

# Numerical Instabilities in Analytical Pipelines Lead to Large and Meaningful Variability in Brain Networks

Gregory Kiar<sup>1</sup>, Yohan Chatelain<sup>2</sup>, Pablo de Oliveira Castro<sup>3</sup>, Eric Petit<sup>4</sup>, Ariel Rokem<sup>5</sup>, Gaël Varoquaux<sup>6</sup>, Bratislav Misic<sup>1</sup>, Alan C. Evans<sup>1†</sup>, Tristan Glatard<sup>2†</sup>

The analysis of brain-imaging data requires complex and often non-linear transformations to support findings on brain function or pathologies. And yet, recent work has shown that variability in the choices that one makes when analyzing data can lead to quantitatively and qualitatively different results, endangering the trust in conclusions<sup>1-3</sup>. Even within a given method or analytical technique, numerical instabilities could compromise findings<sup>4-7</sup>. We instrumented a structural-connectome estimation pipeline with Monte Carlo Arithmetic<sup>8,9</sup>, a technique to introduce random noise in floating-point computations, and evaluated the stability of the derived connectomes, their features<sup>10,11</sup>, and the impact on a downstream analysis<sup>12,13</sup>. The stability of results was found to be highly dependent upon which features of the connectomes were evaluated, and ranged from perfectly stable (i.e. no observed variability across executions) to highly unstable (i.e. the results contained no trustworthy significant information). While the extreme range and variability in results presented here could severely hamper our understanding of brain organization in brain-imaging studies, it also leads to an increase in the reliability of datasets. This paper highlights the potential of leveraging the induced variance in estimates of brain connectivity to reduce the bias in networks alongside increasing the robustness of their applications in the detection or classification of individual differences. This paper demonstrates that stability evaluations are necessary for understanding error and bias inherent to scientific computing, and that they should be a component of typical analytical workflows.

## Keywords

Stability — Reproducibility — Network Neuroscience — Neuroimaging

<sup>1</sup>Montréal Neurological Institute, McGill University, Montréal, QC, Canada; <sup>2</sup>Department of Computer Science and Software Engineering, Concordia University, Montréal, QC, Canada; <sup>3</sup>Department of Computer Science, Université of Versailles, Versailles, France; <sup>4</sup>Exascale Computing Lab, Intel, Paris, France; <sup>5</sup>Department of Psychology and eScience Institute, University of Washington, Seattle, WA, USA; <sup>6</sup>Parietal project-team, INRIA Saclay-ile de France, France; †Authors contributed equally.

1 The modelling of brain networks, called connectomics, 7 This can not only improve understanding of so-called “connec-  
2 has shaped our understanding of the structure and function 8 topathies”, such as Alzheimer’s Disease and Schizophrenia,  
3 of the brain across a variety of organisms and scales over 9 but potentially pave the way for therapeutics<sup>19-23</sup>.  
4 the last decade<sup>11, 14-18</sup>. In humans, these wiring diagrams are 10 However, the analysis of brain imaging data relies on com-  
5 obtained *in vivo* through Magnetic Resonance Imaging (MRI), 11 plex computational methods and software. Tools are trusted to  
6 and show promise towards identifying biomarkers of disease. 12 perform everything from pre-processing tasks to downstream

13 statistical evaluation. While these tools undoubtedly undergo 49  
14 rigorous evaluation on bespoke datasets, in the absence of 50  
15 ground-truth this is often evaluated through measures of re- 51  
16 liability<sup>24–27</sup>, proxy outcome statistics, or agreement with 52  
17 existing theory. Importantly, this means that tools are not 53  
18 necessarily of known or consistent quality, and it is not un- 54  
19 common that equivalent experiments may lead to diverging 55  
20 conclusions<sup>1,5–7</sup>. While many scientific disciplines suffer 56  
21 from a lack of reproducibility<sup>28</sup>, this was recently explored 57  
22 in brain imaging by a 70 team consortium which performed 58  
23 equivalent analyses and found widely inconsistent results<sup>1</sup>, 59  
24 and it is likely that software instabilities played a role.

25 The present study approached evaluating reproducibility 61  
26 from a computational perspective in which a series of brain 62  
27 imaging studies were numerically perturbed such that the 63  
28 plausibility of results was not affected, and the biological 64  
29 implications of the observed instabilities were quantified. We 65  
30 accomplished this through the use of Monte Carlo Arithmetic 66  
31 (MCA)<sup>8</sup>, a technique which enables characterization of the 67  
32 sensitivity of a system to small perturbations. We explored 68  
33 the impact of perturbations through the direct comparison 69  
34 of structural connectomes, the consistency of their features, 70  
35 and their eventual application in a neuroscience study. Finally 71  
36 we conclude on the consequences and opportunities afforded 72  
37 by the observed instabilities and make recommendations for 73  
38 the roles stability analyses may play towards increasing the 74  
39 reliability of brain imaging research.

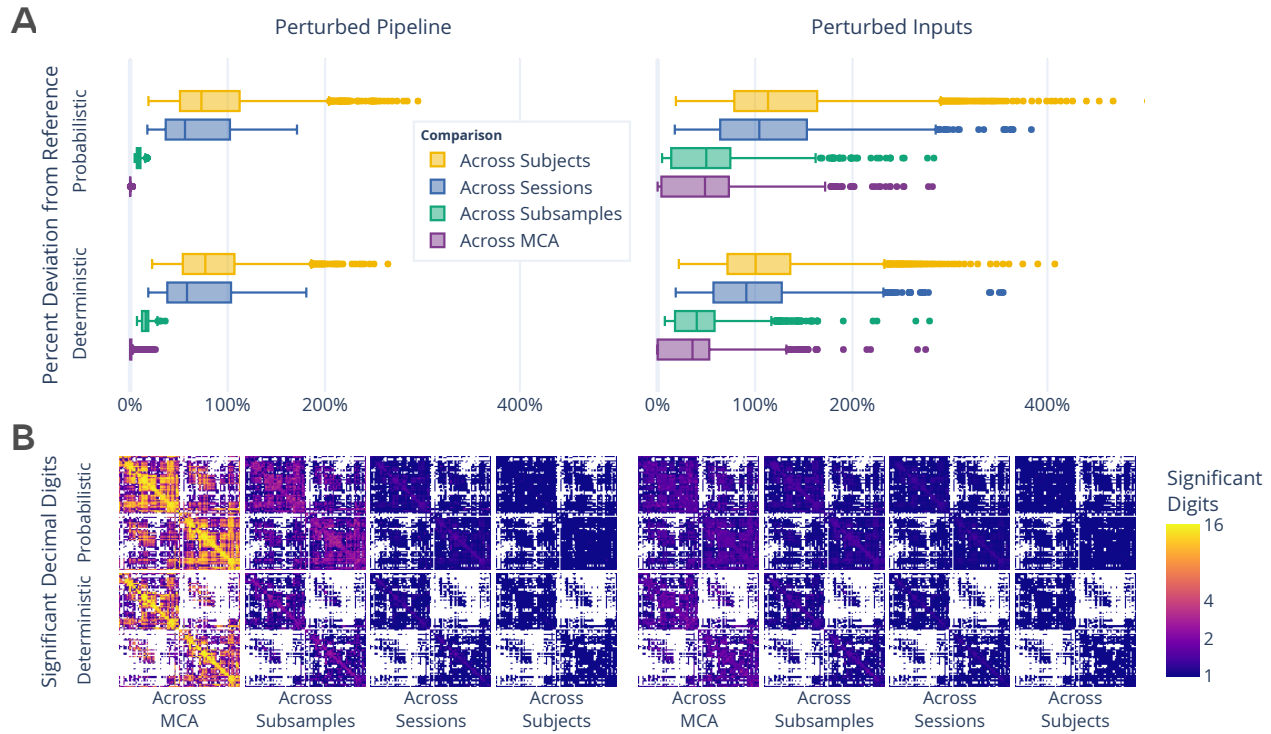
#### 40 **Graphs Vary Widely With Perturbations**

41 Prior to exploring the analytic impact of instabilities, a direct 78  
42 understanding of the induced variability was required. A sub- 79  
43 set of the Nathan Kline Institute Rockland Sample (NKIRS) 80  
44 dataset<sup>29</sup> was randomly selected to contain 25 individuals with 81  
45 two sessions of imaging data, each of which was subsampled 82  
46 into two components, resulting in four collections per individ- 83  
47 ual. Structural connectomes were generated with canonical 84  
48 deterministic and probabilistic pipelines<sup>30,31</sup> which were in- 85

strumented with MCA, replicating computational noise at  
either the inputs or throughout the pipelines<sup>4,9</sup>. The pipelines  
were sampled 20 times per collection and once without per-  
turbations, resulting in a total of 4,200 connectomes.

