tend to be rather high, with the average value within each bundle in the range 0.81 ± 0.04 to 0.95 ± 0.01 and a median across bundles of 0.89 (Figure 1B). We find similar results for mean diffusivity (MD; Fig. S4) and replicate similar results in a second dataset (Fig. 3B).

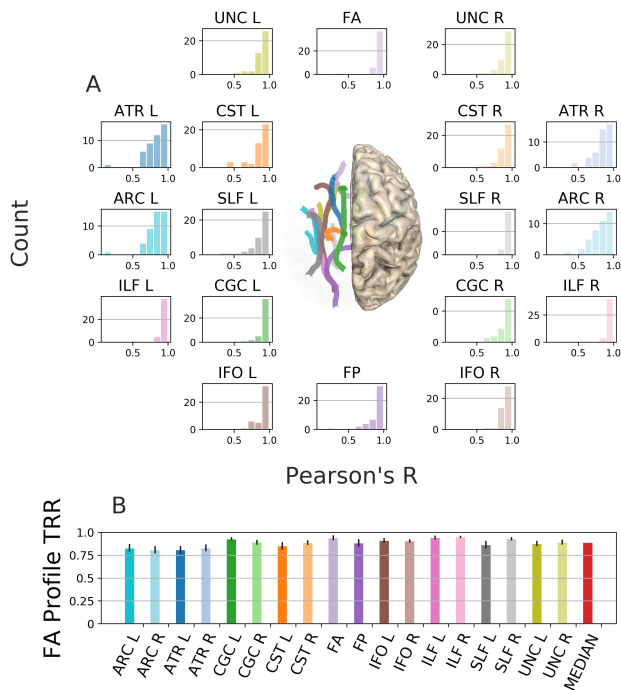


Fig. 1. FA profile test-retest reliability **A:** Histograms of individual subject Pearson’s r between the FA tract profiles across sessions for a given bundle. Colors encode the bundles, matching the diagram showing the rough anatomical positions of the bundles for the left side of the brain (center). **B:** Mean (\pm 95% confidence interval) TRR for each bundle, color-coded to match the histograms and the bundles diagram, with median across bundles in red.

Subject reliability assesses the reliability of mean tract profiles across individuals. Subject FA TRR in the HCP-TR also tends to be high, but the values vary more across bundles with a range of 0.64 ± 0.21 to 0.90 ± 0.08 and a median across bundles of 0.75. We can see that subject TRR is lower than profile TRR (Figure 2). This trend is consistent for MD (Fig. S5) as well as for another dataset (Fig. 3C).

One way of benchmarking whether these results are meaningful as an indication of high reliability is to compare them to the TRR of the images themselves. We assess measure-

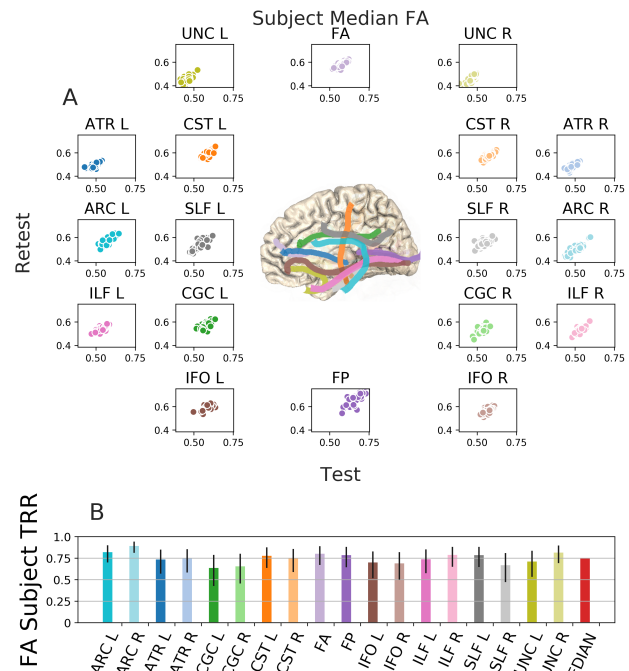


Fig. 2. Subject test-retest reliability **A:** Mean tract profiles for a given bundle and the FA scalar for each subject using the first and second session of HCP-TR. Colors encode bundle information, matching the core of the bundles (center). **B:** subject reliability is calculated from the Pearson’s correlation coefficient of these distributions, with median across bundles in red.

ment TRR by calculating the Pearson correlation coefficient between the diffusion-weighted measurements in each voxel in the white matter for each subject between the first and second measurement. For each subject, we extract the median value and then calculate the median across subjects. This is 0.94 on HCP in volumes where $b = 1000$, and 0.92 on HCP in volumes where $b > 0$ (including $b = 1000$, $b = 2000$, and $b = 3000$). This suggests that though there is some loss of information from the original measurement to tractometry within and individual, and further when examined across individuals, these are still highly reliable measures. This pattern is replicated in another dataset, which has a measurement TRR of 0.74 and is described in the “University of Washington Pre-K (UW-PREK)” section of the Supplement.

Test-retest reliability of tractometry in different implementations, datasets, and tractography methods.

We compared TRR across datasets and implementations. In both datasets, we found high TRR in the results of tractography and bundle recognition: wDSC was larger than 0.7 for all but one bundle (Fig. 3A): the delineation of the anterior forceps (FA) seems relatively unreliable using pyAFQ in the UW-PREK dataset (in FA, pyAFQ subject TRR is only 0.34 ± 0.28 compared to mAFQ’s 0.81 ± 0.12). We found overall high profile TRR that did not always translate to high subject TRR (Fig. 3B-G). For example, for FA in UW-PREK, median profile TRRs are 0.79 for pyAFQ and 0.81 for mAFQ while median subject TRRs are both only 0.73. mAFQ is one of the most popular software pipelines currently available for tractometry analysis, so it provides an important point for

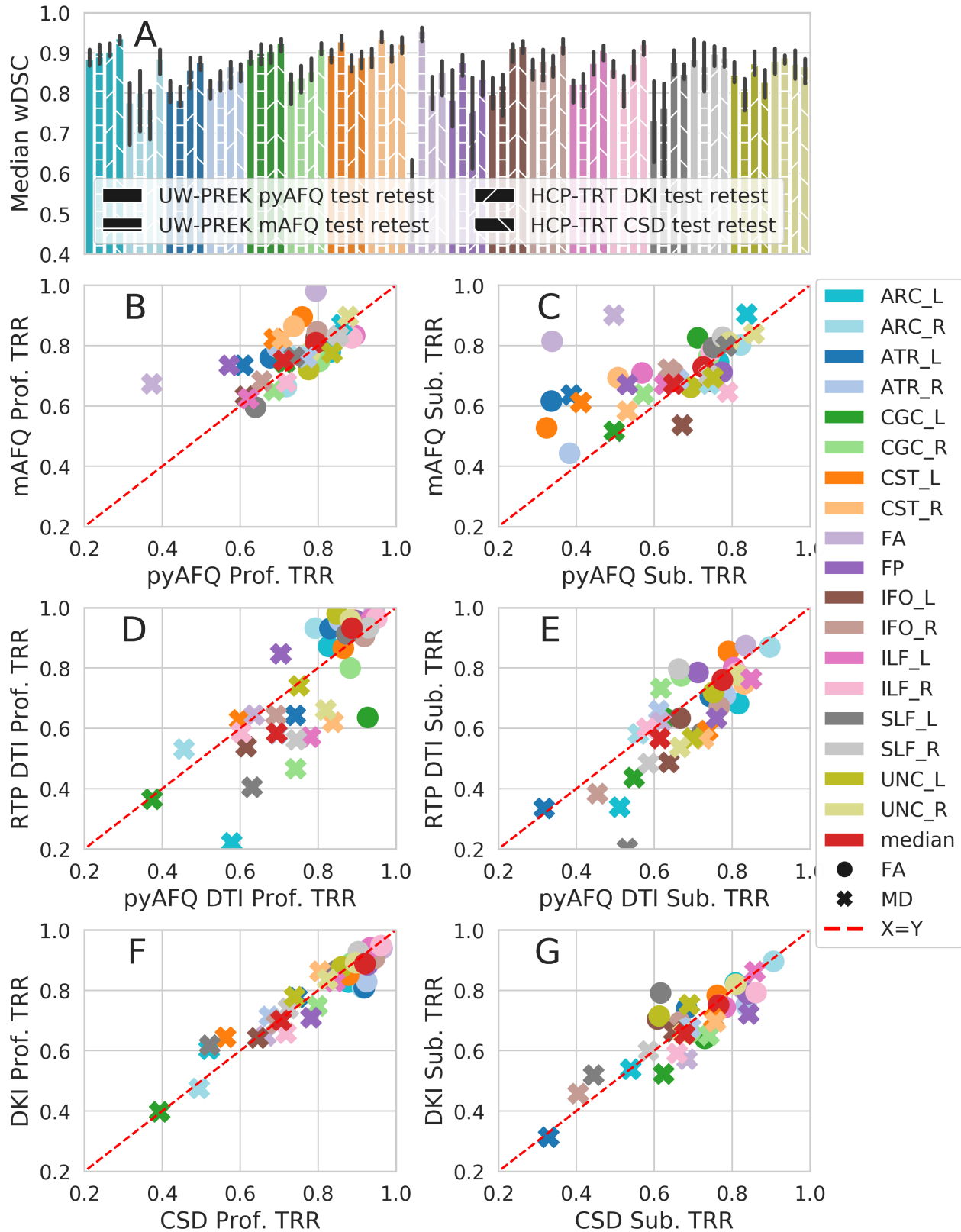


Fig. 3. wDSC, profile, and subject TRR of: pyAFQ and mAFQ on UW-PREK; pyAFQ on HCP-TR using different ODF models; and RTP on HCP-TR. Colors indicate bundle. In **A**: texture indicates the dataset and methods being compared. Error bars show the 95% confidence interval. **B**, **D**, and **F** show profile TRR and **C**, **E**, and **G** show subject TRR. In **B** and **C**, we compare the TRR of mAFQ and pyAFQ on UW-PREK. In **D** and **E**, we compare pyAFQ and RTP on HCP-TR using only single shell data. In **F** and **G**, we compare DKI and CSD TRR on HCP-TR. Point shapes indicate the extracted scalar. The red dotted line is equal TRR between methods.

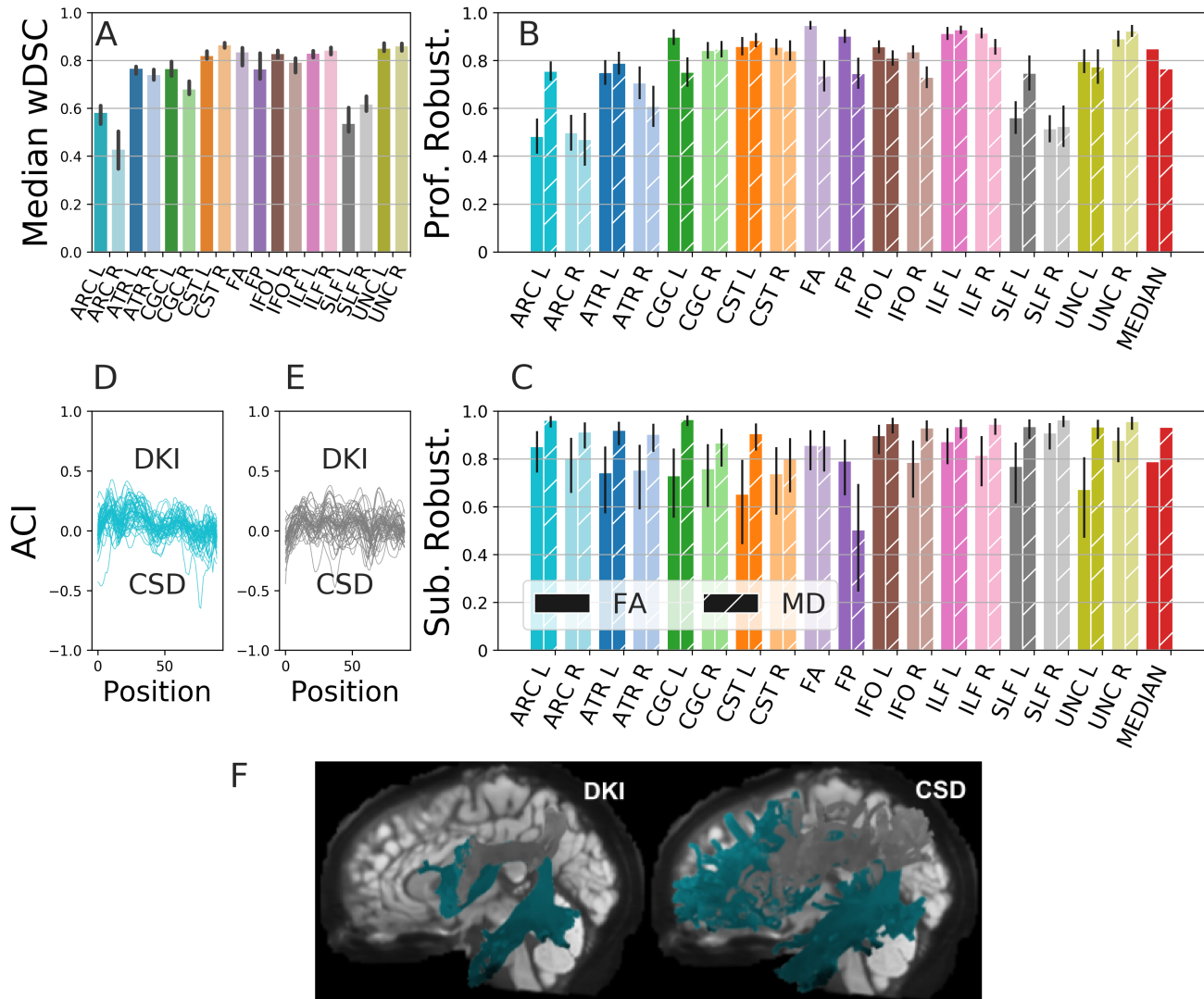


Fig. 4. ODF model robustness. We compared DKI- and CSD-derived tractography. Colors encode bundle information as in Figures 1 and 2. Textured hatching encodes FA/MD information. **A** wDSC robustness. **B** Profile robustness. **C** Subject robustness. **D, E** Adjusted contrast index profile (ACIP) between ARC L and SLF L tract profiles of each algorithm. Individual ACIP are plotted with thin lines, and means are plotted with a thick line. Positive ACI indicates DKI found a higher value of FA than CSD at that node. The 95% confidence interval on the mean is shaded. **F** Tractography and bundle recognition results for ARC L and SLF L respectively for one example subject. Error bars represent 95% confidence interval.

comparison. In comparing different software implementations, we found that mAFQ has higher subject TRR relative to pyAFQ in the UW-PREK dataset, when TRR is relatively low for pyAFQ (see the FA, CST L, and ATR L in Fig. 3C). On the other hand, in the HCP-TR dataset pyAFQ we used the RTP pipeline (36, 37), which is an extension of mAFQ, and found that pyAFQ tends to have slightly higher profile TRR than RTP for MD, but slightly lower profile TRR for FA (Fig. 3D). The pyAFQ and RTP subject TRR are highly comparable (Fig. 3E). In FA, the median pyAFQ subject TRR for FA is 0.77 while the median RTP subject TRR is 0.76. Comparing different ODF models in pyAFQ, we found that the DKI and CSD ODF models have highly similar TRR, both at the level of wDSC (Fig. 3A), as well as at the level of profile and subject TRR (Fig. 3F-G).

Robustness: comparison between distinct tractography models and bundles recognition algorithms. To assess the robustness of tractometry results to different models and algorithms, we used the same measures that were used to calculate TRR.

