

BIDSonym: a BIDS App for the pseudo-anonymization of neuroimaging datasets

Peer Herholz¹, Rita Marie Ludwig², and Jean-Baptiste Poline¹

¹ NeuroDataScience-Origami lab, McConnell Brain Imaging Centre, The Neuro (Montreal Neurological Institute-Hospital), Faculty of Medicine, McGill University, Montreal, Quebec, Canada
² University of Oregon, Eugene, United States

DOI: [10.21105/joss.03169](https://doi.org/10.21105/joss.03169)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Olivia Guest](#) ↗

Reviewers:

- [@deep-introspection](#)
- [@chrisgorgo](#)
- [@neuromusic](#)

Submitted: 11 March 2021

Published: 07 June 2021

License

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

Statement of Need

Due to the evolution of research incentives, technical advancements and the development of new standards ([Eickhoff et al., 2016](#); [Gorgolewski et al., 2016](#); [Nichols et al., 2017](#); [Poldrack et al., 2013](#); [Poldrack & Gorgolewski, 2017, 2014](#)), increasingly greater amounts of neuroimaging data are being shared either publicly or made available through data user agreements. These datasets originate from small samples of participants collected by individual research groups, as well as from “Big Data” samples including thousands of participants collected by large research consortia (UK Biobank ([Sudlow et al., 2015](#)), HCP ([Van Essen et al., 2013](#)), ABIDE ([Di Martino et al., 2014](#)), ADNI ([Mueller et al., 2005](#)), etc.). While data sharing is important and beneficial ([Eickhoff et al., 2016](#); [Nichols et al., 2017](#); [Poldrack & Gorgolewski, 2014](#); [Poline et al., 2012](#)), privacy of participant data must be protected ([Bannier et al., 2020](#); [Brakewood & Poldrack, 2013](#)). To that end, Ethic Review Boards and data sharing platforms typically require that uploaded datasets are provided in anonymized or pseudo-anonymized form, limiting participant reidentification. However, the (pseudo-) anonymization process is deceptively complex; attempts at ensuring data privacy must take into consideration all dataset components, including imaging modalities, as well as national legal and ethical frameworks. Several algorithms have been developed to (pseudo-) anonymize imaging datasets but they offer limited solutions. Some are attached to specific software, some are limited to specific computing environments; most miss an in-depth assessment and treatment of the metadata attached to the dataset, or lack the capacity to automatize (pseudo-) anonymization across large datasets. BIDSonym was created to address these points in one simple, flexible, and general tool that offers users an array of automated (pseudo-)anonymization options to augment participant privacy in neuroimaging datasets. There are two components of neuroimaging datasets that arguably pose the largest risk to maintaining participant privacy: the structural images and accompanying metadata (e.g. metadata text files or information embedded in image file headers). Structural images contain visible identifiable participant information via facial features like the eyes, nose, and mouth, and privacy is usually addressed through a process called “defacing” within which all or a subset of these features are removed from the final structural data files. The metadata text files may additionally contain identifiable participant data through the recording of acquisition time and location, and personal details such as date of birth, height, and weight. Here, privacy is maintained by removing or blurring this information from the final dataset. BIDSonym addresses both vulnerabilities in neuroimaging datasets, obviating the need for multiple steps within a data sharing pipeline to ensure participant privacy.

Summary

In concordance with the BIDS-App template ([Gorgolewski et al., 2017](#)), BIDSonym operates as a command line tool written in Python ([Rossum, 1995](#)) and is intended to run in its containerized version (either using Docker (<https://www.docker.com>) or Singularity (<https://sylabs.io>),

45 providing all necessary software dependencies. However, it is also available as a Python pack-
46 age via PyPi (<https://pypi.org>) to facilitate reuse in a development environment. BIDSonym
47 expects BIDS datasets (Gorgolewski et al., 2016) and provides three core functionalities as
48 depicted in Figure 1: defacing of structural (i.e. T1 and T2 weighted images, adaptation of
49 potentially sensitive metadata information, and evaluation of (pseudo-) anonymization results.