The stability of connectomes was evaluated through the  
deviation from reference and the number of significant digits  
(Figure 1). The comparisons were grouped according to dif-  
ferences across simulations, subsampling of data, sessions of  
acquisition, or subjects. While the similarity of connectomes  
decreases as the collections become more distinct, connec-  
tomes generated with input perturbations show considerable  
variability, often reaching deviations equal to or greater than  
those observed across individuals or sessions (Figure 1A;  
right). This finding suggests that instabilities inherent to  
these pipelines may mask session or individual differences,  
limiting the trustworthiness of derived connectomes. While  
both pipelines show similar performance, the probabilistic  
pipeline was more stable in the face of pipeline perturbations  
whereas the deterministic was more stable to input pertur-  
bations ( $p < 0.0001$  for all; exploratory). The stability of  
correlations can be found in Supplemental Section S1.

The number of significant digits per edge across connec-  
tomes (Figure 1B) similarly decreases across groups. While  
the cross-MCA comparison of connectomes generated with  
pipeline perturbations show nearly perfect precision for many  
edges (approaching the maximum of 15.7 digits for 64-bit  
data), this evaluation uniquely shows considerable drop off  
in performance across data subsampling (average of  $< 4$  dig-  
its). In addition, input perturbations show no more than an  
average of 3 significant digits across all groups, demonstrat-  
ing a significant limitation in the reliability independent edge  
weights. Significance across individuals did not exceed a  
single digit per edge in any case, indicating that only the  
magnitude of edges in naively computed groupwise average  
connectomes can be trusted. The combination of these results  
with those presented in Figure 1A suggests that while specific  
edge weights are largely affected by instabilities, macro-scale



**Figure 1.** Exploration of perturbation-induced deviations from reference connectomes. **(A)** The absolute deviations, in the form of normalized percent deviation from reference, shown as the across MCA series relative to Across Subsample, Across Session, and Across Subject variations. **(B)** The number of significant decimal digits in each set of connectomes as obtained after evaluating the effect of perturbations. In the case of 16, values can be fully relied upon, whereas in the case of 1 only the first digit of a value can be trusted. Pipeline- and input-perturbations are shown on the left and right, respectively.

86 network topology is stable.

### 87 **Subject-Specific Signal is Amplified While Off-Target** 88 **Biases Are Reduced**

89 We assessed the reproducibility of the dataset through mimick-  
90 ing and extending a typical test-retest experiment<sup>26</sup> in which  
91 the similarity of samples across multiple measurements were  
92 compared to distinct samples in the dataset (Table 1, with  
93 additional experiments and explanation in Supplemental Sec-  
94 tion S2). The ability to separate connectomes across subjects  
95 (Hypothesis 1) is an essential prerequisite for the application  
96 of brain imaging towards identifying individual differences<sup>18</sup>.  
97 In testing hypothesis 1, we observe that the dataset is sep-  
98 arable with a score of 0.64 and 0.65 ( $p < 0.001$ ; optimal  
99 score: 1.0; chance: 0.04) without any instrumentation. How-

100 ever, we can see that inducing instabilities through MCA  
101 improves the reliability of the dataset to over 0.75 in each  
102 case ( $p < 0.001$  for all), significantly higher than without  
103 instrumentation ( $p < 0.005$  for all). This result impactfully  
104 suggests the utility of perturbation methods for synthesizing  
105 robust and reliable individual estimates of connectivity, serv-  
106 ing as a cost effective and context-agnostic method for dataset  
107 augmentation.

108 While the separability of individuals is essential for the  
109 identification of brain networks, it is similarly reliant on net-  
110 work similarity across equivalent acquisitions (Hypothesis 2).  
111 In this case, connectomes were grouped based upon session,  
112 rather than subject, and the ability to distinguish one session  
113 from another was computed within-individual and aggregated.  
114 Both the unperturbed and pipeline perturbation settings per-

**Table 1.** The impact of instabilities as evaluated through the separability of the dataset based on individual (or subject) differences, session, and subsample. The performance is reported as mean Discriminability. While a perfectly separable dataset would be represented by a score of 1.0, the chance performance, indicating minimal separability, is  $1/\text{the number of classes}$ .  $H_3$  could not be tested using the reference executions due to too few possible comparisons. The alternative hypothesis, indicating significant separation, was accepted for all experiments, with  $p < 0.005$ .

| Comparison                | Chance | Target | Reference Execution |       | Perturbed Pipeline |       | Perturbed Inputs |       |
|---------------------------|--------|--------|---------------------|-------|--------------------|-------|------------------|-------|
|                           |        |        | Det.                | Prob. | Det.               | Prob. | Det.             | Prob. |
| $H_1$ : Across Subjects   | 0.04   | 1.0    | 0.64                | 0.65  | 0.82               | 0.82  | 0.77             | 0.75  |
| $H_2$ : Across Sessions   | 0.5    | 0.5    | 1.00                | 1.00  | 1.00               | 1.00  | 0.88             | 0.85  |
| $H_3$ : Across Subsamples | 0.5    | 0.5    |                     |       | 0.99               | 1.00  | 0.71             | 0.61  |

fectly preserved differences between cross-sectional sessions with a score of 1.0 ( $p < 0.005$ ; optimal score: 0.5; chance:

0.5), indicating a dominant session-dependent signal for all individuals despite no intended biological differences. However, while still significant relative to chance (score: 0.85 and 0.88;  $p < 0.005$  for both), input perturbations lead to significantly lower separability of the dataset ( $p < 0.005$  for all). This reduction of the difference between sessions of data within individuals suggests that increased variance caused by input perturbations reduces the impact of non-biological acquisition-dependent bias inherent in the brain graphs.

Though the previous sets of experiments inextricably evaluate the interaction between the dataset and tool, the use of

subsampling allowed for characterizing the separability of networks sampled from within a single acquisition (Hypothesis 3). While this experiment could not be evaluated using

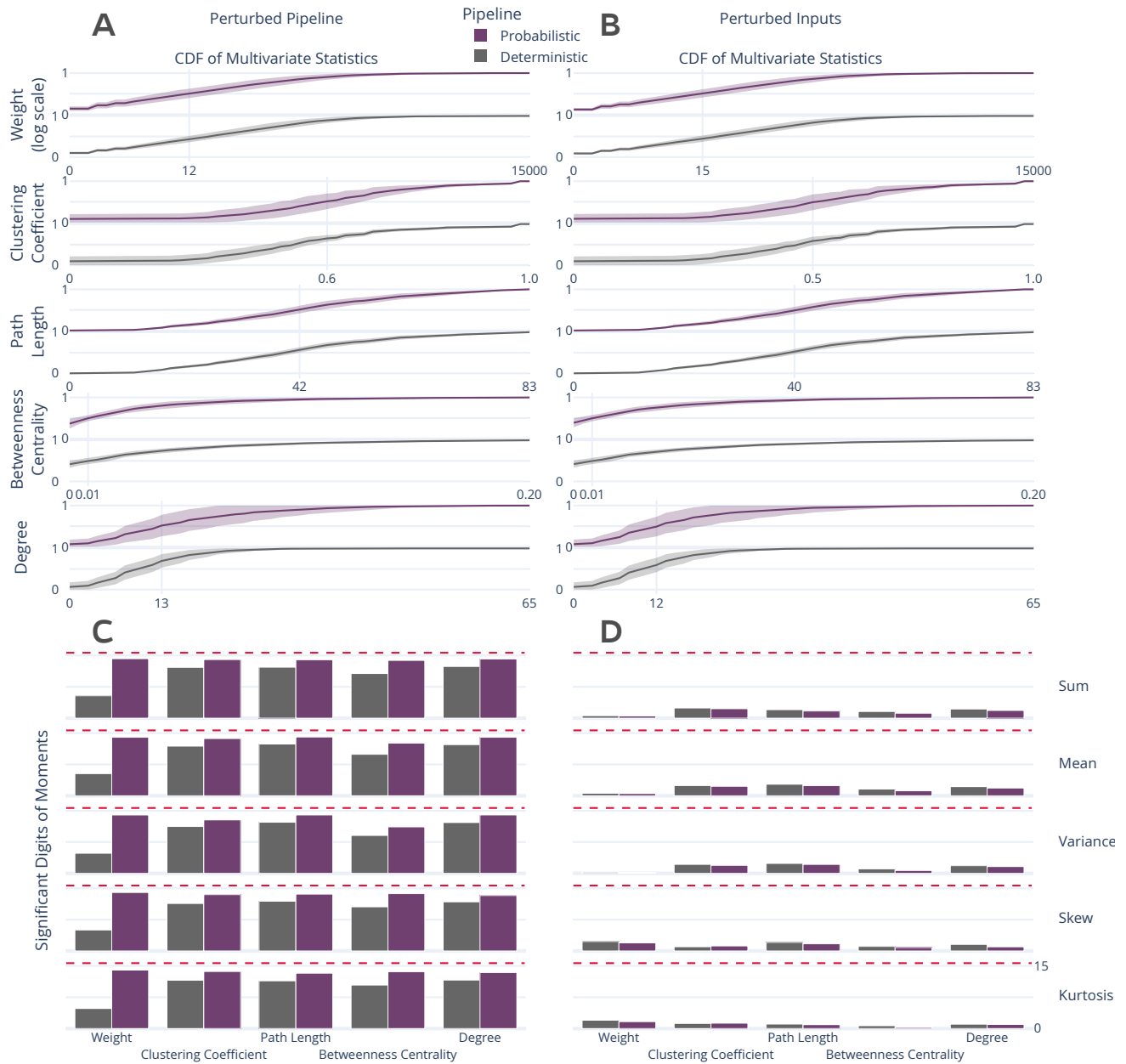
reference executions, the executions performed with pipeline perturbations showed near perfect separation between subsamples, with scores of 0.99 and 1.0 ( $p < 0.005$ ; optimal: 0.5; chance: 0.5). Given that there is no variability in data acquisition or preprocessing that contributes to this reliable identification of scans, the separability observed in this experiment may only be due to instability or bias inherent to the pipelines. The high variability introduced through input perturbations considerably lowered the reliability towards chance (score: 0.71 and 0.61;  $p < 0.005$  for all), further supporting