Tractometry results can be robust to differences in ODF models used in tractography. We compared two algorithms: tractography using DKI- and CSD-derived ODFs. The weighted Dice similarity coefficient (wDSC) for this comparison can be rather high in some cases (e.g., the uncinate and corticospinal tracts, Figure 4A), but produce results that appear very different for some bundles, such as the arcuate and superior longitudinal fasciculi (ARC and SLF) (see also Figure 4D). Despite these discrepancies, profile and subject robustness are high for most bundles (median FA of 0.85 and 0.79, respectively) (Figure 4B,C). In contrast to the re-

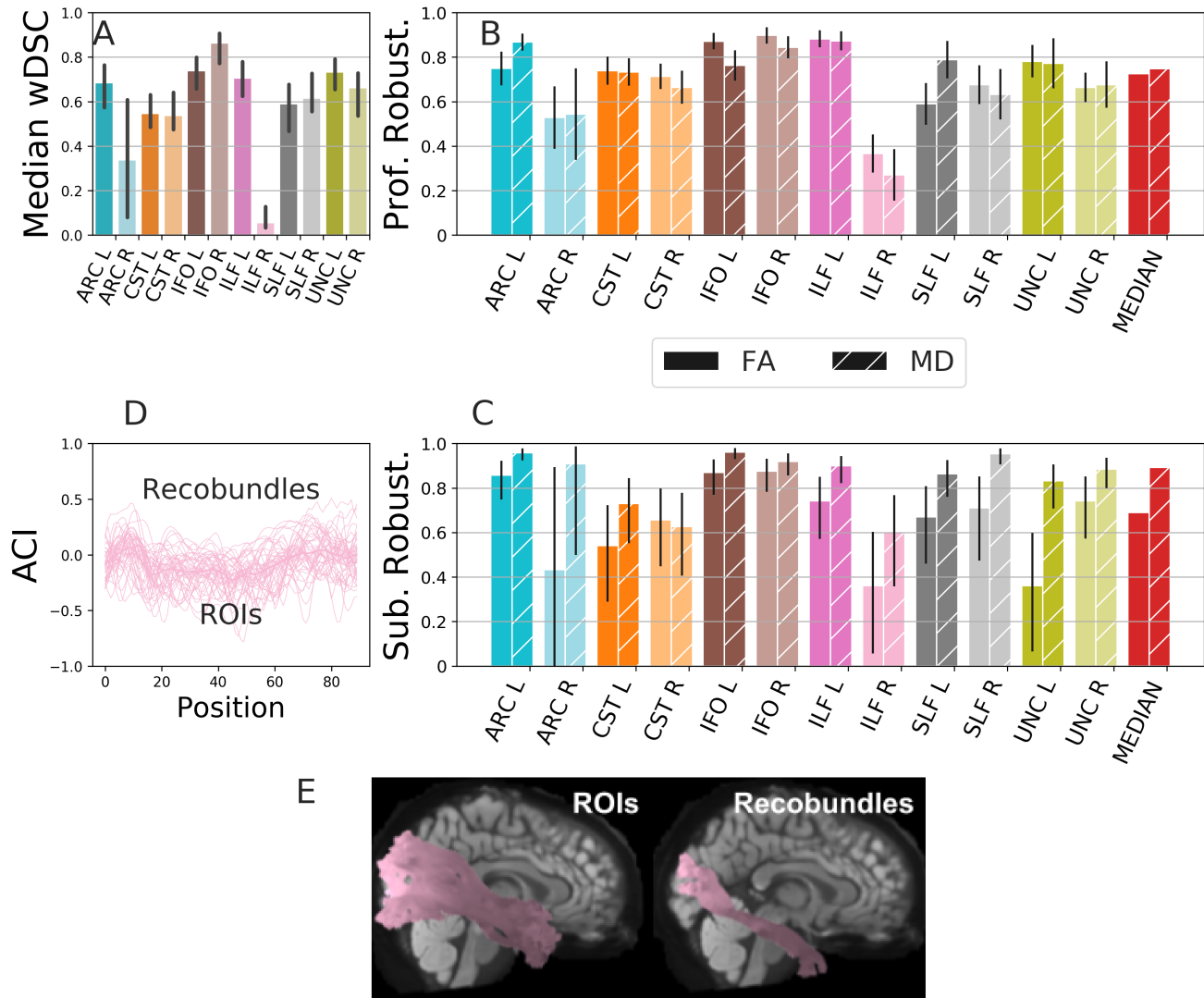


Fig. 5. Segmentation algorithm robustness. **A** wDSC. **B** Profile robustness. **C** Subject robustness. The right arcuate was not found by our Recobundles implementation in 37 out of the 44 subjects, which causes the large error bars for that tract. **D** The ILF R FA ACIP, where positive ACIP indicates Recobundles found a higher value of FA than the waypoint ROIs approach at that node. **E** shows the ILF R found by each algorithm for an example subject with wDSC= 0.11.

sults found in TRR, MD subject robustness is consistently higher than FA subject robustness. The two bundles with the most marked differences between the two ODF models are the SLF and ARC (Figure 4D). These bundles have low wDSC and profile robustness, yet their subject robustness remains remarkably high (In FA, 0.80 ± 0.14 for ARC R and 0.91 ± 0.16 for SLF R) (Figure 4C). These differences are partially explained due to the fact that there are systematic biases in the sampling of white matter by bundles generated with these two ODF models, as demonstrated by the non-zero adjusted contrast index profile (ACIP) between the two models (Figure 4E).

Most white matter bundles are highly robust across bundle recognition methods. We compared bundle recognition with the same tractography results using two different approaches: the default waypoint ROI approach (9), and an alternative approach (Recobundles) that uses atlas templates in the space

of the streamlines (38). Between these algorithms, wDSC is around or above 0.5 for all but two bundles, (1) ARC R and (2) ILF R (Figure 5). There are particular reasons why wDSC is low for these two bundles: (1) RecoBundles often did not recognize any streamlines as belonging to the right arcuate fasciculus and (2) there is an asymmetry in the ILF atlas bundle(7), which results in discrepancies between ILF R recognized with waypoint ROIs and with RecoBundles. Despite these two bundles, we find high robustness overall. For MD, the first quartile subject robustness is 0.83 (Figure 5C, D).

Note the RecoBundles algorithm is somewhat sensitive to choice of parameters, and we always used only one parameter setting (39). It is possible that these results could be improved by using multiple parameter settings and atlases, as in Ocampo-Pineda *et al* (40).

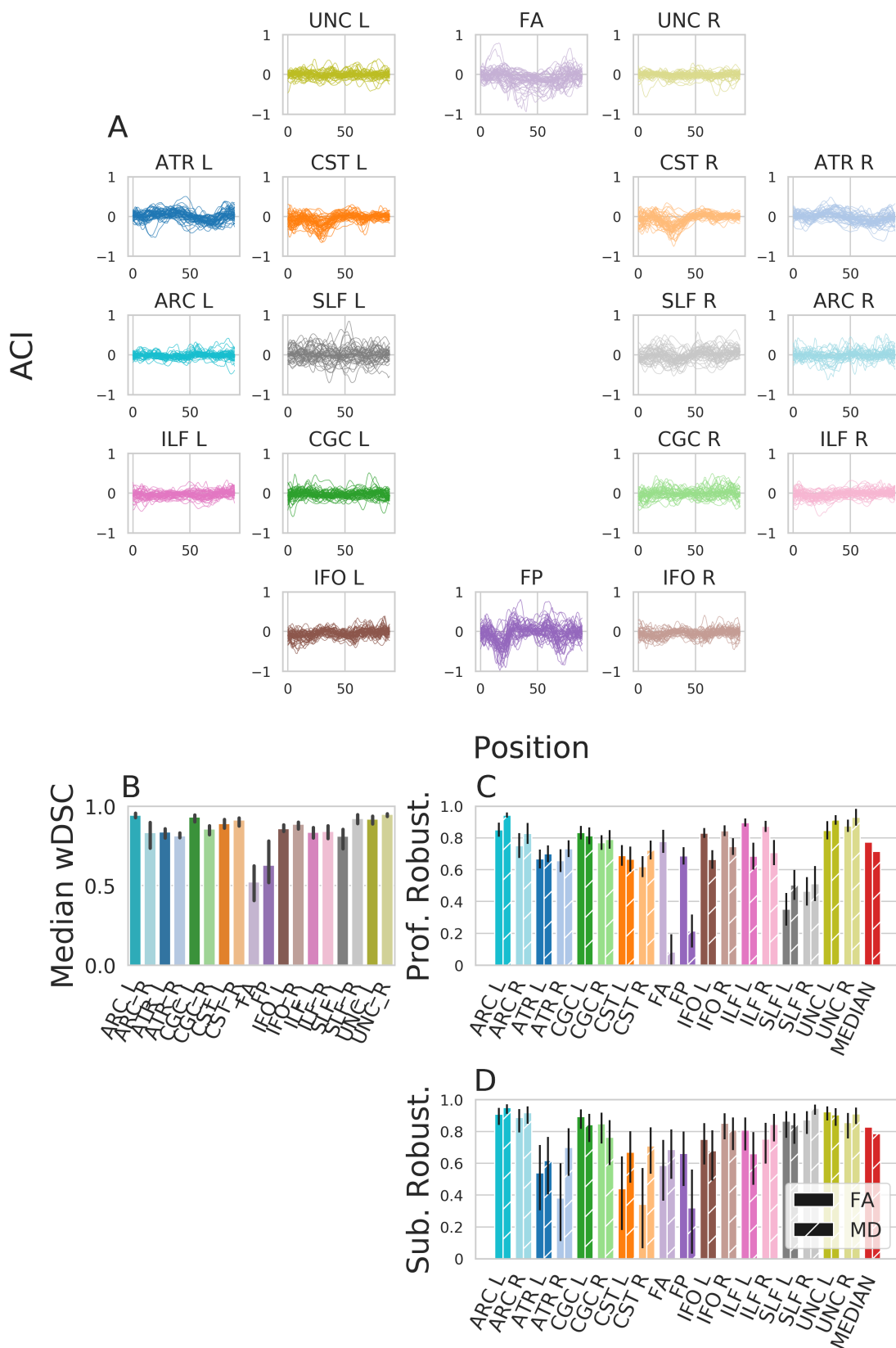


Fig. 6. Robustness between pyAFQ and mAFQ on UW-PREK session #1 data. **A** ACIP between the tract profiles from UW-PREK using pyAFQ and mAFQ (mean in thick line, and individual data in thin lines). Positive ACI indicates pyAFQ found a higher value than mAFQ at that node. The 95% confidence interval on the mean is shaded. Robustness in wDSC (**A**) bundle profiles (**B**) and across subjects (**C**). Error bars show the 95% confidence interval.

Tractometry results are robust to differences in software implementation. Overall, we found that robustness of tractometry across these different software implementations is high in most white matter bundles. In the mAFQ/pyAFQ comparison, most bundles have a wDSC around or above 0.8, except the two callosal bundles (FA and FP), which have a much lower overlap (Fig. 6A). Consistent with this pattern, profile and subject robustness is also overall rather high (Fig. 6B, C). The median values across bundles are 0.72 and 0.79 for MD profile and subject robustness, respectively. For some bundles, like the right and left uncinate, there is large agreement between pyAFQ and mAFQ (for subject MD: UNC L $r = 0.90 \pm 0.04$, UNC R $r = 0.91 \pm 0.06$). However, the callosal bundles also have particularly low mean diffusivity (MD) profile robustness (Fig. 6B) (FP = $r = 0.08 \pm 0.11$, FA $r = 0.21 \pm 0.10$).

The robustness of tractometry to the differences between the pyAFQ and mAFQ implementation depends on the bundle, scalar, and reliability metric. In addition, for many bundles, the ACIP between mAFQ and pyAFQ results is very close to 0, indicating no systematic differences (Fig. 6D). In some bundles – the corticospinal tract (CST) and the anterior thalamic radiations (ATR) – there are small systematic differences between mAFQ and pyAFQ. In the Forceps Posterior (FP), pyAFQ consistently finds smaller FA values than mAFQ in a section on the left side. Notice that the forceps anterior has an ACIP that deviates only slightly from 0, even though the forceps recognitions did not have as much overlap as other bundle recognitions (see Fig. 6A).

Discussion

Previous work has called into question the reliability of neuroimaging analysis (e.g., (24, 41, 42)). We assessed the reliability of a specific approach, tractometry, which is grounded in decades of anatomical knowledge, and we demonstrate that this approach is reproducible, reliable and robust. A tractometry analysis typically combines the outputs of tractography with diffusion reconstruction at the level of the individual voxels within each bundle. One of the major challenges facing researchers who use tractometry is that there are many ways to analyze diffusion data, including different models of diffusion at the level of individual voxels; techniques to connect voxels through tractography; and approaches to classify tractography results into major white matter bundles. Here, we analyzed the reliability of tractometry analysis at several different levels. We analyzed both test-retest reliability of tractometry results and their robustness to changes in analytic details, such as choice of tractography method, bundle recognition algorithm, and software implementation (Fig 6).

Test-retest reliability of tractometry. Test-retest reliability (TRR) of tractometry is usually rather high, comparable in some tracts and measurements to the TRR of the measurement. In comparing the HCP-TR analysis and UW-PREK analysis, we note that higher measurement reliability goes hand in hand with tractometry reliability.

In terms of the anatomical definitions of the bundles, quantified as the TRR wDSC, we find reliable results in both datasets and with both software implementations and both tractography methods that we tested. With pyAFQ we found a relatively low TRR in the frontal callosal bundle (FA) in the UW-PREK dataset. This could be due to the sensitivity of the definition of this bundle to susceptibility distortion artifacts in the frontal poles of the two hemispheres. This low TRR was not found with mAFQ, suggesting that this low TRR is not a necessary feature of the analysis, and is a potential avenue for improvement to pyAFQ. While the two implementations were created by teams with partial overlap and despite the fact that pyAFQ implementation drew both inspiration as well as specific implementation details from mAFQ, many details of implementation still differ substantially. For example, the implementations of tractography algorithms are quite different – pyAFQ relies on DIPY (27) for its tractography, while mAFQ uses implementations provided in Vistasoft (43). The two pipelines also use different registration algorithms, with pyAFQ relying on the SyN algorithm (31), while mAFQ relies on registration methods implemented as part of the Statistical Parametric Mapping (SPM) software (44). These differences may explain the discrepancies observed.

We also find that TRR is high at the level of profiles within subjects and mean tract profiles across subjects. This is generally observed in both datasets that we examined, and using different analysis methods and software implementations. For the UW-PREK dataset, subject TRR tends to be higher in mAFQ than in pyAFQ. On the other hand, for the HCP-TR dataset, pyAFQ subject TRR tends to be higher than that obtained with RTP, which is a fork and extension of mAFQ (36, 37). Generally, TRR of FA profiles and also TRR of mean FA across subjects tend to be higher than those of MD. This could be because the assessment of MD is more sensitive to partial volume effects. In contrast to FA, MD is also not bounded, which means that extreme values at the boundaries of tissue types can have a substantial effect on TRR.

Robustness of tractometry. As highlighted in the recent work by Botvinik-Nezer *et al* (24) and in parallel by Schilling *et al* (41), inferences from even a single dataset can vary significantly, depending on the decisions and analysis pipelines that are used. The analysis approaches used in tractometry embody many assumptions made at the different stages of analysis: the model of the signal in each individual voxel, the manner in which streamlines are generated in tractography, the definition of bundles, and the extraction of tract profiles. While TRR is important, it does not guard against systematic errors in the analysis approach. One way to test model assumptions and software failures is to create ground truth data against which different methods and implementations can be tested (13, 45, 46). However, this approach also relies on certain assumptions about the mechanisms that generate the data that is considered ground truth, making this approach more straightforward for some methods than others. Here, we instead assessed the robustness of tractometry results to perturbations of analytic components, focusing on the mod-

elling of ODFs in individual voxels and the approach taken to bundle recognition.

Subject robustness remains high despite differences in the spatial extent of bundles. We replicated previous findings that the definition of major bundles can vary in terms of their spatial extent (quantified via wDSC) (13, 35, 41, 47), depending on the software implementation or the ODF model used. As we show, low wDSC robustness often corresponds to low profile robustness, and vice versa (Fig 6B,C, Fig 4A,B, and Fig 5A,B). That is, when two algorithms detect bundles with small spatial overlap, the shape of the resulting tract profiles are also different from each other. However, low wDSC and profile robustness does not always translate to low subject robustness. Algorithms can detect bundles with low spatial overlap and of different shapes yet still agree on the ordering of the mean of the profiles, i.e., which subjects have high or low FA in a given bundle. A clear example of this is the SLF and ARC in Fig 4 (wDSC and profile robustness are low, yet subject robustness is very high). This suggests that tractometry can overcome failures in precise delineation of the major bundles by averaging tissue properties within the core of the white matter. Conversely, important details that are sensitive to these choices may be missed when averaging along the length of the tracts.

