

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, Bramsh Q. Chandio^j, Eleftherios Garyfallidisⁱ, and Ariel Rokem^{1a,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

^jDepartment of Intelligent Systems Engineering, Luddy School of Informatics, Computing and Engineering, Indiana University Bloomington, Bloomington, IN, 47408, 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

1 Introduction

2 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

15 these are models of the anatomy, we refer to these estimates as “bundles” to distinguish them from the anatomical pathways 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).

16 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 recognition 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, 18–21). For the inferences made from tractometry to be valid and useful, tract profiles need to be reliable.

17 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 (22, 23). Here, we use the definitions proposed in (24). *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* (25) show that even when full computational reproducibility is achieved, the re-

sults 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 (26)), 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 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 (28). Our implementation was also guided by a previous MATLAB implementation of tractometry (mAFQ) (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 (28), are used to generate streamlines throughout the brain white matter (Fig. S1A). Next, the T1-weighted MNI template (29, 30) is registered to the anisotropic power map (APM) (31, 32) computed from the diffusion data, that has a T1-like contrast (Fig. S1B) using the symmetric image normalization method (33) implemented in DIPY (28). 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 (34). 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 (35) 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., fractional anisotropy (FA) and mean diffusivity (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). To make sure that profiles represent properties of the core white matter, we remove the first and last 5 nodes of the profile, then further remove any nodes where either the FA is less than 0.2 or the MD is greater than 0.002. This removes nodes that contain partial volume artifacts (16).

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) (36). 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 the diffusion kurtosis model (DKI) as the orientation distribution function (ODF) model. In the second configuration, we used constrained spherical deconvolution (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 recognition of each bundle was compared across measurements and methods using the Dice coefficient, weighted by streamline count (wDSC) (37). Tract profiles were compared with three measures: (1) Profile reliability: mean intraclass correlation coefficient (ICC) across points in different tract profiles for different data, which quantifies the *agreement* of tract profiles (38, 39); (2) Subject reliability: Spearman’s rank correlation coefficient (Spearman’s ρ) between the mean of the tract profiles across individuals, which quantifies the *consistency* of the mean of tract profiles; (3) an adjusted contrast index profile (ACIP) 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. For example, if we calculate the

164 subject reliability across analysis methods, we would call 219
 165 that “subject robustness”. If we calculated subject reliability 220
 166 across measurements, we would call that “subject TRR”. We 221
 167 explain profile and subject reliability in more detail below; 222
 168 we explain wDSC and ACIP in more detail in the ‘Measures
 169 of Reliability’ section of Supplementary Methods 223

224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
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 ICC 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 ICC. Consistent profile shapes are important for
 distinguishing bundles. Profile reliability provides an assess-
 ment of the overall reliability of the tract profiles, summariz-
 ing over the full length of the bundle, for a particular scalar.
 We calculate the 95% confidence interval on profile reliabili-
 ties using the standard error of the measurement.

184 In some cases, there is low between-subject variance in
 185 tract profile shape (for example, this is often the case in
 186 CST). We use ICC to account for this, as ICC will penal-
 187 ize low between-subject variance in addition to rewarding
 188 high within-subject variance. Profile reliability is a way of
 189 quantifying the *agreement* between profiles. Qualitatively,
 190 we use four descriptions for profile reliability: excellent (ICC
 191 > 0.75), good (ICC = 0.60 to 0.74), fair (ICC = 0.40 to 0.59),
 192 and poor (ICC < 0.40) (40).

193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
Subject reliability. We calculate subject reliability to compare
 individual differences in profiles, per bundle and per scalar,
 following (41). Given two measurements for each subject,
 we first take the mean of each profile within each individ-
 ual, measurement and scalar. Then we calculate Spearman’s
 ρ from the means from different subjects for a given bundle
 and scalar across the measurements. High subject reliabil-
 ity means the ordering of an individual’s tract profile mean
 among other individuals is consistent across measurements
 or methods. This is akin to test reliability which is computed
 for any clinical measure.

204 One downside of subject reliability is that the shape of the
 205 extracted profile is not considered. Additionally, if one mea-
 206 surement or method produces higher values for all subjects
 207 uniformly, subject reliability would not be affected. Instead,
 208 the intent of subject reliability is to well summarize the
 209 preservation of relative differences between individuals for
 210 mean tract profiles. In other words, subject reliability quan-
 211 tifies the *consistency* of mean profiles. The 95% confidence
 212 interval on subject reliabilities are parametric.

213 Results

214 Tractometry using pyAFQ classifies streamlines into bundles 242
 215 that represent major anatomical pathways. The streamlines 243
 216 are used to sample dMRI-derived scalars into bundle profiles
 217 that are calculated for every individual and can be summa- 244
 218 rized for a group of subjects. An example of the process and 245

result of the tract profile extraction process is shown in Sup-
 plementary 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 reli-
 ability (TRR) at several different levels of tractometry. For ex-
 ample, the correlation between two profiles provides a mea-
 sure of the reliability of the overall tract profile in that sub-
 ject. Analyzing the Human Connectome Project’s test-retest
 dataset (HCP-TR), we find that for fractional anisotropy (FA)
 calculated using 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
 0.77 ± 0.05 to 0.92 ± 0.02 and a median across bundles of
 0.86 (Figure 1B). We find similar results for mean diffusivity
 (MD; Fig. S4) and replicate similar results in a second dataset
 (Fig. 3B).

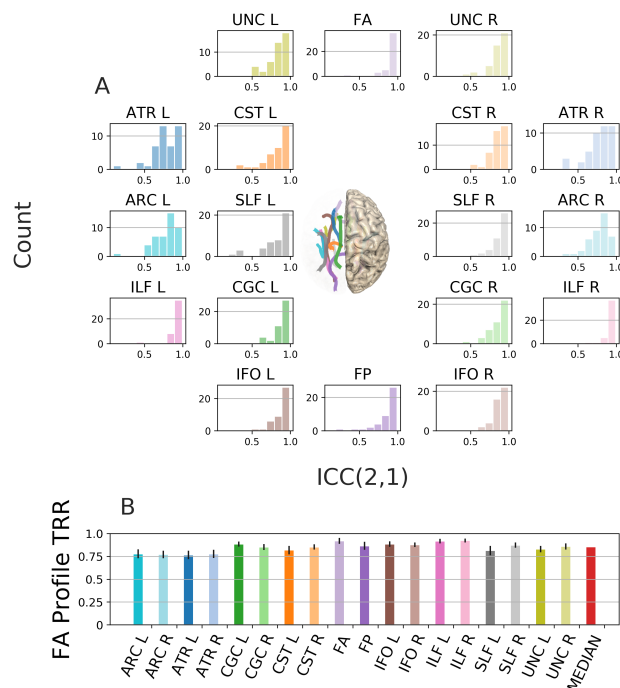


Fig. 1. FA profile test-retest reliability A: Histograms of individual subject ICC between the FA tract profiles across sessions for a given bundle. Colors encode the bundles, matching the diagram showing the rough anatomical positions of the bundles for the left side of the brain (center). B: Mean (\pm 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 0.57 ± 0.24 to 0.85 ± 0.12 and a median across bundles of 0.73. 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).

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

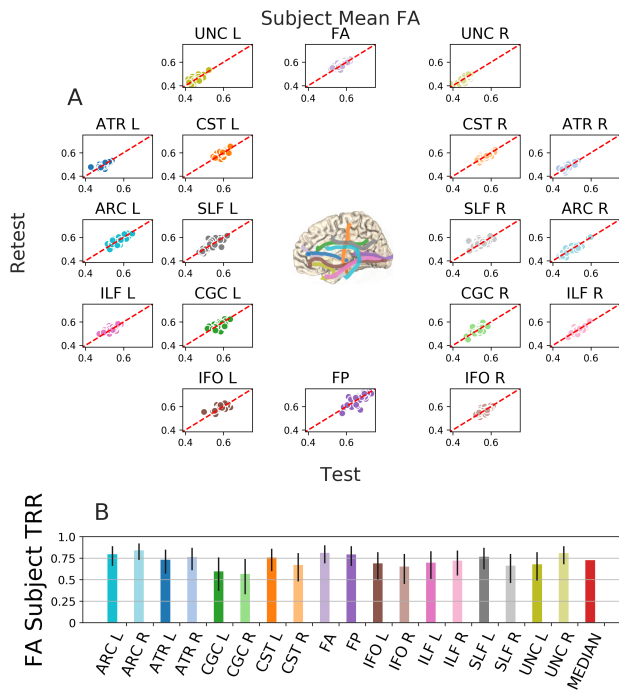


Fig. 2. Subject test-retest reliability **A:** Mean tract profiles for a given bundle and the FA scalar for each subject using the first and second session of HCP-TR. Colors encode bundle information, matching the core of the bundles (center). **B:** subject reliability is calculated from the Spearman's ρ of these distributions, with median across bundles in red (\pm 95% confidence interval).

246 compared TRR across datasets and implementations. In both
 247 datasets, we found high TRR in the results of tractography
 248 and bundle recognition: wDSC was larger than 0.7 for all
 249 but one bundle (Fig. 3A): the delineation of the anterior
 250 forceps (FA bundle) seems relatively unreliable using pyAFQ
 251 in the UW-PREK dataset (using the FA scalar, pyAFQ
 252 subject TRR is only 0.37 ± 0.28 compared to mAFQ's $0.84 \pm$
 253 0.10). We found overall high profile TRR that did not always
 254 translate to high subject TRR (Fig. 3B-G). For example, for
 255 FA in UW-PREK, median profile TRRs are 0.75 for pyAFQ
 256 and 0.77 for mAFQ while median subject TRRs are 0.70 for
 257 pyAFQ and 0.75 for mAFQ. Note that profile and subject
 258 TRR have different denominators (for example, subjects that
 259 have similar mean profiles to each other would have low sub-
 260 ject TRR, even if the profiles are reliable, because it is harder
 261 to distinguish between subjects in this case). mAFQ is one of
 262 the most popular software pipelines currently available for
 263 tractometry analysis, so it provides an important point for
 264 comparison. In comparing different software implemen-
 265 tations, we found that mAFQ has higher subject TRR relative
 266 to pyAFQ in the UW-PREK dataset, when TRR is relatively
 267 low for pyAFQ (see the FA bundle, CST L, and ATR L in
 268 Fig. 3C). On the other hand, in the HCP-TR dataset pyAFQ
 269 we used the RTP pipeline (42, 43), which is an extension of
 270 mAFQ, and found that pyAFQ tends to have slightly higher
 271 profile TRR than RTP for MD, but slightly lower profile TRR
 272 for FA (Fig. 3D). The pyAFQ and RTP subject TRR are
 273 highly comparable (Fig. 3E). In FA, the median pyAFQ sub-
 274 ject TRR for FA is 0.76 while the median RTP subject TRR is

275 0.74. Comparing different ODF models in pyAFQ, we found
 276 that the DKI and CSD ODF models have highly similar TRR,
 277 both at the level of wDSC (Fig. 3A), as well as at the level of
 278 profile and subject TRR (Fig. 3F-G).

**Robustness: comparison between distinct tractogra-
 phy 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 compar-
 ison 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 0.77
 and 0.75, respectively) (Figure 4B,C). In contrast to the re-
 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 re-
 mains remarkably high (In FA, 0.75 ± 0.17 for ARC R and
 0.88 ± 0.09 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 ap-
 proach (RecoBundles) that uses atlas templates in the space
 of the streamlines (44). Between these algorithms, wDSC is
 around or above 0.6 for all but one bundle, ILF R (Figure 5).
 There is an asymmetry in the ILF atlas bundle(7), which re-
 sults in discrepancies between ILF R recognized with way-
 point ROIs and with RecoBundles. Despite this bundle, we
 find high robustness overall. For MD, the first quartile subject
 robustness is 0.82 (Figure 5C, D).

**Tractometry results are robust to differences in software im-
 plementation.** Overall, we found that robustness of tractom-
 etry across these different software implementations is high
 in most white matter bundles. In the mAFQ/pyAFQ compar-
 ison, most bundles have a wDSC around or above 0.8,
 except the two callosal bundles (FA bundle 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 0.71
 and 0.77 for FA profile and subject robustness, respectively.

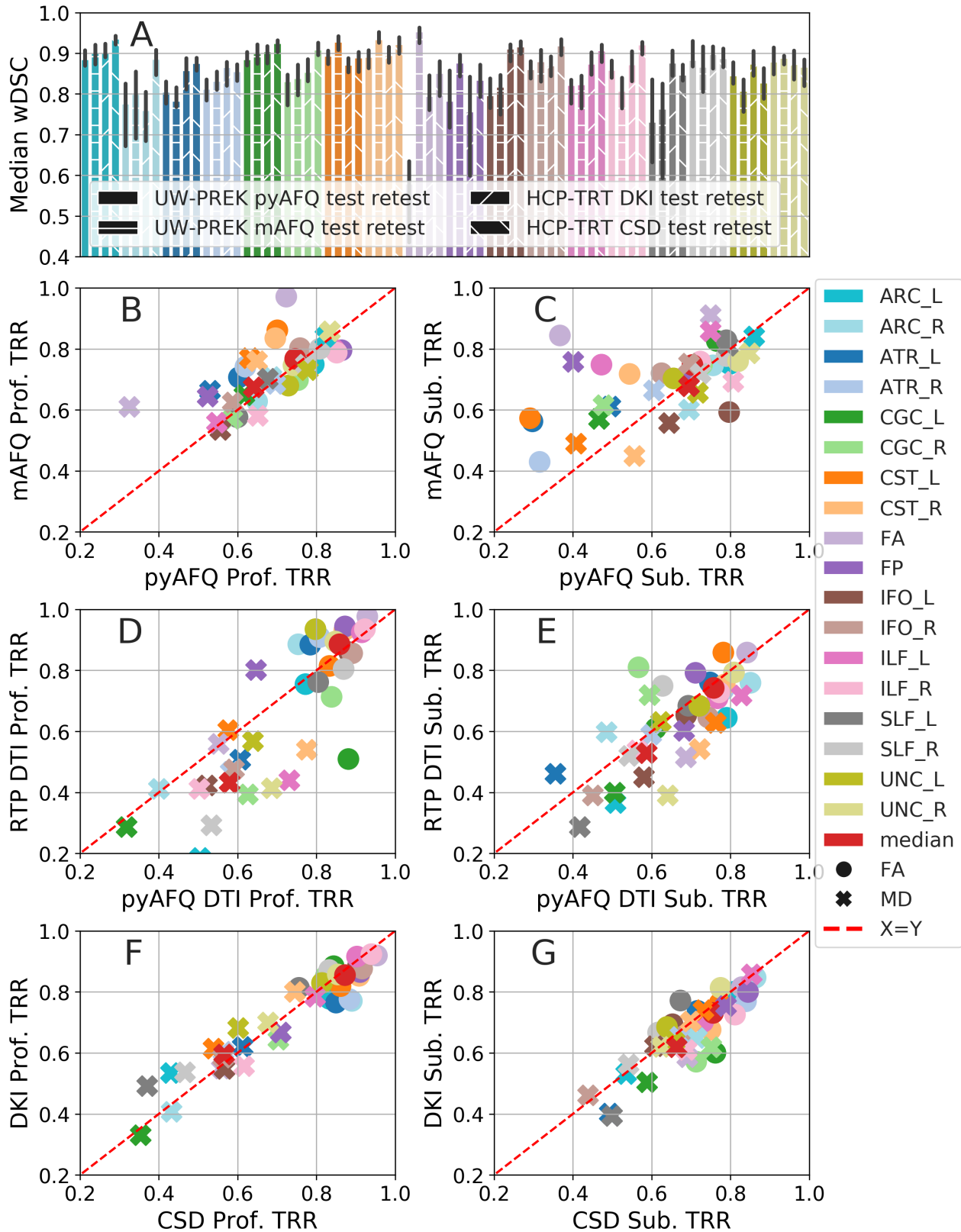


Fig. 3. wDSC, profile, and subject TRR of: pyAFQ and mAFQ on UW-PREK; pyAFQ on HCP-TR using different ODF models; and RTP on HCP-TR. Colors indicate bundle. In **A**: texture indicates the dataset and methods being compared. Error bars show the 95% confidence interval. **B**, **D**, and **F** show profile TRR and **C**, **E**, and **G** show subject TRR. Profile and subject TRR calculations are demonstrated with HCP-TR using DKI in figures 1 and 2 respectively. 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.