this as an effective method for obtaining lower-bias estimates of individual connectivity. Across all cases, the induced perturbations showed an amplification of meaningful biological signal alongside a reduction of off-target signal. This result appears strikingly like a manifestation of the well-known bias-variance tradeoff<sup>32</sup> in machine learning, a concept which observes a decrease in bias as variance is favoured by a model. In particular, this highlights that numerical perturbations can be used to not only evaluate the stability of pipelines, but that the induced variance may be leveraged for the interpretation as a robust distributions of possible results.

### Distributions of Graph Statistics Are Reliable, But Individual Statistics Are Not

Exploring the stability of topological features of connectomes is relevant for typical analyses, as low dimensional features are often more suitable than full connectomes for many analytical methods in practice<sup>11</sup>. A separate subset of the NKIRS dataset was randomly selected to contain a single non-sampled session for 100 individuals, and connectomes were generated as above.

The stability of several commonly-used multivariate graph features<sup>10</sup> was explored in Figure 2. The cumulative density of the features was computed within individuals and the mean density and associated standard error were computed



**Figure 2.** Distribution and stability assessment of multivariate graph statistics. (A, B) The cumulative distribution functions of multivariate statistics across all subjects and perturbation settings. There was no significant difference between the distributions in A and B. (C, D) The number of significant digits in the first 5 five moments of each statistic across perturbations. The dashed red line refers to the maximum possible number of significant digits.

166 for across individuals (Figures 2A and 2B). There was no sig- 171 In addition to the comparison of distributions, the stabil-  
 167 nificant difference between the distributions for each feature 172 ity of the first 5 moments of these features was evaluated  
 168 across the two perturbation settings, suggesting that the topo- 173 (Figures 2C and 2D). In the face of pipeline perturbations,  
 169 logical features summarized by these multivariate features are 174 the feature-moments were stable with more than 10 signifi-  
 170 robust across both perturbation modes. 175 cant digits with the exception of edge weight when using the

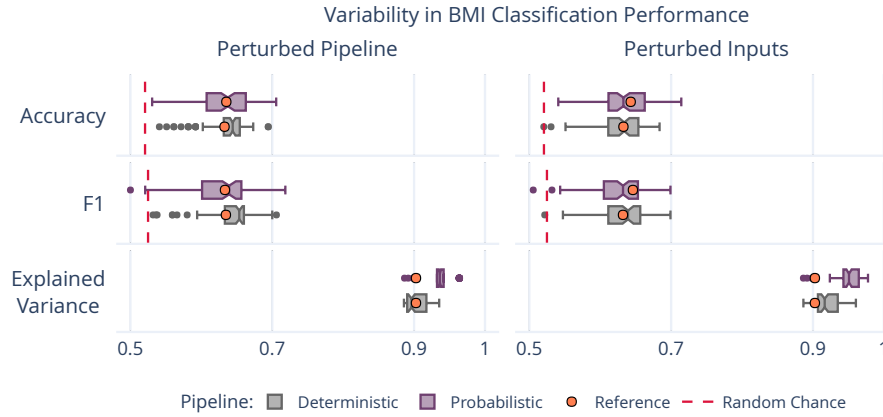


176 deterministic pipeline, though the probabilistic pipeline was 212 outcome. Importantly, this finding does not suggest that mod-  
177 more stable for all comparisons ( $p < 0.0001$ ; exploratory). 213 elling brain-phenotype relationships is not possible, but rather  
178 In stark contrast, input perturbations led to highly unstable 214 it sheds light on impactful uncertainty that must be accounted  
179 feature-moments (Figure 2D), such that none contained more 215 for in this process, and supports the use of ensemble modeling  
180 than 5 significant digits of information and several contained 216 techniques.  
181 less than a single significant digit, indicating a complete lack  
182 of reliability. This dramatic degradation in stability for in- 217 **Discussion**  
183 dividual measures strongly suggests that these features may 218 The perturbation of structural connectome estimation pipelines  
184 be unreliable as individual biomarkers when derived from a 219 with small amounts of noise, on the order of machine error,  
185 single pipeline evaluation, though their reliability may be in- 220 led to considerable variability in derived brain graphs. Across  
186 creased when studying their distributions across perturbations. 221 all analyses the stability of results ranged from nearly per-  
187 A similar analysis was performed for univariate statistics and 222 fectly trustworthy (i.e. no variation) to completely unreliable  
188 can be found in Supplemental Section S3. 223 (i.e. containing no trustworthy information). Given that the  
224 magnitude of introduced numerical noise is to be expected

### 189 **Uncertainty in Brain-Phenotype Relationships**

190 While the variability of connectomes and their features was 225 in typical settings, this finding has potentially significant im-  
191 summarized above, networks are commonly used as inputs to 226 plications for inferences in brain imaging as it is currently  
192 machine learning models tasked with learning brain-phenotype 227 performed. In particular, this bounds the success of studying  
193 relationships<sup>18</sup>. To explore the stability of these analyses, we 228 individual differences, a central objective in brain imaging<sup>18</sup>,  
194 modelled the relationship between high- or low- Body Mass 229 given that the quality of relationships between phenotypic  
195 Index (BMI) groups and brain connectivity<sup>12,13</sup>, using stan- 230 data and brain networks will be limited by the stability of the  
196 dard dimensionality reduction and classification tools, and 231 connectomes themselves. This issue was accentuated through  
197 compared this to reference and random performance (Fig- 232 the crucial finding that individually derived network features  
198 ure 3). 233 were unreliable despite there being no significant difference  
234 in their aggregated distributions. This finding is not damn-

199 The analysis was perturbed through distinct samplings of 235 ing for the study of brain networks as a whole, but rather is  
200 the dataset across both pipelines and perturbation methods. 236 strong support for the aggregation of networks, either across  
201 The accuracy and F1 score for the perturbed models varied 237 perturbations for an individual or across groups, over the use  
202 from 0.520 – 0.716 and 0.510 – 0.725, respectively, rang- 238 of individual estimates.  
203 ing from at or below random performance to outperforming 239 **Underestimated False Positive Rates** While the instabil-  
204 performance on the reference dataset. This large variability 240 ity of brain networks was used here to demonstrate the lim-  
205 illustrates a previously uncharacterized margin of uncertainty 241 itations of modelling brain-phenotype relationships in the  
206 in the modelling of this relationship, and limits confidence in 242 context of machine learning, this limitation extends to classi-  
207 reported accuracy scores on singly processed datasets. The 243 cal hypothesis testing, as well. Though performing individual  
208 portion of explained variance in these samples ranged from 244 comparisons in a hypothesis testing framework will be accom-  
209 88.6% – 97.8%, similar to the reference, suggesting that the 245 panied by reported false positive rates, the accuracy of these  
210 range in performance was not due to a gain or loss of mean- 246 rates is critically dependent upon the reliability of the samples  
211 ingful signal, but rather the reduction of bias towards specific 247 used. In reality, the true false positive rate for a test would be



**Figure 3.** Variability in BMI classification across the sampling of an MCA-perturbed dataset. The dashed red lines indicate random-chance performance, and the orange dots show the performance using the reference executions.