Our high subject-level robustness results (Fig 6C, Fig 4C, and Fig 5C) dovetail with the results of a recently-published study that used tractometry in a sample of 45 participants (48), and found high subject-level correlations between the mean tract values of FA and MD for two different pipelines: deterministic tractography using DTI as the ODF model (essentially identical to a pipeline used in our supplementary analysis, described in “DTI Configuration”), and probabilistic tractography using CSD as the ODF model. Consistent with our results on the HCP-TR dataset, slightly higher subject robustness was found for MD than for FA.

Exceptions & Limitations. High profile robustness did not always imply high subject robustness (e.g., the FP in Fig 4 has high profile robustness, but low subject robustness). This suggests that there are other sources of between-subject variance that do not correspond directly to profile robustness within an individual.

There are still significant challenges to robustness that arise from the way in which the major bundles are defined. This problem was highlighted in recent work that demonstrated that different researchers use different criteria to define bundles of streamlines that represent the same tract (41). In our case, this challenge is represented by the relatively low robustness between the waypoint ROI algorithm for bundle definition and the RecoBundles algorithm. In this comparison, the wDSC exceeds 0.8 in only one bundle and is below 0.4 in two cases. While both algorithms identify a bundle of streamlines that represents the right ILF, this bundle differs substantially between the two algorithms. Even so, profile and subject robustness can still be rather high, even in some cases in which rather middling overlap is found between the anatomical extent of the bundles. This challenge highlights

the need for more precise definitions of the models of brain tracts that are derived from dMRI, but also highlights the need for clear, automated and reproducible software to perform bundle recognition.

In addition to decisions about analysis approach, which may be theoretically motivated, software implementations may contain systematic errors in executing the different steps and different software may be prone to different kinds of failure modes. Since other software implementations (9, 36) of the AFQ approach have been in widespread use in multiple different datasets and research settings, we also compared the results across different software implementations (Fig. 6). While there are some systematic differences between implementations, tractometry is overall quite robust to differences between software implementations.

Another important limitation of this work is that we have only analyzed samples of healthy individuals. Where brains are severely deformed (e.g., in TBI, brain tumors and so forth), particular care would be needed to check the results of bundle recognition, and separate considerations would be needed in order to reach conclusions about the reliability of the inferences made.

Computational reproducibility via open-source software.

Reproducibility is a bedrock of science, but achieving full computational reproducibility is a high bar that requires access to the software, data and computational environment that a researcher uses (21). One of the goals of pyAFQ is to provide a platform for reproducible tractometry. It is embedded in an ecosystem of tools for reproducible neuroimaging and is extensible. This is shown in Fig. S6 and Fig S2 and is further discussed in “Supplementary Discussion of pyAFQ”. Results from the present article and supplements can be reproduced using a set of Jupyter notebooks provided here: https://github.com/36000/Tractometry_TRR_and_robustness. After installing the version of pyAFQ that we used (0.6), reproduction should be straightforward on standard operating systems and architectures, or in cloud computing systems (see code and Supplementary Methods). But making the results fully reproducible demonstrates some of the challenges of reproducibility (see also (49) for an extensive discussion of these issues). For example, the datasets that we used cannot be fully shared. The UW-PREK dataset cannot be fully shared, though we are able to share an extract of the data that should be sufficient to reproduce results based on tract profiles provided here. In addition, we provide web-based visualizations of the data using a tool that we have previously developed for transparent data sharing of tractometry data (50): https://yeatmanlab.github.io/UW_PREK_pyAFQ_pre_browser and https://yeatmanlab.github.io/UW_PREK_pyAFQ_post_browser. The HCP-TR dataset is relatively straightforward for others to access in its preprocessed form through the HCP, and because the study IDs can be openly shared in our code, anyone with such access should be able to reproduce the figures in full.

Nevertheless, further exploration and extension of these results is limited. For example, if other researchers would be interested in comparing our TRR results to another tractometry pipeline (e.g., TRACULA (11), another popular tractometry pipeline) or another bundle recognition algorithm (e.g., TractSeg (51), which uses a neural network to recognize bundles, or Classifyber (52), which uses a linear classifier), their ability to do so would be limited, in the absence of the original raw and/or processed data. Using these resources, it should be possible to re-execute our workflows and replicate most of our results (53).

ACKNOWLEDGEMENTS

This work was supported through grant 1RF1MH121868-01 from the National Institute of Mental Health/The BRAIN Initiative and through grant 5R01EB027585-02 to Eleftherios Garyfallidis (Indiana University) from the National Institute of Biomedical Imaging and Bioengineering. We are also grateful for support from the Gordon and Betty Moore Foundation and the Alfred P. Sloan Foundation to the University of Washington eScience Institute Data Science Environment, as well as support from the Washington Research Foundation to eScience and to the University of Washington Institute for Neuroengineering. Thanks to Andreas Neef for feedback on the pyAFQ software. Data were provided in part by the Human Connectome Project, WU-Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Washington University.

Bibliography

1. Steven E Petersen and Olaf Sporns. Brain Networks and Cognitive Architectures. *Neuron*, 88(1):207–219, October 2015. Publisher: Elsevier.
2. Danielle S Bassett and Olaf Sporns. Network neuroscience. *Nat. Neurosci.*, 20(3):353–364, February 2017.
3. T E Conturo, N F Lori, T S Cull, E Akbudak, A Z Snyder, J S Shimony, R C McKinstry, H Burton, and M E Raichle. Tracking neuronal fiber pathways in the living human brain. *Proc. Natl. Acad. Sci. U. S. A.*, 96(18):10422–10427, August 1999.
4. Susumu Mori and Peter C M Van Zijl. Fiber tracking: principles and strategies—a technical review. *NMR in Biomedicine: An International Journal Devoted to the Development and Application of Magnetic Resonance In Vivo*, 15(7-8):468–480, 2002. Publisher: Wiley Online Library.
5. Setsu Wakana, Hangyi Jiang, Lidia M Nagae-Poetscher, Peter C M van Zijl, and Susumu Mori. Fiber tract-based atlas of human white matter anatomy. *Radiology*, 230(1):77–87, January 2004.
6. Kenichi Oishi, Karl Zilles, Katrin Amunts, Andreia Faria, Hangyi Jiang, Xin Li, Kazi Akhter, Kegang Hua, Roger Woods, Arthur W Toga, G Bruce Pike, Pedro Rosa-Neto, Alan Evans, Jiayang Zhang, Hao Huang, Michael I Miller, Peter C M van Zijl, John Mazziotta, and Susumu Mori. Human brain white matter atlas: identification and assignment of common anatomical structures in superficial white matter. *Neuroimage*, 43(3):447–457, November 2008.
7. Fang-Cheng Yeh, Sandip Panesar, David Fernandes, Antonio Meola, Masanori Yoshino, Juan C. Fernandez-Miranda, Jean M. Vettel, and Timothy Verstynen. Population-averaged atlas of the macroscale human structural connectome and its network topology. *NeuroImage*, 178:57–68, 2018. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2018.05.027.
8. Eleftherios Garyfallidis, Marc-Alexandre Côté, Francois Rheault, Jasmeen Sidhu, Janice Hau, Laurent Petit, David Fortin, Stephen Cunanne, and Maxime Descoteaux. Recognition of white matter bundles using local and global streamline-based registration and clustering. *Neuroimage*, July 2017. doi: 10.1016/j.neuroimage.2017.07.015.
9. Jason D. Yeatman, Robert F. Dougherty, Nathaniel J. Myall, Brian A. Wandell, and Heidi M. Feldman. Tract Profiles of White Matter Properties: Automating Fiber-Tract Quantification. *PLOS ONE*, 7(11):e49790, November 2012. ISSN 1932-6203. doi: 10.1371/journal.pone.0049790. Publisher: Public Library of Science.
10. Marco Catani and Michel Thiebaut de Schotten. A diffusion tensor imaging tractography atlas for virtual in vivo dissections. *Cortex*, 44(8):1105–1132, September 2008. Publisher: Elsevier.
11. Anastasia Yendiki, Patricia Panneck, Priti Srinivasan, Allison Stevens, Lilla Zöllei, Jean Augustinack, Ruopeng Wang, David Salat, Stefan Ehrlich, Tim Behrens, Saad Jbabdi, Randy Gollub, and Bruce Fischl. Automated probabilistic reconstruction of white-matter pathways in health and disease using an atlas of the underlying anatomy. *Front. Neuroinform.*, 5:23, October 2011.
12. Demian Wassermann, Nikos Makris, Yogesh Rathi, Martha Shenton, Ron Kikinis, Marek Kubicki, and Carl-Fredrik Westin. The white matter query language: a novel approach for describing human white matter anatomy. *Brain Struct. Funct.*, 221(9):4705–4721, December 2016.
13. Klaus H. Maier-Hein, Peter F. Neher, Jean-Christophe Houde, Marc-Alexandre Côté, Eleftherios Garyfallidis, Jidan Zhong, Maxime Chamberland, Fang-Cheng Yeh, Ying-Chia Lin, Qing Ji, Wilburn E. Reddick, John O. Glass, David Qixiang Chen, Yuanjing Feng, Chengfeng Gao, Ye Wu, Jieyan Ma, Renjie He, Qiang Li, Carl-Fredrik Westin, Samuel Deslauriers-Gauthier, J. Omar Ocegueda González, Michael Paquette, Samuel St-Jean, Gabriel Girard, François Rheault, Jasmeen Sidhu, Chantal M. W. Tax, Fenghua Guo, Hamed Y. Mesri, Szabolcs Dávid, Martijn Froeling, Anneriet M. Heemskerk, Alexander Leemans, Arnaud Boré, Basile Pinsard, Christophe Bedetti, Matthieu Desrosiers, Simona Brambati, Julien Doyon, Alessia Sarica, Roberta Vasta, Antonio Cerasa, Aldo Quattrone, Jason Yeatman, Ali R. Khan, Wes Hodges, Simon Alexander, David Romascano, Muhamed Barakovic, Anna Aurià, Oscar Esteban, Alia Lemkaddem, Jean-Philippe Thiran, H. Ertan Cetingul, Benjamin L. Odry, Boris Mailhe, Mariappan S. Nadar, Fabrizio Pizzagalli, Gautam Prasad, Julio E. Villalón-Reina, Justin Galvis, Paul M. Thompson, Francisco De Santiago Requejo, Pedro Luque Laguna, Luis Miguel Lacerda, Rachel Barrett, Flavio Dell’Acqua, Marco Catani, Laurent Petit, Emmanuel Caruyer, Alessandro Daducci, Tim B. Dyrby, Tim Holland-Letz, Claus C. Hilgetag, Bram Stieltjes, and Maxime Descoteaux. The challenge of mapping the human connectome based on diffusion tractography. *Nature Communications*, 8(1):1349, November 2017. ISSN 2041-1723. doi: 10.1038/s41467-017-01285-x. Number: 1 Publisher: Nature Publishing Group.
14. Cibu Thomas, Frank Q Ye, M Okan Irfanoglu, Pooja Modi, Kadharbatcha S Saleem, David A Leopold, and Carlo Pierpaoli. Anatomical accuracy of brain connections derived from diffusion MRI tractography is inherently limited. *Proc. Natl. Acad. Sci. U. S. A.*, 111(46):16574–16579, November 2014.
15. Kurt G Schilling, Laurent Petit, Francois Rheault, Samuel Remedios, Carlo Pierpaoli, Adam W Anderson, Bennett A Landman, and Maxime Descoteaux. Brain connections derived from diffusion mri tractography can be highly anatomically accurate—if we know where white matter pathways start, where they end, and where they do not go. *Brain Structure and Function*, 225(8):2387–2402, 2020.
16. Ariel Rokem, Jason D Yeatman, Franco Pestilli, Kendrick N Kay, Aviv Mezer, Stefan van der Walt, and Brian A Wandell. Evaluating the accuracy of diffusion MRI models in white matter. November 2014. _eprint: 1411.0721.
17. Dmitry S Novikov, Valerij G Kiselev, and Sune N Jespersen. On modeling. *Magn. Reson. Med.*, 79(6):3172–3193, June 2018.
18. John B Colby, Lindsay Soderberg, Catherine Lebel, Ivo D Dinov, Paul M Thompson, and Elizabeth R Sowell. Along-tract statistics allow for enhanced tractography analysis. *Neuroimage*, 59(4):3227–3242, February 2012.
19. Adam Richie-Halford, Jason Yeatman, Noah Simon, and Ariel Rokem. Multidimensional analysis and detection of informative features in diffusion MRI measurements of human white matter. December 2019.
20. Michael Dayan, Elizabeth Monohan, Sneha Pandya, Amy Kuceyeski, Thanh D Nguyen, Ashish Raj, and Susan A Gauthier. Profilmetry: A new statistical framework for the characterization of white matter pathways, with application to multiple sclerosis. *Hum. Brain Mapp.*, December 2015.
21. David L Donoho. An invitation to reproducible computational research. *Biostatistics*, 11(3):385–388, July 2010.
22. Peter Ivie and Douglas Thain. Reproducibility in Scientific Computing. *ACM Comput. Surv.*, 51(3):1–36, July 2018. Place: New York, NY, USA Publisher: Association for Computing Machinery.
23. The Turing Way Community, Becky Arnold, Louise Bowler, Sarah Gibson, Patricia Hererich, Rosie Higman, Anna Krystalli, Alexander Morley, Martin O’Reilly, and Kirstie Whitaker. *The Turing Way: A Handbook for Reproducible Data Science*. March 2019.
24. Rotem Botvinik-Nezer, Felix Holzmeister, Colin F. Camerer, Anna Dreber, Juergen Huber, Magnus Johannesson, Michael Kirchlner, Roni Iwanir, Jeanette A. Mumford, R. Alison Adcock, Paolo Avesani, Blazej M. Baczkowski, Aahana Bajracharya, Leah Bakst, Sheryl Bail, Marco Barilari, Nadège Bault, Derek Beaton, Julia Beilner, Roland G. Benoit, Ruud M. W. J. Berkers, Jamil P. Bhanji, Bharat B. Biswal, Sebastian Bobadilla-Suarez, Tiago Bortoloni, Katherine L. Bottenhorn, Alexander Bowring, Senne Braem, Hayley R. Brooks, Emily G. Brudner, Cristian B. Calderon, Julia A. Camilleri, Jaime J. Castrellon, Luca Cecchetti, Edna C. Cieslik, Zachary J. Cole, Olivier Collignon, Robert W. Cox, William A. Cunningham, Stefan Czochke, Kamalaker Dadi, Charles P. Davis, Alberto De Luca, Mauricio R. Delgado, Lysia Demetriou, Jeffrey B. Dennison, Xin Di, Erin W. Dickie, Ekaterina Dobryakova, Claire L. Donnat, Juergen Dukart, Niall W. Duncan, Joke Durnez, Amr Eed, Simon B. Eickhoff, Andrew Erhart, Laura Fontanesi, G. Matthew Fricke, Shiguang Fu, Adriana Galván, Remi Gau, Sarah Genon, Tristan Glatard, Enrico Glerean, Jelle J. Goeman, Sergej A. E. Golowin, Carlos González-García, Krzysztof J. Gorgolewski, Cheryl L. Grady, Mikella A. Green, João F. Guassi Moreira, Olivia Guest, Shabnam Hakimi, J. Paul Hamilton, Roeland Hancock, Giacomo Handjaras, Bronson B. Harry, Colin Hawco, Peer Herholz, Gabrielle Herman, Stephan Heunis, Felix Hoffstaedter, Jeremy Hogeveen, Susan Holmes, Chuan-Peng Hu, Scott A. Huettel, Matthew E. Hughes, Vittorio Iacovella, Alexandru D. Iordan, Peder M. Isager, Ayse I. Isik, Andrew Jahn, Matthew R. Johnson, Tom Johnstone, Michael J. E. Joseph, Anthony C. Juliano, Joseph W. Kable, Michalis Kassinosopoulos, Cemal Koba, Xiang-Zhen Kong, Timothy R. Kosciak, Nuri Erkut Kucukboyaci, Brice A. Kuhl, Sebastian Kupek, Angela R. Laird, Claus Lamm, Robert Langner, Nina Lauharatanahirun, Hongmi Lee, Sangil Lee, Alexander Leemans, Andrea Leo, Elise Lesage, Flora Li, Monica Y. C. Li, Phui Cheng Lim, Evan N. Lintz, Schuyler W. Liphardt, Annabel B. Losecaat Vermeer, Bradley C. Love, Michael L. Mack, Norberto Malpica, Theo Marins, Camille Maumet, Kelsey McDonald, Joseph T. McGuire, Helena Melero, Adriana S. Méndez Leal, Benjamin Meyer, Kristin N. Meyer, Glad Mihai, Georgios D. Mitsis, Jorge Moll, Dylan M. Nielson, Gustav Nilsson, Michael P. Notter, Emanuele Olivetti, Adrian I. Onicasc, Paolo Papale, Kaustubh R. Patil, Jonathan E. Peelle, Alexandre Pérez, Doris Pischke, Jean-Baptiste Poline, Yanina Prystauka, Shruti Ray, Patricia A. Reuter-Lorenz, Richard C. Reynolds, Emiliano Ricciardi, Jenny R. Rieck, Anais M. Rodriguez-Thompson, Anthony Romy, Taylor Salo, Gregory R. Samanez-Larkin, Emilio Sanz-Morales, Margaret L. Schlichting, Douglas H. Schultz, Qiang Shen, Margaret A. Sheridan, Jennifer A. Silvers, Kenny Skagerlund, Alec Smith, David V. Smith, Peter Sokol-Hessner, Simon R. Steinkamp, Sarah M. Tashjian, Bertrand Thirion, John N. Thorp, Gustav Tinghög, Loreen Tisdall, Steven H. Tompson, Claudio Toro-Serey, Juan Jesus Torre Tresols, Leonardo Tozzi, Vuong Truong, Luca Turella, Anna E. van ’t Veer, Tom Verguts, Jean M. Vettel, Sagana Vijayarajah, Khoi Vo, Matthew B. Wall, Wouter D. Weeda, Susanne Weis, David J. White, David Wisniewski, Alba Xifra-Porxas, Emily A. Yearling, Sanguk Yoon, Rui Yuan, Kenneth S. L. Yuen, Lei Zhang, Xu Zhang, Joshua E. Zosky, Thomas E. Nichols, Russell A. Poldrack, and Tom Schonberg. Variability in the analysis of a single neuroimaging dataset by many teams. *Nature*, 582(7810):84–88, June 2020. ISSN