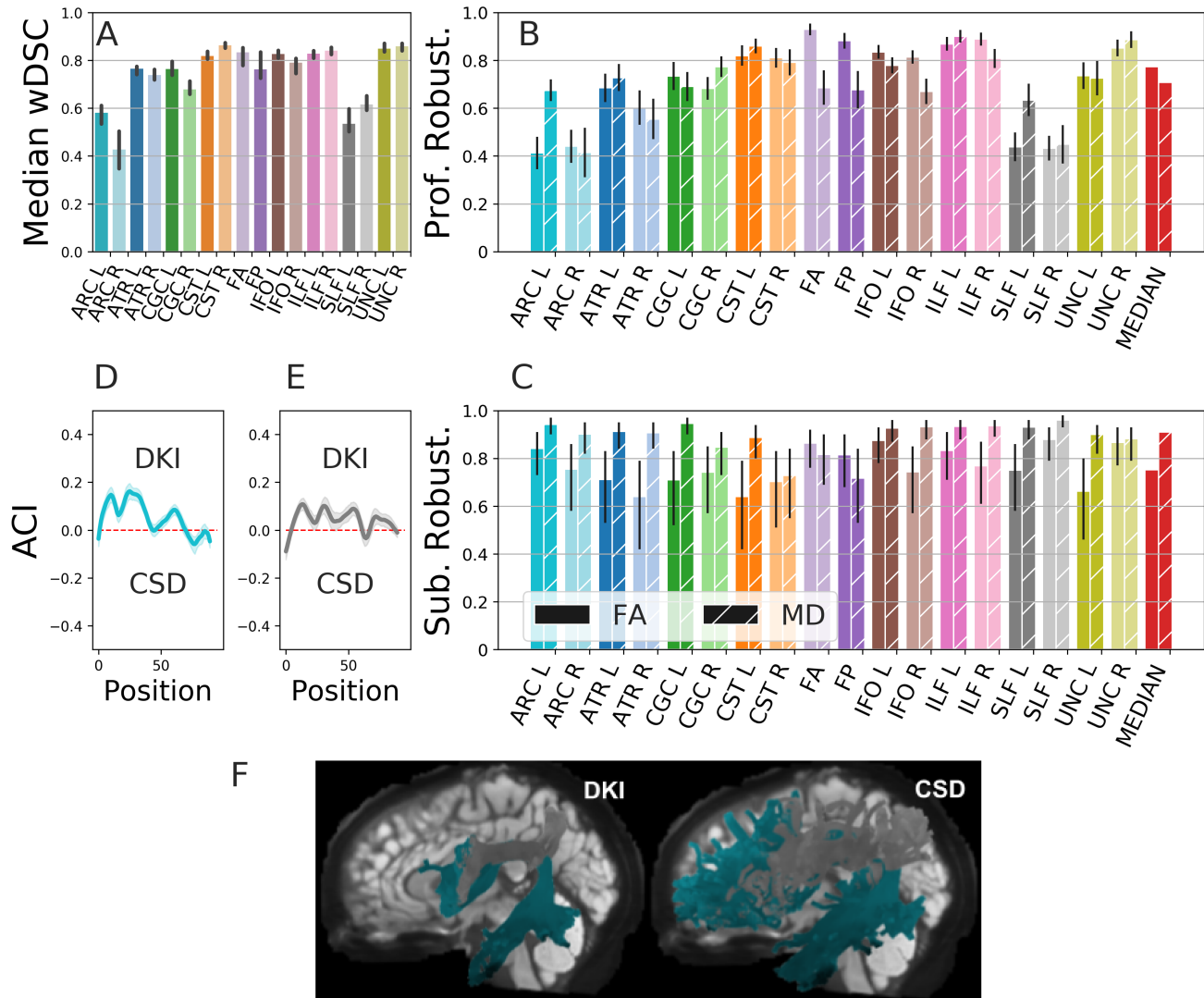


Fig. 4. ODF model robustness. We compared DKI- and CSD-derived tractography. Colors encode bundle information as in Figures 1 and 2. Textured hatching encodes FA/MD information. **A** wDSC robustness. **B** Profile robustness. **C** Subject robustness. Error bars represent 95% confidence interval. **D, E** Adjusted contrast index profile (ACIP) between ARC L and SLF L tract profiles of each algorithm. 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.

329 For some bundles, like the right and left uncinate, there is 346
 330 large agreement between pyAFQ and mAFQ (for subject FA: 347
 331 UNC L $\rho = 0.90 \pm 0.07$, UNC R $\rho = 0.89 \pm 0.08$). How-
 332 ever, the callosal bundles have particularly low mean diffu- 348
 333 sivity (MD) profile robustness (Fig. 6B) (0.07 ± 0.09 for FP,
 334 0.18 ± 0.09 for FA).

335 The robustness of tractometry to the differences between the 351
 336 pyAFQ and mAFQ implementation depends on the bundle, 352
 337 scalar, and reliability metric. In addition, for many bundles, 353
 338 the ACIP between mAFQ and pyAFQ results is very close 354
 339 to 0, indicating no systematic differences (Fig. 6D). In some 355
 340 bundles – the corticospinal tract (CST) and the anterior thala- 356
 341 mic radiations (ATR) – there are small systematic differences 357
 342 between mAFQ and pyAFQ. In the Forceps Posterior (FP), 358
 343 pyAFQ consistently finds smaller FA values than mAFQ in a 359
 344 section on the left side. Notice that the forceps anterior has 360
 345 an ACIP that deviates only slightly from 0, even though the 361

forceps recognitions did not have as much overlap as other bundle recognitions (see Fig. 6A).

Discussion

349 Previous work has called into question the the reliability
 350 of neuroimaging analysis (e.g., (25, 45, 46)). 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 out-
 puts 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 dif-
 ferent models of diffusion at the level of individual voxels;
 techniques to connect voxels through tractography; and ap-
 proaches to classify tractography results into major white

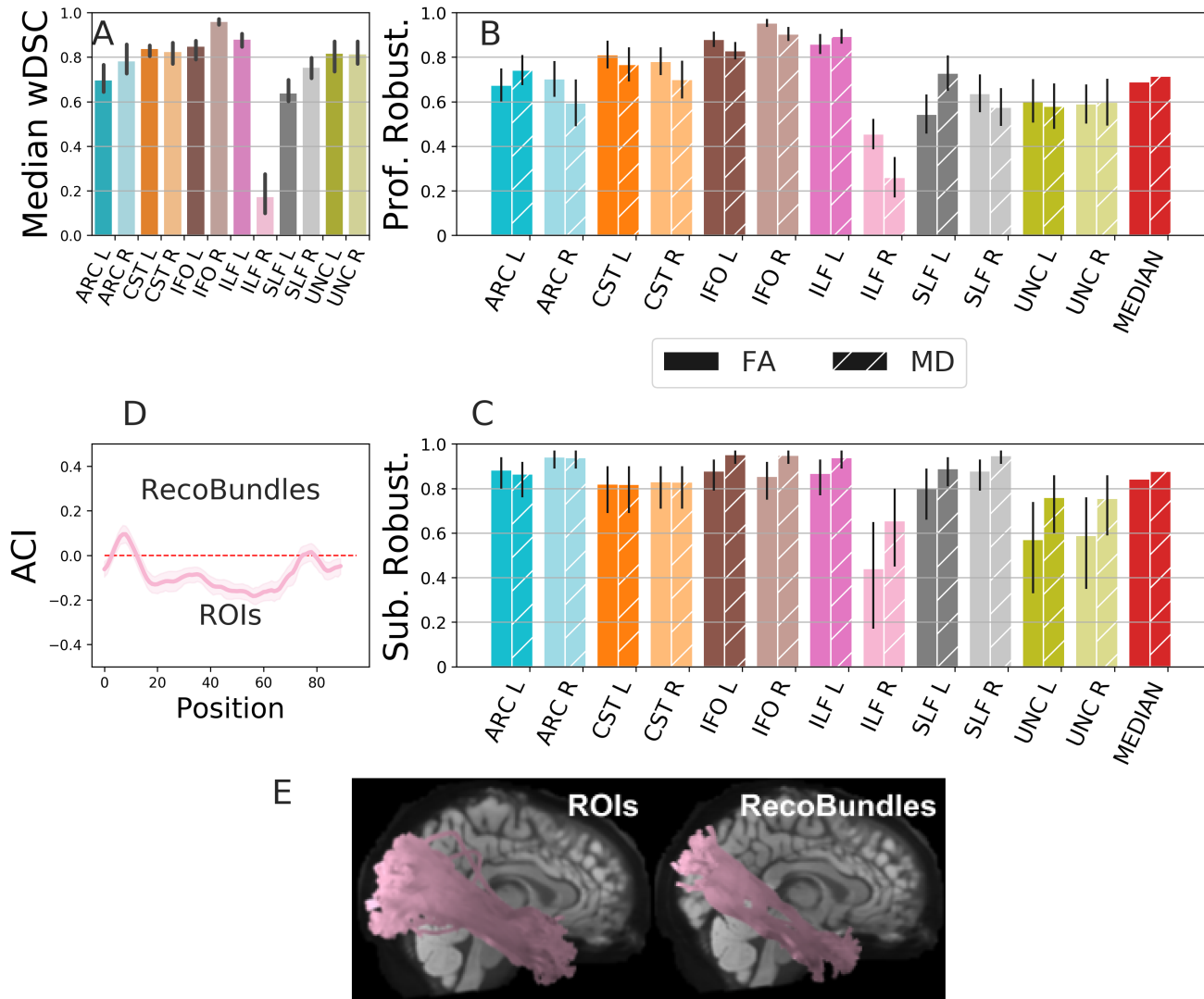


Fig. 5. Recognition algorithm robustness. **A** wDSC. **B** Profile robustness. **C** Subject robustness. Error bars show the 95% confidence interval. **D** The ILF R FA ACIP, where positive ACI indicates RecoBundles found a higher value of FA than the waypoint ROIs approach at that node. **E** shows the ILF R found by each algorithm for an example subject.

362 matter bundles. Here, we analyzed the reliability of tractome- 379
 363 try analysis at several different levels. We analyzed both test- 380
 364 retest reliability of tractometry results and their robustness to 381
 365 changes in analytic details, such as choice of tractography 382
 366 method, bundle recognition algorithm, and software imple- 383
 367 mentation (Fig 6).

368 **Test-retest reliability of tractometry.** Test-retest reliabil- 386
 369 ity (TRR) of tractometry is usually rather high, comparable 387
 370 in some tracts and measurements to the TRR of the measure- 388
 371 ment. In comparing the HCP-TR analysis and UW-PREK 389
 372 analysis, we note that higher measurement reliability goes 390
 373 hand in hand with tractometry reliability. 391

374 In terms of the anatomical definitions of the bundles, quan- 392
 375 tified as the TRR wDSC, we find reliable results in both 393
 376 datasets and with both software implementations and both 394
 377 tractography methods that we tested. With pyAFQ we found 395
 378 a relatively low TRR in the frontal callosal bundle (FA bun- 396

379 dle) in the UW-PREK dataset. This could be due to the sen- 380
 381 sitivity of the definition of this bundle to susceptibility dis- 382
 383 tortion artifacts in the frontal poles of the two hemispheres. 384
 385 This low TRR was not found with mAFQ, suggesting that 386
 387 this low TRR is not a necessary feature of the analysis, and is 388
 389 a potential avenue for improvement to pyAFQ. While the two 390
 391 implementations were created by teams with partial overlap 392
 393 and despite the fact that pyAFQ implementation drew both 394
 395 inspiration as well as specific implementation details from 396
 397 mAFQ, many details of implementation still differ substan- 398
 399 tially. For example, the implementations of tractography al- 400
 401 gorithms are quite different – pyAFQ relies on DIPY (28) 402
 403 for its tractography, while mAFQ uses implementations pro- 404
 405 vided in Vistasoft (47). The two pipelines also use differ- 406
 407 ent registration algorithms, with pyAFQ relying on the SyN 408
 409 algorithm (33), while mAFQ relies on registration methods 410
 411 implemented as part of the Statistical Parametric Mapping 412
 413 (SPM) software (48). These differences may explain the dis-