248 a combination of the reported confidence and the underlying 272 pervasive collection of repeated measurements choreographed  
 249 variability in the results, a typically unknown quantity. 273 by massive cross-institutional consortia<sup>34,35</sup>. The finding that

250 When performing these experiments outside of a repeated- 274 perturbing experiments using MCA both increased the relia-  
 251 measure context, such as that afforded here through MCA, it 275 bility of the dataset and decreased off-target differences across  
 252 is impossible to empirically estimate the reliability of samples. 276 acquisitions opens the door for a promising paradigm shift.  
 253 This means that the reliability of accepted hypotheses is also 277 Given that MCA is data-agnostic, this technique could be used  
 254 unknown, regardless of the reported false positive rate. In 278 effectively in conjunction with, or in lieu of, realistic noise  
 255 fact, it is a virtual certainty that the true false positive rate 279 models to augment existing datasets. While this of course  
 256 for a given hypothesis exceeds the reported value simply as 280 would not replace the need for repeated measurements when  
 257 a result of numerical instabilities. This uncertainty inherent 281 exploring the effect of data collection paradigm or study lon-  
 258 to derived data is compounded with traditional arguments 282 gitudinal progressions of development or disease, it could be  
 259 limiting the trustworthiness of claims<sup>33</sup>, and hampers the 283 used in conjunction with these efforts to increase the reliabil-  
 260 ability of researchers to evaluate the quality of results. The 284 ity of each distinct sample within a dataset. In contexts where  
 261 accompaniment of brain imaging experiments with direct 285 repeated measurements are collected to increase the fidelity of  
 262 evaluations of their stability, as was done here, would allow 286 the dataset, MCA could potentially be employed to increase  
 263 researchers to simultaneously improve the numerical stability 287 the reliability of the dataset and save millions of dollars on  
 264 of their analyses and accurately gauge confidence in them. 288 data collection. This technique also opens the door for the  
 265 The induced variability in derived brain networks may be 289 characterization of reliability across axes which have been  
 266 leveraged to estimate aggregate connectomes with lower bias 290 traditionally inaccessible. For instance, in the absence of a  
 267 than any single independent observation, leading to learned 291 realistic noise model or simulation technique similar to MCA,  
 268 relationships that are more generalizable and ultimately more 292 the evaluation of network stability across data subsampling  
 269 useful. 293 would not have been possible.

270 **Cost-Effective Data Augmentation** The evaluation of reli- 294 **Shortcomings and Future Questions** Given the complex-  
 271 ability in brain imaging has historically relied upon the ex- 295 ity of recompiling complex software libraries, pre-processing

was not perturbed in these experiments. Other work has shown that linear registration, a core piece of many elements of pre-processing such as motion correction and alignment, is sensitive to minor perturbations<sup>7</sup>. It is likely that the instabilities across the entire processing workflow would be compounded with one another, resulting in even greater variability. While the analyses performed in this paper evaluated a single dataset and set of pipelines, extending this work to other modalities and analyses is of interest for future projects.

This paper does not explore methodological flexibility or compare this to numerical instability. Recently, the nearly boundless space of analysis pipelines and their impact on outcomes in brain imaging has been clearly demonstrated<sup>1</sup>. The approach taken in these studies complement one another and explore instability at the opposite ends of the spectrum, with human variability in the construction of an analysis workflow on one end and the unavoidable error implicit in the digital representation of data on the other. It is of extreme interest to combine these approaches and explore the interaction of these scientific degrees of freedom with effects from software implementations, libraries, and parametric choices.

Finally, it is important to state explicitly that the work presented here does not invalidate analytical pipelines used in brain imaging, but merely sheds light on the fact that many studies are accompanied by an unknown degree of uncertainty due to machine-introduced errors. The presence of unknown error-bars associated with experimental findings limits the impact of results due to increased uncertainty. The desired outcome of this paper is to motivate a shift in scientific computing – particularly in neuroimaging – towards a paradigm which favours the explicit evaluation of the trustworthiness of claims alongside the claims themselves.

## References

[1] R. Botvinik-Nezer, F. Holzmeister, C. F. Camerer, A. Dreber, J. Huber, M. Johannesson, M. Kirchler, R. Iwanir, J. A. Mumford, R. A. Adcock *et al.*, “Variability in the analysis of a single neuroimaging dataset by many teams,” *Nature*, pp. 1–7, 2020.

- [2] C. M. Bennett, M. B. Miller, and G. L. Wolford, “Neural correlates of interspecies perspective taking in the post-mortem Atlantic salmon: An argument for multiple comparisons correction,” *Neuroimage*, vol. 47, no. Suppl 1, p. S125, 2009.
- [3] A. Eklund, T. E. Nichols, and H. Knutsson, “Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates,” *Proceedings of the national academy of sciences*, vol. 113, no. 28, pp. 7900–7905, 2016.
- [4] G. Kiar, P. de Oliveira Castro, P. Rioux, E. Petit, S. T. Brown, A. C. Evans, and T. Glatard, “Comparing perturbation models for evaluating stability of neuroimaging pipelines,” *The International Journal of High Performance Computing Applications*, 2020.
- [5] A. Salari, G. Kiar, L. Lewis, A. C. Evans, and T. Glatard, “File-based localization of numerical perturbations in data analysis pipelines,” *arXiv preprint arXiv:2006.04684*, 2020.
- [6] L. B. Lewis, C. Y. Lepage, N. Khalili-Mahani, M. Omidyeganeh, S. Jeon, P. Bermudez, A. Zijdenbos, R. Vincent, R. Adalat, and A. C. Evans, “Robustness and reliability of cortical surface reconstruction in CIVET and FreeSurfer,” *Annual Meeting of the Organization for Human Brain Mapping*, 2017.
- [7] T. Glatard, L. B. Lewis, R. Ferreira da Silva, R. Adalat, N. Beck, C. Lepage, P. Rioux, M.-E. Rousseau, T. Sherif, E. Deelman, N. Khalili-Mahani, and A. C. Evans, “Reproducibility of neuroimaging analyses across operating systems,” *Front. Neuroinform.*, vol. 9, p. 12, Apr. 2015.
- [8] D. S. Parker, *Monte Carlo Arithmetic: exploiting randomness in floating-point arithmetic*. University of California (Los Angeles). Computer Science Department, 1997.
- [9] C. Denis, P. de Oliveira Castro, and E. Petit, “Verificarlo: Checking floating point accuracy through monte carlo arithmetic,” *2016 IEEE 23rd Symposium on Computer Arithmetic (ARITH)*, 2016.
- [10] R. F. Betzel, A. Griffa, P. Hagmann, and B. Mišić, “Distance-dependent consensus thresholds for generating group-representative structural brain networks,” *Network neuroscience*, vol. 3, no. 2, pp. 475–496, 2019.
- [11] M. Rubinov and O. Sporns, “Complex network measures of brain connectivity: uses and interpretations,” *Neuroimage*, vol. 52, no. 3, pp. 1059–1069, Sep. 2010.
- [12] B.-Y. Park, J. Seo, J. Yi, and H. Park, “Structural and functional brain connectivity of people with obesity and prediction of body mass index using connectivity,” *PLoS One*, vol. 10, no. 11, p. e0141376, Nov. 2015.
- [13] A. Gupta, E. A. Mayer, C. P. Sanmiguel, J. D. Van Horn, D. Woodworth, B. M. Ellingson, C. Fling, A. Love, K. Tillisch, and J. S. Labus, “Patterns of brain structural connectivity differentiate normal weight from overweight subjects,” *Neuroimage Clin*, vol. 7, pp. 506–517, Jan. 2015.
- [14] T. E. Behrens and O. Sporns, “Human connectomics,” *Current opinion in neurobiology*, vol. 22, no. 1, pp. 144–153, 2012.