- 1476-4687. doi: 10.1038/s41586-020-2314-9. Number: 7810 Publisher: Nature Publishing Group.
25. Matthew Cieslak, Philip A. Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez Adebimpe, Geoffrey K. Aguirre, Danielle S. Bassett, Richard F. Betzel, Josiane Bourque, Laura M. Cabral, Christos Davatzikos, John Detre, Eric Earl, Mark A. Elliott, Shreyas Fadnavis, Damien A. Fair, Will Foran, Panagiotis Fotiadis, Eleftherios Garyfallidis, Barry Giesbrecht, Ruben C. Gur, Raquel E. Gur, Max Kelz, Anisha Keshavan, Bart S. Larsen, Beatriz Luna, Allyson P. Mackey, Michael Milham, Desmond J. Oathes, Anders Perrone, Adam R. Pines, David R. Roalf, Adam Richie-Halford, Ariel Rokem, Valerie J. Synhor, Tinashe M. Tapera, Ursula A. Tooley, Jean M. Vettel, Jason D. Yeatman, Scott T. Grafton, and Theodore D. Satterthwaite. QSIprep: An integrative platform for preprocessing and reconstructing diffusion MRI. *bioRxiv*, page 2020.09.04.282269, September 2020. doi: 10.1101/2020.09.04.282269. Publisher: Cold Spring Harbor Laboratory Section: New Results.
 26. Hadley Wickham. Tidy data. *J. Stat. Softw.*, 59(10), 2014.
 27. Eleftherios Garyfallidis, Matthew Brett, Bagrat Amirbekian, Ariel Rokem, Stefan Van Der Walt, Maxime Descoteaux, and Ian Nimmo-Smith. Dipy, a library for the analysis of diffusion MRI data. *Frontiers in Neuroinformatics*, 8, 2014. ISSN 1662-5196. doi: 10.3389/fninf.2014.00008. Publisher: Frontiers.
 28. Vladimir Fonov, Alan C. Evans, Kelly Botteron, C. Robert Almli, Robert C. McKinstry, D. Louis Collins, and Brain Development Cooperative Group. Unbiased average age-appropriate atlases for pediatric studies. *NeuroImage*, 54(1):313–327, January 2011. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2010.07.033.
 29. VS Fonov, AC Evans, RC McKinstry, CR Almlí, and DL Collins. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. *NeuroImage*, 47:S102, July 2009. ISSN 1053-8119. doi: 10.1016/S1053-8119(09)70884-5.
 30. Flavio Dell'Acqua, Luis Lacerda, Marco Catani, and Andrew Simmons. Anisotropic Power Maps: A diffusion contrast to reveal low anisotropy tissues from HARDI data. page 1.
 31. B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee. Symmetric Diffeomorphic Image Registration with Cross-Correlation: Evaluating Automated Labeling of Elderly and Neurodegenerative Brain. *Medical image analysis*, 12(1):26–41, February 2008. ISSN 1361-8415. doi: 10.1016/j.media.2007.06.004.
 32. Marco Catani, Robert J. Howard, Sinisa Pajevic, and Derek K. Jones. Virtual in vivo interactive dissection of white matter fasciculi in the human brain. *NeuroImage*, 17(1):77–94, September 2002. ISSN 1053-8119. doi: 10.1006/nimg.2002.1136.
 33. Kegang Hua, Jiangyang Zhang, Setsu Wakana, Hangyi Jiang, Xin Li, Daniel S. Reich, Peter A. Calabresi, James J. Pekar, Peter C. M. van Zijl, and Susumu Mori. Tract probability maps in stereotaxic spaces: analyses of white matter anatomy and tract-specific quantification. *NeuroImage*, 39(1):336–347, January 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.07.053.
 34. Stamatiou N Sotiropoulos, Saad Jbabdi, Junqian Xu, Jesper L Andersson, Steen Moeller, Edward J Auerbach, Matthew F Glasser, Moises Hernandez, Guillermo Sapiro, Mark Jenkinson, David A Feinberg, Essa Yacoub, Christophe Lenglet, David C Van Essen, Kamil Ugurbil, Timothy E J Behrens, and WU-Minn HCP Consortium. Advances in diffusion MRI acquisition and processing in the human connectome project. *NeuroImage*, 80:125–143, October 2013. doi: 10.1016/j.neuroimage.2013.05.057.
 35. Martin Cousineau, Pierre-Marc Jodoin, Eleftherios Garyfallidis, Marc-Alexandre Côté, Félix C. Morency, Verena Rozanski, Marilyn Grand'Maison, Barry J. Bedell, and Maxime Descoteaux. A test-retest study on Parkinson's PPMI dataset yields statistically significant white matter fascicles. *NeuroImage : Clinical*, 16:222, 2017. doi: 10.1016/j.nicl.2017.07.020. Publisher: Elsevier.
 36. Garikoitz Lerma-Usabiaga, Michael L Perry, and Brian A Wandell. Reproducible tract profiles (rtp): from diffusion mri acquisition to publication. *bioRxiv*, page 680173, 2019.
 37. Garikoitz Lerma-Usabiaga, Pratik Mukherjee, Michael L. Perry, and Brian A. Wandell. Data-science ready, multisite, human diffusion MRI white-matter-tract statistics. *Scientific Data*, 7, 2020. doi: 10.1038/s41597-020-00760-3. Publisher: Nature Publishing Group.
 38. Eleftherios Garyfallidis, Marc-Alexandre Côté, Francois Rheault, Jasmeen Sidhu, Janice Hau, Laurent Petit, David Fortin, Stephen Cunanne, and Maxime Descoteaux. Recognition of white matter bundles using local and global streamline-based registration and clustering. *NeuroImage*, 170:283–295, 2018. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2017.07.015.
 39. Francois Rheault, Philippe Poulin, Alex Valcourt Caron, Etienne St-Onge, and Maxime Descoteaux. Common misconceptions, hidden biases and modern challenges of dMRI tractography. *J. Neural Eng.*, January 2020.
 40. Mario Ocampo-Pineda, Simona Schiavi, Francois Rheault, Gabriel Girard, Laurent Petit, Maxime Descoteaux, and Alessandro Daducci. Hierarchical microstructure informed tractography. *Brain Connect.*, January 2021.
 41. Kurt G Schilling, Francois Rheault, Laurent Petit, Colin B Hansen, Vishwesh Nath, Fang-Cheng Yeh, Gabriel Girard, Muhamed Barakovic, Jonathan Rafael-Patino, Thomas Yu, Elda Fisch-Gomez, Marco Pizzolato, Mario Ocampo-Pineda, Simona Schiavi, Erick J Canales-Rodriguez, Alessandro Daducci, Cristina Granziera, Giorgio Innocenti, Jean-Philippe Thiran, Laura Mancini, Stephen Wastling, Sirio Cocozza, Maria Petracca, Giuseppe Pontillo, Matteo Mancini, Sjoerd B Vos, Vejay N Vakharía, John S Duncan, Helena Melero, Lidia Manzanedo, Emilio Sanz-Morales, Ángel Peña-Melián, Fernando Calamante, Arnaud Attyé, Ryan P Cabeen, Laura Korobova, Arthur W Toga, Anupa Ambili Vijayakumari, Drew Parker, Ragini Verma, Ahmed Radwan, Stefan Sunaert, Louise Emsell, Alberto De Luca, Alexander Leemans, Claude J Bajada, Hamied Haroon, Hojjatollah Azadbakht, Maxime Chamberland, Sila Genc, Chantal M W Tax, Ping-Hong Yeh, Rujirutana Srikanchana, Colin Mcknight, Joseph Yuan-Mou Yang, Jian Chen, Claire E Kelly, Chun-Hung Yeh, Jerome Cocheureau, Jerome J Maller, Thomas Welton, Fabien Almirac, Kiran K Seunarine, Chris A Clark, Ramón Zhang, Nikos Makris, Alexandra Golby, Yogesh Rath, Lauren J O'Donnell, Yihao Xia, Dogu Baran Aydogan, Yonggang Shi, Francisco Guerreiro Fernandes, Mathijs Raemaekers, Shaun Warrington, Stijn Michielse, Alonso Ramirez-Manzanera, Luis Concha, Ramón Aranda, Mariano Rivera Meraz, Garikoitz Lerma-Usabiaga, Lucas Roitman, Lucius S Fekonja, Navona Calarco, Michael Joseph, Hajer Nakua, Aristotle N Voineskos, Philippe Karan, Gabrielle Grenier, Jon Haitz Legarreta, Nagesh Adluru, Veena A Nair, Vivek Prabhakaran, Andrew L Alexander, Koji Kamagata, Yuya Saito, Wataru Uchida, Christina Andica, Abe Masahiro, Roza G Bayrak, Claudia A Gandini, Egidio D'Angelo, Fulvia Palesi, Giovanni Savini, Nicolò Rolandi, Pamela Claudio Guevara, Josselin Hoenou, Narciso López-López, Jean-François Mangin, Cyril Poupon, Gaudiero Román, Andrea Vázquez, Chiara Maffei, Mavilde Arantes, José Paulo Andrade, Susana Maria Silva, Rajikha Raja, Vince D Calhoun, Eduardo Caverzasi, Simone Sacco, Michael Lauricella, Franco Pestilli, Daniel Bullock, Yang Zhan, Edith Brignoni-Perez, Catherine Label, Jess E Reynolds, Igor Nestrasil, René Labounek, Christophe Lenglet, Amy Paulson, Stefania Aulicka, Sarah Heilbronner, Katja Heuer, Adam W Anderson, Bennett A Landman, and Maxime Descoteaux. Tractography dissection variability: what happens when 42 groups dissect 14 white matter bundles on the same dataset? October 2020. doi: 10.1101/2020.10.07.321083.
 42. Gregory Kiar, Yohan Chatelain, Pablo de Oliveira Castro, Eric Petit, Ariel Rokem, Gaël Varoquaux, Bratislav Misic, Alan C. Evans, and Tristan Glatard. Numerical Instabilities in Analytical Pipelines Lead to Large and Meaningful Variability in Brain Networks. *bioRxiv*, page 2020.10.15.341495, October 2020. doi: 10.1101/2020.10.15.341495. Publisher: Cold Spring Harbor Laboratory Section: New Results.
 43. Robert F Dougherty, Michal Ben-Shachar, Roland Bammer, Alyssa A Brewer, and Brian A Wandell. Functional organization of human occipital-callosal fiber tracts. *Proc. Natl. Acad. Sci. U. S. A.*, 102(20):7350–7355, May 2005.
 44. Karl J. Friston. Statistical Parametric Mapping. In Rolf Kötter, editor, *Neuroscience Databases: A Practical Guide*, pages 237–250. Springer US, Boston, MA, 2003. ISBN 978-1-4615-1079-6. doi: 10.1007/978-1-4615-1079-6_16.
 45. Garikoitz Lerma-Usabiaga, Noah Benson, Jonathan Winawer, and Brian A Wandell. A validation framework for neuroimaging software: The case of population receptive fields. *PLoS Comput. Biol.*, 16(6):e1007924, June 2020.
 46. Peter F Neher, Frederik B Laun, Bram Stieltjes, and Klaus H Maier-Hein. Fiberfox: facilitating the creation of realistic white matter software phantoms. *Magn. Reson. Med.*, 72(5): 1460–1470, November 2014.
 47. Mariem Boukadi, Karine Marcotte, Christophe Bedetti, Jean-Christophe Houde, Alex Desautels, Samuel Deslauriers-Gauthier, Marianne Chapleau, Arnaud Boré, Maxime Descoteaux, and Simona M. Brambati. Test-Retest Reliability of Diffusion Measures Extracted Along White Matter Language Fiber Bundles Using HARDI-Based Tractography. *Frontiers in Neurosciences*, 12, January 2019. ISSN 1662-4548. doi: 10.3389/fnins.2018.01055.
 48. Maya Yablonski, Benjamin Menashe, and Michal Ben-Shachar. A general role for ventral white matter pathways in morphological processing: Going beyond reading. *NeuroImage*, 226:117577, November 2020.
 49. Justin Kitzes, Daniel Turek, and Fatma Deniz. *The practice of reproducible research: case studies and lessons from the data-intensive sciences*. Univ of California Press, 2017.
 50. Jason D Yeatman, Adam Richie-Halford, Josh K Smith, Anisha Keshavan, and Ariel Rokem. A browser-based tool for visualization and analysis of diffusion MRI data. *Nat. Commun.*, 9(1):940, March 2018.
 51. Jakob Wasserthal, Peter Neher, and Klaus H Maier-Hein. Tractseq-fast and accurate white matter tract segmentation. *NeuroImage*, 183:239–253, 2018.
 52. Giulia Bertò, Daniel Bullock, Pietro Astolfi, Siuchi Hayashi, Luca Zigiotta, Luciano Annicchiarico, Francesco Corsini, Alessandro De Benedictis, Silvio Sarubbo, Franco Pestilli, et al. Classifyber, a robust streamline-based linear classifier for white matter bundle segmentation. *bioRxiv*, 2020.
 53. Satrajit S. Ghosh, Jean-Baptiste Poline, David B. Keator, Yaroslav O. Halchenko, Adam G. Thomas, Daniel A. Kessler, and David N. Kennedy. A very simple, re-executable neuroimaging publication. *F1000Research*, 6:124, June 2017. ISSN 2046-1402. doi: 10.12688/f1000research.10783.2.
 54. Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J van der Walt, Matthew Brett, Joshua Wilson, K Jarrod Millman, Nikolay Mayorov, Andrew R J Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E A Quinterro, Charles R Harris, Anne M Archibald, António H Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: fundamental algorithms for scientific computing in python. *Nat. Methods*, 17(3):261–272, March 2020.
 55. Sphinx, November 2020.
 56. Sphinx-Gallery, November 2020.
 57. Ariel Rokem, Jason D. Yeatman, Franco Pestilli, Kendrick N. Kay, Aviv Mezer, Stefan van der Walt, and Brian A. Wandell. Evaluating the Accuracy of Diffusion MRI Models in White Matter. *PLOS ONE*, 10(4):e0123272, April 2015. ISSN 1932-6203. doi: 10.1371/journal.pone.0123272. Publisher: Public Library of Science.
 58. Brian Hansen and Sune Norhøj Jespersen. Data for evaluation of fast kurtosis strategies, b-value optimization and exploration of diffusion MRI contrast. *Scientific Data*, 3(1):160072, August 2016. ISSN 2052-4463. doi: 10.1038/sdata.2016.72. Number: 1 Publisher: Nature Publishing Group.
 59. Matthew Rocklin. Dask: Parallel Computation with Blocked algorithms and Task Scheduling. pages 126–132, Austin, Texas, 2015. doi: 10.25080/Majora-7b98c3cd-013.
 60. Adam Richie-Halford and Ariel Rokem. Cloudknot: A Python Library to Run your Existing Code on AWS Batch. *Proceedings of the 17th Python in Science Conference*, pages 8–14, 2018. doi: 10.25080/Majora-4af1f417-001. Conference Name: Proceedings of the 17th Python in Science Conference.
 61. Tom Preston-Werner. toml, January 2021. original-date: 2013-02-24T03:03:57Z.
 62. Tristan Glatard, Gregory Kiar, Tristan Aumentado-Armstrong, Natacha Beck, Pierre Bellec, Rémi Bernard, Axel Bonnet, Shawn T. Brown, Sorina Camarasu-Pop, Frédéric Cervenansky, Samir Das, Rafael Ferreira da Silva, Guillaume Flandin, Pascal Girard, Krzysztof J. Gorgolewski, Charles R. G. Guttman, Valérie Hayot-Sasson, Pierre-Olivier Quirin, Pierre Rioux, Marc-Étienne Rousseau, and Alan C. Evans. Boutiques: a flexible framework to integrate command-line applications in computing platforms. *GigaScience*, 7(5), May 2018. doi: 10.1093/gigascience/giy016. Publisher: Oxford Academic.
 63. Tal Yarkoni, Christopher J. Markiewicz, Alejandro de la Vega, Krzysztof J. Gorgolewski, Taylor Salo, Yaroslav O. Halchenko, Quinten McNamara, Krista DeStasio, Jean-Baptiste Poline, Hans Johnson, Oscar Esteban, Dmitry Petrov, James D. Kent, Stefan Appelhoff, Valérie Hayot-Sasson, Dylan M. Nielson, Johan Carlin, Gregory Kiar, Kirstie Whitaker,

