



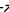
# datashuttle: automated data management for experimental neuroscience

Joseph J. Ziminski<sup>1</sup>, Nikoloz Sirmipilatze<sup>1</sup>, Brandon D. Peri<sup>2</sup>, Shrey Singh<sup>3</sup>, Sepiedeh Keshavarzi<sup>2</sup>, and Adam L. Tyson<sup>1</sup>

<sup>1</sup> Neuroinformatics Unit, Sainsbury Wellcome Centre & Gatsby Computational Neuroscience Unit, University College London, London, U.K <sup>2</sup> Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge, United Kingdom <sup>3</sup> Netaji Subhas University of Technology, New Delhi, India

DOI: [DOIunavailable](#)

## Software

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

Editor: [Pending Editor](#) 

Reviewers:

- [@Pending Reviewers](#)

Submitted: N/A

Published: N/A

## License

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

## Summary

Datashuttle is a Python package that facilitates data standardisation in neuroscience. Experimental data are often stored using custom folder structures and naming conventions, which hinders data sharing, reproducibility, and the development of community tools. Datashuttle addresses this by providing user-friendly tools to create, validate, and transfer standardised experimental data folders. The package can be used programmatically—integrated into existing Python scripts for data acquisition—or via a graphical user interface.

## Statement of Need

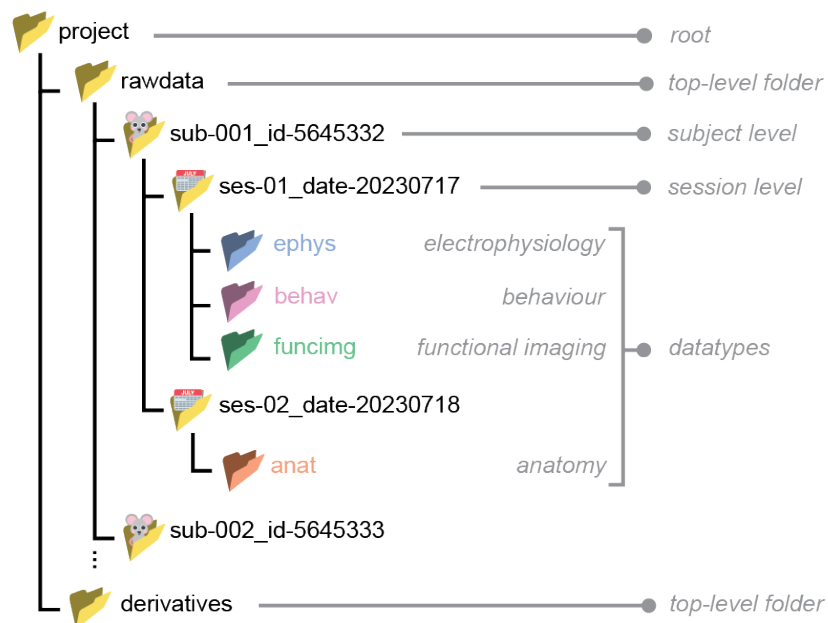
The past decade has seen significant progress in the development of neuroscience data standards. Experimental datasets have become increasingly complex, with multiple modalities (e.g. behaviour, electrophysiology and imaging) often collected from a single subject. At its core, standardisation facilitates reproducibility by ensuring these complex datasets are well organised, accessible, machine-readable and sufficiently documented (Martone, 2024). This standardisation permits robust, automated project management including the transfer of experimental data between machines and validation of project contents. Detailed specifications covering folder, file and metadata naming and structural conventions are required to achieve this goal.

Development and dissemination of comprehensive standards has been driven by community organisations such as the International Neuroinformatics Coordinating Facility (INCF) (Abrams et al., 2022). This includes adoption of the FAIR principles (Abrams et al., 2022; Wilkinson et al., 2016) ensuring data are Findable, Accessible, Interoperable and Reusable. Important standardisation initiatives include the Brain Imaging Dataset Structure (BIDS) (Gorgolewski et al., 2016), a file, folder and metadata standard widely used in neuroimaging, and the open file format Neurodata Without Borders (Rübel et al., 2022). These initiatives aim to achieve ‘full’ standardisation that enables automated analysis of machine-readable experimental datasets.