- [15] M. Xia, Q. Lin, Y. Bi, and Y. He, “Connectomic insights into topologically centralized network edges and relevant motifs in the human brain,” *Frontiers in human neuroscience*, vol. 10, p. 158, 2016.
- [16] J. L. Morgan and J. W. Lichtman, “Why not connectomics?” *Nature methods*, vol. 10, no. 6, p. 494, 2013.
- [17] M. P. Van den Heuvel, E. T. Bullmore, and O. Sporns, “Comparative connectomics,” *Trends in cognitive sciences*, vol. 20, no. 5, pp. 345–361, 2016.
- [18] J. Dubois and R. Adolphs, “Building a science of individual differences from fMRI,” *Trends Cogn. Sci.*, vol. 20, no. 6, pp. 425–443, Jun. 2016.
- [19] A. Fornito and E. T. Bullmore, “Connectomics: a new paradigm for understanding brain disease,” *European Neuropsychopharmacology*, vol. 25, no. 5, pp. 733–748, 2015.
- [20] G. Deco and M. L. Kringelbach, “Great expectations: using whole-brain computational connectomics for understanding neuropsychiatric disorders,” *Neuron*, vol. 84, no. 5, pp. 892–905, 2014.
- [21] T. Xie and Y. He, “Mapping the alzheimer’s brain with connectomics,” *Frontiers in psychiatry*, vol. 2, p. 77, 2012.
- [22] M. Filippi, M. P. van den Heuvel, A. Fornito, Y. He, H. E. H. Pol, F. Agosta, G. Comi, and M. A. Rocca, “Assessment of system dysfunction in the brain through mri-based connectomics,” *The Lancet Neurology*, vol. 12, no. 12, pp. 1189–1199, 2013.
- [23] M. P. Van Den Heuvel and A. Fornito, “Brain networks in schizophrenia,” *Neuropsychology review*, vol. 24, no. 1, pp. 32–48, 2014.
- [24] J. J. Bartko, “The intraclass correlation coefficient as a measure of reliability,” *Psychol. Rep.*, vol. 19, no. 1, pp. 3–11, Aug. 1966.
- [25] A. M. Brandmaier, E. Wenger, N. C. Bodammer, S. Kühn, N. Raz, and U. Lindenberger, “Assessing reliability in neuroimaging research through intra-class effect decomposition (ICED),” *Elife*, vol. 7, Jul. 2018.
- [26] E. W. Bridgeford, S. Wang, Z. Yang, Z. Wang, T. Xu, C. Craddock, J. Dey, G. Kiar, W. Gray-Roncal, C. Coulantoni *et al.*, “Eliminating accidental deviations to minimize generalization error: applications in connectomics and genomics,” *bioRxiv*, p. 802629, 2020.
- [27] G. Kiar, E. Bridgeford, W. G. Roncal, V. Chandrashekar, and others, “A High-Throughput pipeline identifies robust connectomes but troublesome variability,” *bioRxiv*, 2018.
- [28] M. Baker, “1,500 scientists lift the lid on reproducibility,” *Nature*, 2016.
- [29] K. B. Nooner, S. J. Colcombe, R. H. Tobe, M. Mennes *et al.*, “The NKI-Rockland sample: A model for accelerating the pace of discovery science in psychiatry,” *Front. Neurosci.*, vol. 6, p. 152, Oct. 2012.
- [30] E. Garyfallidis, M. Brett, B. Amirbekian, A. Rokem, S. van der Walt, M. Descoteaux, I. Nimmo-Smith, and Dipy Contributors, “Dipy, a library for the analysis of diffusion MRI data,” *Front. Neuroinform.*, vol. 8, p. 8, Feb. 2014.
- [31] E. Garyfallidis, M. Brett, M. M. Correia, G. B. Williams, and I. Nimmo-Smith, “QuickBundles, a method for tractography simplification,” *Front. Neurosci.*, vol. 6, p. 175, Dec. 2012.
- [32] S. Geman, E. Bienenstock, and R. Doursat, “Neural networks and the bias/variance dilemma,” *Neural computation*, vol. 4, no. 1, pp. 1–58, 1992.
- [33] J. P. Ioannidis, “Why most published research findings are false,” *PLoS medicine*, vol. 2, no. 8, p. e124, 2005.
- [34] D. C. Van Essen, S. M. Smith, D. M. Barch, T. E. Behrens, E. Yacoub, K. Ugurbil, W.-M. H. Consortium *et al.*, “The WU-Minn human connectome project: an overview,” *Neuroimage*, vol. 80, pp. 62–79, 2013.
- [35] X.-N. Zuo, J. S. Anderson, P. Bellec, R. M. Birn, B. B. Biswal, J. Blautzik, J. C. Breitner, R. L. Buckner, V. D. Calhoun, F. X. Castellanos *et al.*, “An open science resource for establishing reliability and reproducibility in functional connectomics,” *Scientific data*, vol. 1, no. 1, pp. 1–13, 2014.
- [36] M. Jenkinson, C. F. Beckmann, T. E. J. Behrens, M. W. Woolrich, and S. M. Smith, “FSL,” *Neuroimage*, vol. 62, no. 2, pp. 782–790, Aug. 2012.
- [37] J. L. Lancaster, D. Tordesillas-Gutiérrez, M. Martínez, F. Salinas, A. Evans, K. Zilles, J. C. Mazziotta, and P. T. Fox, “Bias between mni and talairach coordinates analyzed using the icbm-152 brain template,” *Human brain mapping*, vol. 28, no. 11, pp. 1194–1205, 2007.
- [38] A. Klein and J. Tourville, “101 labeled brain images and a consistent human cortical labeling protocol,” *Front. Neurosci.*, vol. 6, p. 171, Dec. 2012.
- [39] D. Sohier, P. De Oliveira Castro, F. Févotte, B. Lathuilière, E. Petit, and O. Jamond, “Confidence intervals for stochastic arithmetic,” Jul. 2018.
- [40] J. Benesty, J. Chen, Y. Huang, and I. Cohen, “Pearson correlation coefficient,” in *Noise Reduction in Speech Processing*, I. Cohen, Y. Huang, J. Chen, and J. Benesty, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 1–4.
- [41] C. A. Raji, A. J. Ho, N. N. Parikhshak, J. T. Becker, O. L. Lopez, L. H. Kuller, X. Hua, A. D. Leow, A. W. Toga, and P. M. Thompson, “Brain structure and obesity,” *Hum. Brain Mapp.*, vol. 31, no. 3, pp. 353–364, Mar. 2010.
- [42] T. Glatard, G. Kiar, T. Aumentado-Armstrong, N. Beck, P. Bellec, R. Bernard, A. Bonnet, S. T. Brown, S. Camarasu-Pop, F. Cervenansky, S. Das, R. Ferreira da Silva, G. Flandin, P. Girard, K. J. Gorgolewski, C. R. G. Guttman, V. Hayot-Sasson, P.-O. Quirion, P. Rioux, M.-É. Rousseau, and A. C. Evans, “Boutiques: a flexible framework to integrate command-line applications in computing platforms,” *Gigascience*, vol. 7, no. 5, May 2018.
- [43] G. Kiar, S. T. Brown, T. Glatard, and A. C. Evans, “A serverless tool for platform agnostic computational experiment management,” *Front. Neuroinform.*, vol. 13, p. 12, Mar. 2019.

468 [44] H. Huang and M. Ding, “Linking functional connectivity and structural  
469 connectivity quantitatively: a comparison of methods,” *Brain connectiv-*  
470 *ity*, vol. 6, no. 2, pp. 99–108, 2016.

## Methods

### Dataset

The Nathan Kline Institute Rockland Sample (NKI-RS)<sup>29</sup> dataset contains high-fidelity imaging and phenotypic data from over 1,000 individuals spread across the lifespan. A subset of this dataset was chosen for each experiment to both match sample sizes presented in the original analyses and to minimize the computational burden of performing MCA. The selected subset comprises 100 individuals ranging in age from 6 – 79 with a mean of 36.8 (original: 6 – 81, mean 37.8), 60% female (original: 60%), with 52% having a BMI over 25 (original: 54%).

Each selected individual had at least a single session of both structural T1-weighted (MPRAGE) and diffusion-weighted (DWI) MR imaging data. DWI data was acquired with 137 diffusion directions; more information regarding the acquisition of this dataset can be found in the NKI-RS data release<sup>29</sup>.

In addition to the 100 sessions mentioned above, 25 individuals had a second session to be used in a test-retest analysis. Two additional copies of the data for these individuals were generated, including only the odd or even diffusion directions (64 + 9 B0 volumes = 73 in either case). This allowed for an extra level of stability evaluation to be performed between the levels of MCA and session-level variation.

In total, the dataset is composed of 100 downsampled sessions of data originating from 50 acquisitions and 25 individuals for in depth stability analysis, and an additional 100 sessions of full-resolution data from 100 individuals for subsequent analyses.

