

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Reproducible Preprint

- [Jupyter Book](#)

Code

- [Technical Screening](#)
- [Submitted Repository](#)

Reproducibility Assets

- [Repository](#)
- [Dataset](#)
- [Jupyter Book](#)
- [Container](#)

Moderator: [NeuroLibre](#)

Screener(s):

- [@roboneuro](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

# 1 Parcellating the parcellation issue - a proof of concept 2 for reproducible analyses using Neurolibre

3 **Pierre Bellec** <sup>1,2</sup>, **Saâd Jbabdi** <sup>3</sup>, and **R. Cameron Craddock** <sup>4</sup>

4 **1** Université de Montréal, Montréal, Canada **2** Centre de recherche de l'université de Montréal, Montréal, CA **3** University of Oxford, Oxford, UK **4** [brainhack.org](http://brainhack.org)



**THIS PDF IS INTENDED FOR CONTENT REGISTRATION PURPOSES ONLY! FOR FULL ACCESS AND INTERACTIVE READING OF THIS PUBLICATION, PLEASE VISIT [THE REPRODUCIBLE PREPRINT](#).**

## 7 Summary

8 Back in 2017, a special issue on the topic of **brain parcellation and segmentation** was published  
9 in the journal Neuroimage. We acted as guest editors for this special issue, and wrote an  
10 editorial ([Craddock et al., 2018](#)) providing an overview of all papers, sorted into categories.  
11 The categories were generated using a data-driven parcellation analysis, based on the words  
12 contained in the abstract of the articles. This jupyter book will allow interested readers to  
13 reproduce this analysis, as a proof of concept for reproducible publications using **jupyter books**  
14 and the **Neurolibre** preprint server.

## 15 Acknowledgements

16 NeuroLibre is sponsored by the Canadian Open Neuroscience Platform (CONP), Brain Canada,  
17 Cancer Computers, the Courtois foundation, the Quebec Bioimaging Network, and Healthy  
18 Brains for Healthy Life.



## 19 NOTE

20 **NOTE:** The following section in this document repeats the narrative content  
21 exactly as found in the **corresponding NeuroLibre Reproducible Preprint (NRP)**.  
22 The content was automatically incorporated into this PDF using the NeuroLibre  
23 publication workflow ([Karakuzu et al., 2022](#)) to credit the referenced resources.  
24 The submitting author of the preprint has verified and approved the inclusion of  
25 this section through a GitHub pull request made to the **source repository** from  
26 which this document was built. Please note that the figures and tables have been  
27 excluded from this (static) document. **To interactively explore such outputs and**  
28 **re-generate them, please visit the corresponding NRP.** For more information on  
29 integrated research objects (e.g., NRPs) that bundle narrative and executable  
30 content for reproducible and transparent publications, please refer to DuPre et al.  
31 ([2022](#)). NeuroLibre is sponsored by the Canadian Open Neuroscience Platform  
32 (CONP) ([Harding et al., 2023](#)).

## 33 **Text mining**

### 34 **List of papers**

35 We first assembled the title, the name of the corresponding author, and the abstract for all the  
36 articles into a tabular-separated values (tsv) file, which we publicly archived on [Figshare](#). We  
37 use the [Repo2Data](#) tool developed by the NeuroLibre team to collect these data and include  
38 them in our reproducible computational environment.

### 39 **Word features**

40 For each paper, we used [scikit-learn](#) ([Kramer, 2016](#)) to extract a bag of words representation  
41 for each abstract, picking on the 300 most important terms seen across all articles based on  
42 a term frequency-inverse document frequency ([tf-idf](#)) [index](#). Following that, a special value  
43 decomposition was used to further reduce the dimensionality of the abstracts to 10 components.  
44 We ended up with a component matrix of dimension 38 (articles) times 10 (abstract text  
45 components). The distribution of each of the 38 articles across the 10 components is represented  
46 below. Note how some articles have particular high loadings on specific components, suggesting  
47 these may capture particular topics. Rather than visually inspect the component loadings  
48 to group paper ourselves, we are going to resort to an automated parcellation (clustering)  
49 technique.