- Satrjit Ghosh, Adina Wagner, Elizabeth DuPre, Andrew Janke, Alexander Ivanov, Ashley Gillman, Johannes Wennberg, Lee S. Tirrell, Steven Tilley II, Adam Li, Jon Haitz Legarreta, Mainak Jas, Michael Hanke, Russell Poldrack, Chadwick Boulay, Chris Holdgraf, Evgenii Kalenkovich, Isla Staden, Remi Gau, Ariel Rokem, Bertrand Thirion, Dave F. Kleinschmidt, Erin W. Dickie, John A. Lee, Mathias Goncalves, Matteo Visconti di Oleggio Castello, Michael Philipp Notter, Pauline Roca, and Ross Blair. PyBIDS: Python tools for BIDS datasets, July 2020.
64. Tal Yarkoni, Christopher J. Markiewicz, Alejandro de la Vega, Krzysztof J. Gorgolewski, Taylor Salo, Yaroslav O. Halchenko, Quinten McNamara, Krista DeStasio, Jean-Baptiste Poline, Dmitry Petrov, Valérie Hayot-Sasson, Dylan M. Nielson, Johan Carlin, Gregory Kiar, Kirstie Whitaker, Elizabeth DuPre, Adina Wagner, Lee S. Tirrell, Mainak Jas, Michael Hanke, Russell A. Poldrack, Oscar Esteban, Stefan Appelhoff, Chris Holdgraf, Isla Staden, Bertrand Thirion, Dave F. Kleinschmidt, John A. Lee, Matteo Visconti di Oleggio Castello, Michael P. Notter, and Ross Blair. PyBIDS: Python tools for BIDS datasets. *Journal of Open Source Software*, 4(40):1294, August 2019. ISSN 2475-9066. doi: 10.21105/joss.01294.
65. Krzysztof J. Gorgolewski, Tibor Auer, Vince D. Calhoun, R. Cameron Craddock, Samir Das, Eugene P. Duff, Guillaume Flandin, Satrajit S. Ghosh, Tristan Glatard, Yaroslav O. Halchenko, Daniel A. Handwerker, Michael Hanke, David Keator, Xiangrui Li, Zachary Michael, Camille Maumet, B. Nolan Nichols, Thomas E. Nichols, John Pellman, Jean-Baptiste Poline, Ariel Rokem, Gunnar Schaefer, Vanessa Sochat, William Triplett, Jessica A. Turner, Gaël Varoquaux, and Russell A. Poldrack. The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Scientific Data*, 3(1):160044, June 2016. ISSN 2052-4463. doi: 10.1038/sdata.2016.44. Number: 1 Publisher: Nature Publishing Group.
66. Matthew Brett, Christopher J. Markiewicz, Michael Hanke, Marc-Alexandre Côté, Ben Cipollini, Paul McCarthy, Dorota Jarecka, Christopher P. Cheng, Yaroslav O. Halchenko, Michiel Cottaar, Eric Larson, Satrajit Ghosh, Demian Wassermann, Stephan Gerhard, Gregory R. Lee, Hao-Ting Wang, Erik Kastman, Jakub Kaczmarzyk, Roberto Guidotti, Or Duek, Jonathan Daniel, Ariel Rokem, Cindee Madison, Brendan Moloney, Félix C. Morency, Mathias Goncalves, Ross Markello, Cameron Riddell, Christopher Burns, Jarrod Millman, Alexandre Gramfort, Jaakko Leppäkangas, Anibal Solón, Jasper J.F. van den Bosch, Robert D. Vincent, Henry Krav, Krish Subramanian, Krzysztof J. Gorgolewski, Pradeep Reddy Raamana, Julian Blug, B. Nolan Nichols, Eric M. Baker, Soichi Hayashi, Basile Pinsard, Christian Haselgrove, Mark Hymers, Oscar Esteban, Serge Koudoro, Fernando Pérez-García, Nikolaos N. Oosterhof, Bago Amirbekian, Ian Nimmo-Smith, Ly Nguyen, Samir Reddigari, Samuel St-Jean, Egor Panfilov, Eleftherios Garyfallidis, Gael Varoquaux, Jon Haitz Legarreta, Kevin S. Hahn, Oliver P. Hinds, Bennet Fauber, Jean-Baptiste Poline, Jon Stutters, Kesshi Jordan, Matthew Cieslak, Miguel Estevan Moreno, Valentin Haenel, Yannick Schwartz, Zvi Baratz, Benjamin C Darwin, Bertrand Thirion, Carl Gauthier, Dimitri Papadopoulos Orfanos, Igor Solovey, Ivan Gonzalez, Jath Palasubramaniam, Justin Lecher, Katrin Leinweber, Konstantinos Raktivan, Markéta Calábková, Peter Fischer, Philippe Gervais, Syam Gadde, Thomas Ballinger, Thomas Roos, Venkateswara Reddy Reddam, and freec84. nipy/nibabel: 3.2.0, October 2020.
67. Maxime Descoteaux, Rachid Deriche, Thomas R. Knösche, and Alfred Anwander. Deterministic and probabilistic tractography based on complex fibre orientation distributions. *IEEE transactions on medical imaging*, 28(2):269–286, February 2009. ISSN 1558-254X. doi: 10.1109/TMI.2008.2004424.
68. P. J. Basser, J. Mattiello, and D. LeBihan. Estimation of the effective self-diffusion tensor from the NMR spin echo. *Journal of Magnetic Resonance. Series B*, 103(3):247–254, March 1994. ISSN 1064-1866. doi: 10.1006/jmrb.1994.1037.
69. Peter J. Basser and Carlo Pierpaoli. Microstructural and physiological features of tissues elucidated by quantitative-diffusion-tensor MRI. 1996. *Journal of Magnetic Resonance (San Diego, Calif. : 1997)*, 213(2):560–570, December 2011. ISSN 1096-0856. doi: 10.1016/j.jmr.2011.09.022.
70. Ali Tabesh, Jens H. Jensen, Babak A. Ardekani, and Joseph A. Hergl. Estimation of tensors and tensor-derived measures in diffusional kurtosis imaging. *Magnetic Resonance in Medicine*, 65(3):823–836, March 2011. ISSN 1522-2594. doi: 10.1002/mrm.22655.
71. J.-Donald Tournier, Fernando Calamante, David G. Gadian, and Alan Connelly. Direct estimation of the fiber orientation density function from diffusion-weighted MRI data using spherical deconvolution. *NeuroImage*, 23(3):1176–1185, November 2004. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2004.07.037.
72. J.-Donald Tournier, Fernando Calamante, and Alan Connelly. Robust determination of the fibre orientation distribution in diffusion MRI: non-negativity constrained super-resolved spherical deconvolution. *NeuroImage*, 35(4):1459–1472, May 2007. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.02.016.
73. Ben Jeurissen, Jacques-Donald Tournier, Thijs Dhollander, Alan Connelly, and Jan Sijbers. Multi-tissue constrained spherical deconvolution for improved analysis of multi-shell diffusion MRI data. *NeuroImage*, 103:411–426, December 2014. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2014.07.061.
74. Gabriel Girard, Kevin Whittingstall, Rachid Deriche, and Maxime Descoteaux. Towards quantitative connectivity analysis: reducing tractography biases. *NeuroImage*, 98:266–278, September 2014. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2014.04.074.
75. Robert E. Smith, Jacques-Donald Tournier, Fernando Calamante, and Alan Connelly. Anatomically-constrained tractography: improved diffusion MRI streamlines tractography through effective use of anatomical information. *NeuroImage*, 62(3):1924–1938, September 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.06.005.
76. Marc-Alexandre Côté, Gabriel Girard, Arnaud Boré, Eleftherios Garyfallidis, Jean-Christophe Houde, and Maxime Descoteaux. Tractometer: towards validation of tractography pipelines. *Medical Image Analysis*, 17(7):844–857, October 2013. ISSN 1361-8423. doi: 10.1016/j.media.2013.03.009.
77. Fidel Alfaro-Almagro, Mark Jenkinson, Neal K. Bangerter, Jesper L. R. Andersson, Ludovica Griffanti, Gwenaëlle Douaud, Stamatios N. Sotiropoulos, Saad Jbabdi, Moises Hernandez-Fernandez, Emmanuel Vallee, Diego Vidaurre, Matthew Webster, Paul McCarthy, Christopher Rorden, Alessandro Daducci, Daniel C. Alexander, Hui Zhang, Iulius Dragonu, Paul M. Matthews, Karla L. Miller, and Stephen M. Smith. Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank. *NeuroImage*, 166:400–424, February 2018. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2017.10.034.
78. Karla L. Miller, Fidel Alfaro-Almagro, Neal K. Bangerter, David L. Thomas, Essa Yacoub, Junqian Xu, Andreas J. Bartsch, Saad Jbabdi, Stamatios N. Sotiropoulos, Jesper L. R. Andersson, Ludovica Griffanti, Gwenaëlle Douaud, Thomas W. Okell, Peter Weale, Iulius Dragonu, Steve Garratt, Sarah Hudson, Rory Collins, Mark Jenkinson, Paul M. Matthews, and Stephen M. Smith. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nature Neuroscience*, 19(11):1523–1536, November 2016. ISSN 1546-1726. doi: 10.1038/nn.4393. Number: 11 Publisher: Nature Publishing Group.
79. Eleftherios Garyfallidis, Omar Ocegueda, Demian Wassermann, and Maxime Descoteaux. Robust and efficient linear registration of white-matter fascicles in the space of streamlines. *NeuroImage*, 117:124–140, August 2015. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2015.05.016.
80. Oscar Esteban, Rastko Ciric, Christopher J. Markiewicz, Yaroslav O. Halchenko, Mathias Goncalves, Satrajit S. Ghosh, Russell A. Poldrack, and Krzysztof J. Gorgolewski. TemplateFlow: Standardizing standard 3D spaces in neuroimaging, November 2019.
81. Nobuyuki Otsu. A Threshold Selection Method from Gray-Level Histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, 9(1):62–66, January 1979. ISSN 2168-2909. doi: 10.1109/TSMC.1979.4310076. Conference Name: IEEE Transactions on Systems, Man, and Cybernetics.
82. David Chen, Flavio Dell’Acqua, Ariel Rokem, Eleftherios Garyfallidis, Jidan Zhong, and Mojgan Hodaie. Diffusion Weighted Image Co-registration: Investigation of Best Practices. December 2019. doi: 10.1101/864108.
83. Setsu Wakana, Arvind Caprihan, Martina M. Panzenboeck, James H. Fallon, Michele Perry, Randy L. Gollub, Kegang Hua, Jiangyang Zhang, Hangyi Jiang, Prachi Dubey, Ari Blitz, Peter van Zijl, and Susumu Mori. Reproducibility of Quantitative Tractography Methods Applied to Cerebral White Matter. *NeuroImage*, 36(3):630–644, July 2007. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.02.049.
84. N. Tzourio-Mazoyer, B. Landeau, D. Papathanassiou, F. Crivello, O. Etard, N. Delcroix, B. Mazoyer, and M. Joliot. Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *NeuroImage*, 15(1):273–289, January 2002. ISSN 1053-8119. doi: 10.1006/nimg.2001.0978.
85. C. Bradford Barber, David P. Dobkin, and Hannu Huuhdanpaa. The quickhull algorithm for convex hulls. *ACM Transactions on Mathematical Software*, 22(4):469–483, December 1996. ISSN 0098-3500, 1557-7295. doi: 10.1145/235815.235821.
86. FURY, October 2020.
87. Plotly Python Graphing Library, October 2020.
88. David C. Van Essen, Stephen M. Smith, Deanna M. Barch, Timothy E. J. Behrens, Essa Yacoub, and Kamil Ugurbil. The WU-Minn Human Connectome Project: An overview. *NeuroImage*, 80:62–79, October 2013. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2013.05.041.
89. Lin-Ching Chang, Derek K. Jones, and Carlo Pierpaoli. RESTORE: robust estimation of tensors by outlier rejection. *Magnetic Resonance in Medicine*, 53(5):1088–1095, May 2005. ISSN 0740-3194. doi: 10.1002/mrm.20426.
90. J.-Donald Tournier, Robert Smith, David Raffelt, Rami Tabbara, Thijs Dhollander, Maximilian Pietsch, Daan Christiaens, Ben Jeurissen, Chun-Hung Yeh, and Alan Connelly. MRtrix3: A fast, flexible and open software framework for medical image processing and visualisation. *NeuroImage*, 202:116137, November 2019.
91. Lee R. Dice. Measures of the Amount of Ecologic Association Between Species. *Ecology*, 26(3):297–302, 1945. ISSN 00129658, 19399170. doi: 10.2307/1932409. Publisher: Ecological Society of America.
92. Elizabeth Huber, Rafael Neto Henriques, Julia P. Owen, Ariel Rokem, and Jason D. Yeatman. Applying microstructural models to understand the role of white matter in cognitive development. *Developmental Cognitive Neuroscience*, 36, February 2019. ISSN 1878-9293. doi: 10.1016/j.dcn.2019.100624.
93. Matthew Cieslak, Philip A Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez Adebimpe, Geoffrey K Aguirre, Danielle S Bassett, Richard F Betzel, Josiane Bourque, et al. Qsiprep: An integrative platform for preprocessing and reconstructing diffusion mri. *bioRxiv*, 2020.
94. J.-Donald Tournier, Robert Smith, David Raffelt, Rami Tabbara, Thijs Dhollander, Maximilian Pietsch, Daan Christiaens, Ben Jeurissen, Chun-Hung Yeh, and Alan Connelly. MRtrix3: A fast, flexible and open software framework for medical image processing and visualisation. *NeuroImage*, 202:116137, 2019.
95. Lindsay M Alexander, Jasmine Escalera, Lei Ai, Charissa Andreotti, Karina Febre, Alexander Mangone, Natan Vega-Potler, Nicolas Langer, Alexis Alexander, Meagan Kovacs, Shannon Litke, Bridget O’Hagan, Jennifer Andersen, Batya Bronstein, Anastasia Bui, Marijanje Bushey, Henry Butler, Victoria Castagna, Nicolas Camacho, Elisha Chan, Danielle Citera, Jon Clucas, Samantha Cohen, Sarah Dufek, Megan Eaves, Brian Fradera, Judith Gardner, Natalie Grant-Villegas, Gabriella Green, Camille Gregory, Emily Hart, Shana Harris, Megan Horton, Danielle Kahn, Katherine Kabotyanski, Bernard Karmel, Simon P Kelly, Kayla Kleinman, Bonhwang Koo, Eliza Kramer, Elizabeth Lennon, Catherine Lord, Ginny Mantello, Amy Margolis, Kathleen R Merikangas, Judith Milham, Giuseppe Minniti, Rebecca Neuhaus, Alexandra Levine, Yael Osman, Lucas C Parra, Ken R Pugh, Amy Racanello, Anita Restrepo, Tian Saltzman, Batya Septimus, Russell Tobe, Rachel Waltz, Anna Williams, Anna Yeo, Francisco X Castellanos, Arno Klein, Tomas Paus, Bennett L Leventhal, R Cameron Craddock, Harold S Koplewicz, and Michael P Milham. An open resource for transdiagnostic research in pediatric mental health and learning disorders. *Sci Data*, 4:170181, December 2017.
96. Martin Lindquist. Neuroimaging results altered by varying analysis pipelines. *Nature*, 582(7810):36–37, June 2020. doi: 10.1038/d41586-020-01282-z. Number: 7810 Publisher: Nature Publishing Group.
97. Robert F Dougherty, Michal Ben-Shachar, Gayle K Deutsch, Arvel Hernandez, Glenn R Fox, and Brian A Wandell. Temporal-callosal pathway diffusivity predicts phonological skills in children. *Proc. Natl. Acad. Sci. U. S. A.*, 104(20):8556–8561, May 2007.