### Processing

The dataset was preprocessed using a standard FSL<sup>36</sup> workflow consisting of eddy-current correction and alignment. The MNI152 atlas<sup>37</sup> was aligned to each session of data, and the resulting transformation was applied to the DKT parcellation<sup>38</sup>. Downsampling the diffusion data took place after preprocess-

ing was performed on full-resolution sessions, ensuring that an additional confound was not introduced in this process when comparing between downsampled sessions. The preprocessing described here was performed once without MCA, and thus is not being evaluated.

Structural connectomes were generated from preprocessed data using two canonical pipelines from Dipy<sup>30</sup>: deterministic and probabilistic. In the deterministic pipeline, a constant solid angle model was used to estimate tensors at each voxel and streamlines were then generated using the EuDX algorithm<sup>31</sup>. In the probabilistic pipeline, a constrained spherical deconvolution model was fit at each voxel and streamlines were generated by iteratively sampling the resulting fiber orientation distributions. In both cases tracking occurred with 8 seeds per 3D voxel and edges were added to the graph based on the location of terminal nodes with weight determined by fiber count.

The random state of the probabilistic pipeline was fixed for all analyses. Fixing this random seed allowed for explicit attribution of observed variability to Monte Carlo simulations rather than internal state of the algorithm.

### Perturbations

All connectomes were generated with one reference execution where no perturbation was introduced in the processing. For all other executions, all floating point operations were instrumented with Monte Carlo Arithmetic (MCA)<sup>8</sup> through Verificarlo<sup>9</sup>. MCA simulates the distribution of errors implicit to all instrumented floating point operations (flop). This rounding is performed on a value  $x$  at precision  $t$  by:

$$\text{inexact}(x) = x + 2^{e_x - t} \xi \quad (1)$$

where  $e_x$  is the exponent value of  $x$  and  $\xi$  is a uniform random variable in the range  $(-\frac{1}{2}, \frac{1}{2})$ . MCA can be introduced in two places for each flop: before or after evaluation. Performing MCA on the inputs of an operation limits its precision, while performing MCA on the output of an operation high-

lights round-off errors that may be introduced. The former is referred to as Precision Bounding (PB) and the latter is called Random Rounding (RR).

Using MCA, the execution of a pipeline may be performed many times to produce a distribution of results. Studying the distribution of these results can then lead to insights on the stability of the instrumented tools or functions. To this end, a complete software stack was instrumented with MCA and is made available on GitHub at <https://github.com/gkiar/fuzzy>.

Both the RR and PB variants of MCA were used independently for all experiments. As was presented in<sup>4</sup>, both the degree of instrumentation (i.e. number of affected libraries) and the perturbation mode have an effect on the distribution of observed results. For this work, the RR-MCA was applied across the bulk of the relevant libraries and is referred to as Pipeline Perturbation. In this case the bulk of numerical operations were affected by MCA.

Conversely, the case in which PB-MCA was applied across the operations in a small subset of libraries is here referred to as Input Perturbation. In this case, the inputs to operations within the instrumented libraries (namely, Python and Cython) were perturbed, resulting in less frequent, data-centric perturbations. Alongside the stated theoretical differences, Input Perturbation is considerably less computationally expensive than Pipeline Perturbation.

All perturbations targeted the least-significant-bit for all data ( $t = 24$  and  $t = 53$  in float32 and float64, respectively<sup>9</sup>). Simulations were performed 20 times for each pipeline execution. A detailed motivation for the number of simulations can be found in<sup>39</sup>.

## Evaluation

The magnitude and importance of instabilities in pipelines can be considered at a number of analytical levels, namely: the induced variability of derivatives directly, the resulting downstream impact on summary statistics or features, or the

ultimate change in analyses or findings. We explore the nature and severity of instabilities through each of these lenses. Unless otherwise stated, all p-values were computed using Wilcoxon signed-rank tests.

## Direct Evaluation of the Graphs

The differences between simulated graphs was measured directly through both a direct variance quantification and a comparison to other sources of variance such as individual- and session-level differences.

**Quantification of Variability** Graphs, in the form of adjacency matrices, were compared to one another using three metrics: normalized percent deviation, Pearson correlation, and edgewise significant digits. The normalized percent deviation measure, defined in<sup>4</sup>, scales the norm of the difference between a simulated graph and the reference execution (that without intentional perturbation) with respect to the norm of the reference graph. The purpose of this comparison is to provide insight on the scale of differences in observed graphs relative to the original signal intensity. A Pearson correlation coefficient<sup>40</sup> was computed in complement to normalized percent deviation to identify the consistency of structure and not just intensity between observed graphs.

Finally, the estimated number of significant digits,  $s'$ , for each edge in the graph is calculated as:

$$s' = -\log_{10} \frac{\sigma}{|\mu|} \quad (2)$$

where  $\mu$  and  $\sigma$  are the mean and unbiased estimator of standard deviation across graphs, respectively. The upper bound on significant digits is 15.7 for 64-bit floating point data.

The percent deviation, correlation, and number of significant digits were each calculated within a single session of data, thereby removing any subject- and session-effects and providing a direct measure of the tool-introduced variability across perturbations. A distribution was formed by aggregating these individual results.

611 **Class-based Variability Evaluation** To gain a concrete un-  
612 derstanding of the significance of observed variations we ex-  
613 plore the separability of our results with respect to under-  
614 stood sources of variability, such as subject-, session-, and pipeline-  
615 level effects. This can be probed through Discriminability<sup>26</sup>,  
616 a technique similar to ICC<sup>24</sup> which relies on the mean of a  
617 ranked distribution of distances between observations belong-  
618 ing to a defined set of classes. The discriminability statistic is  
619 formalized as follows:

$$Disc. = Pr(\|g_{ij} - g_{i'j'}\| \leq \|g_{ij} - g_{i'j}\|) \quad (3)$$

620 where  $g_{ij}$  is a graph belonging to class  $i$  that was measured  
621 at observation  $j$ , where  $i \neq i'$  and  $j \neq j'$ .

622 Discriminability can then be read as the probability that an  
623 observation belonging to a given class will be more similar to  
624 other observations within that class than observations of a dif-  
625 ferent class. It is a measure of reproducibility, and is discussed  
626 in detail in<sup>26</sup>. This definition allows for the exploration of  
627 deviations across arbitrarily defined classes which in practice  
628 can be any of those listed above. We combine this statistic  
629 with permutation testing to test hypotheses on whether differ-  
630 ences between classes are statistically significant in each of  
631 these settings.

632 With this in mind, three hypotheses were defined. For  
633 each setting, we state the alternate hypotheses, the variable(s)  
634 which were used to determine class membership, and the  
635 remaining variables which may be sampled when obtaining  
636 multiple observations. Each hypothesis was tested indepen-  
637 dently for each pipeline and perturbation mode, and in every  
638 case where it was possible the hypotheses were tested using  
639 the reference executions alongside using MCA.

640  $H_{A1}$ : Individuals are distinct from one another

641 Class definition: *Subject ID*

642 Comparisons: *Session (1 subsample), Subsample (1*  
643 *session), MCA (1 subsample, 1 session)*

644  $H_{A2}$ : Sessions within an individual are distinct

645 Class definition: *Session ID | Subject ID*

646 Comparisons: *Subsample, MCA (1 subsample)*

647  $H_{A3}$ : Subsamples are distinct

648 Class definition: *Subsample | Subject ID, Session ID*

649 Comparisons: *MCA*

650 As a result, we tested 3 hypotheses across 6 MCA ex-  
651 periments and 3 reference experiments on 2 pipelines and 2  
652 perturbation modes, resulting in a total of 30 distinct tests.

### 653 Evaluating Graph-Theoretical Metrics

654 While connectomes may be used directly for some analyses,  
655 it is common practice to summarize them with structural mea-  
656 sures, which can then be used as lower-dimensional proxies  
657 of connectivity in so-called graph-theoretical studies<sup>11</sup>. We  
658 explored the stability of several commonly-used univariate  
659 (graphwise) and multivariate (nodewise or edgewise) features.  
660 The features computed and subsequent methods for compari-  
661 son in this section were selected to closely match those com-  
662 puted in<sup>10</sup>.

663 **Univariate Differences** For each univariate statistic (edge  
664 count, mean clustering coefficient, global efficiency, modu-  
665 larity of the largest connected component, assortativity, and  
666 mean path length) a distribution of values across all perturba-  
667 tions within subjects was observed. A Z-score was computed  
668 for each sample with respect to the distribution of feature  
669 values within an individual, and the proportion of "classically  
670 significant" Z-scores, i.e. corresponding to  $p < 0.05$ , was  
671 reported and aggregated across all subjects. The number of  
672 significant digits contained within an estimate derived from a  
673 single subject were calculated and aggregated.