### 50 **Parcellate the papers**

51 Now that the content of each paper has been condensed into only 10 (hopefully informative)  
52 numbers, we can run these features into a trusted, classic parcellation algorithm: Ward's  
53 agglomerative hierarchical clustering, as implemented in the [scipy](#) library. We cut the hierarchy  
54 to extract 7 "paper parcels", and also use the hierarchy to re-order the papers, such that  
55 similar papers are close in order, as illustrated in a dendrogram representation.

### 56 **Similarity matrix**

57 So, to get a better feel of the similarity between papers that was fed into the clustering  
58 procedure, we extracted the 38x38 (papers x papers) correlation matrix across features. Papers  
59 are re-ordered in the matrix according to the above hierarchy. Each "paper parcel" has been  
60 indicated by a white square along the diagonal, which represents the similarity measures  
61 between papers falling into the same parcel.

### 62 **Word cloud**

63 Now, each paper of the special issue has been assigned to one and only one out of 7 possible  
64 "paper parcel". For each paper parcel, we can evaluate which words contribute more to the  
65 dominant component associated with that parcel.

### 66 **Categories**

67 Thanks to the word clouds, these simple data-driven categories turned out to be fairly easily  
68 interpretable. For example, the word cloud of the category number 4 features prominently  
69 words like "white", "matter" and "bundles". If we examine the exact list of papers included in  
70 this category, we see that it is composed of four papers, which all considered parcels derived  
71 from white matter bundles with diffusion imaging. We can also check the distribution of  
72 component loadings for this category alone. As expected, there is a certain similarity in the  
73 component loadings for these papers, in particular along component 4:

## 74 References

- 75 Craddock, R. C., Bellec, P., & Jbabdi, S. (2018). Neuroimage special issue on brain seg-  
76 mentation and parcellation - editorial. *Neuroimage*, *170*, 1–4. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.neuroimage.2017.11.063)  
77 [neuroimage.2017.11.063](https://doi.org/10.1016/j.neuroimage.2017.11.063)
- 78 DuPre, E., Holdgraf, C., Karakuzu, A., Tetrel, L., Bellec, P., Stikov, N., & Poline, J.-B. (2022).  
79 Beyond advertising: New infrastructures for publishing integrated research objects. *PLOS*  
80 *Computational Biology*, *18*(1), e1009651. <https://doi.org/10.1371/journal.pcbi.1009651>
- 81 Harding, R. J., Bermudez, P., Bernier, A., Beauvais, M., Bellec, P., Hill, S., Karakuzu, A.,  
82 Knoppers, B. M., Pavlidis, P., Poline, J.-B., Roskams, J., Stikov, N., Stone, J., Strother, S.,  
83 Consortium, C., & Evans, A. C. (2023). The Canadian Open Neuroscience Platform—An  
84 open science framework for the neuroscience community. *PLOS Computational Biology*,  
85 *19*(7), 1–14. <https://doi.org/10.1371/journal.pcbi.1011230>
- 86 Karakuzu, A., DuPre, E., Tetrel, L., Bermudez, P., Boudreau, M., Chin, M., Poline, J.-B.,  
87 Das, S., Bellec, P., & Stikov, N. (2022). *NeuroLibre : A preprint server for full-fledged*  
88 *reproducible neuroscience*. OSF Preprints. <https://doi.org/10.31219/osf.io/h89js>
- 89 Kramer, O. (2016). Scikit-learn. In *Machine learning for evolution strategies* (pp. 45–53).  
90 Springer International Publishing. [https://doi.org/10.1007/978-3-319-33383-0\\_5](https://doi.org/10.1007/978-3-319-33383-0_5)

DRAFT