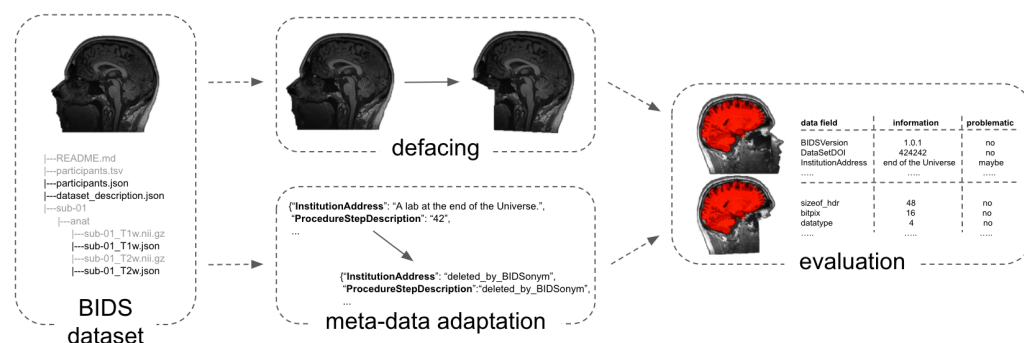


Figure 1: Overview of BIDSonym's functionality | Providing a dataset in BIDS as input, structural images are defaced, meta-data fields adapted as requested and the performance of the defacing, as well as all meta-data fields (in both the json sidcar files and image headers) evaluated.

50 Focusing on the first aspect, BIDSonym supports a multitude of commonly used defacing
51 algorithms and tools, including pydeface (Gulban et al., 2019), mri_deface (Bischoff-Grethe
52 et al., 2007), quickshear (Schimke & Hale, 2011), and mrdefacer (Hanke & Halchenko, 2018).
53 Based on the chosen tool, facial features of the structural images of either specified participants
54 or the whole data set are then removed. Figure 2 provides an example for the different
55 algorithms and tools available for defacing a structural T1 weighted image. Furthermore,
56 structural images from other modalities (e.g. T2 weighted) can also be defaced in which case
57 the defaced T1 weighted structural image will be utilized as a deface mask. In order to account
58 for possible errors during the defacing process, BIDSonym moves the non-defaced original
59 structural images into a distinct directory before defacing, allowing users to test multiple
60 defacing settings and/or options. Following BIDS (Gorgolewski et al., 2016), those files are
61 moved to a different directory ('/sourcedata') and a description identifier ('*_desc-nondeid*')
62 is added to the filenames.

63 A comparable behavior is implemented with regard to metadata information in BIDSonym's
64 second core functionality. The information in the metadata files that accompany neuroimaging
65 data and in the headers of neuroimaging data files will be gathered and listed within a tabular
66 file ('*.tsv'). If specified, the extracted information will then be queried for potentially sensitive
67 information (e.g. name, date or place of birth, etc.) and marked accordingly. Additionally,
68 users can specify that certain information should be deleted from metadata files, in which case
69 they will be moved and renamed as described for the neuroimaging data.

70 BIDSonym's third core functionality implements quality control assessments of the defacing
71 results and of the information present in the data. Concerning the first, this includes (inter-
72 active) plots that allow to evaluate if the applied algorithm and settings were too stringent
73 (e.g. removing voxels belonging to the brain). Regarding the second, information present in
74 the meta-data files and image headers are gathered in tabular format within respective files
75 ('*.tsv'). Each table contains all key-value pairs present in a given file and aims to provide an
76 assessment of potentially sensitive information.