674 **Multivariate Differences** In the case of both nodewise (de-  
675 gree distribution, clustering coefficient, betweenness central-  
676 ity) and edgewise (weight distribution, connection length) fea-  
677 tures, the cumulative density functions of their distributions  
678 were evaluated over a fixed range and subsequently aggre-



679 gated across individuals. The number of significant digits  
680 for each moment of these distributions (sum, mean, variance,  
681 skew, and kurtosis) were calculated across observations within  
682 a sample and aggregated.

### 683 **Evaluating A Brain-Phenotype Analysis**

684 Though each of the above approaches explores the instabil-  
685 ity of derived connectomes and their features, many modern  
686 studies employ modeling or machine-learning approaches, for  
687 instance to learn brain-phenotype relationships or identify dif-  
688 ferences across groups. We carried out one such study and ex-  
689 plored the instability of its results with respect to the upstream  
690 variability of connectomes characterized in the previous sec-  
691 tions. We performed the modeling task with a single sampled  
692 connectome per individual and repeated this sampling and  
693 modelling 20 times. We report the model performance for  
694 each sampling of the dataset and summarize its variance.

695 **BMI Classification** Structural changes have been linked to  
696 obesity in adolescents and adults<sup>41</sup>. We classified normal-  
697 weight and overweight individuals from their structural net-  
698 works (using for overweight a cutoff of  $BMI > 25^{13}$ ). We  
699 reduced the dimensionality of the connectomes through prin-  
700 cipal component analysis (PCA), and provided the first N-  
701 components to a logistic regression classifier for predicting  
702 BMI class membership, similar to methods shown in<sup>12,13</sup>.  
703 The number of components was selected as the minimum set  
704 which explained  $> 90\%$  of the variance when averaged across  
705 the training set for each fold within the cross validation of  
706 the original graphs; this resulted in a feature of 20 compo-  
707 nents. We trained the model using  $k$ -fold cross validation,  
708 with  $k = 2, 5, 10$ , and  $N$  (equivalent to leave-one-out; LOO).

### 709 **Data Availability**

710 The unprocessed dataset is available through The Consortium  
711 of Reliability and Reproducibility ([http://fcon\\_1000.projects.nitrc.org/indi/enhanced/](http://fcon_1000.projects.nitrc.org/indi/enhanced/)), including  
712 both the imaging data as well as phenotypic data which may  
713 be obtained upon submission and compliance with a Data Us-

715 age Agreement. The connectomes generated through simula-  
716 tions have been bundled and stored permanently (<https://doi.org/10.5281/zenodo.4041549>), and are made  
717 available through The Canadian Open Neuroscience Platform  
718 (<https://portal.conp.ca/search>, search term "Kiar").

### 720 **Code Availability**

721 All software developed for processing or evaluation is publicly  
722 available on GitHub at [https://github.com/gkpapers/](https://github.com/gkpapers/2020ImpactOfInstability)  
723 `2020ImpactOfInstability`. Experiments were launched  
724 using Boutiques<sup>42</sup> and Clowdr<sup>43</sup> in Compute Canada's HPC  
725 cluster environment. MCA instrumentation was achieved  
726 through Verificarlo<sup>9</sup> available on Github at <https://github.com/verificarlo/verificarlo>. A set of MCA in-  
727 strumented software containers is available on Github at <https://github.com/gkiar/fuzzy>.

### 730 **Author Contributions**

731 GK was responsible for the experimental design, data pro-  
732 cessing, analysis, interpretation, and the majority of writing.  
733 All authors contributed to the revision of the manuscript. YC,  
734 POC, and EP were responsible for MCA tool development and  
735 software testing. AR, GV, and BM contributed to experimen-  
736 tal design and interpretation. TG contributed to experimental  
737 design, analysis, and interpretation. TG and ACE were respon-  
738 sible for supervising and supporting all contributions made by  
739 GK. The authors declare no competing interests for this work.

### 740 **Acknowledgments**

741 This research was financially supported by the Natural Sci-  
742 ences and Engineering Research Council of Canada (NSERC)  
743 (award no. CGSD3-519497-2018). This work was also sup-  
744 ported in part by funding provided by Brain Canada, in partner-  
745 ship with Health Canada, for the Canadian Open Neuroscience  
746 Platform initiative.

### 747 **Additional Information**

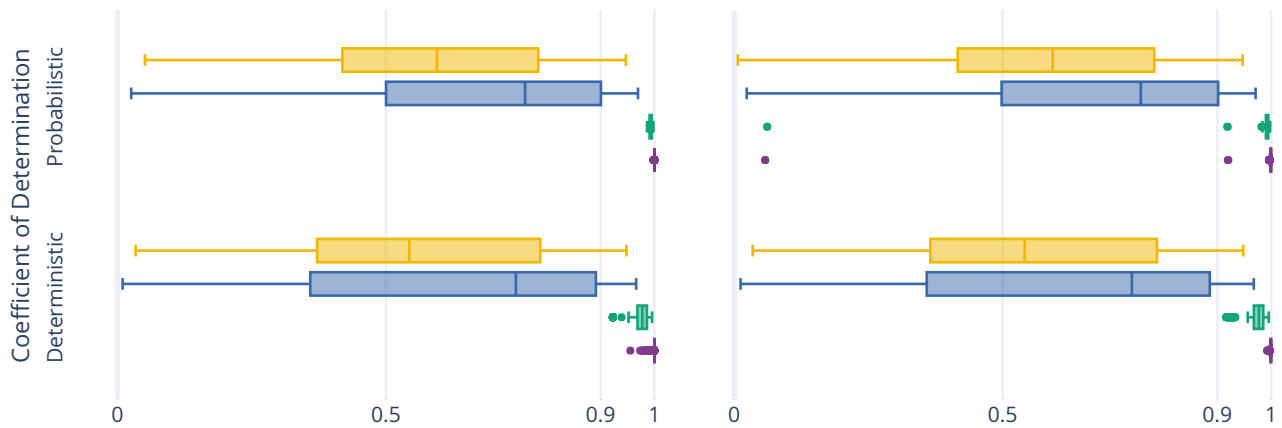
748 Supplementary Information is available for this paper. Corre-  
749 spondence and requests for materials should be addressed to

<sup>750</sup> Tristan Glatard at [tristan.glatard@concordia.ca](mailto:tristan.glatard@concordia.ca).

## S1. Graph Correlation

751  
752 The correlations between observed graphs (Figure S1) across each grouping follow the same trend to as percent deviation, as  
753 shown in Figure 1. However, notably different from percent deviation, there is no significant difference in the correlations  
754 between pipeline or input instrumentations. By this measure, the probabilistic pipeline is more stable in all cross-MCA and  
755 cross-directions except for the combination of input perturbation and cross-MCA ( $p < 0.0001$  for all; exploratory).

756 The marked lack in drop-off of performance across these settings, inconsistent with the measures show in Figure 1 is due  
757 to the nature of the measure and the graphs. Given that structural graphs are sparse and contain considerable numbers of  
758 zero-weighted edges, the presence or absence of an edge dominated the correlation measure where it was less impactful for the  
759 others. For this reason and others<sup>44</sup>, correlation is not a commonly used measure in the context of structural connectivity.