Supplementary Methods

Automated Fiber Quantification in Python (pyAFQ). Based on a previous MATLAB implementation (9), we developed a software library that automates dMRI-based tractometry analysis. The library is called pyAFQ (Python Automated Fiber Quantification), and it is implemented as open-source software here: <https://github.com/yeatmanlab/pyAFQ>. The software is developed under the permissive OSI-approved BSD license. It allows users to specify the methods and parameters they want to use for tractometry. pyAFQ uses many components of the scientific Python ecosystem (54). In particular, it relies heavily on implementations of algorithms for diffusion reconstruction, orientation determination, tractography and image registration implemented in Diffusion Imaging in Python (DIPY), an open-source, Python library for computational neuroanatomy (27). The pyAFQ software implements extensive documentation with Sphinx (55), including a gallery of executable examples, implemented using Sphinx Gallery (56). Unit testing is implemented using pytest, with continuous integration implemented to test proposed changes to the library, as well as longer nightly tests that check that pipelines of operations are not adversely affected by changes that are introduced in developing the software. pyAFQ's test suite uses the HARDI data collected for (57), CFIN (58), and data from the Human Connectome Project. pyAFQ can be parallelized across subjects and sessions using dask (59). The analysis performed in this paper primarily used pyAFQ run using Cloudknot (60) on Amazon Web Services (AWS).

There are many ways to analyze dMRI data and to estimate tractometry-based tract-profiles. For example, many different models are used to determine the directions of tracking within each voxel and to connect different voxels with a variety of tractography algorithms. Similarly, different models can be used to determine the tissue properties within a voxel. However, it is hard to determine which methods to use, because different methods may be appropriate for different datasets, depending on their characteristics: the measurements conducted, the signal to noise ratio (SNR) of the data and so forth. Software to support analysis of a variety of datasets should make it easy to use many different methods and to compare results between methods. All of the choices the user can make in each of the steps of pyAFQ are delineated below and summarized in Fig. S2. The software implements a library with an object-oriented application programming interface (API), as well as a command-line interface (CLI). Using pyAFQ's API, pyAFQ can be run with only a few lines of code. The API is also flexible, giving the user the ability to choose which algorithms and parameters to use. For users unfamiliar with python, pyAFQ has a command line interface (CLI) which uses a configuration file written in TOML (61). pyAFQ also has a Boutiques configuration file and can be executed using Boutiques (62).

Locating and mapping data (BIDS). The first step in analysis is to find the files that the software will use. pyAFQ relies on pyBIDS (63, 64) to query data that is provided in the BIDS format (65). It looks for dMRI, b-value, and b-vector files stored in standard formats (see <https://yeatmanlab.github.io/pyAFQ/usage/data.html> for details). Additionally, the user can provide files from other processing pipelines to be used as a brain mask during registration or as start or stop masks during tractography, as well as completed tractography results. We typically use the Nibabel software library to interact with neuroimaging files (66). Following the BIDS standard, the outputs of pyAFQ are put in the BIDS derivatives folder, in a pipeline directory labelled as "afq". The derivative BIDS format follows as much as possible the draft implementation of the BIDS derivatives for dMRI data.

Tractography. There are several methods for computational tractography. The pyAFQ software exposes many of these as options. It allows users to choose from multiple fiber orientation distribution functions (67) that determine the direction of tracking in each step of the process: based on Diffusion Tensor Imaging (DTI) (68, 69), Diffusion Kurtosis Imaging (DKI) (70), Constrained Spherical Deconvolution (CSD) (71, 72), and Multi-Shell Multi-Tissue Constrained Spherical Deconvolution (MSMT-CSD) (73). Deterministic and probabilistic tractography algorithms can be used and stopping criteria can be implemented for particle filtering tractography, using the continuous map criterion (74) or anatomically-constrained tractography (75). The default tractography setting uses DTI, deterministic direction finding, a max turning angle per step of 30°, one seed per voxel, and retains only streamlines between 10 and 1000mm long. Many of our tractography defaults are inspired by the results of (76) and (9). The default seed and stop masks are created by thresholding FA at 0.2. All of these parameters can be customized using pyAFQ's API or CLI.

Template registration. The user can specify their own template and subject image to register, however pyAFQ also provides four builtin options: register subject non-diffusion weighted image (also known as b0) to the Montreal Neurological Institute (MNI) T2 template (28, 29); register subject FA to a group mean fractional anisotropy (FA) template from the UK Biobank (77, 78); register a subject's anisotropic power map (APM) (30?) to the MNI T1 template; and register subject streamlines to the 16 bundles human connectome project (HCP) atlas (7) using streamline registration (SLR) (79). The first three of these builtin techniques use the nonlinear Symmetric Diffeomorphic Registration (SyN) (31) after an optional linear preregistration, both implemented in DIPY. pyAFQ uses Templateflow (80) to get MNI T1/T2 templates for registration. The default registration behavior is to consider all b-values under 50 to be b0, mask the subject's APM using DIPY's median_otsu image recognition algorithm (81) on the subject b0, and register the masked power map to the masked MNI T1 template. Per default, we chose to

use the APM for registration based on previous findings that show this is a good choice (82) and based on our own experience. All of these parameters can be customized using pyAFQ's API and CLI.

Bundle recognition and cleaning. To identify the streamlines that best represent a particular anatomical pathway, we perform bundle recognition. The default behavior is to perform the initial classification using probability maps, and then segment with waypoint ROIs defined in (83), then filter the classified streamlines by their termination locations, using the AAL atlas (84), where streamlines must be within 4mm of the expected endpoint region. Waypoint ROIs are moved into the subject space and then patched up using the Quickhull Algorithm (85). There is also an option, turned off by default, to clip streamline edges at the ROIs (83).

In addition to the waypoint-based recognition described above, pyAFQ also allows the user to choose to use a streamline atlas based bundle recognition method, called RecoBundles (38). Parameters for either algorithm can be customized using pyAFQ's API and CLI.

After recognition, cleaning is performed based on the Mahalanobis distance of each streamline from the mean in each node. This process was originally described in (9). By default, pyAFQ resamples streamlines to 100 points (nodes) and performs 5 rounds of cleaning with a distance threshold of 5 standard deviations from the mean of the node coordinates at each point, and a length threshold of 4 standard deviations from the mean length. Cleaning is also stopped if a bundle has less than 20 streamlines. All of these parameters can be customized using pyAFQ's API and CLI.

Tract Profile Extraction. After cleaning, pyAFQ computes and visualizes tract profiles. The mean profile (called a "tract profile") is calculated using the same Mahalanobis distance-based weighting strategy as in Yeatman et al. (9), implemented in DIPY. Visualization can be performed using one of two backends: fury (86) or plotly (87), which create either animated gifs or interactive html files respectively. Visualizations are created for the whole brain tractometry and for each individual bundle.

Data. We measured the reliability of tractometry using two datasets with contrasting characteristics.

Human Connectome Project (HCP-TR). The WU-Minn Human Connectome Project (HCP) (88) includes measurements of diffusion MRI data from almost all of the 1,200 participants. Here, we focus our analysis on a subset of these subjects for which test-retest data are available. We refer to this data as HCP-TR. This dataset contains dMRI data from 44 individuals. This represents a relatively high-quality, high-resolution dataset, with multiple diffusion directions and multiple b-values. The acquisition parameters of HCP-TR are described in detail elsewhere (34). We used data that had been preprocessed through the HCP pipelines, as provided through the AWS Open Data program (<https://registry.opendata.aws/hcp-openaccess/>).

University of Washington Pre-K (UW-PREK). Two measurements were conducted in each participant 1 day apart. These were acquired with 32 directions, $b=1,500$ s/mm², 2 mm³ isotropic resolution, TR/TE=7200/83 msec. Data were preprocessed using FSL for eddy current, motion correction, and susceptibility distortion correction. Analysis using the mAFQ was conducted as previously described (9). We converted UW-PREK to BIDS format (65) for input into pyAFQ's API.

We attempted to configure pyAFQ to most closely match the mAFQ configuration. We used robust estimation of tensors by outlier rejection (RESTORE) (89) to fit the DTI model. In tractography, we used 160,000 seeds randomly distributed wherever DTI FA is higher than 0.3. We used only 1 round of cleaning. We ran this on both the UW-PREK pre and post sessions, and compared its reproducibility to the results on the same datasets with mAFQ. We also compared the robustness of the results between the pyAFQ and mAFQ algorithms on the pre-session data only.

Configurations. For all configurations, we used the Freesurfer brain segmentation provided by HCP to calculate a permissive brain mask, with all portions of the image not labelled as 0, considered part of the brain. The brain mask is used when fitting the ODF models. We compared the TRR of each configuration, as well as the robustness of the results across configurations. We also compared the TRR of these configurations to the TRR of results published by Lerma-Usabiaga and colleagues (37), denoted RTP.

DTI Configuration. In addition to the three configurations enumerated in the present paper, we processed HCP-TR with a fourth configuration. We used only measurements with b-values between 990 and 1010 s/mm². We used DTI as the ODF model for tractography and profile extraction.

Recobundles Configuration. One of the configurations we ran on the HCP-TR data used Recobundles (8). pyAFQ provides programmatic access to two atlases, one being the full 80 bundles human connectome project (HCP) atlas (7), and other being a 16 bundle subset of that atlas. We ran Recobundles on HCP-TR using the full 80 bundles atlas. We tried only one Recobundles parameter configuration: a model cluster threshold of 5, a reduction threshold of 20, no refinement, a pruning threshold of 5, local streamline-based linear registration on with an asymmetric metric. We used this configuration for all 80 bundles. Multi-shell data and the DKI ODF model were used. We used nonlinear symmetric diffeomorphic registration and a brain mask based on the HCP-provided segmentation. We were not attempting to reproduce the results of (8), as we only tried one configuration, and RecoBundles performance is known to be dependent on the fine-tuning of parameters (39).

RTP. As a point of comparison, we used an open dataset of HCP-TR derivatives that was published by Lerma-Usabiaga and colleagues (37). They processed HCP-TR using the Reproducible Tract Profiles (RTP) pipeline (36). This pipeline is a full end-to-end pipeline and system for deployment of analysis that receives as input raw MRI data as acquired on the scanner. While it applies different preprocessing steps and uses different tractography algorithms than mAFQ, relying on MRTRIX for many of these steps (90), the bundle recognition steps closely resemble the ones used in mAFQ, relying on functions that stem from the same MATLAB codebase as mAFQ. The end result of RTP are tract profiles in an easy-to-use and data-science ready JSON format. We denote their results as RTP and compare them to the HCP-TR results computed with pyAFQ.

Measures of reliability. pyAFQ gives the user the choice of which underlying algorithms to use when performing tractometry, as shown in Fig. S2. We use this feature of pyAFQ to run multiple analyses on HCP-TR and UW-PREK, which both have test-retest data. The analyses we selected represent only a small subset of the possible configurations of pyAFQ. However, because the software is freely available and easily configurable with the API or CLI, it would be straightforward to test other analyses. To compare the results on test-retest data (TRR) and compare results across analyses (robustness), we use four different measures of reliability. Each one of these measures emphasizes different aspects of reliability.

Weighted Dice similarity coefficient (wDSC). The anatomical reliability of bundle recognition solutions is assessed by comparing their spatial overlap in the white matter volume. First, for every voxel in the white matter, we count the number of streamlines that pass through that voxel for a given bundle, then divide by the total number of streamlines. This creates what we call a streamline density map (27). We could compare streamline density maps using a Dice similarity coefficient (91), but that would require applying a threshold to the density maps, and could give a few streamlines a large influence on the calculation. Instead, we use the weighted Dice similarity coefficient (wDSC) (35):

$$D(i, j) = \frac{\sum_{v \in \mathcal{V}_i \cap \mathcal{V}_j} W_{i,v} + W_{j,v}}{\sum_{v \in \mathcal{V}_i} W_{i,v} + \sum_{v \in \mathcal{V}_j} W_{j,v}} \quad (1)$$

where v is a voxel index, $W_{i,v}$ is the streamline density for a bundle i in voxel v , and v' are voxels where the two bundles i and j intersect. wDSC provides a measure of the reliability in the spatial extent of bundles, in a manner that is independent from the assessment of tract profiles.

Adjusted contrast index profile (ACIP). We use an adjusted contrast index to directly compare the values of individual nodes in the tract profiles in different measurements. For two values (V_1, V_2) in different profiles, the adjusted contrast index (ACI) is calculated using Eq (2).

$$ACI(V1, V2) = 2 \frac{V_2 - V_1}{V_2 + V_1} \quad (2)$$

We multiply by 2 to make the contrast index have comparable values to percent difference. In contrast to percent difference, the ACI does not require one of the variables to be a reference, and $ACI(V1, V2) = -ACI(V2, V1)$. Calculating and then plotting the ACI for each point between two profiles highlights the differences between profiles, producing the adjusted contrast index profile (ACIP). ACIP emphasizes discrepancies in estimates along the length of the tract in a manner that does not depend on the scale of the measurement (e.g., the different scales of FA and MD).

Profile reliability. We use profile reliability to compare the shapes of profiles per bundle and per scalar. Given two sets of data (either test-retest or from different analyses), we first calculate the Pearson's r between tract profiles for each subject in a given bundle and scalar. Then, we take the mean of those correlations. We do this for every bundle and for every scalar. We call this profile reliability because larger differences in the overall values along the profiles will result in a smaller mean of the Pearson's r . Consistent profile shapes are important for distinguishing bundles. Profile reliability provides an assessment of the overall reliability of the tract profiles, summarizing over the full length of the bundle, for a particular scalar. We calculate the 95% confidence interval on profile reliabilities using the standard error of the measurement.