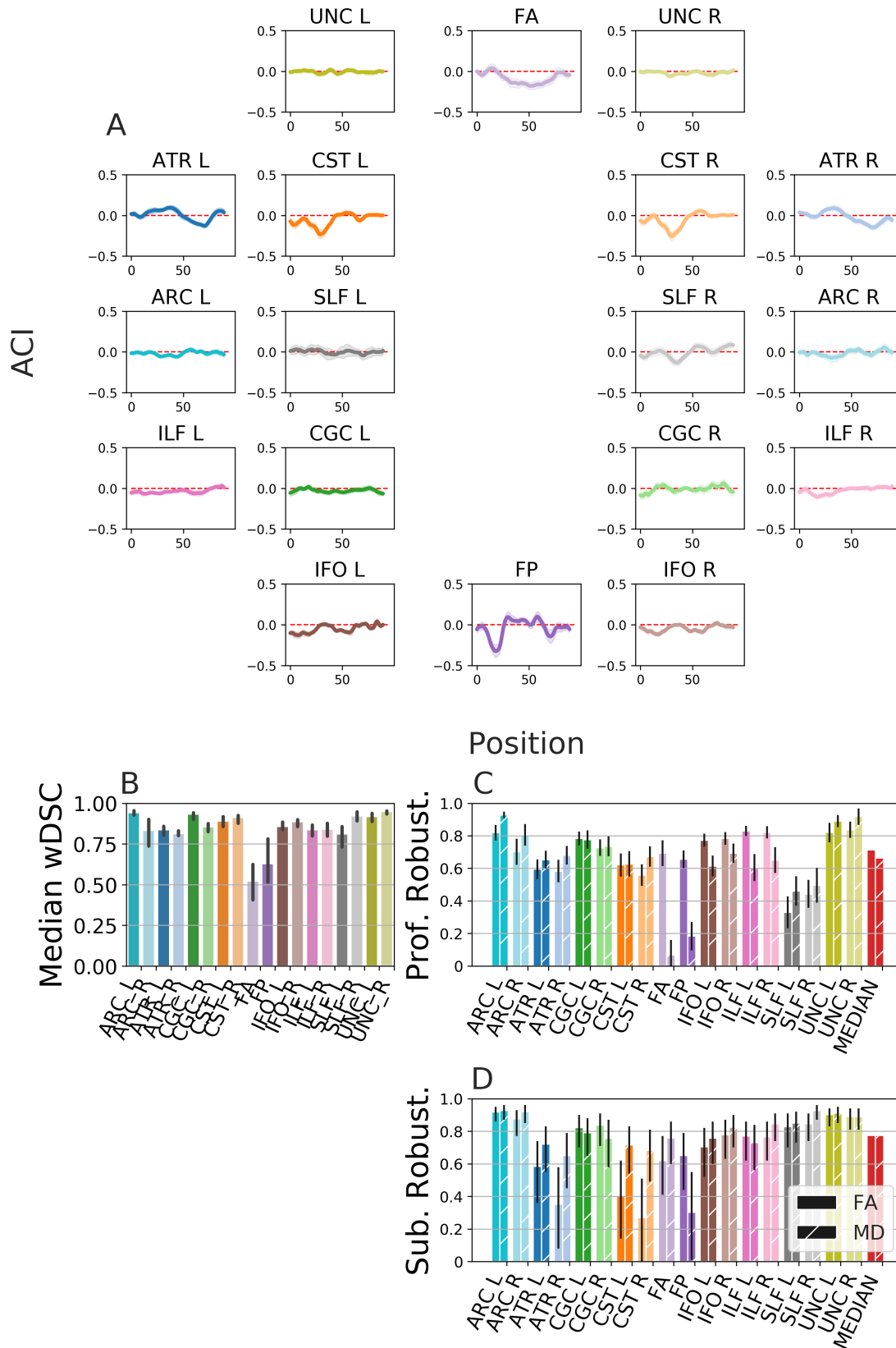


Fig. 6. Robustness between pyAFQ and mAFQ on UW-PREK session # 1 data. **A** ACIP between the FA tract profiles from UW-PREK using pyAFQ and mAFQ. Positive ACI indicates pyAFQ found a higher value than mAFQ at that node. The 95% confidence interval on the mean is shaded. Robustness in wDSC (**B**) bundle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

397 discrepancies observed. 453
398 We also find that TRR is high at the level of profiles within 454
399 subjects and mean tract profiles across subjects. This is gen- 455
400 erally observed in both datasets that we examined, and us- 456
401 ing different analysis methods and software implementations. 457
402 For the UW-PREK dataset, subject TRR tends to be higher 458
403 in mAFQ than in pyAFQ. On the other hand, for the HCP- 459
404 TR dataset, pyAFQ subject TRR tends to be higher than that 460
405 obtained with RTP, which is a fork and extension of mAFQ 461
406 (42, 43). Generally, TRR of FA profiles and also TRR of 462
407 mean FA across subjects tend to be higher than those of MD. 463
408 This could be because the assessment of MD is more sensi- 464
409 tive to partial volume effects. In contrast to FA, MD is also 465
410 not bounded, which means that extreme values at the bound- 466
411 aries of tissue types can have a substantial effect on TRR. 467

412 **Robustness of tractometry.** As highlighted in the recent 468
413 work by Botvinik-Nezer *et al* (25) and in parallel by Schilling 469
414 *et al* (45), inferences from even a single dataset can vary sig- 470
415 nificantly, depending on the decisions and analysis pipelines 471
416 that are used. The analysis approaches used in tractometry 472
417 embody many assumptions made at the different stages of 473
418 analysis: the model of the signal in each individual voxel, the 474
419 manner in which streamlines are generated in tractography, 475
420 the definition of bundles, and the extraction of tract profiles. 476
421 While TRR is important, it does not guard against systematic 477
422 errors in the analysis approach. One way to test model as- 478
423 sumptions and software failures is to create ground truth data 479
424 against which different methods and implementations can be 480
425 tested (13, 49, 50). However, this approach also relies on 481
426 certain assumptions about the mechanisms that generate the 482
427 data that is considered ground truth, making this approach 483
428 more straightforward for some methods than others. Here, 484
429 we instead assessed the robustness of tractometry results to 485
430 perturbations of analytic components, focusing on the mod- 486
431 elling of ODFs in individual voxels and the approach taken 487
432 to bundle recognition. 488
489

433 **Subject robustness remains high despite differences in the 490**
434 **spatial extent of bundles.** We replicated previous findings 491
435 that the definition of major bundles can vary in terms of their 492
436 spatial extent (quantified via wDSC) (13, 37, 40, 45), depend- 493
437 ing on the software implementation or the ODF model used. 494
438 As we show, low wDSC robustness often corresponds to low 495
439 profile robustness, and vice versa (Fig 6B,C, Fig 4A,B, and 496
440 Fig 5A,B). That is, when two algorithms detect bundles with 497
441 small spatial overlap, the shape of the resulting tract profiles 498
442 are also different from each other. However, low wDSC and 499
443 profile robustness does not always translate to low subject 500
444 robustness. Algorithms can detect bundles with low spatial 501
445 overlap and of different shapes yet still agree on the ordering 502
446 of the mean of the profiles, i.e., which subjects have high or 503
447 low FA in a given bundle. A clear example of this is the SLF 504
448 and ARC in Fig 4 (wDSC and profile robustness are low, yet 505
449 subject robustness is very high). This suggests that tractome- 506
450 try can overcome failures in precise delineation of the major 507
451 bundles by averaging tissue properties within the core of the 508
452 white matter. Conversely, important details that are sensitive 509

to these choices may be missed when averaging along the length of the tracts. Moreover, this may also reflect biases in the measurement that cannot be overcome at either stage of the analysis: tractography or bundle recognition.

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 (51), and found high subject-level correlations between the mean tract values of FA and MD for two different pipelines: deterministic tractography using the diffusion tensor model (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 (45). 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, 42) 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 infer-

ences 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 (22). 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). In the UW-PREK dataset, we shared the tract profiles and we provide web-based visualizations using a tool that previously developed for transparent data sharing of tractometry data (52): 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. Using these resources, it should be possible to re-execute our workflows and replicate most of our results (53). 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 (54), which uses a neural network to recognize bundles, or Classifyer (55), which uses a linear classifier), they could do so with the HCP-TR dataset, inspired by our scripts, and the visualization tools in the pyAFQ software.

Future Work. There are many aspects of reliability that could be further explored. We explored robustness with respect to ODF models and bundle recognition algorithms; robustness could also be explored with respect to: data acquisition parameters within the same subject; preprocessing methods; profile extraction method (for example, comparing our current approach with the BUndle ANalytics (BUAN) (56)); and the effects of profile realignment on tract profile reliability (57). Another possibility for teasing apart measurement and tractography effects would be to test profile TRR using the streamline of one scan on the results of the second scan (by registering the streamline themselves, to avoid data interpolation in volume registration). This could tease apart the effects of tractography from the voxel-level models of tissue properties, because it is not necessary that these would be sensitive to the same constraints (e.g., different sensitivity to noise). The methods we demonstrate and resources we

provide in this paper should be useful for anyone wishing to further explore reliability in tractometry.

ACKNOWLEDGEMENTS

This work was supported through grant 1RF1MH121868-01 from the National Institute of Mental Health/The BRAIN Initiative, through grant 5R01EB027585-02 to Eleftherios Garyfallidis (Indiana University) from the National Institute of Biomedical Imaging and Bioengineering and through Azure Cloud Computing Credits for Research & Teaching provided through University of Washington Research Computing and the University of Washington eScience Institute. We are also grateful for support from the Gordon & 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 the eScience Institute 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, Jiaqiang 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, Wilburr 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 Pávid, Martijn Froeling, Anneriet M. Heemskerk, Alexander Leemans, Arnaud Boré, Basile Dinsard, 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 Auria, Oscar Esteban, Alia Lemkaddem, Jean-Philippe Thiran, H. Ertan Cetingul, Benjamin L. Odry, Boris Mailhe, Mariappan S. Nadar, Fabrizio Pizzagalli, Gautam Prasad, Julio E. Villalon-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,