However, the adoption of data standards in systems neuroscience is not yet widespread (Klingner et al., 2023). This is due in part to the inherent, and necessary, complexity required to achieve full standardisation (Pierré et al., 2024) and lack of user-friendly tools to automate management of standardised projects without requiring coding experience. Further, not all systems neuroscience methods have a corresponding data specification (e.g. fibre photometry). Researchers often default to custom folder structures in lieu of full standardisation, leading to inconsistencies both across and within laboratories.

Datashuttle aims to bridge the gap between ‘no standardisation’ and full standardisation by implementing a simple specification called ‘NeuroBlueprint’ (Ziminski et al., 2025) ([Figure 1](#)).

NeuroBlueprint mandates only folder naming and structure conventions, while recommending file and metadata-naming schemes. It is designed to be easy to adopt, meaning it is suitable for the busy data-acquisition stages of a project in which applying full standardisation is often too onerous. The structure and format are heavily inspired by BIDS, in order to reduce redundancy across specifications and facilitate later transition to this more comprehensive schema. This means that while NeuroBlueprint is not sufficiently standardised to ensure data are FAIR, it provides an easy-to-use starting point that requires relatively little effort to adopt.



**Figure 1:** The NeuroBlueprint specification. Raw data (i.e. as collected from acquisition machines) are organised hierarchically by subject, session, and datatype. Subject and session names consist of key-value pairs. Only the sub- and ses-keys are required and others are optional. Acquired data are placed in the datatype folder, with valid datatype names defined in the specification. Derived data are stored in the top-level derivatives folder, and while not mandated, it is advised to organise these similar to the rawdata directory.

Datashuttle automates the creation, validation and transfer of experimental folders in NeuroBlueprint standard. It is designed to drop into existing scripted or manual data-acquisition pipelines, ensuring standardisation at the point of data collection. Datashuttle offers flexible data transfer capabilities that make standardisation practical and convenient, rather than an added burden. With minimal dependencies and no lock-in, it provides a lightweight, adaptable solution for managing neuroscience project workflows.

## Features

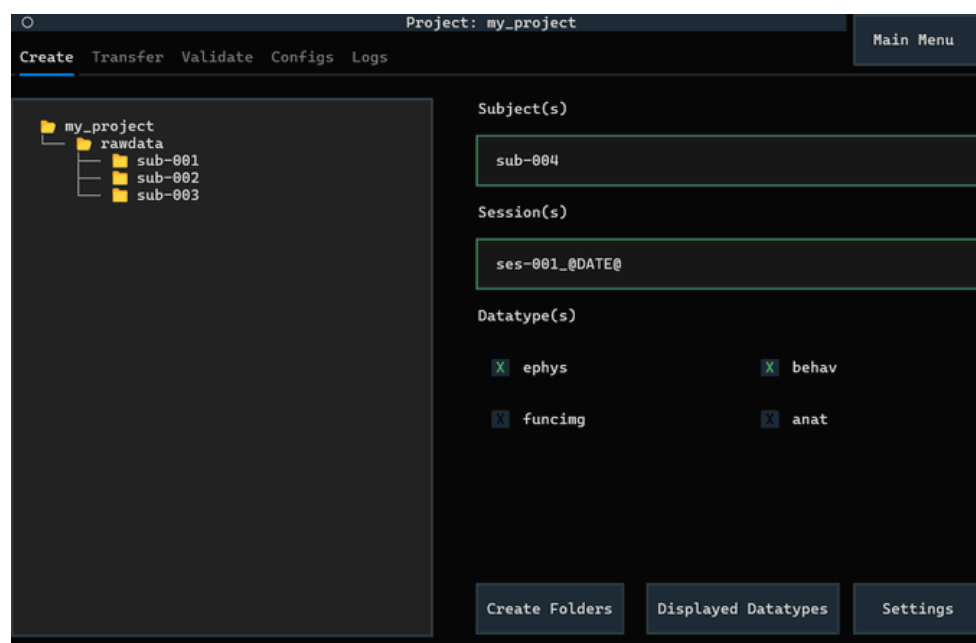
Datashuttle can be installed via the package manager conda. While pip installation is supported, the non-Python dependency RClone (used to manage data transfers) must be installed separately. The cross-platform terminal user interface (TUI) is built with Textual (McGugan, 2021) and can be used in the system terminal.