Subject reliability. We calculate subject reliability to compare individual differences in profiles, per bundle and per scalar, following (92). Given two measurements for each subject, we first take the mean of each profile within each individual, measurement and scalar. Then we correlate the means from different subjects for a given bundle and scalar across the measurements. High subject reliability means the individual differences in tract profile means are consistent across measurements. This is akin to test reliability which is computed for any clinical measure.

One downside of subject reliability is that the shape of the extracted profile is not considered. However, it well summarizes the preservation of relative differences between individuals for mean tract profiles. We use the Fisher transformation to calculate

the 95% confidence interval for subject reliabilities, because when the transformation is applied to Pearson's r , the result is approximately normal.

Supplementary Discussion of pyAFQ

pyAFQ is embedded in an ecosystem of tools for reproducible neuroimaging. The wider ecosystem of tools and standards surrounding pyAFQ is shown in Fig. S6. Each tool has its own place in the ecosystem. We rely heavily on implementations of dMRI analysis algorithms implemented in DIPY (27). Reproducibility and interoperability are also facilitated by relying on the BIDS format (65) and the pyBIDS software (63, 64). Requiring a BIDS-like input makes integration with other software in the ecosystem easier. For example, it is fairly straightforward to use the outputs of BIDS-compatible preprocessing pipelines, such as qsbprep (93), as inputs to pyAFQ. Furthermore, the modularity of the pyAFQ pipeline means that outputs of other tractography software (e.g., MRTRIX (94)) can be used as inputs to bundle recognition, with BIDS filters as the metadata that allows finding and incorporating through the right data.

Cloud-based processing is going to be more important as large datasets are processed. pyAFQ does not depend on proprietary software and can be scaled to large datasets using cloud computing platforms. In this paper, we used Cloudknot (60) to scale pyAFQ across subjects and methods on AWS. However, because pyAFQ is a Python package, it can easily be run on any cloud computing platform. Computing in the public cloud also supports reproducible research, as computations conducted on the public cloud are perfectly portable to other users of the software. Our software is written with that in mind, including functions that know how to easily access datasets that are already stored in the cloud (e.g., HCP and Healthy Brain Network (95) datasets). We know that one of the most important ways in which users can diagnose whether processing worked as expected is by visually inspecting the results. Thus, we provide several different visualization methods, relying on the VTK-derived FURY library, or on browser-friendly visualizations with Plotly. pyAFQ outputs are also fully compatible with AFQ-Browser, a browser-based tool for interactive visualization and exploration of tractometry results (50).

Finally, beyond visualization and summary of the results, and tools for analysis of reliability presented in this work, pyAFQ does not provide a substantial set of tools for statistical analysis of tractometry results. Instead, the outputs of pyAFQ are provided as “tidy” CSV tables (26). This means that it is compatible as inputs to the AFQ Insight tool for statistical analysis (19), but also amenable to many other statistical analysis approaches. This output should facilitate interdisciplinary use of dMRI data, as it is provided in a format that is widely used in statistics and machine learning.

pyAFQ is extensible. In general, variability in results would be reduced with a standard pipeline that could be used across all studies and datasets. However, as noted by Lindquist, “studies tend to be too varied for one pipeline to always be appropriate” (96). This is particularly true as new measurement techniques, new processing methods and new analysis approaches for dMRI are evolving. Therefore, the pyAFQ pipeline was designed to be flexible, making it easier to reproduce results, while providing researchers with many choices for the appropriate analysis, depending on their data and questions. pyAFQ allows the user to make many decisions (Fig S2), and all of those decisions can be encoded in a configuration file. That configuration file can be used to reproduce the same analysis pipeline given the same version of pyAFQ is used. By providing the configuration file or the arguments passed to the main API, one can clearly satisfy the requirement for a re-executable workflow outlined in (53).

To extend to new bundles, pyAFQ allows users to define new queries that recognize bundles that are not part of the set of 18 detected by the original mAFQ software. For a simple example, we use a set of alternative waypoint ROIs to detect different portions of the corpus callosum (97) (Fig S7A). These alternative ROIs are included in pyAFQ but not used by default. In more complicated example, another set of ROIs is used to recognize the location of the optic radiations (OR; Fig S7). Because these are relatively small and winding, their delineation requires additional components: it requires several waypoint ROIs used not only as inclusion criteria, but also as exclusion criteria, and it requires delineation of endpoints in the cortex that are not part of the AAL atlas, which is used in the standard set of bundles. It also requires oversampling of streamlines, so in order to obtain a proper definition of the OR, tractography is configured to use 125 seeds per voxel (instead of the default 8). All of these components can be integrated into calls to the software API, without needing to change any of its internals. This includes any custom waypoint ROIs, inclusive or exclusive, as well as probability maps, endpoint locations, and whether the bundle crosses the midline.

Supplementary Figures and Tables

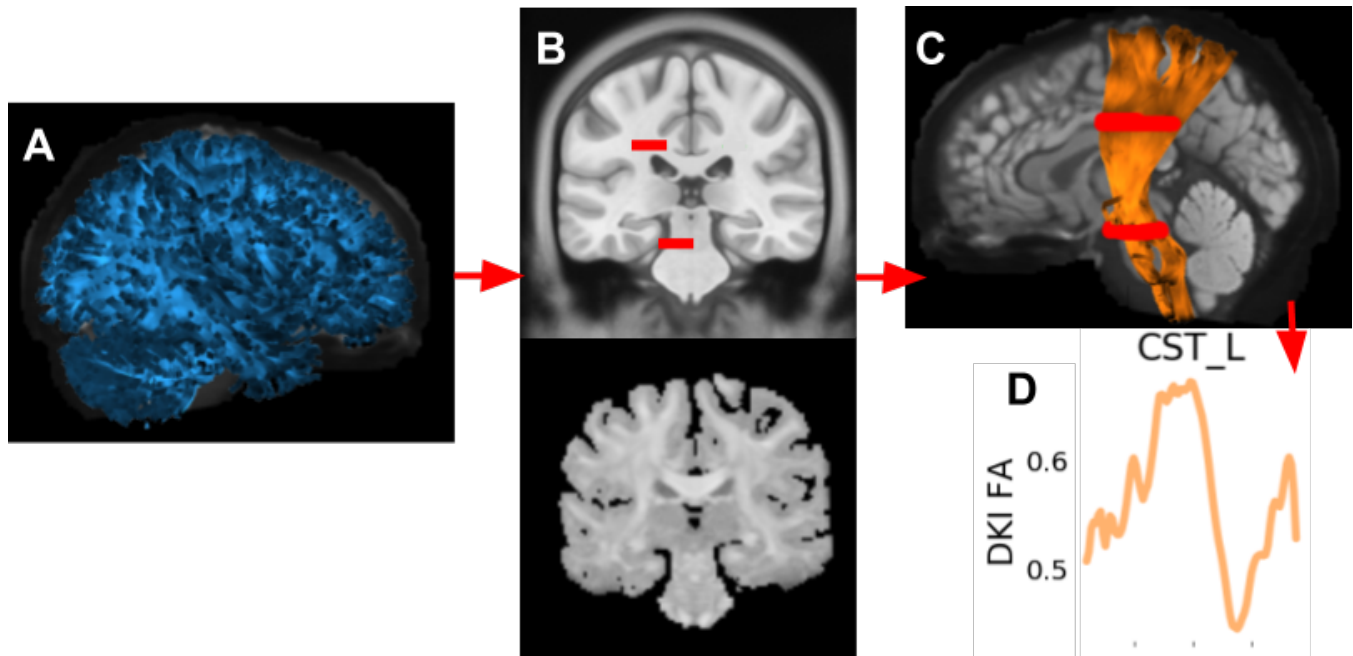


Fig. S1. The stages of tractometry. **A** Computational tractography generates streamlines estimating the trajectories of white matter connections. **B** An anatomical template is registered to each subject's individual brain. Here, in a mid-coronal view, the MNI T1-weighted template (28, 29), shown with the locations of waypoint ROIs for classification of the left corticospinal tract (5) (slightly enlarged for visualization purposes). The subject's anisotropic power map (APM) (30) is used as the target for registration, due to its similarity to the T1 contrast. **C** Classification of the streamlines. Here, in a lateral view, the streamlines classified as belonging to the left corticospinal tract (CST L), overlaid on a mid-sagittal slice of the subject's non-diffusion-weighted (b0) image. The streamlines are shaded by the subject's fractional anisotropy (FA) along their length. **D**, Tract profiles are extracted from the bundles. Here, the FA profile for CST L.

ARC L	Left Arcuate
ARC R	Right Arcuate
ATR L	Left Thalamic Radiation
ATR R	Right Thalamic Radiation
CGC L	Left Cingulum Cingulate
CGC R	Right Cingulum Cingulate
CST L	Left Corticospinal
CST R	Right Corticospinal
FA	Callosum Forceps Minor
FP	Callosum Forceps Major
IFO L	Left Inferior Fronto-occipital Fasciculus
IFO R	Right Inferior Fronto-occipital Fasciculus
ILF L	Left Inferior Longitudinal Fasciculus
ILF R	Right Inferior Longitudinal Fasciculus
SLF L	Left Superior Longitudinal Fasciculus
SLF R	Right Superior Longitudinal Fasciculus
UNC L	Left Uncinate
UNC R	Right Uncinate

Table S1. Abbreviations of the major white matter pathways recognized by pyAFQ.

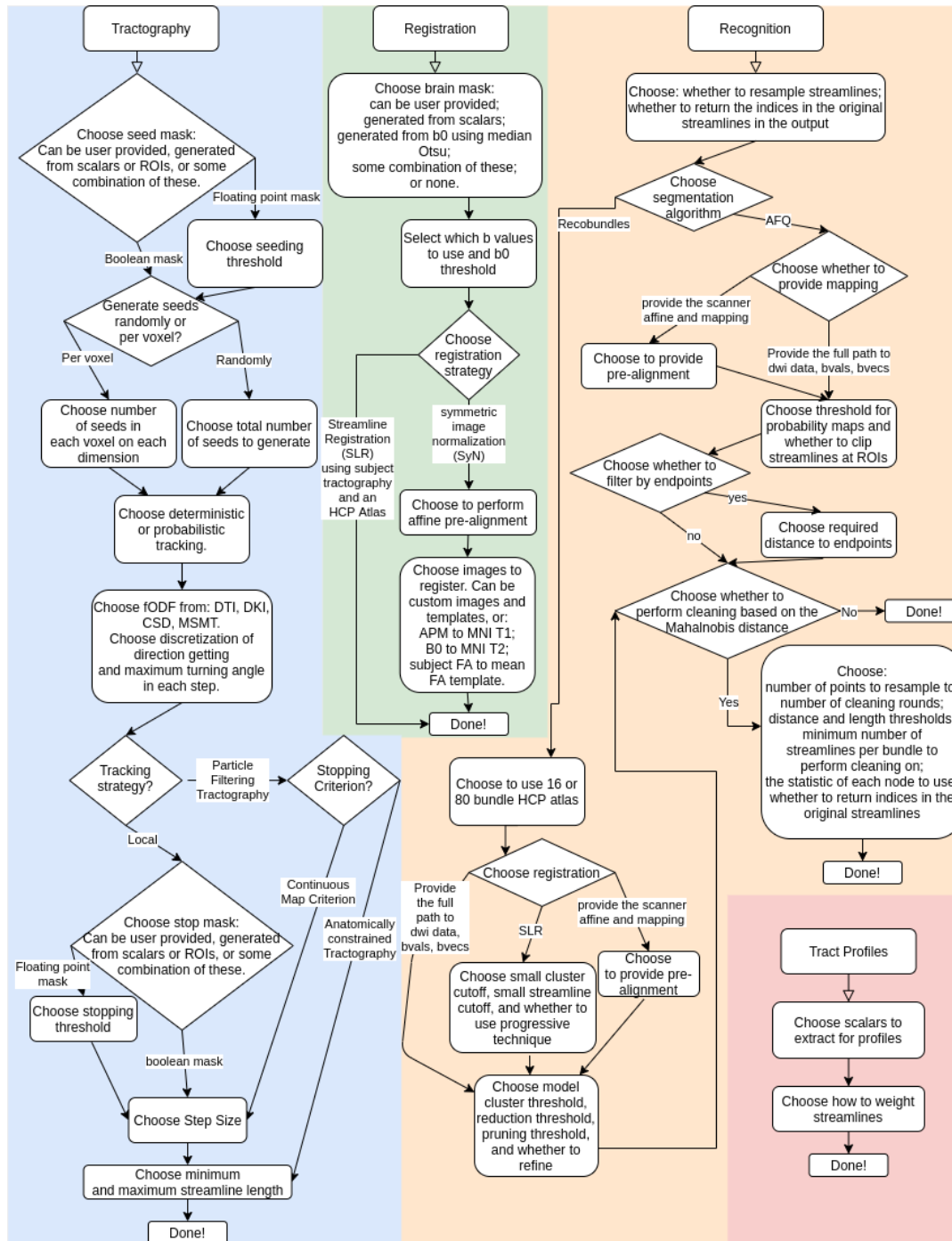


Fig. S2. Choices the user can make for how to run pyAFQ. The colors represent different steps of tractometry. Tractography is shaded blue, registration is shaded green, recognition is shaded orange, and tract profiles is shaded red. Every rounded box and diamond contains one or more choices, except for the rounded boxes marked "Done!", which indicates all choices have been made. Diamonds indicate the path you take depends on the choice in the diamond. pyAFQ has reasonable defaults for all of these decisions; however it also makes it simple for the user to customize their tractometry pipeline according to their needs.

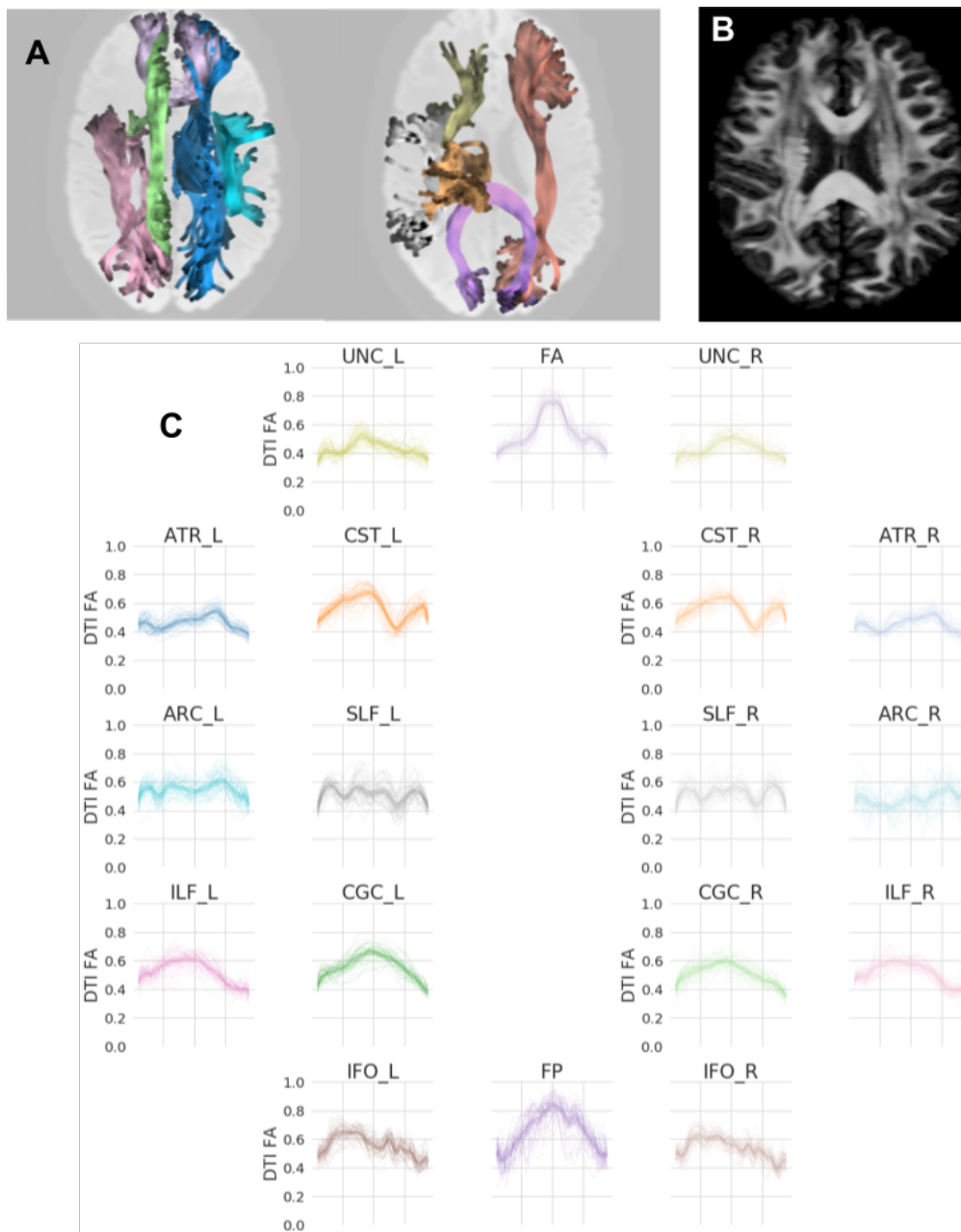


Fig. S3. Extraction of tract profiles from the recognition of white matter into major bundles of streamlines. **A** Representative bundles from an example subject in the HCP-TR dataset. Streamlines are colored by bundle, and are shaded by the interpolated FA value at each point. The background is the mean non diffusion-weighted image (b_0). **B** The same subject's fractional anisotropy (FA). **C** extracting FA along each bundle and plotting the FA in a tract profile. Individual tract profiles are plotted with thin lines and the mean tract profile is plotted with a thick line. The tract profiles are colored according to their bundle are laid out in positions that reflect their anatomical positions (compare **A** and **C**).