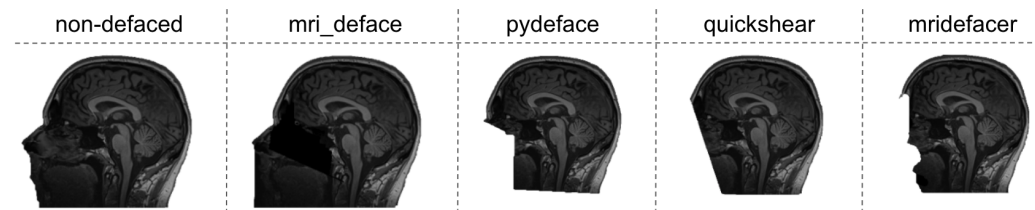


Figure 2: Defacing examples | Results of the different algorithms and tools (columns) included in BIDSonym, displayed in comparison to the corresponding original structural image (most left).

77 More information on BIDSonym's workflow, the corresponding processing steps and outcomes,
78 as well as installation instructions can be found in the respective documentation ([https://](https://peerherholz.github.io/BIDSonym)
79 peerherholz.github.io/BIDSonym) and GitHub repository ([https://github.com/PeerHerholz/](https://github.com/PeerHerholz/BIDSonym)
80 [BIDSonym](https://github.com/PeerHerholz/BIDSonym)). BIDSonym provides a straightforward and flexible way to pseudo-anonymize
81 neuroimaging datasets by a variety of means, operating on both small and large datasets
82 through its implementation following the BIDS-App template (Gorgolewski et al., 2017).
83 BIDSonym depends on the nibabel (Brett et al., 2020), nipy (Gorgolewski et al., 2011),
84 Nilearn (Abraham et al., 2014), pybids (Yarkoni et al., 2019) and pandas (McKinney, 2010)
85 python packages (all are well maintained and tested) and is licensed under the BSD-3 li-
86 cense (<https://opensource.org/licenses/BSD-3-Clause>). As data sharing becomes more widely
87 adopted, BIDSonym fills an important gap for the neuroimaging community.

88 Acknowledgements

89 P.H. and J.B.P. were supported in parts by funding from the Canada First Research Excellence
90 Fund, awarded to McGill University for the Healthy Brains for Healthy Lives initiative, the
91 National Institutes of Health (NIH) NIH-NIBIB P41 EB019936 (ReproNim), as well as the
92 National Institute of Mental Health of the NIH under Award Number R01MH096906. P.H. was
93 additionally supported by research scholar award from Brain Canada, in partnership with Health
94 Canada, for the Canadian Open Neuroscience Platform initiative. This project originated as
95 part of Neurohackademy which is funded by the National Institute of Mental Health through
96 a grant to Ariel Rokem and Tal Yarkoni (R25MH112480). Finally, all the contributors listed
97 in the project's Zenodo and GitHub repository have contributed code and intellectual labor to
98 further improve BIDSonym. The same holds true for users that reported issues and continue
99 to do so.

100 References

- 101 Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort,
102 A., Thirion, B., & Varoquaux, G. (2014). Machine learning for neuroimaging with scikit-
103 learn. *Frontiers in Neuroinformatics*, 8. <https://doi.org/10.3389/fninf.2014.00014>
- 104 Bannier, E., Barker, G., Borghesani, V., Broeckx, N., Clement, P., Vaya, M. de la I., Emblem,
105 K. E., Ghosh, S., Glerean, E., Gorgolewski, K. J., Havu, M., Halchenko, Y. O., Herholz,
106 P., Hespel, A., Heunis, S., Hu, Y., Chuan-Peng, H., Huijser, D., Jancalek, R., ... Zhu, H.
107 (2020). *The Open Brain Consent: Informing research participants and obtaining consent*
108 *to share brain imaging data*. PsyArXiv. <https://doi.org/10.31234/osf.io/f6mnp>
- 109 Bischoff-Grethe, A., Ozyurt, I. B., Busa, E., Quinn, B. T., Fennema-Notestine, C., Clark, C.
110 P., Morris, S., Bondi, M. W., Jernigan, T. L., Dale, A. M., Brown, G. G., & Fischl, B.
111 (2007). A Technique for the Deidentification of Structural Brain MR Images. *Human*
112 *Brain Mapping*, 28(9), 892–903. <https://doi.org/10.1002/hbm.20312>