**Figure S1.** The correlation between perturbed connectomes and their reference.

760

## S2. Complete Discriminability Analysis

**Table S1.** The complete results from the Discriminability analysis, with results reported as mean  $\pm$  standard deviation Discriminability. As was the case in the condensed table, the alternative hypothesis, indicating significant separation across groups, was accepted for all experiments, with  $p < 0.005$ .

| Exp. | Subj. | Sess. | Samp. | Reference Execution |                 | Perturbed Pipeline |                 | Perturbed Inputs |                 |
|------|-------|-------|-------|---------------------|-----------------|--------------------|-----------------|------------------|-----------------|
|      |       |       |       | Det.                | Prob.           | Det.               | Prob.           | Det.             | Prob.           |
| 1.1  | All   | All   | 1     | 0.64 $\pm$ 0.00     | 0.65 $\pm$ 0.00 | 0.82 $\pm$ 0.00    | 0.82 $\pm$ 0.00 | 0.77 $\pm$ 0.00  | 0.75 $\pm$ 0.00 |
| 1.2  | All   | 1     | All   | 1.00 $\pm$ 0.00     | 1.00 $\pm$ 0.00 | 1.00 $\pm$ 0.00    | 1.00 $\pm$ 0.00 | 0.93 $\pm$ 0.02  | 0.90 $\pm$ 0.02 |
| 1.3  | All   | 1     | 1     |                     |                 | 1.00 $\pm$ 0.00    | 1.00 $\pm$ 0.00 | 0.94 $\pm$ 0.02  | 0.90 $\pm$ 0.02 |
| 2.4  | 1     | All   | All   | 1.00 $\pm$ 0.00     | 1.00 $\pm$ 0.00 | 1.00 $\pm$ 0.00    | 1.00 $\pm$ 0.00 | 0.88 $\pm$ 0.12  | 0.85 $\pm$ 0.12 |
| 2.5  | 1     | All   | 1     |                     |                 | 1.00 $\pm$ 0.00    | 1.00 $\pm$ 0.00 | 0.89 $\pm$ 0.11  | 0.84 $\pm$ 0.12 |
| 3.6  | 1     | 1     | All   |                     |                 | 0.99 $\pm$ 0.03    | 1.00 $\pm$ 0.00 | 0.71 $\pm$ 0.07  | 0.61 $\pm$ 0.05 |

761 The complete discriminability analysis includes comparisons across more axes of variability than the condensed version.  
 762 The reduction in the main body was such that only axes which would be relevant for a typical analysis were presented. Here,  
 763 each of Hypothesis 1, testing the difference across subjects, and 2, testing the difference across sessions, were accompanied  
 764 with additional comparisons to those shown in the main body.

765 **Subject Variation** Alongside experiment 1.1, that which mimicked a typical test-retest scenario, experiments 1.2 and 1.3  
 766 could be considered a test-retest with a handicap, given a single acquisition per individual was compared either across  
 767 subsamples or simulations, respectively. For this reason, it is unsurprising that the dataset achieved considerably higher  
 768 discriminability scores.

769 **Session Variation** Similar to subject variation, the session variation was also modelled across either both or a single  
 770 subsample. In both of these cases the performance was similar, and the finding that input perturbation reduced the off-target  
 771 signal was consistent.

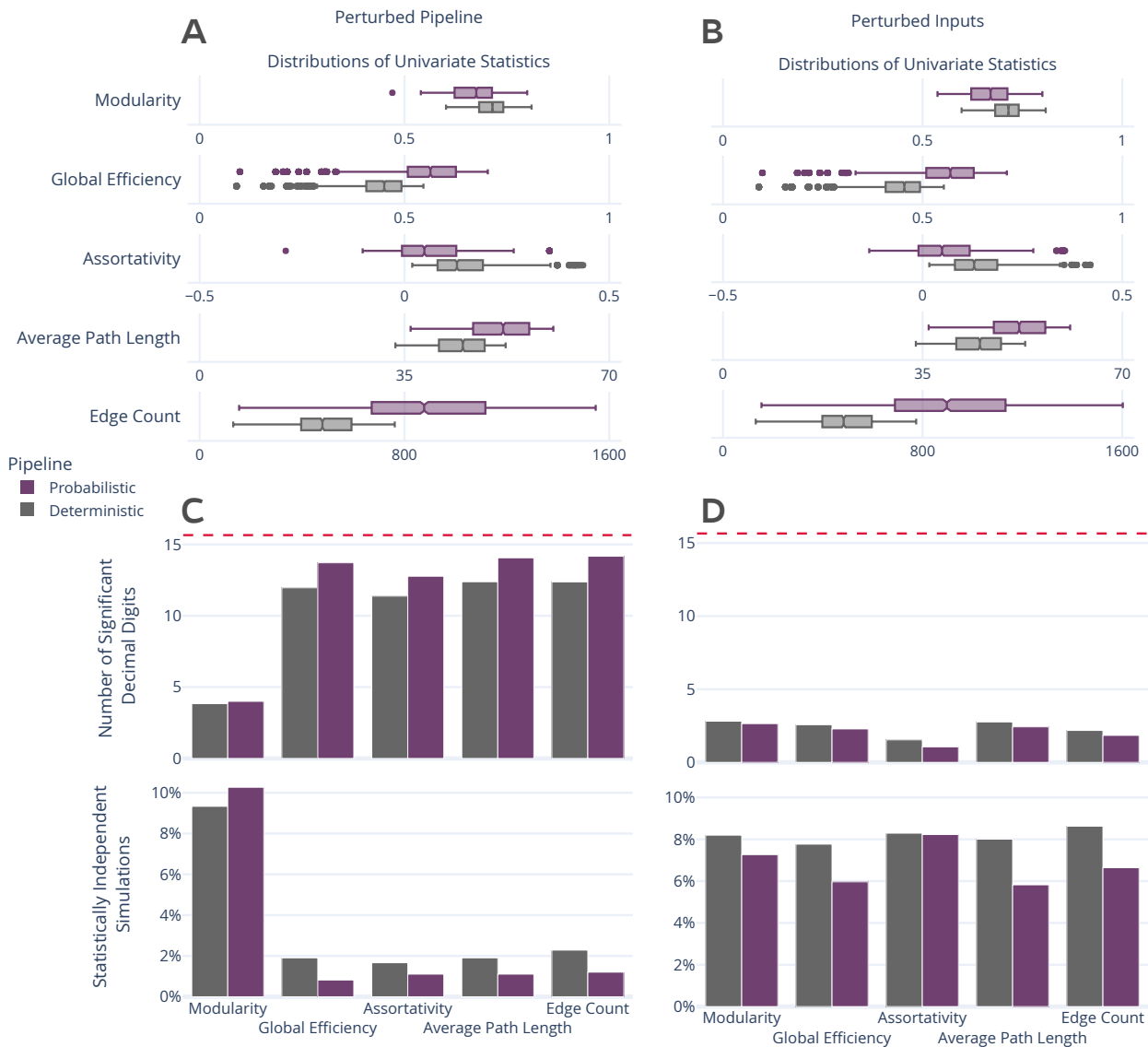
### S3. Univariate Graph Statistics

772  
773 Figure S2 explores the stability of univariate graph-theoretical metrics computed from the perturbed graphs, including modularity,  
774 global efficiency, assortativity, average path length, and edge count. When aggregated across individuals and perturbations, the  
775 distributions of these statistics (Figures S2A and S22B) showed no significant differences between perturbation methods for  
776 either deterministic or probabilistic pipelines.

777 However, when quantifying the stability of these measures across connectomes derived from a single session of data, the  
778 two perturbation methods show considerable differences. The number of significant digits in univariate statistics for Pipeline  
779 Perturbation instrumented connectome generation exceeded 11 digits for all measures except modularity, which contained  
780 more than 4 significant digits of information (Figure S2C). When detecting outliers from the distributions of observed statistics  
781 for a given session, the false positive rate (using a threshold of  $p = 0.05$ ) was approximately 2% for all statistics with the  
782 exception of modularity which again was less stable with an approximately 10% false positive rate. The probabilistic pipeline  
783 is significantly more stable than the deterministic pipeline ( $p < 0.0001$ ; exploratory) for all features except modularity. When  
784 similarly evaluating these features from connectomes generated in the input perturbation setting, no statistic was stable with  
785 more than 3 significant digits or a false positive rate lower than nearly 6% (Figure S2D). The deterministic pipeline was more  
786 stable than the probabilistic pipeline in this setting ( $p < 0.0001$ ; exploratory).

787 Two notable differences between the two perturbation methods are, first, the uniformity in the stability of the statistics,  
788 and second, the dramatic decline in stability of individual statistics in the input perturbation setting despite the consistency in  
789 the overall distribution of values. It is unclear at present if the discrepancy between the stability of modularity in the pipeline  
790 perturbation context versus the other statistics suggests the implementation of this measure is the source of instability or if it is  
791 implicit to the measure itself. The dramatic decline in the stability of features derived from input perturbed graphs despite no  
792 difference in their overall distribution both shows that while individual estimates may be unstable the comparison between  
793 aggregates or groups may be considered much more reliable; this finding is consistent with that presented for multivariate  
794 statistics.





**Figure S2.** Distribution and stability assessment of univariate graph statistics. **(A, B)** The distributions of each computed univariate statistic across all subjects and perturbations for Pipeline and Input settings, respectively. There was no significant difference between the distributions in A and B. **(C, D; top)** The number of significant decimal digits in each statistic across perturbations, averaged across individuals. The dashed red line refers to the maximum possible number of significant digits. **(C, D; bottom)** The percentage of connectomes which were deemed significantly different ( $p < 0.05$ ) from the others obtained for an individual.