The typical workflow begins with researchers creating standardised folders at the start of each experimental session. Data generated during acquisition (e.g. from cameras, behaviour-monitoring devices or electrophysiology probes) are saved into the created folders. Real-time validation features ensure that common errors such as duplicate subject or session IDs

are caught immediately. At the end of the experimental session, data are transferred to the laboratory's central storage. Transfers can be made to a remote server either via a mounted drive or SSH, while cloud services such as Google Drive and AWS S3 Buckets are also supported.

## Folder Creation

NeuroBlueprint-formatted folder trees can be created for a given subject, session and datatype (e.g. 'behav' for behaviour), with online validation to reduce the likelihood of errors in user input (Figure 2).



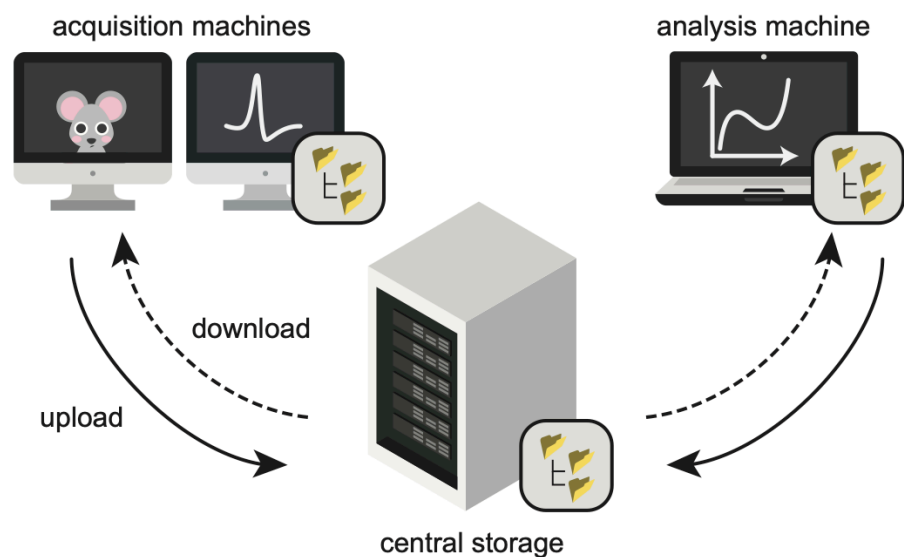
**Figure 2:** The Create Folders screen. Subject, session and datatype folders can be created through this interface. The text input border provides real-time validation results, while tags such as @DATE@ can be used to auto-format the system date. The current project is displayed on the left-hand directory tree, which can be used to copy file-paths and open the operating filesystem.

## Validation

Validation catches issues such as duplicate subject or session IDs, inconsistent number of leading zeros, bad key-value pair formatting and other common typographical errors. Validation can be performed on an entire project, listing any formatting errors that are discovered. Additionally, real-time validation during folder creation ensures no non-NeuroBlueprint format folders can be made. Custom extensions to the validation can be added, with subject and session names validated against user-defined regular expression templates.

## Data Transfer

Datashuttle uses the open source tool RClone (Craig-Wood, 2014) to perform data transfers. Experimental data can be 'uploaded' (from the local machine to central storage) or 'downloaded' (from the central storage to the local machine) (Figure 3). A benefit of standardisation is machine-readable folder names—meaning it is simple to select arbitrary subsets of data for transfer e.g. only the first five subjects.



**Figure 3:** Data transfers in datashuttle. A typical workflow involves transferring data from an acquisition machine to a central laboratory storage. Later, the entire dataset or subsets of it (e.g. only electrophysiology data) may be downloaded to an analysis machine for processing.

## Logging

In order to track full provenance of the project, datashuttle operations are logged to file with `fancylog` (Ziminski & Tyson, 2025). Logs can be accessed directly from disk or displayed in the TUI.

## Future Directions

Datashuttle will continue to evolve alongside the NeuroBlueprint specification, implementing upcoming extensions as they emerge. While Datashuttle does not currently support a meta-data standard, this will be a key focus for future development to enable improved validation and automation.

Currently, NeuroBlueprint is designed for experiments in which subjects go through the experimental procedures individually. However multi-animal experiments investigating social behaviours, in which animals interact during experimental sessions, are a growing area of neuroscience research. Future updates to both NeuroBlueprint and datashuttle will aim to support this use case.