- 645 November 2017. ISSN 2041-1723. doi: 10.1038/s41467-017-01285-x. Number: 1 Pub- 731
646 lisher: Nature Publishing Group. 732
- 647 14. Cibu Thomas, Frank Q Ye, M Okan Irfanoglu, Pooja Modi, Kadharbatha S Saleem, David A 733
648 Leopold, and Carlo Pierpaoli. Anatomical accuracy of brain connections derived from diffu- 734
649 sion MRI tractography is inherently limited. *Proc. Natl. Acad. Sci. U. S. A.*, 111(46):16574– 735
650 16579, November 2014. 736
- 651 15. Kurt G Schilling, Laurent Petit, Francois Rheault, Samuel Remedios, Carlo Pierpaoli, 737
652 Adam W Anderson, Bennett A Landman, and Maxime Descoteaux. Brain connections de- 738
653 rived from diffusion mri tractography can be highly anatomically accurate—if we know where 739
654 white matter pathways start, where they end, and where they do not go. *Brain Structure and* 740
655 *Function*, 225(8):2387–2402, 2020. 741
- 656 16. Ariel Rokem, Jason D. Yeatman, Franco Pestilli, Kendrick N. Kay, Aviv Mezer, Stefan van der 742
657 Walt, and Brian A. Wandell. Evaluating the Accuracy of Diffusion MRI Models in White 743
658 Matter. *PLOS ONE*, 10(4):e0123272, April 2015. ISSN 1932-6203. doi: 10.1371/journal. 744
659 pone.0123272. Publisher: Public Library of Science. 745
- 660 17. Dmitry S Novikov, Valerij G Kiselev, and Sune N Jespersen. On modeling. *Magn. Reson.* 746
661 *Med.*, 79(6):3172–3193, June 2018. 747
- 662 18. Derek K Jones, Adam R Travis, Greg Eden, Carlo Pierpaoli, and Peter J Basser. PASTA: A 748
663 pointwise assessment of streamline tractography attributes. *Magn. Reson. Med.*, 53(6): 749
664 1462–1467, June 2005. 750
- 665 19. John B Colby, Lindsay Soderberg, Catherine Lebel, Ivo D Dinov, Paul M Thompson, and 751
666 Elizabeth R Sowell. Along-tract statistics allow for enhanced tractography analysis. *Neu-* 752
667 *roimage*, 59(4):3227–3242, February 2012. 753
- 668 20. Adam Richie-Halford, Jason Yeatman, Noah Simon, and Ariel Rokem. Multidimensional 754
669 analysis and detection of informative features in diffusion MRI measurements of human 755
670 white matter. *PLoS Computational Biology*, in press, 2021. doi: [https://doi.org/10.1371/](https://doi.org/10.1371/journal.pcbi.1009136) 756
671 [journal.pcbi.1009136](https://doi.org/10.1371/journal.pcbi.1009136). 757
- 672 21. Michael Dayan, Elizabeth Monohan, Sneha Pandya, Amy Kuceyeski, Thanh D Nguyen, 758
673 Ashish Raj, and Susan A Gauthier. Profilmetry: A new statistical framework for the char- 759
674 acterization of white matter pathways, with application to multiple sclerosis. *Hum. Brain* 760
675 *Mapp.*, December 2015. 761
- 676 22. David L Donoho. An invitation to reproducible computational research. *Biostatistics*, 11(3): 762
677 385–388, July 2010. 763
- 678 23. Peter Ivić and Douglas Thain. Reproducibility in Scientific Computing. *ACM Comput. Surv.*, 764
679 51(3):1–36, July 2018. Place: New York, NY, USA Publisher: Association for Computing 765
680 Machinery. 766
- 681 24. The Turing Way Community, Becky Arnold, Louise Bowler, Sarah Gibson, Patricia Herterich, 767
682 Rosie Higman, Anna Krystallidi, Alexander Morley, Martin O'Reilly, and Kirstie Whitaker. *The* 768
683 *Turing Way: A Handbook for Reproducible Data Science*. March 2019. 769
- 684 25. Rotem Botvinnik-Nezer, Felix Holzmeister, Colin F. Camerer, Anna Dreber, Juergen Huber, 770
685 Magnus Johannesson, Michael Kircher, Roni Iwanir, Jeanette A. Mumford, R. Alison Ad- 771
686 cock, Paolo Avesani, Blazej M. Baczkowski, Aahana Bajracharya, Leah Bakst, Sheryl Ball, 772
687 Marco Barilari, Nadège Bault, Derek Beaton, Julia Beintner, Roland G. Benoit, Ruud M. W. J. 773
688 Berkers, Jamil P. Bhanji, Bharat B. Biswal, Sebastian Bobadilla-Suarez, Tiago Bortolini, 774
689 Katherine L. Bottenhorn, Alexander Bowring, Senne Braem, Hayley R. Brooks, Emily G. 775
690 Brudner, Cristian B. Calderon, Julia A. Camilleri, Jaime J. Castellon, Luca Cecchetti, 776
691 Edna C. Cieslik, Zachary J. Cole, Olivier Collignon, Robert W. Cox, William A. Cunningham, 777
692 Stefan Czoschke, Kamalaker Dadi, Charles P. Davis, Alberto De Luca, Mauricio R. Del- 778
693 gado, Lysia Demetriou, Jeffrey B. Dennison, Xin Di, Erin W. Dickie, Ekaterina Dobryakova, 779
694 Claire L. Donnat, Juergen Funkert, Niall W. Duncan, Joke Durnez, Amr Ed, Simon B. Eick- 780
695 hoff, Andrew Erhart, Laura Fontanesi, G. Matthew Fricke, Shiguang Fu, Adriana Galván, 781
696 Remi Gau, Sarah Genon, Tristan Glatar, Enrico Glerean, Jelle J. Goeman, Sergej A. E. 782
697 Golowin, Carlos González-García, Krzysztof J. Gorgolewski, Cheryl L. Grady, Mikella A. 783
698 Green, João F. Guassi Moreira, Olivia Guest, Shabnam Hakimi, J. Paul Hamilton, Roeland 784
699 Hancock, Giacomo Handjaras, Bronson B. Harry, Colin Hawco, Peer Herholz, Gabrielle 785
700 Herman, Stephan Heunis, Felix Hoffstaedter, Jeremy Hogeveen, Susan Holmes, Chuan- 786
701 Peng Hu, Scott A. Huettel, Matthew E. Hughes, Vittorio Iacovella, Alexandru D. Iordan, 787
702 Peder M. Isager, Ayse I. Isik, Andrew Jahn, Matthew R. Johnson, Tom Johnstone, Michael 788
703 J. E. Joseph, Anthony C. Juliano, Joseph W. Kable, Michalis Kassinosopoulos, Cemal Koba, 789
704 Xiang-Zhen Kong, Timothy R. Kosciak, Nuri Erkut Kucukboyaci, Brice A. Kuhl, Sebastian 790
705 Kupek, Angela R. Laird, Claus Lamm, Robert Langner, Nina Lauharatanahirun, Hongmi 791
706 Lee, Sangil Lee, Alexander Leemans, Andrea Leo, Elise Lesage, Flora Li, Monica Y. C. 792
707 Li, Phui Cheng Lim, Evan N. Lintz, Schuyler W. Liphardt, Annabel B. Losecaat Vermeer, 793
708 Bradley C. Love, Michael L. Mack, Norberto Malpica, Theo Marins, Camille Maumet, Kelsey 794
709 McDonald, Joseph T. McGuire, Helena Melero, Adriana S. Méndez Leal, Benjamin Meyer, 795
710 Kristin N. Meyer, Glad Mihai, Georgios D. Mitsis, Jorge Moll, Dylan M. Nielson, Gustav Nil- 796
711 sonne, Michael P. Nottter, Emanuele Olivetti, Adrian I. Onicas, Paolo Papale, Kaustubh R. 797
712 Patil, Jonathan E. Peelle, Alexandre Pérez, Doris Pischke, Jean-Baptiste Poline, Yanina 798
713 Prystauka, Shruti Ray, Patricia A. Reuter-Lorenz, Richard C. Reynolds, Emiliano Ricciardi, 799
714 Jenny R. Rieck, Anais M. Rodriguez-Thompson, Anthony Romyon, Taylor Salo, Giorgio R. 800
715 Samanez-Larkin, Emilio Sanz-Morales, Margaret L. Schlichting, Douglas H. Schultz, Oyang 801
716 Shen, Margaret A. Sheridan, Jennifer A. Silvers, Kenny Skagerlund, Alec Smith, David V. 802
717 Smith, Peter Sokol-Hessner, Simon R. Steinkamp, Sarah M. Tashjian, Bertrand Thieroy, 803
718 John N. Thorp, Gustav Tinghög, Loreen Tisdall, Steven H. Tompson, Claudio Toro-Sherin, 804
719 Juan Jesus Torre Tresols, Leonardo Tozzi, Vuong Truong, Luca Turella, Anna E. van 't Veer, 805
720 Tom Verguts, Jean M. Vettel, Sagana Vijayarajah, Khoi Vo, Matthew B. Wall, Wouter D. 806
721 Weeda, Susanne Weis, David J. White, David Wisniewski, Alba Xifra-Porxas, Emily A. Year- 807
722 ling, Sangsuk Yoon, Rui Yuan, Kenneth S. L. Yuen, Lei Zhang, Xu Zhang, Joshua E. Zosky, 808
723 Thomas E. Nichols, Russell A. Poldrack, and Tom Schonberg. Variability in the analysis of 809
724 a single neuroimaging dataset by many teams. *Nature*, 582(7810):84–88, June 2020. ISSN 810
725 1476-4687. doi: 10.1038/s41586-020-2314-9. Number: 7810 Publisher: Nature Publishing 811
726 Group. 812
- 727 26. Matthew Cieslak, Philip A. Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez 813
728 Adebimpe, Geoffrey K. Aguirre, Danielle S. Bassett, Richard F. Betzel, Josiane Bourque, 814
729 Laura M. Cabral, Christos Davatzikos, John Detre, Eric Earl, Mark A. Elliott, Shreyas 815
730 Fadnavis, Damien A. Fair, Will Foran, Panagiotis Fotiadis, Eleftherios Garyfallidis, Barry 816
- Giesbrecht, Ruben C. Gur, Raquel E. Gur, Max Kelz, Anisha Keshavan, Bart S. Larsen, 817
818 Beatriz Luna, Allyson P. Mackey, Michael Milham, Desmond J. Oathes, Anders Perrone, 819
820 Adam R. Pines, David R. Roalf, Adam Richie-Halford, Ariel Rokem, Valerie J. Snyder, 821
822 Tinashe M. Taper, Ursula A. Tooley, Jean M. Vettel, Jason D. Yeatman, Scott T. Grafton, 823
824 and Theodore D. Satterthwaite. QSIprep: An integrative platform for preprocessing and 825
826 reconstructing diffusion MRI. *bioRxiv*, page 2020.09.04.282269, September 2020. doi: 827
828 10.1101/2020.09.04.282269. Publisher: Cold Spring Harbor Laboratory Section: New Re- 829
830 sults. 831
- 832 27. Hadley Wickham. Tidy data. *J. Stat. Softw.*, 59(10), 2014. 833
- 834 28. Eleftherios Garyfallidis, Matthew Brett, Bagrat Amirbekian, Ariel Rokem, Stefan Van 835
836 Der Walt, Maxime Descoteaux, and Ian Nimmo-Smith. Dipy, a library for the analysis 837
838 of diffusion MRI data. *Frontiers in Neuroinformatics*, 8, 2014. ISSN 1662-5196. doi: 839
840 10.3389/fninf.2014.00008. Publisher: Frontiers. 841
- 842 29. Vladimir Fonov, Alan C. Evans, Kelly Botteron, C. Robert Almli, Robert C. McKinsty, 843
844 D. Louis Collins, and Brain Development Cooperative Group. Unbiased average age- 845
846 appropriate atlases for pediatric studies. *NeuroImage*, 54(1):313–327, January 2011. ISSN 847
848 1095-9572. doi: 10.1016/j.neuroimage.2010.07.033. 849
- 850 30. VS Fonov, AC Evans, RC McKinsty, CR Almli, and DL Collins. Unbiased nonlinear average 851
852 age-appropriate brain templates from birth to adulthood. *NeuroImage*, 47:S102, July 2009. 853
854 ISSN 1053-8119. doi: 10.1016/S1053-8119(09)70884-5. 855
- 856 31. Flavio Dell'Acqua, Luis Lacerda, Marco Catani, and Andrew Simmons. Anisotropic Power 857
858 Maps: A diffusion contrast to reveal low anisotropy tissues from HARDI data. page 1. 859
- 860 32. David Qixiang Chen, Flavio Dell'Acqua, Ariel Rokem, Eleftherios Garyfallidis, David J. 861
862 Hayes, Jidan Zhong, and Mojgan Hodaie. Diffusion weighted image co-registration: In- 863
864 vestigation of best practices. *bioRxiv*, 2019. doi: 10.1101/864108. 865
- 866 33. B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee. Symmetric Diffeomorphic Image 867
868 Registration with Cross-Correlation: Evaluating Automated Labeling of Elderly and Neuro- 869
870 degenerative Brain. *Medical image analysis*, 12(1):26–41, February 2008. ISSN 1361- 871
872 8415. doi: 10.1016/j.media.2007.06.004. 873
- 874 34. Marco Catani, Robert J. Howard, Sinisa Pajevic, and Derek K. Jones. Virtual in vivo inter- 875
876 active dissection of white matter fasciculi in the human brain. *NeuroImage*, 17(1):77–94, 877
878 September 2002. ISSN 1053-8119. doi: 10.1006/nimg.2002.1136. 879
- 880 35. Kegang Hua, JIANGYANG ZHANG, Setsu Wakana, Hangyi Jiang, Xin Li, Daniel S. Re- 881
882 ich, Peter A. Calabresi, James J. Pekar, Peter C. M. van Zijl, and Susumu Mori. Tract 883
884 probability maps in stereotaxic spaces: analyses of white matter anatomy and tract- 885
886 specific quantification. *NeuroImage*, 39(1):336–347, January 2008. ISSN 1053-8119. doi: 887
888 10.1016/j.neuroimage.2007.07.053. 889
- 890 36. Stamatios N Sotiropoulos, Saad Jbabdi, Junqian Xu, Jesper L Andersson, Steen Moeller, 891
892 Edward J Auerbach, Matthew F Glasser, Moises Hernandez, Guillermo Sapiro, Mark Jenk- 893
894 inson, David A Feinberg, Essa Yacoub, Christophe Lenglet, David C Van Essen, Kamil 895
896 Ugurbil, Timothy E J Behrens, and WU-Minn HCP Consortium. Advances in diffusion MRI 897
898 acquisition and processing in the human connectome project. *Neuroimage*, 80:125–143, 899
900 October 2013. doi: 10.1016/j.neuroimage.2013.05.057. 901
- 902 37. Martin Cousineau, Pierre-Marc Jodoin, Eleftherios Garyfallidis, Marc-Alexandre Côté, 903
904 Félix C. Morency, Verena Rozanski, Marilyn Grand'Maison, Barry J. Bedell, and Maxime 905
906 Descoteaux. A test-retest study on Parkinson's PPMI dataset yields statistically significant 907
908 white matter fascicles. *NeuroImage: Clinical*, 16:222, 2017. doi: 10.1016/j.nicl.2017.07.020. 909
910 Publisher: Elsevier. 911
- 912 38. Kenneth O. McGraw and S. P. Wong. Forming inferences about some intraclass correlation 913
914 coefficients. *Psychological Methods*, 1(1):30–46, 1996. ISSN 1939-1463(Electronic),1082- 915
916 989X(Print). doi: 10.1037/1082-989X.1.1.30. Place: US Publisher: American Psychological 917
918 Association. 919
- 920 39. Mariem Boukadi, Karine Marcotte, Christophe Bedetti, Jean-Christophe Houde, Alex 921
922 Desautels, Samuel Deslauriers-Gauthier, Marianne Chapleau, Arnaud Boré, Maxime De- 923
924 scoteaux, and Simona M Brambati. Test-Retest reliability of diffusion measures extracted 925
926 along white matter language fiber bundles using HARDI-Based tractography. *Front. Neuro-* 927
928 *roscl.*, 12:1055, 2018. 929
- 930 40. Mariem Boukadi, Karine Marcotte, Christophe Bedetti, Jean-Christophe Houde, Alex 931
932 Desautels, Samuel Deslauriers-Gauthier, Marianne Chapleau, Arnaud Boré, Maxime De- 933
934 scoteaux, and Simona M. Brambati. Test-Retest Reliability of Diffusion Measures Extracted 935
936 Along White Matter Language Fiber Bundles Using HARDI-Based Tractography. *Frontiers* 937
938 *in Neuroscience*, 12, January 2019. ISSN 1662-4548. doi: 10.3389/fnins.2018.01055. 939
- 940 41. Elizabeth Huber, Rafael Neto Henriques, Julia P. Owen, Ariel Rokem, and Jason D. Yeat- 941
942 man. Applying microstructural models to understand the role of white matter in cognitive 943
944 development. *Developmental Cognitive Neuroscience*, 36, February 2019. ISSN 1878- 945
946 9293. doi: 10.1016/j.dcn.2019.100624. 947
- 948 42. Garikoitz Lerma-Usabiaga, Michael L. Perry, and Brian A. Wandell. Reproducible tract pro- 949
950 files (rtp): from diffusion mri acquisition to publication. *bioRxiv*, page 680173, 2019. 951
- 952 43. Garikoitz Lerma-Usabiaga, Pratik Mukherjee, Michael L. Perry, and Brian A. Wandell. Data- 953
954 science ready, multisite, human diffusion MRI white-matter-tract statistics. *Scientific Data*, 7, 955
956 2020. doi: 10.1038/s41597-020-00760-3. Publisher: Nature Publishing Group. 957
- 958 44. Eleftherios Garyfallidis, Marc-Alexandre Côté, Francois Rheault, Jasmeen Sidhu, Janice 959
960 Hau, Laurent Petit, David Fortin, Stephen Cunnane, and Maxime Descoteaux. Recognition 961
962 of white matter bundles using local and global streamline-based registration and clustering. 963
964 *NeuroImage*, 170:283–295, 2018. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2017.07.015. 965
- 966 45. Kurt G Schilling, Francois Rheault, Laurent Petit, Colin B Hansen, Vishvesh Nath, Fang- 967
968 Cheng Yeh, Gabriel Girard, Muhamed Barakovic, Jonathan Rafael-Patino, Thomas Yu, Elda 969
970 Fisch-Gomez, Marco Pizzolato, Mario Ocampo-Pineda, Simona Schiavi, Erick J Canales- 971
972 Rodríguez, Alessandro Daducci, Cristina Gancziera, Giorgio Innocenti, Jean-Philippe Thi- 973
974 ran, Laura Mancini, Stephen Wastling, Sirio Crocotta, Maria Petracca, Giuseppe Pontillo, 975
976 Matteo Mancini, Sjoerd B Vos, Vejaj N Vakharia, John S Duncan, Helena Melero, Lidia 977
978 Manzanedo, Emilio Sanz-Morales, Ángel Peña-Melán, Fernando Calamante, Arnaud Ad- 979
980 rnaud, Ryan P Cabeen, Laura Korobova, Arthur W Toga, Anupa Ambili Vijayakumari, Drew 981
982 Parker, Ragini Verma, Ahmed Radwan, Stefan Sunaert, Louise Emsell, Alberto De Luca, 983
984 Alexander Leemans, Claude J Bajada, Hamed Haroon, Hojjatollah Azadbakht, Maxime 985
986 Chamberland, Sila Genc, Chantal M W Tax, Ping-Hong Yeh, Rujirutana Srikanhchana, Colin 987