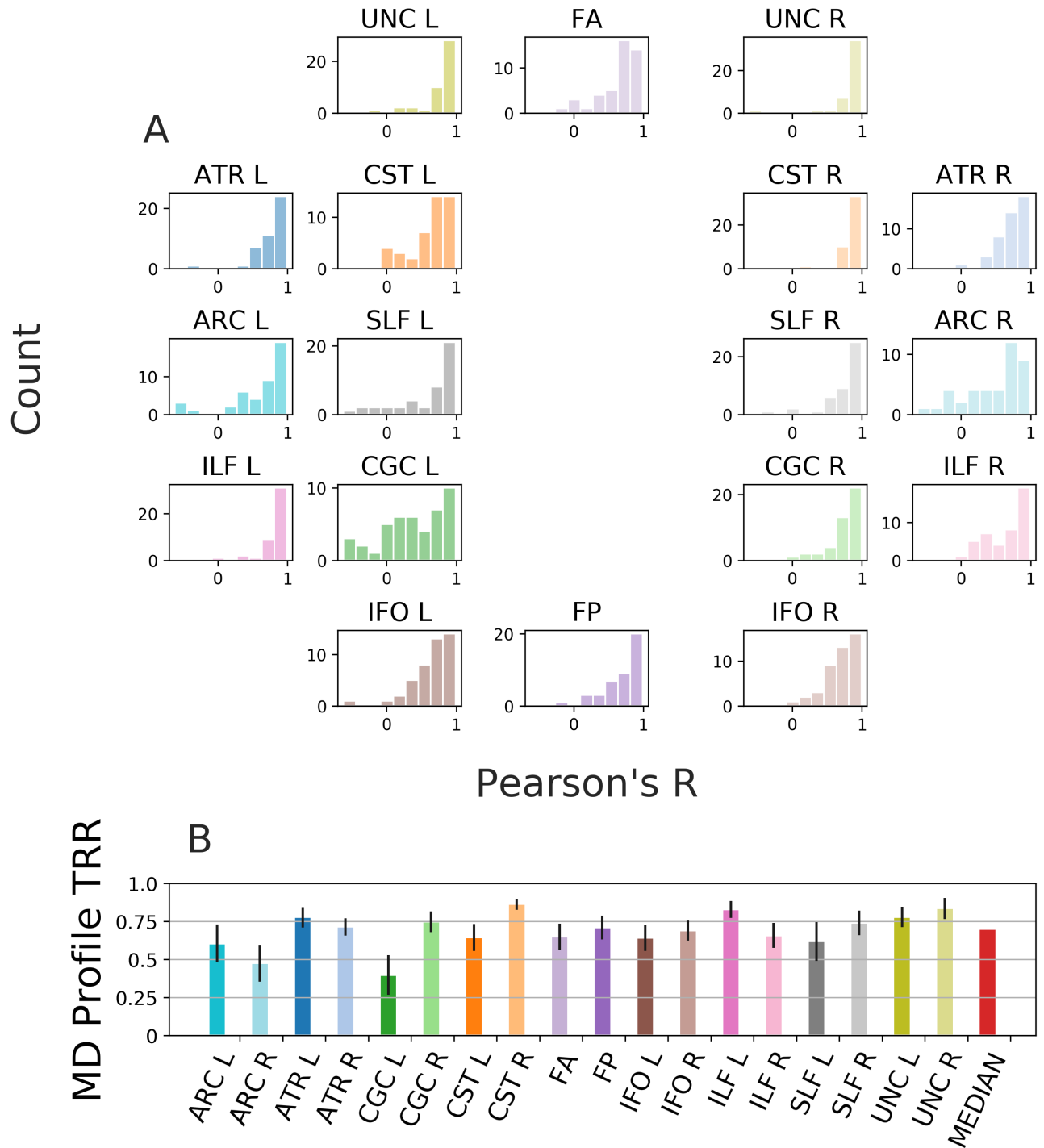


Fig. S4. MD profile test-retest reliability **A:** Histograms of individual subject Pearson's r between the MD tract profiles across sessions for a given bundle. Colors encode the bundles, matching the diagram showing the rough anatomical positions of the bundles for the left side of the brain (center). **B:** Mean (\pm 95% confidence interval) TRR for each bundle, color-coded to match the histograms and the bundles diagram, with median across bundles in red.

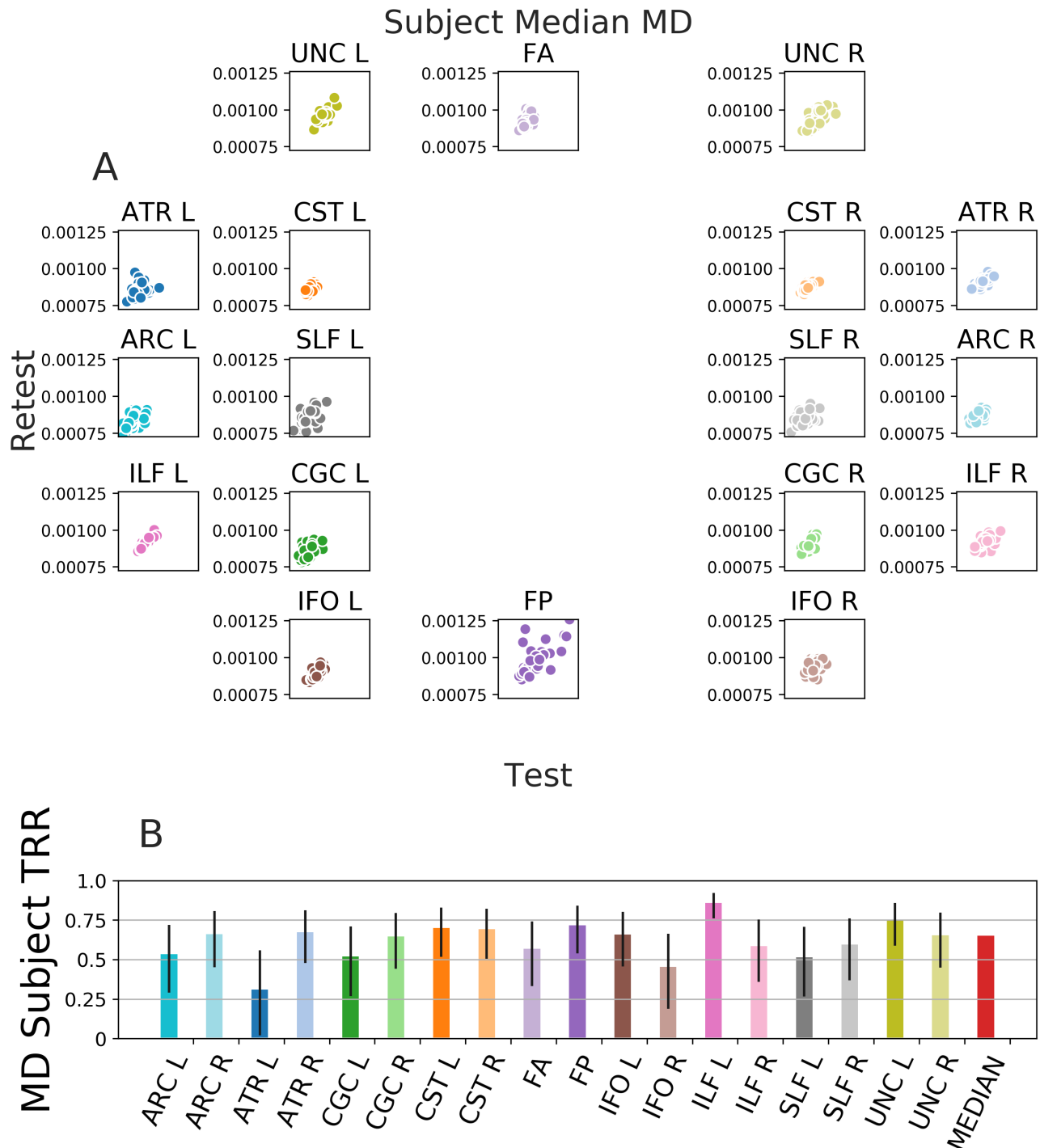


Fig. S5. Subject test-retest reliability **A:** Mean tract profiles for a given bundle and the MD scalar for each subject using the first and second session of HCP-TR. Colors encode bundle information, matching the core of the bundles (center). **B:** subject reliability is calculated from the Pearson's correlation coefficient of these distributions, with median across bundles in red.

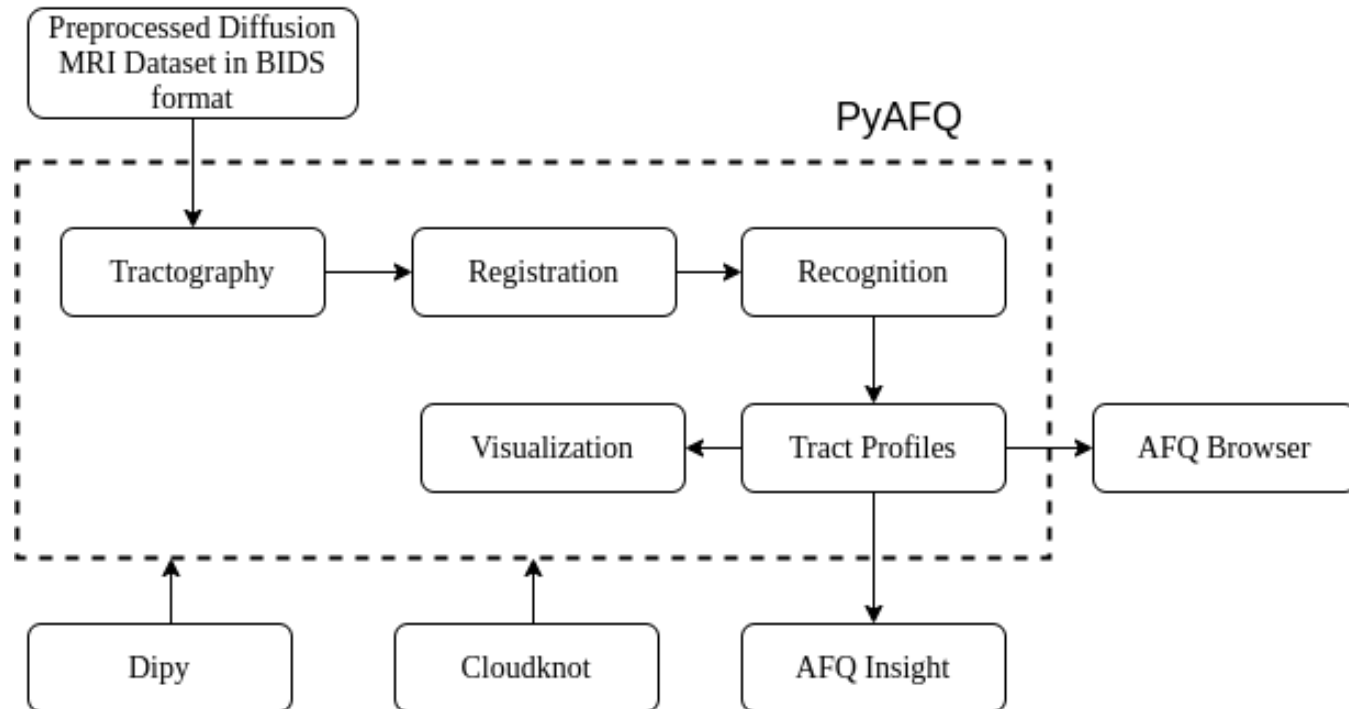


Fig. S6. The pyAFQ software is intergrated into an ecosystem for reproducible tractometry Steps performed by pyAFQ are enclosed in the dotted rectangle, whereas steps outside that rectangle are performed by other software. Upper left: pyAFQ requires preprocessed diffusion MRI data in BIDS format. This could be from QSIprep (25) or dMRIprep (<https://github.com/nipreps/dmriprep>). Bottom right: pyAFQ outputs can serve as inputs to AFQ Browser for further interaction and visualization (50) or AFQ Insight for statistical analysis (19). Bottom left: pyAFQ uses DIPY (27) for the implementation of dMRI algorithms. pyAFQ uses Cloudknot (60) to scale processing by parallelizing across subjects in AWS.

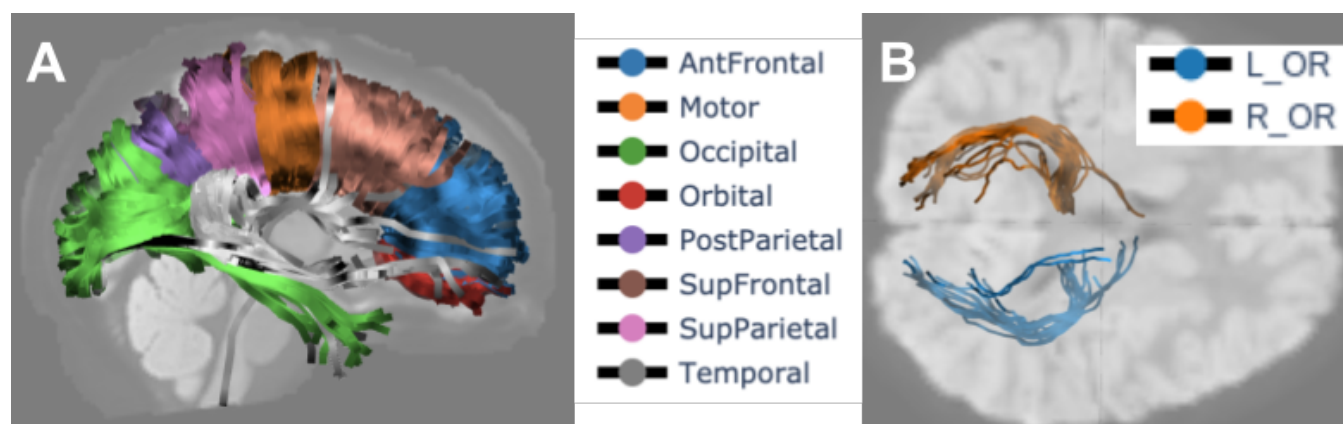


Fig. S7. Callosal bundles from HCP-TR, optic radiations from UW-PREK, found by pyAFQ. Streamlines are colored according to their bundles and shaded according to FA. The background images are each a b0 slice. **A** callosal bundles found by pyAFQ on an example subject from HCP-TR. **B** optic radiations found by pyAFQ on an example subject from UW-PREK.