- 113 Brakewood, B., & Poldrack, R. A. (2013). The ethics of secondary data analysis: Considering
114 the application of Belmont principles to the sharing of neuroimaging data. *NeuroImage*,
115 82, 671–676. <https://doi.org/10.1016/j.neuroimage.2013.02.040>
- 116 Brett, M., Markiewicz, C. J., Hanke, M., Côté, M.-A., Cipollini, B., McCarthy, P., Jarecka,
117 D., Cheng, C. P., Halchenko, Y. O., Cottaar, M., Ghosh, S., Larson, E., Wassermann,
118 D., Gerhard, S., Lee, G. R., Wang, H.-T., Kastman, E., Kaczmarzyk, J., Guidotti, R., ...
119 freec84. (2020). *Nibabel*. Zenodo. <https://doi.org/10.5281/zenodo.3757992>
- 120 Di Martino, A., Yan, C.-G., Li, Q., Denio, E., Castellanos, F. X., Alaerts, K., Anderson, J.
121 S., Assaf, M., Bookheimer, S. Y., Dapretto, M., Deen, B., Delmonte, S., Dinstein, I.,
122 Ertl-Wagner, B., Fair, D. A., Gallagher, L., Kennedy, D. P., Keown, C. L., Keysers, C.,
123 ... Milham, M. P. (2014). The autism brain imaging data exchange: Towards a large-
124 scale evaluation of the intrinsic brain architecture in autism. *Molecular Psychiatry*, 19(6),
125 659–667. <https://doi.org/10.1038/mp.2013.78>
- 126 Eickhoff, S., Nichols, T. E., Van Horn, J. D., & Turner, J. A. (2016). Sharing the wealth:
127 Neuroimaging data repositories. *NeuroImage*, 124(Pt B), 1065–1068. <https://doi.org/10.1016/j.neuroimage.2015.10.079>
- 129 Gorgolewski, K. J., Alfaro-Almagro, F., Auer, T., Bellec, P., Capotă, M., Chakravarty, M. M.,
130 Churchill, N. W., Cohen, A. L., Craddock, R. C., Devenyi, G. A., Eklund, A., Esteban,
131 O., Flandin, G., Ghosh, S. S., Guntupalli, J. S., Jenkinson, M., Keshavan, A., Kiar, G.,
132 Liem, F., ... Poldrack, R. A. (2017). BIDS apps: Improving ease of use, accessibility, and
133 reproducibility of neuroimaging data analysis methods. *PLOS Computational Biology*,
134 13(3), e1005209. <https://doi.org/10.1371/journal.pcbi.1005209>
- 135 Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P., Flandin,
136 G., Ghosh, S. S., Glatard, T., Halchenko, Y. O., Handwerker, D. A., Hanke, M., Keator,
137 D., Li, X., Michael, Z., Maumet, C., Nichols, B. N., Nichols, T. E., Pellman, J., ...
138 Poldrack, R. A. (2016). The brain imaging data structure, a format for organizing and
139 describing outputs of neuroimaging experiments. *Scientific Data*, 3(1), 1–9. <https://doi.org/10.1038/sdata.2016.44>
- 141 Gorgolewski, K. J., Burns, C. D., Madison, C., Clark, D., Halchenko, Y. O., Waskom, M.
142 L., & Ghosh, S. S. (2011). Nipype: A Flexible, Lightweight and Extensible Neuroimaging
143 Data Processing Framework in Python. *Frontiers in Neuroinformatics*, 5. <https://doi.org/10.3389/fninf.2011.00013>
- 145 Gulban, O. F., Nielson, D., Poldrack, R., Lee, J., Gorgolewski, K. J., Sochat, V., & Ghosh,
146 S. (2019). *Pydeface*. <http://doi.org/10.5281/zenodo.3524401>
- 147 Hanke, M., & Halchenko, Y. (2018). *Mridefacer*. <https://github.com/mih/mridefacer>
- 148 McKinney, W. (2010). Data structures for statistical computing in python. *Proceedings of*
149 *the 9th Python in Science Conference*, 445, 51–56.
- 150 Mueller, S. G., Weiner, M. W., Thal, L. J., Petersen, R. C., Jack, C., Jagust, W., Tro-
151 janowski, J. Q., Toga, A. W., & Beckett, L. (2005). The Alzheimer's Disease Neu-
152 roimaging Initiative. *Neuroimaging Clinics of North America*, 15(4), 869–877. <https://doi.org/10.1016/j.nic.2005.09.008>
- 154 Nichols, T. E., Das, S., Eickhoff, S. B., Evans, A. C., Glatard, T., Hanke, M., Kriegeskorte,
155 N., Milham, M. P., Poldrack, R. A., Poline, J.-B., Proal, E., Thirion, B., Van Essen,
156 D. C., White, T., & Yeo, B. T. T. (2017). Best practices in data analysis and sharing
157 in neuroimaging using MRI. *Nature Neuroscience*, 20(3), 299–303. <https://doi.org/10.1038/nn.4500>
- 159 Poldrack, R. A., Barch, D. M., Mitchell, J., Wager, T., Wagner, A. D., Devlin, J. T., Cumba,
160 C., Koyejo, O., & Milham, M. (2013). Toward open sharing of task-based fMRI data: The