- 817 Mcknight, Joseph Yuan-Mou Yang, Jian Chen, Claire E Kelly, Chun-Hung Yeh, Jerome 903
818 Cochereau, Jerome J Maller, Thomas Welton, Fabien Almaric, Kiran K Seunarine, Chris A 904
819 Clark, Fan Zhang, Nikos Makris, Alexandra Golby, Yogesh Rath, Lauren J O'Donnell, Yi- 905
820 hao Xia, Dogu Baran Aydogan, Yijiang Gong Shi, Francisco Guerreiro Fernandes, Mathijs 906
821 Raemaekers, Shaun Warrington, Stijn Michiels, Alonso Ramirez-Manzanares, Luis Con- 907
822 cha, Ramón Aranda, Mariano Rivera Meraz, Garikoitz Lerma-Usabiaga, Lucas Roitman, 908
823 Lucius S Fekonja, Navona Calarco, Michael Joseph, Hajer Nakua, Aristotle N Voinoskos, 909
824 Philippe Karan, Gabrielle Grenier, Jon Haitz Legarreta, Nagesh Adluru, Veena A Nair, 910
825 Vivek Prabhakaran, Andrew L Alexander, Koji Kamagata, Yuya Saito, Wataru Uchida, 911
826 Christina Andica, Abe Masahiro, Roza G Bayrak, Claudia A Gandini, Egidio D'Angelo, 912
827 Fulvia Palesi, Giovanni Savini, Nicolò Rolandi, Pamela Guevara, Josselin Houenou, Nar- 913
828 ciso López-López, Jean-François Mangin, Cyril Poupon, Claudio Román, Andrea Vázquez, 914
829 Chiara Maffei, Mavilde Arantes, José Paulo Andrade, Susana Maria Silva, Rajikha Raja, 915
830 Vince D Calhoun, Eduardo Caverzas, Simone Sacco, Michael Lauricella, Franco Pestilli, 916
831 Daniel Bullock, Yang Zhan, Edith Brignoni-Perez, Catherine Label, Jess E Reynolds, Igor 917
832 Nestrasil, René Labounek, Christophe Lenglet, Amy Paulson, Stefania Aulicka, Sarah Heil- 918
833 bronner, Katja Heuer, Adam W Anderson, Bennett A Landman, and Maxime Descoteaux. 919
834 Tractography dissection variability: what happens when 42 groups dissect 14 white matter 920
835 bundles on the same dataset? October 2020. doi: 10.1101/2020.10.07.321083. 921
836 46. Gregory Kiar, Yohan Chatelain, Pablo de Oliveira Castro, Eric Petit, Ariel Rokem, Gaël 922
837 Varoquaux, Bratislav Misic, Alan C. Evans, and Tristan Glatard. Numerical Instabilities in 923
838 Analytical Pipelines Lead to Large and Meaningful Variability in Brain Networks. *bioRxiv*, 924
839 page 2020.10.15.341495, October 2020. doi: 10.1101/2020.10.15.341495. Publisher: Cold 925
840 Spring Harbor Laboratory Section: New Results. 926
841 47. Robert F Dougherty, Michal Ben-Shachar, Roland Bammer, Alyssa A Brewer, and Brian A 927
842 Wandell. Functional organization of human occipital-callosal fiber tracts. *Proc. Natl. Acad. Sci. U. S. A.*, 102(20):7350–7355, May 2005. 928
843 48. Karl J. Friston. Statistical Parametric Mapping. In Rolf Kötter, editor, *Neuroscience* 930
844 *Databases: A Practical Guide*, pages 237–250. Springer US, Boston, MA, 2003. ISBN 931
845 978-1-4615-1079-6. doi: 10.1007/978-1-4615-1079-6_16. 932
846 49. Garikoitz Lerma-Usabiaga, Noah Benson, Jonathan Winawer, and Brian A Wandell. A 933
847 validation framework for neuroimaging software: The case of population receptive fields. 934
848 *PLoS Comput. Biol.*, 16(6):e1007924, June 2020. 935
849 50. Peter F Neher, Frederik B Laun, Bram Stieltjes, and Klaus H Maier-Hein. Fiberfox: 936
850 tating the creation of realistic white matter software phantoms. *Magn. Reson. Med.*, 72(5): 937
851 1460–1470, November 2014. 938
852 51. Maya Yablonski, Benjamin Menashe, and Michal Ben-Shachar. A general role for ventral 939
853 white matter pathways in morphological processing: Going beyond reading. *NeuroImage*, 940
854 226:117577, November 2020. 941
855 52. Jason D Yeatman, Adam Richie-Halford, Josh K Smith, Anisha Keshavan, and Ariel Rokem. 942
856 A browser-based tool for visualization and analysis of diffusion MRI data. *Nat. Commun.*, 943
857 (1):940, March 2018. 944
858 53. Satrajit S. Ghosh, Jean-Baptiste Poline, David B. Keator, Yaroslav O. Halchenko, Adam G. 945
859 Thomas, Daniel A. Kessler, and David N. Kennedy. A very simple, re-executable neu- 946
860 roimaging publication. *F1000Research*, 6:124, June 2017. ISSN 2046-1402. doi: 947
861 10.12688/f1000research.10783.2. 948
862 54. Jakob Wasserthal, Peter Neher, and Klaus H Maier-Hein. Tractseg-fast and accurate white 949
863 matter tract segmentation. *NeuroImage*, 183:239–253, 2018. 950
864 55. Giulia Berò, Daniel Bullock, Pietro Astolfi, Siocai Hayashi, Luca Zigiotta, Luciano Annic- 951
865 chiarico, Francesco Corsini, Alessandro De Benedictis, Silvio Sarubbo, Franco Pestilli, et al. 952
866 Classifyer, a robust streamline-based linear classifier for white matter bundle segmenta- 953
867 tion. *BioRxiv*, 2020. 954
868 56. Bramsh Qamar Chandio, Shannon Leigh Risacher, Franco Pestilli, Daniel Bullock, Fang- 955
869 Cheng Yeh, Serge Koudoro, Ariel Rokem, Jaroslaw Harezlak, and Eleftherios Garyfallidis. 956
870 Bundle analytics, a computational framework for investigating the shapes and profiles of 957
871 brain pathways across populations. *Scientific Reports*, 10(1):17149, October 2020. ISSN 958
872 2045-2322. doi: 10.1038/s41598-020-74054-4. Number: 1 Publisher: Nature Publishing 959
873 Group. 960
874 57. Samuel St-Jean, Maxime Chamberland, Max A. Viergever, and Alexander Leemans. Re- 961
875 ducing variability in along-tract analysis with diffusion profile realignment. *NeuroImage*, 199: 962
876 663–679, October 2019. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2019.06.016. 963
877 58. Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cour- 964
878 napeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J 965
879 van der Walt, Matthew Brett, Joshua Wilson, K Jarrod Millman, Nikolay Mayorov, Andrew 966
880 R J Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, Ilhan Polat, Yu Feng, Eric W 967
881 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, 968
882 E A Quintero, Charles R Harris, Anne M Archibald, Antônio H Ribeiro, Fabian Pedregosa, 969
883 Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: fundamental algorithms for sci- 970
884 entific computing in python. *Nat. Methods*, 17(3):261–272, March 2020. 971
885 59. Sphinx, November 2020. 972
886 60. Sphinx-Gallery, November 2020. 973
887 61. Brian Hansen and Sune Nørhøj Jespersen. Data for evaluation of fast kurtosis strategies, 974
888 b-value optimization and exploration of diffusion MRI contrast. *Scientific Data*, 3(1):160072, 975
889 August 2016. ISSN 2052-4463. doi: 10.1038/sdata.2016.72. Number: 1 Publisher: Nature 976
890 Publishing Group. 977
891 62. Matthew Rocklin. Dask: Parallel Computation with Blocked Algorithms and Task Scheduling. 978
892 pages 126–132, Austin, Texas, 2015. doi: 10.25080/Majora-7b98c3ed-013. 979
893 63. Adam Richie-Halford and Ariel Rokem. Cloudknot: A Python Library to Run your Existing 980
894 Code on AWS Batch. *Proceedings of the 17th Python in Science Conference*, pages 8–14, 981
895 2018. doi: 10.25080/Majora-4af1f417-001. Conference Name: Proceedings of the 17th 982
896 Python in Science Conference. 983
897 64. Tom Preston-Werner. tom, January 2021. original-date: 2013-02-24T03:03:57Z. 984
898 65. Tristan Glatard, Gregory Kiar, Tristan Aumentado-Armstrong, Natacha Beck, Pierre Bellec, 985
899 Rémi Bernard, Axel Bonnet, Shawn T. Brown, Sorina Camarasu-Pop, Frédéric Cervenac, 986
900 sky, Samir Das, Rafael Ferreira da Silva, Guillaume Flandin, Pascal Girard, Krzysztof J. 987
901 Gorgolewski, Charles R. G. Guttman, Valérie Hayot-Sasson, Pierre-Olivier Quirion, Pierre 988
902 Rioux, Marc-Étienne Rousseau, and Alan C. Evans. Boutiques: a flexible framework to in-
tegrate command-line applications in computing platforms. *GigaScience*, 7(5), May 2018.
doi: 10.1093/gigascience/giy016. Publisher: Oxford Academic.
66. 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,
Satrajit 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 Hold-
graf, Evgenii Kalenkovich, Isla Staden, Remi Gau, Ariel Rokem, Bertrand Thirion, Dave F.
Kleinschmidt, Erin W. Dickie, John A. Lee, Mathias Goncalves, Matteo Visconti di Olegio-
Castello, Michael Philipp Notter, Pauline Roca, and Ross Blair. PyBIDS: Python tools
for BIDS datasets, July 2020.
67. 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 Oleggio di 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.
68. 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 Schafer, Vanessa Sochat, William Triplett, Jessica A.
Turner, Gaël Varoquaux, and Russell A. Poldrack. The brain imaging data structure, a for-
mat 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.
69. Matthew Brett, Christopher J. Markiewicz, Michael Hanke, Marc-Alexandre Côté, Ben
Cipollini, Paul McCarthy, Dorota Jarecka, Christopher P. Cheng, Yaroslav O. Halchenko,
Michiel Cottar, Eric Larson, Satrajit Ghosh, Demian Wassermann, Stephan Gerhard,
Gregory R. Lee, Hao-Ting Wang, Erik Kastman, Jakub Kaczmarzyk, Roberto Giorditi,
Or Duek, Jonathan Daniel, Ariel Rokem, Cindee Madison, Brendan Moloney, Félix C.
Morency, Mathias Goncalves, Ross Markello, Cameron Riddell, Christopher Burns, Jar-
rod Millman, Alexandre Gramfort, Jaakko Leppäakangas, Anibal Sólón, Jasper J.F. van den
Bosch, Robert D. Vincent, Henry Braun, Krish Subramaniam, Krzysztof J. Gorgolewski,
Pradeep Reddy Raamana, Julian Klug, B. Nolan Nichols, Eric M. Baker, Soichi Hayashi,
Basile Pinsard, Christian Haselgrove, Mark Hymers, Oscar Esteban, Serge Koudoro,
Fernando Pérez-García, Nikolaas 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, Keshi 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 freee83a. nipy/nibabel: 3.2.0, October 2020.
70. Maxime Descoteaux, Rachid Deriche, Thomas R. Knösche, and Alfred Anwander. De-
terministic 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.
71. 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.
72. 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.
73. Ali Tabesh, Jens H. Jensen, Babak A. Ardekani, and Joseph A. Helpert. 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.
74. J.-Donald Tournier, Fernando Calamante, David G. Gadian, and Alan Connelly. Direct es-
timation 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.
75. J.-Donald Tournier, Fernando Calamante, and Alan Connelly. Robust determination of
the fiber 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.
76. Ben Jeurissen, Jacques-Donald Tournier, Thijs Dhollander, Alan Connelly, and Jan Sijbers.
Multi-tissue constrained spherical deconvolution for improved analysis of multi-shell dif-
fusion MRI data. *NeuroImage*, 103:411–426, December 2014. ISSN 1095-9572. doi:
10.1016/j.neuroimage.2014.07.061.
77. 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.
78. 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, Septem-
ber 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.06.005.
79. Marc-Alexandre Côté, Gabriel Girard, Arnaud Boré, Eleftherios Garyfallidis, Jean-
Christophe Houde, and Maxime Descoteaux. Tractometer: towards validation of tractog-
raphy pipelines. *Medical Image Analysis*, 17(7):844–857, October 2013. ISSN 1361-8423.