## Availability

Datashuttle's source code is available at <https://github.com/neuroinformatics-unit/datashuttle> and documentation published at <https://datashuttle.neuroinformatics.dev>.

## Acknowledgements

J.J.Z., N.S. and A.L.T. were funded by the core grant to the Sainsbury Wellcome Centre (Wellcome - 219627/Z/19/Z, Gatsby Charitable Foundation - GAT3755). A.L.T. was funded

by the core grant to the Gatsby Computational Neuroscience Unit (Gatsby Charitable Foundation - GAT3850). B.P. and S.K. are funded by a Wellcome Trust Career Development Award (226039/Z/22/Z to S.K.). We also thank Laura Porta, Alessandro Felder and Igor Tatarnikov for comments on the manuscript and codebase and all contributors to the datashuttle project.

## References

- Abrams, M. B., Bjaalie, J. G., Das, S., Egan, G. F., Ghosh, S. S., Goscinski, W. J., Grethe, J. S., Koteleski, J. H., Ho, E. T. W., Kennedy, D. N., Lanyon, L. J., Leergaard, T. B., Mayberg, H. S., Milanese, L., Mouček, R., Poline, J. B., Roy, P. K., Strother, S. C., Tang, T. B., ... Martone, M. E. (2022). A standards organization for open and FAIR neuroscience: The international neuroinformatics coordinating facility. *Neuroinformatics*, 20. <https://doi.org/10.1007/s12021-020-09509-0>
- Craig-Wood, N. (2014). *RClone [software]*. <https://rclone.org/>
- Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P., Flandin, G., Ghosh, S. S., Glatard, T., Halchenko, Y. O., Handwerker, D. A., Hanke, M., Keator, D., Li, X., Michael, Z., Maumet, C., Nichols, B. N., Nichols, T. E., Pellman, J., ... Poldrack, R. A. (2016). The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Scientific Data*, 3. <https://doi.org/10.1038/sdata.2016.44>
- Klingner, C. M., Denker, M., Grün, S., Hanke, M., Oeltze-Jafra, S., Ohl, F. W., Radny, J., Rotter, S., Scherberger, H., Stein, A., Wachtler, T., Witte, O. W., & Ritter, P. (2023). Results of a community survey of the german national research data infrastructure initiative neuroscience. *eNeuro*, 10, ENEURO.0215–22.2023. <https://doi.org/10.1523/ENEURO.0215-22.2023>
- Martone, M. E. (2024). The past, present and future of neuroscience data sharing: A perspective on the state of practices and infrastructure for FAIR. *Frontiers in Neuroinformatics*, 17. <https://doi.org/10.3389/fninf.2023.1276407>
- McGugan, W. (2021). *Textual [software]*. <https://github.com/Textualize/textual>
- Pierré, A., Pham, T., Pearl, J., Datta, S. R., Ritt, J. T., & Fleischmann, A. (2024). A perspective on neuroscience data standardization with neurodata without borders. *The Journal of Neuroscience*, 44. <https://doi.org/10.1523/JNEUROSCI.0381-24.2024>
- Rübel, O., Tritt, A., Ly, R., Dichter, B. K., Ghosh, S., Niu, L., Baker, P., Soltesz, I., Ng, L., Svoboda, K., Frank, L., & Bouchard, K. E. (2022). The neurodata without borders ecosystem for neurophysiological data science. *eLife*, 11. <https://doi.org/10.7554/eLife.78362>
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., Silva Santos, L. B. da, Bourne, P. E., Bouwman, J., Brookes, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T., Finkers, R., ... Mons, B. (2016). The FAIR guiding principles for scientific data management and stewardship. *Scientific Data*, 3. <https://doi.org/10.1038/sdata.2016.18>
- Ziminski, J. J., Sirmpilatze, N., Porta, L., Plattner, V., Peri, B. D., Keshavarzi, S., & Tyson, A. L. (2025). *NeuroBlueprint*. Zenodo. <https://doi.org/10.5281/zenodo.15720970>
- Ziminski, J. J., & Tyson, A. L. (2025). *Fancylog [software]*. <https://doi.org/10.5281/zenodo.15776028>