- 161 OpenfMRI project. *Frontiers in Neuroinformatics*, 7. <https://doi.org/10.3389/fninf.2013.00012>
 162
- 163 Poldrack, R. A., & Gorgolewski, K. J. (2017). OpenfMRI: Open sharing of task fMRI data.
 164 *NeuroImage*, 144, 259–261. <https://doi.org/10.1016/j.neuroimage.2015.05.073>
- 165 Poldrack, R. A., & Gorgolewski, K. J. (2014). Making big data open: Data sharing in
 166 neuroimaging. *Nature Neuroscience*, 17(11), 1510–1517. <https://doi.org/10.1038/nn.3818>
 167
- 168 Poline, J.-B., Breeze, J. L., Ghosh, S. S., Gorgolewski, K. J., Halchenko, Y. O., Hanke, M.,
 169 Helmer, K. G., Marcus, D. S., Poldrack, R. A., Schwartz, Y., Ashburner, J., & Kennedy,
 170 D. N. (2012). Data sharing in neuroimaging research. *Frontiers in Neuroinformatics*, 6.
 171 <https://doi.org/10.3389/fninf.2012.00009>
- 172 Rossum, G. (1995). *Python reference manual* [Technical Report]. CWI (Centre for Mathe-
 173 matics; Computer Science).
- 174 Schimke, N., & Hale, J. (2011). Quickshear defacing for neuroimages. *Proceedings of the*
 175 *2nd USENIX Conference on Health Security and Privacy*, 11.
- 176 Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P.,
 177 Green, J., Landray, M., Liu, B., Matthews, P., Ong, G., Pell, J., Silman, A., Young, A.,
 178 Sprosen, T., Peakman, T., & Collins, R. (2015). UK Biobank: An Open Access Resource
 179 for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age.
 180 *PLoS Medicine*, 12(3). <https://doi.org/10.1371/journal.pmed.1001779>
- 181 Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E. J., Yacoub, E., & Ugurbil,
 182 K. (2013). The WU-Minn Human Connectome Project: An overview. *NeuroImage*, 80,
 183 62–79. <https://doi.org/10.1016/j.neuroimage.2013.05.041>
- 184 Yarkoni, T., Markiewicz, C. J., Vega, A. de la, Gorgolewski, K. J., Salo, T., Halchenko, Y. O.,
 185 McNamara, Q., DeStasio, K., Poline, J.-B., Petrov, D., Hayot-Sasson, V., Nielson, D. M.,
 186 Carlin, J., Kiar, G., Whitaker, K., DuPre, E., Wagner, A., Tirrell, L. S., Jas, M., ... Blair,
 187 R. (2019). PyBIDS: Python tools for BIDS datasets. *Journal of Open Source Software*,
 188 4(40), 1294. <https://doi.org/10.21105/joss.01294>