- 989 doi: 10.1016/j.media.2013.03.009.
- 990 80. Fidel Alfaro-Almagro, Mark Jenkinson, Neal K. Bangerter, Jesper L. R. Andersson, Ludovica
- 991 Griffanti, Gwenaëlle Douaud, Stamatios N. Sotiropoulos, Saad Jbabdi, Moises Hernandez-
- 992 Fernandez, Emmanuel Vallee, Diego Vidaurre, Matthew Webster, Paul McCarthy, Christo-
- 993 pher Rorden, Alessandro Daducci, Daniel C. Alexander, Hui Zhang, Iulius Dragonu, Paul M.
- 994 Matthews, Karla L. Miller, and Stephen M. Smith. Image processing and Quality Control
- 995 for the first 10,000 brain imaging datasets from UK Biobank. *NeuroImage*, 166:400–424,
- 996 February 2018. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2017.10.034.
- 997 81. Karla L. Miller, Fidel Alfaro-Almagro, Neal K. Bangerter, David L. Thomas, Essa Yacoub,
- 998 Junqian Xu, Andreas J. Bartsch, Saad Jbabdi, Stamatios N. Sotiropoulos, Jesper L. R.
- 999 Andersson, Ludovica Griffanti, Gwenaëlle Douaud, Thomas W. Okell, Peter Weale, Iulius
- 1000 Dragonu, Steve Garratt, Sarah Hudson, Rory Collins, Mark Jenkinson, Paul M. Matthews,
- 1001 and Stephen M. Smith. Multimodal population brain imaging in the UK Biobank prospective
- 1002 epidemiological study. *Nature Neuroscience*, 19(11):1523–1536, November 2016. ISSN
- 1003 1546-1726. doi: 10.1038/nn.4393. Number: 11 Publisher: Nature Publishing Group.
- 1004 82. Eleftherios Garyfallidis, Omar Ocegueda, Demian Wassermann, and Maxime Descoteaux.
- 1005 Robust and efficient linear registration of white-matter fascicles in the space of streamlines.
- 1006 *NeuroImage*, 117:124–140, August 2015. ISSN 1053-8119. doi: 10.1016/j.neuroimage.
- 1007 2015.05.016.
- 1008 83. Oscar Esteban, Rastko Ciric, Christopher J. Markiewicz, Yaroslav O. Halchenko, Mathias
- 1009 Goncalves, Satrajit S. Ghosh, Russell A. Poldrack, and Krzysztof J. Gorgolewski. Template-
- 1010 Flow: Standardizing standard 3D spaces in neuroimaging, November 2019.
- 1011 84. Nobuyuki Otsu. A Threshold Selection Method from Gray-Level Histograms. *IEEE Transac-*
- 1012 *tions on Systems, Man, and Cybernetics*, 9(1):62–66, January 1979. ISSN 2168-2909. doi:
- 1013 10.1109/TSMC.1979.4310076. Conference Name: IEEE Transactions on Systems, Man,
- 1014 and Cybernetics.
- 1015 85. David Chen, Flavio Dell'Acqua, Ariel Rokem, Eleftherios Garyfallidis, Jidan Zhong, and
- 1016 Mojgan Hodaie. Diffusion Weighted Image Co-registration: Investigation of Best Practices.
- 1017 December 2019. doi: 10.1101/864108.
- 1018 86. Setsu Wakana, Arvind Caprihan, Martina M. Panzenboeck, James H. Fallon, Michele Perry,
- 1019 Randy L. Gollub, Kegang Hua, Jiangyang Zhang, Hangyi Jiang, Prachi Dubey, Ari Blitz,
- 1020 Peter van Zijl, and Susumu Mori. Reproducibility of Quantitative Tractography Methods
- 1021 Applied to Cerebral White Matter. *NeuroImage*, 36(3):630–644, July 2007. ISSN 1053-
- 1022 8119. doi: 10.1016/j.neuroimage.2007.02.049.
- 1023 87. N. Tzourio-Mazoyer, B. Landeau, D. Papathanassiou, F. Crivello, O. Etard, N. Delcroix,
- 1024 B. Mazoyer, and M. Joliot. Automated anatomical labeling of activations in SPM using a
- 1025 macroscopic anatomical parcellation of the MNI MRI single-subject brain. *NeuroImage*, 15
- 1026 (1):273–289, January 2002. ISSN 1053-8119. doi: 10.1006/nimg.2001.0978.
- 1027 88. C. Bradford Barber, David P. Dobkin, and Hannu Huhdanpaa. The quickhull algorithm for
- 1028 convex hulls. *ACM Transactions on Mathematical Software*, 22(4):469–483, December
- 1029 1996. ISSN 0098-3500, 1557-7295. doi: 10.1145/235815.235821.
- 1030 89. FURY, October 2020.
- 1031 90. Plotly Python Graphing Library, October 2020.
- 1032 91. David C. Van Essen, Stephen M. Smith, Deanna M. Barch, Timothy E. J. Behrens, Essa
- 1033 Yacoub, and Kamil Ugurbil. The WU-Minn Human Connectome Project: An overview. *Neu-*
- 1034 *roImage*, 80:62–79, October 2013. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2013.05.
- 1035 041.
- 1036 92. Lin-Ching Chang, Derek K. Jones, and Carlo Pierpaoli. RESTORE: robust estimation of
- 1037 tensors by outlier rejection. *Magnetic Resonance in Medicine*, 53(5):1088–1095, May 2005.
- 1038 ISSN 0740-3194. doi: 10.1002/mrm.20426.
- 1039 93. J-Donald Tournier, Robert Smith, David Raffelt, Rami Tabbara, Thijs Dhollander, Maximilian
- 1040 Pietsch, Daan Christiaens, Ben Jeurissen, Chun-Hung Yeh, and Alan Connelly. MRtrix3: A
- 1041 fast, flexible and open software framework for medical image processing and visualisation.
- 1042 *NeuroImage*, 202:116137, November 2019.
- 1043 94. Lee R. Dice. Measures of the Amount of Ecologic Association Between Species. *Ecol-*
- 1044 *ogy*, 26(3):297–302, 1945. ISSN 00129658, 19399170. doi: 10.2307/1932409. Publisher:
- 1045 Ecological Society of America.
- 1046 95. Matthew Cieslak, Philip A Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez
- 1047 Adebimpe, Geoffrey K Aguirre, Danielle S Bassett, Richard F Betzel, Josiane Bourque, et al.
- 1048 Qsirep: An integrative platform for preprocessing and reconstructing diffusion mri. *bioRxiv*,
- 1049 2020.
- 1050 96. J-Donald Tournier, Robert Smith, David Raffelt, Rami Tabbara, Thijs Dhollander, Maximilian
- 1051 Pietsch, Daan Christiaens, Ben Jeurissen, Chun-Hung Yeh, and Alan Connelly. Mrtrix3: A
- 1052 fast, flexible and open software framework for medical image processing and visualisation.
- 1053 *NeuroImage*, 202:116137, 2019.
- 1054 97. Lindsay M Alexander, Jasmine Escalera, Lei Ai, Charissa Andreotti, Karina Febre, Alexan-
- 1055 der Mangone, Natan Vega-Potler, Nicolas Langer, Alexis Alexander, Meagan Kovacs, Shan-
- 1056 non Litke, Bridget O'Hagan, Jennifer Andersen, Batya Bronstein, Anastasia Bui, Marijayne
- 1057 Bushey, Henry Butler, Victoria Castagna, Nicolas Camacho, Elisha Chan, Danielle Citera,
- 1058 Jon Clucas, Samantha Cohen, Sarah Dufek, Megan Eaves, Brian Fradera, Judith Gardner,
- 1059 Natalie Grant-Villegas, Gabriella Green, Camille Gregory, Emily Hart, Shana Harris, Megan
- 1060 Horton, Danielle Kahn, Katherine Kabotyanski, Bernard Karmel, Simon P Kelly, Kayla Klein-
- 1061 man, Bonhwang Koo, Eliza Kramer, Elizabeth Lennon, Catherine Lord, Ginny Mantello, Amy
- 1062 Margolis, Kathleen R Merikangas, Judith Milham, Giuseppe Minniti, Rebecca Neuhaus,
- 1063 Alexandra Levine, Yael Osman, Lucas C Parra, Ken R Pugh, Amy Racanello, Anita Res-
- 1064 strepo, Tian Saltzman, Batya Septimus, Russell Tobe, Rachel Waltz, Anna Williams, Anna
- 1065 Yeo, Francisco X Castellanos, Arno Klein, Tomas Paus, Bennett L Leventhal, R Cameron
- 1066 Craddock, Harold S Koplewicz, and Michael P Milham. An open resource for transdiagnostic
- 1067 research in pediatric mental health and learning disorders. *Sci Data*, 4:170181, December
- 1068 2017.
- 1069 98. Martin Lindquist. Neuroimaging results altered by varying analysis pipelines. *Nature*, 582
- 1070 (7810):36–37, June 2020. doi: 10.1038/d41586-020-01282-z. Number: 7810 Publisher:
- 1071 Nature Publishing Group.
- 1072 99. Robert F Dougherty, Michal Ben-Shachar, Gayle K Deutsch, Arvel Hernandez, Glenn R
- 1073 Fox, and Brian A Wandell. Temporal-callosal pathway diffusivity predicts phonological skills
- 1074 in children. *Proc. Natl. Acad. Sci. U. S. A.*, 104(20):8556–8561, May 2007.

1075 **Supplementary Methods**

1076 **Automated Fiber Quantification in Python (pyAFQ).** Inspired by a previous MATLAB implementation (9), We developed
1077 a software library that automates dMRI-based tractometry analysis. The library is called pyAFQ (Python Automated Fiber
1078 Quantification), and it is implemented as open-source software here: <https://github.com/yeatmanlab/pyAFQ>. The
1079 software is developed under the permissive OSI-approved BSD license. It allows users to specify the methods and param-
1080 eters they want to use for tractometry. pyAFQ uses many components of the scientific Python ecosystem (58). In particular,
1081 it relies heavily on implementations of algorithms for diffusion reconstruction, orientation determination, tractography and
1082 image registration implemented in Diffusion Imaging in Python (DIPY), an open-source, Python library for computational neu-
1083 roanatomy (28). The pyAFQ software implements extensive documentation with Sphinx (59), including a gallery of executable
1084 examples, implemented using Sphinx Gallery (60). Unit testing is implemented using pytest, with continuous integration im-
1085 plemented to test proposed changes to the library, as well as longer nightly tests that check that pipelines of operations are
1086 not adversely affected by changes that are introduced in developing the software. pyAFQ's test suite uses the HARDI data
1087 collected for (16), CFIN (61), and data from the Human Connectome Project. pyAFQ can be parallelized across subjects and
1088 sessions using dask (62). The analysis performed in this paper primarily used pyAFQ run using Cloudknot (63) on Amazon
1089 Web Services (AWS).

1090 There are many ways to analyze dMRI data and to estimate tractometry-based tract-profiles. For example, many different
1091 models are used to determine the directions of tracking within each voxel and to connect different voxels with a variety of
1092 tractography algorithms. Similarly, different models can be used to determine the tissue properties within a voxel. However, it
1093 is hard to determine which methods to use, because different methods may be appropriate for different datasets, depending
1094 on their characteristics: the measurements conducted, the signal to noise ratio (SNR) of the data and so forth. Software to support
1095 analysis of a variety of datasets should make it easy to use many different methods and to compare results between methods.
1096 All of the choices the user can make in each of the steps of pyAFQ are delineated below and summarized in Fig. S2. The
1097 software implements a library with an object-oriented application programming interface (API), as well as a command-line
1098 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
1099 the ability to choose which algorithms and parameters to use. For users unfamiliar with python, pyAFQ has a command line
1100 interface (CLI) which uses a configuration file written in TOML (64). pyAFQ also has a Boutiques configuration file and can
1101 be executed using Boutiques (65).

1102 **Locating and mapping data (BIDS).** The first step in analysis is to find the files that the software will use. pyAFQ relies on
1103 pyBIDS (66, 67) to query data that is provided in the BIDS format (68). It looks for dMRI, b-value, and b-vector files stored
1104 in standard formats (see <https://yeatmanlab.github.io/pyAFQ/usage/data.html> for details). Additionally,
1105 the user can provide files from other processing pipelines to be used as a brain mask during registration or as start or stop
1106 masks during tractography, as well as completed tractography results. We typically use the Nibabel software library to interact
1107 with neuroimaging files (69). Following the BIDS standard, the outputs of pyAFQ are put in the BIDS derivatives folder, in a
1108 pipeline directory labelled as "afq". The derivative BIDS format follows as much as possible the draft implementation of the
1109 BIDS derivatives for dMRI data.

1110 **Tractography.** There are several methods for computational tractography. The pyAFQ software exposes many of these as op-
1111 tions. It allows users to choose from multiple fiber orientation distribution functions (70) that determine the direction of tracking
1112 in each step of the process: based on Diffusion Tensor Imaging (DTI) (71, 72), Diffusion Kurtosis Imaging (DKI) (73), Con-
1113 strained Spherical Deconvolution (CSD) (74, 75), and Multi-Shell Multi-Tissue Constrained Spherical Deconvolution (MSMT-
1114 CSD) (76). Deterministic and probabilistic tractography algorithms can be used and stopping criteria can be implemented for
1115 particle filtering tractography, using the continuous map criterion (77) or anatomically-constrained tractography (78). The de-
1116 fault tractography setting uses DTI, deterministic direction finding, a max turning angle per step of 30° , one seed per voxel, and
1117 retains only streamlines between 10 and 1000mm long. Many of our tractography defaults are inspired by the results of (79)
1118 and (9). The default seed and stop masks are created by thresholding FA at 0.2. All of these parameters can be customized
1119 using pyAFQ's API or CLI.

1120 **Template registration.** The user can specify their own template and subject image to register, however pyAFQ also provides four
1121 builtin options: register subject non-diffusion weighted image (also known as b0) to the Montreal Neurological Institute (MNI)
1122 T2 template (29, 30); register subject FA to a group mean fractional anisotropy (FA) template from the UK Biobank (80, 81);
1123 register a subject's anisotropic power map (APM) (31, 32) to the MNI T1 template; and register subject streamlines to the 16
1124 bundles human connectome project (HCP) atlas (7) using streamline registration (SLR) (82). The first three of these builtin
1125 techniques use the nonlinear Symmetric Diffeomorphic Registration (SyN) (33) after an optional linear preregistration, both
1126 implemented in DIPY. pyAFQ uses Templateflow (83) to get MNI T1/T2 templates for registration. The default registration
1127 behavior is to consider all b-values under 50 to be b0, mask the subject's APM using DIPY's median_otsu image recognition
1128 algorithm (84) on the subject b0, and register the masked power map to the masked MNI T1 template. Per default, we chose to

1129 use the APM for registration based on previous findings that show this is a good choice (85) and based on our own experience.
1130 All of these parameters can be customized using pyAFQ's API and CLI.

1131 **Bundle recognition and cleaning.** To identify the streamlines that best represent a particular anatomical pathway, we perform
1132 bundle recognition. The default behavior is to perform the initial classification using probability maps, and then segment with
1133 waypoint ROIs defined in (86), then filter the classified streamlines by their termination locations, using the AAL atlas (87),
1134 where streamlines must be within 4mm of the expected endpoint region. Waypoint ROIs are moved into the subject space and
1135 then patched up using the Quickhull Algorithm (88). There is also an option, turned off by default, to clip streamline edges at
1136 the ROIs (86).

1137 In addition to the waypoint-based recognition described above, pyAFQ also allows the user to choose to use a streamline atlas
1138 based bundle recognition method, called RecoBundles (44). Parameters for either algorithm can be customized using pyAFQ's
1139 API and CLI.

1140 After recognition, cleaning is performed based on the Mahalanobis distance of each streamline from the mean in each node.
1141 This process was originally described in (9). By default, pyAFQ resamples streamlines to 100 points (nodes) and performs
1142 5 rounds of cleaning with a distance threshold of 5 standard deviations from the mean of the node coordinates at each point,
1143 and a length threshold of 4 standard deviations from the mean length. Cleaning is also stopped if a bundle has less than 20
1144 streamlines. All of these parameters can be customized using pyAFQ's API and CLI.

1145 **Tract Profile Extraction.** After cleaning, pyAFQ computes and visualizes tract profiles. The mean profile (called a "tract profile")
1146 is calculated using the same Mahalanobis distance-based weighting strategy as in Yeatman et al. (9), implemented in DIPY.
1147 Visualization can be performed using one of two backends: fury (89) or plotly (90), which create either animated gifs or
1148 interactive html files respectively. Visualizations are created for the whole brain tractometry and for each individual bundle.

1149 **Data.** We measured the reliability of tractometry using two datasets with contrasting characteristics.

1150 **Human Connectome Project (HCP-TR).** The WU-Minn Human Connectome Project (HCP) (91) includes measurements of
1151 diffusion MRI data from almost all of the 1,200 participants. Here, we focus our analysis on a subset of these subjects for
1152 which test-retest data are available. We refer to this data as HCP-TR. This dataset contains dMRI data from 44 individuals.
1153 This represents a relatively high-quality, high-resolution dataset, with multiple diffusion directions and multiple b-values. The
1154 acquisition parameters of HCP-TR are described in detail elsewhere (36). We used data that had been preprocessed through the
1155 HCP pipelines, as provided through the AWS Open Data program (<https://registry.opendata.aws/hcp-openaccess/>).

1156 **University of Washington Pre-K (UW-PREK).** Two measurements were conducted in each participant 1 day apart. These were
1157 acquired with 32 directions, $b=1,500 \text{ s/mm}^2$, 2 mm^3 isotropic resolution, $TR/TE=7200/83 \text{ msec}$. Data were preprocessed using
1158 FSL for eddy current, motion correction, and susceptibility distortion correction. Analysis using the mAFQ was conducted as
1159 previously described (9). We converted UW-PREK to BIDS format (68) for input into pyAFQ's API.

1160 We attempted to configure pyAFQ to most closely match the mAFQ configuration. We used robust estimation of tensors by
1161 outlier rejection (RESTORE) (92) to fit the DTI model. In tractography, we used 160,000 seeds randomly distributed wherever
1162 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
1163 compared its reproducibility to the results on the same datasets with mAFQ. We also compared the robustness of the results
1164 between the pyAFQ and mAFQ algorithms on the pre-session data only.

1165 **Configurations.** For all configurations, we used the Freesurfer brain segmentation provided by HCP to calculate a permissive
1166 brain mask, with all portions of the image not labelled as 0, considered part of the brain. The brain mask is used when fitting
1167 the ODF models. We compared the TRR of each configuration, as well as the robustness of the results across configurations.
1168 We also compared the TRR of these configurations to the TRR of results published by Lerma-Usabiaga and colleagues (43),
1169 denoted RTP.

1170 **DTI Configuration.** In addition to the three configurations enumerated in the present paper, we processed HCP-TR with a fourth
1171 configuration. We used only measurements with b-values between 990 and 1010 s/mm^2 . We used DTI as the ODF model for
1172 tractography and profile extraction. We compared this configuration to RTP in 3D,E. We also analysed DTI for robustness and
1173 found its results to be nearly identical to DKI.

1174 **RecoBundles Configuration.** One of the configurations we ran on the HCP-TR data used RecoBundles (8). pyAFQ provides
1175 programmatic access to two atlases, one being the full 80 bundles human connectome project (HCP) atlas (7), and other being
1176 a 16 bundle subset of that atlas. We ran RecoBundles on HCP-TR using the full 80 bundles atlas. We use the following
1177 RecoBundles parameter configuration: a model cluster threshold of 1.25, a reduction threshold of 25, no refinement, a pruning
1178 threshold of 12, local streamline-based linear registration on with an asymmetric metric. We used this configuration for all 80
1179 bundles. Multi-shell data and the DKI ODF model were used. We used nonlinear symmetric diffeomorphic registration and a
1180 brain mask based on the HCP-provided segmentation.

1181 **RTP.** As a point of comparison, we used an open dataset of HCP-TR derivatives that was published by Lerma-Usabiaga and
1182 colleagues (43). They processed HCP-TR using the Reproducible Tract Profiles (RTP) pipeline (42). This pipeline is a full
1183 end-to-end pipeline and system for deployment of analysis that receives as input raw MRI data as acquired on the scanner.
1184 While it applies different preprocessing steps and uses different tractography algorithms than mAFQ, relying on MRTRIX for
1185 many of these steps (93), the bundle recognition steps closely resemble the ones used in mAFQ, relying on functions that stem
1186 from the same MATLAB codebase as mAFQ. The end result of RTP are tract profiles in an easy-to-use and data-science ready
1187 JSON format. We denote their results as RTP and compare them to the HCP-TR results computed with pyAFQ.

1188 **Measures of reliability.** pyAFQ gives the user the choice of which underlying algorithms to use when performing tractometry,
1189 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-
1190 retest data. The analyses we selected represent only a small subset of the possible configurations of pyAFQ. However, because
1191 the software is freely available and easily configurable with the API or CLI, it would be straightforward to test other analyses. To
1192 compare the results on test-retest data (TRR) and compare results across analyses (robustness), we use four different measures
1193 of reliability. Each one of these measures emphasizes different aspects of reliability.

1194 **Weighted Dice similarity coefficient (wDSC).** The anatomical reliability of bundle recognition solutions is assessed by compar-
1195 ing their spatial overlap in the white matter volume. First, for every voxel in the white matter, we count the number of
1196 streamlines that pass through that voxel for a given bundle, then divide by the total number of streamlines in that bundle. This
1197 creates what we call a streamline density map (28). We could compare streamline density maps using a Dice similarity coeffi-
1198 cient (94), but that would require applying a threshold to the density maps, and could give a few streamlines a large influence
1199 on the calculation. Instead, we use the weighted Dice similarity coefficient (wDSC) (37):

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

1200 where v is a voxel index, $W_{i,v}$ is the streamline density for a bundle i in voxel v , and v^i are voxels where the two bundles i and
1201 j intersect. wDSC provides a measure of the reliability in the spatial extent of bundles, in a manner that is independent from
1202 the assessment of tract profiles.

1203 **Adjusted contrast index profile (ACIP).** We use an adjusted contrast index to directly compare the values of individual nodes in
1204 the tract profiles in different measurements. For two values (V_1, V_2) in different profiles, the adjusted contrast index (ACI)
1205 is calculated using Eq (2).

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

1206 We multiply by 2 to make the contrast index have comparable values to fractional difference. In contrast to fractional difference,
1207 however, the ACI does not require one of the variables to be a reference, and $ACI(V1, V2) = -ACI(V2, V1)$. Calculating and
1208 then plotting the ACI for each point between two profiles highlights the differences between profiles, producing the adjusted
1209 contrast index profile (ACIP). ACIP emphasizes discrepancies in estimates along the length of the tract in a manner that does
1210 not depend on the scale of the measurement (e.g., the different scales of FA and MD).

1211 **Supplementary Discussion of pyAFQ**

1212 **pyAFQ is embedded in an ecosystem of tools for reproducible neuroimaging.** The wider ecosystem of tools and standards
1213 surrounding pyAFQ is shown in Fig. S6. Each tool has its own place in the ecosystem. We rely heavily on implementations
1214 of dMRI analysis algorithms implemented in DIPY (28). Reproducibility and interoperability are also facilitated by relying on
1215 the BIDS format (68) and the pyBIDS software (66, 67). Requiring a BIDS-like input makes integration with other software in
1216 the ecosystem easier. For example, it is fairly straightforward to use the outputs of BIDS-compatible preprocessing pipelines,
1217 such as qsiprep (95), as inputs to pyAFQ. Furthermore, the modularity of the pyAFQ pipeline means that outputs of other
1218 tractography software (e.g., MRTRIX (96)) can be used as inputs to bundle recognition, with BIDS filters as the metadata that
1219 allows finding and incorporating through the right data.

1220 Cloud-based processing is going to be more important as large datasets are processed. pyAFQ does not depend on proprietary
1221 software and can be scaled to large datasets using cloud computing platforms. In this paper, we used Cloudknot (63) to scale
1222 pyAFQ across subjects and methods on AWS. However, because pyAFQ is a Python package, it can easily be run on any cloud
1223 computing platform. Computing in the public cloud also supports reproducible research, as computations conducted on the
1224 public cloud are perfectly portable to other users of the software. Our software is written with that in mind, including functions
1225 that know how to easily access datasets that are already stored in the cloud (e.g., HCP and Healthy Brain Network (97) datasets).
1226 We know that one of the most important ways in which users can diagnose whether processing worked as expected is by visually
1227 inspecting the results. Thus, we provide several different visualization methods, relying on the VTK-derived FURY library, or
1228 on browser-friendly visualizations with Plotly. pyAFQ outputs are also fully compatible with AFQ-Browser, a browser-based
1229 tool for interactive visualization and exploration of tractometry results (52).

1230 Finally, beyond visualization and summary of the results, and tools for analysis of reliability presented in this work, pyAFQ
1231 does not provide a substantial set of tools for statistical analysis of tractometry results. Instead, the outputs of pyAFQ are
1232 provided as “tidy” CSV tables (27). This means that it is compatible as inputs to the AFQ Insight tool for statistical analysis
1233 (20), but also amenable to many other statistical analysis approaches. This output should facilitate interdisciplinary use of
1234 dMRI data, as it is provided in a format that is widely used in statistics and machine learning.

1235 **pyAFQ is extensible.** In general, variability in results would be reduced with a standard pipeline that could be used across all
1236 studies and datasets. However, as noted by Lindquist, “studies tend to be too varied for one pipeline to always be appropri-
1237 ate” (98). This is particularly true as new measurement techniques, new processing methods and new analysis approaches for
1238 dMRI are evolving. Therefore, the pyAFQ pipeline was designed to be flexible, making it easier to reproduce results, while
1239 providing researchers with many choices for the appropriate analysis, depending on their data and questions. pyAFQ allows the
1240 user to make many decisions (Fig S2), and all of those decisions can be encoded in a configuration file. That configuration file
1241 can be used to reproduce the same analysis pipeline given the same version of pyAFQ is used. By providing the configuration
1242 file or the arguments passed to the main API, one can clearly satisfy the requirement for a re-executable workflow outlined
1243 in (53).

1244 To extend to new bundles, pyAFQ allows users to define new queries that recognize bundles that are not part of the set of 18
1245 detected by the original mAFQ software. For a simple example, we use a set of alternative waypoint ROIs to detect different
1246 portions of the corpus callosum (99) (Fig S7A). These alternative ROIs are included in pyAFQ but not used by default. In more
1247 complicated example, another set of ROIs is used to recognize the location of the optic radiations (OR; Fig S7). Because these
1248 are relatively small and winding, their delineation requires additional components: it requires several waypoint ROIs used not
1249 only as inclusion criteria, but also as exclusion criteria, and it requires delineation of endpoints in the cortex that are not part of
1250 the AAL atlas, which is used in the standard set of bundles. It also requires oversampling of streamlines, so in order to obtain
1251 a proper definition of the OR, tractography is configured to use 125 seeds per voxel (instead of the default 8). All of these
1252 components can be integrated into calls to the software API, without needing to change any of its internals. This includes any
1253 custom waypoint ROIs, inclusive or exclusive, as well as probability maps, endpoint locations, and whether the bundle crosses
1254 the midline.

1255 **Supplementary Figures and Tables**

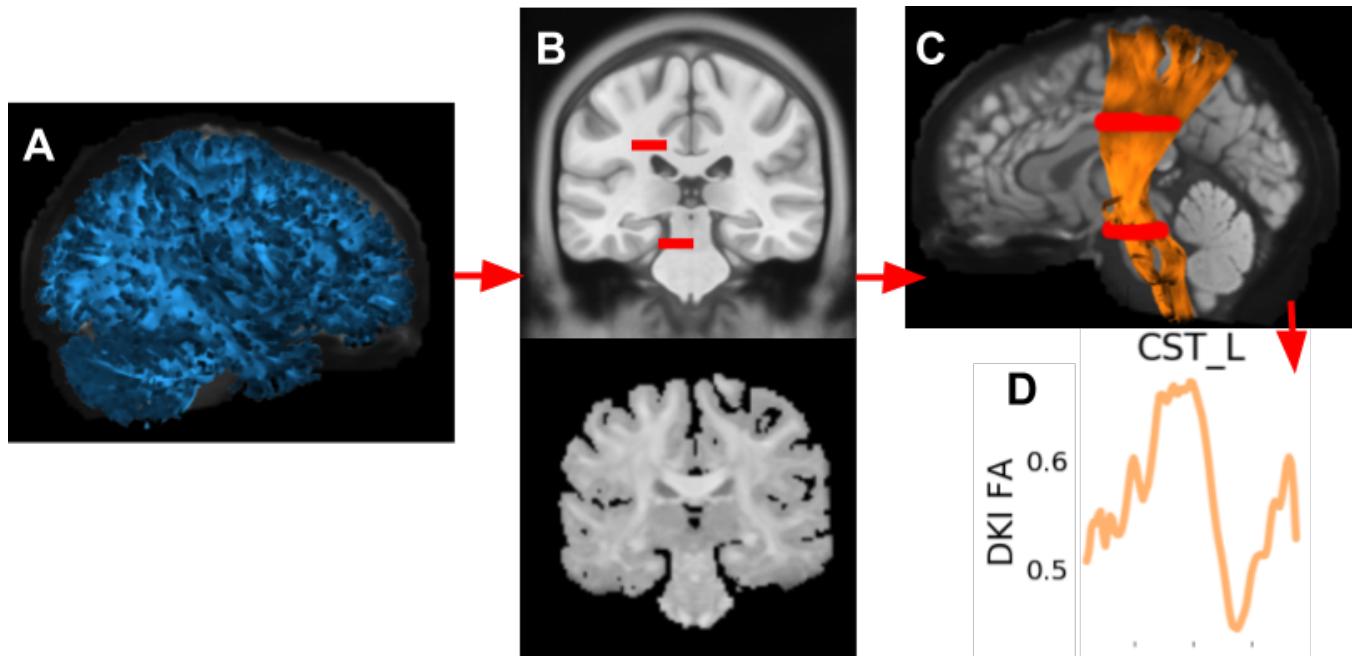


Fig. S1. The stages of tractometry. **A** Computational tractography generates streamlines estimating the trajectories of white matter connections. **B** An anatomical template is registered to each subject's individual brain. Here, in a mid-coronal view, the MNI T1-weighted template (29, 30), shown 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) (31) 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.

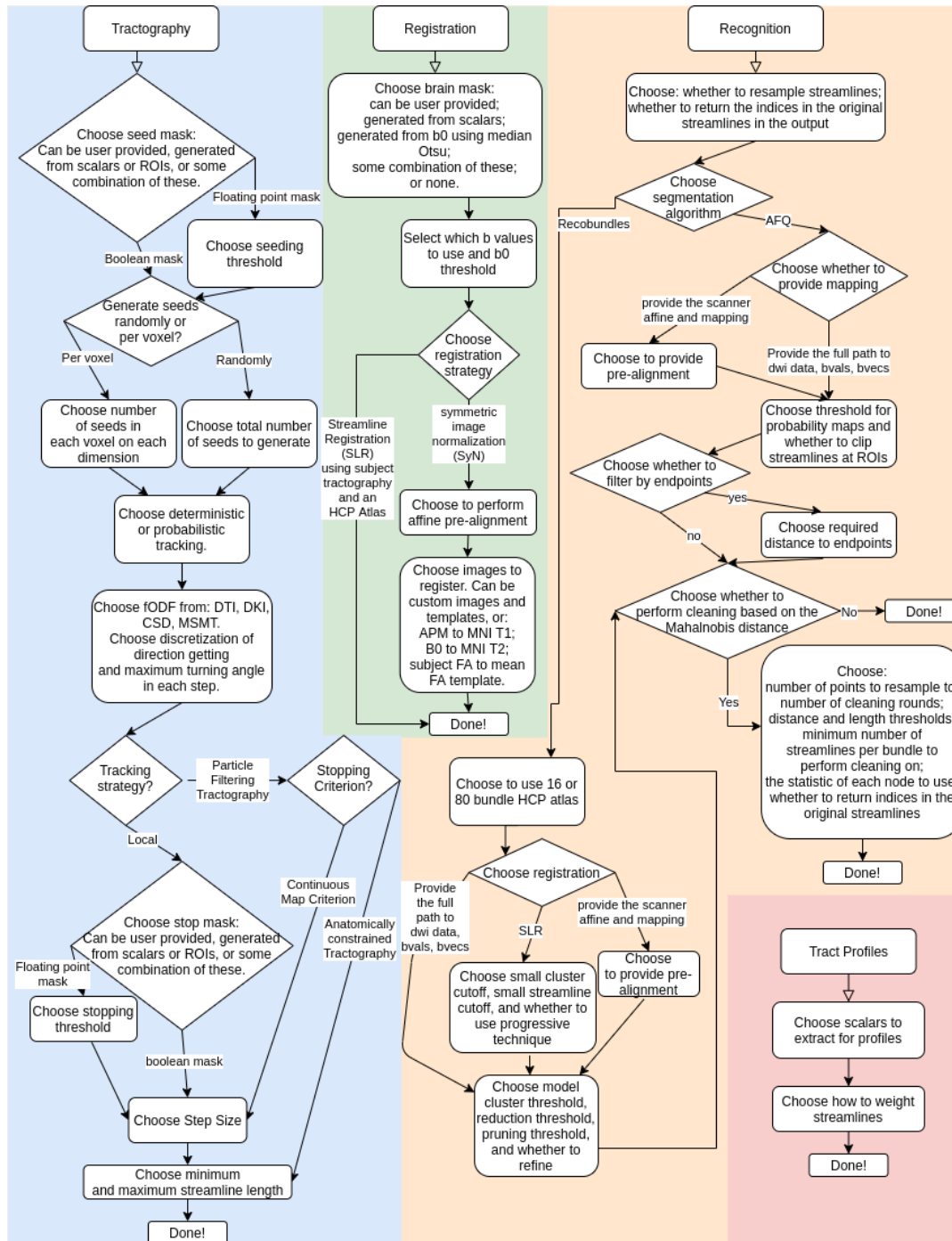


Fig. S2. Choices the user can make for how to run pyAFQ. The colors represent different steps of tractometry. Tractography is shaded blue, registration is shaded green, recognition is shaded orange, and tract profiles is shaded red. Every rounded box and diamond contains one or more choices, except for 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.

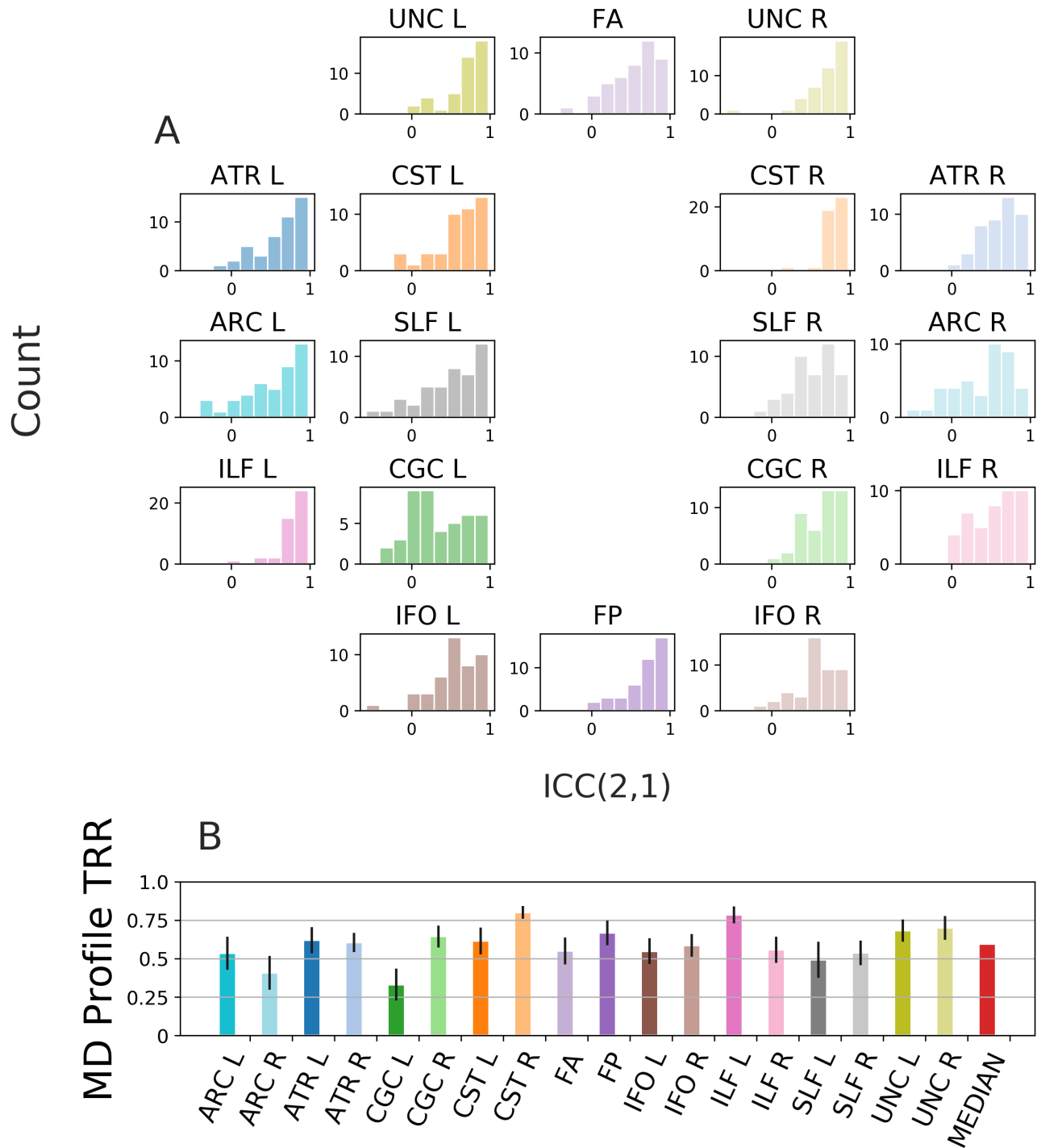


Fig. S4. MD profile test-retest reliability **A:** Histograms of individual subject ICC between the MD tract profiles across sessions for a given bundle. Colors encode the bundles, matching the diagram showing the rough anatomical positions of the bundles for the left side of the brain (center). **B:** 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.

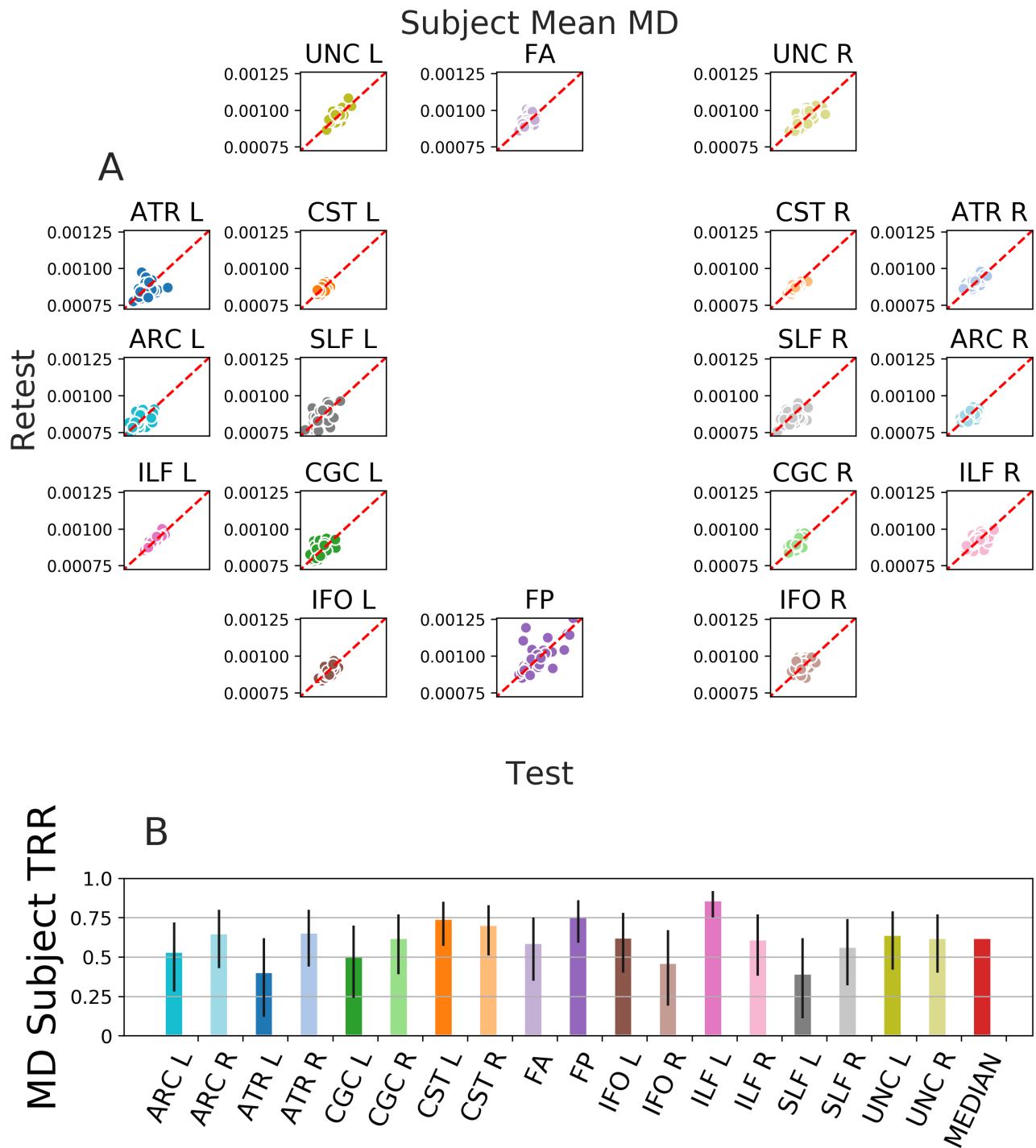


Fig. S5. Subject test-retest reliability **A:** Mean tract profiles for a given bundle and the MD scalar for each subject using the first and second session of HCP-TR. Colors encode bundle information, matching the core of the bundles (center). **B:** subject reliability is calculated from the Spearman's ρ of these distributions, with median across bundles in red. Error bars show the 95% confidence interval.

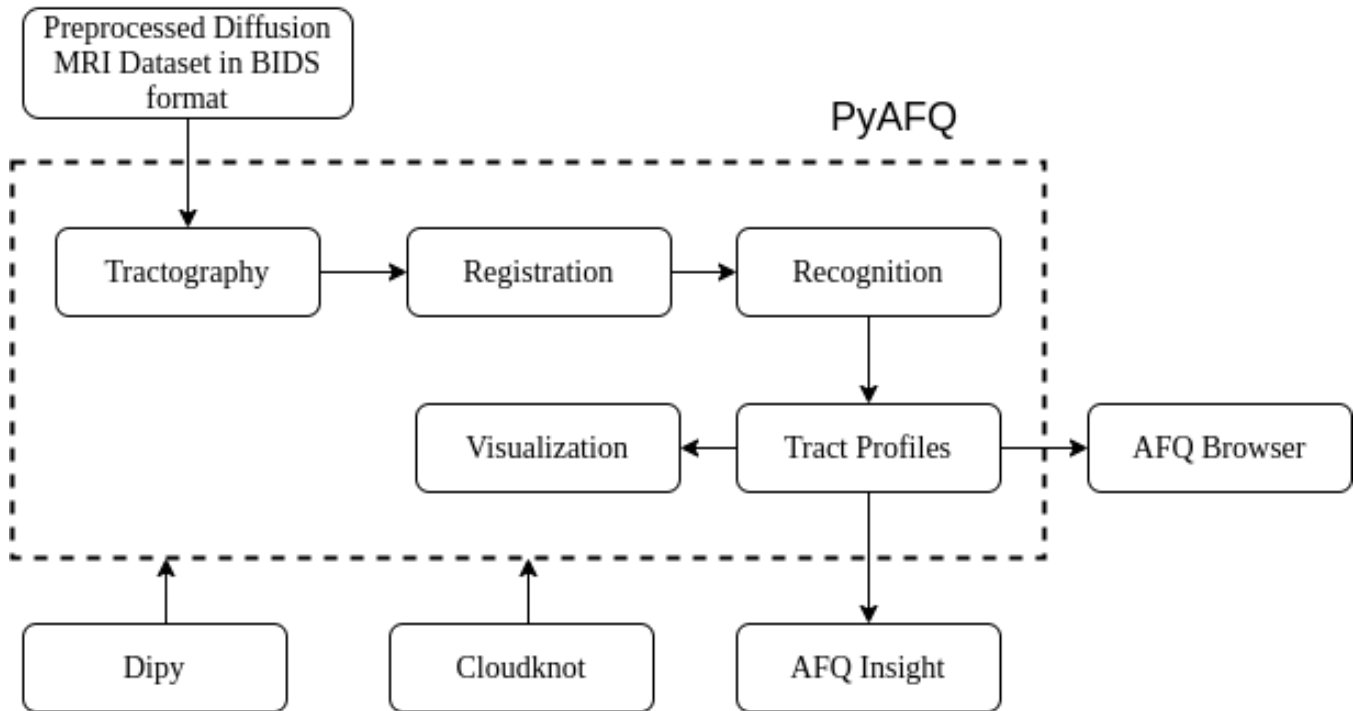


Fig. S6. The pyAFQ software is intergrated into an ecosystem for reproducible tractometry Steps performed by pyAFQ are enclosed in the dotted rectangle, whereas steps outside that rectangle are performed by other software. Upper left: pyAFQ requires preprocessed diffusion MRI data in BIDS format. This could be from QSIprep (26) or dMRIprep (<https://github.com/nipreps/dmriprep>). Bottom right: pyAFQ outputs can serve as inputs to AFQ Browser for further interaction and visualization (52) or AFQ Insight for statistical analysis (20). Bottom left: pyAFQ uses DIPY (28) for the implementation of dMRI algorithms. pyAFQ uses Cloudknot (63) to scale processing by parallelizing across subjects in AWS.

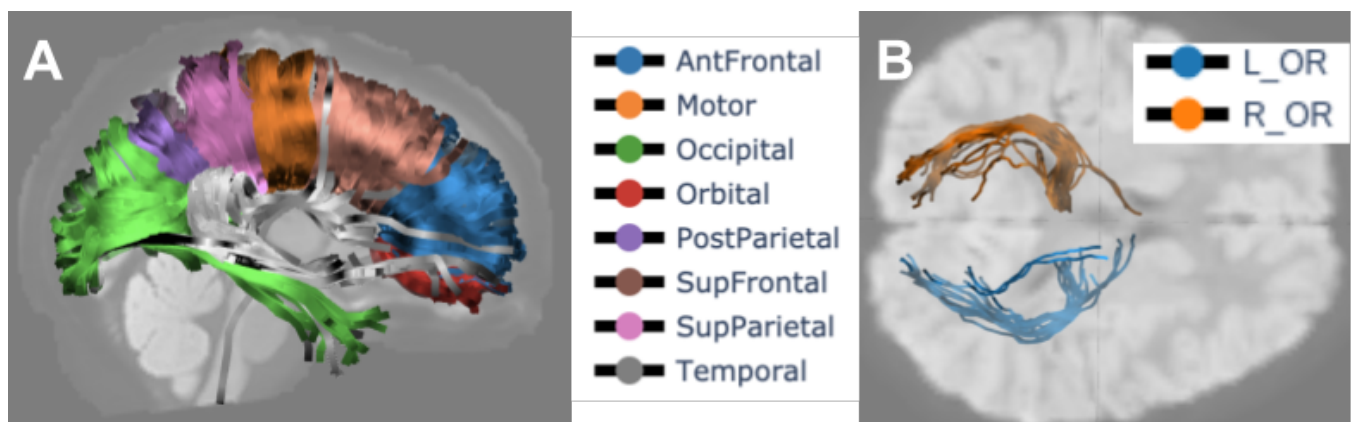


Fig. S7. Callosal bundles from HCP-TR, optic radiations from UW-PREK, found by pyAFQ. Streamlines are colored according to their bundles and shaded according to FA. The background images are each a b0 slice. **A** callosal bundles found by pyAFQ on an example subject from HCP-TR. **B** optic radiations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