

# GigaScience

## The genetics-BIDS extension: Easing the search for genetic data associated with human brain imaging --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-20-00216
<b>Full Title:</b>	The genetics-BIDS extension: Easing the search for genetic data associated with human brain imaging
<b>Article Type:</b>	Commentary
<b>Funding Information:</b>	
<b>Abstract:</b>	Metadata are key in our ability to search databases. Without them, researchers would spend hours examining datasets in the hope to find data with features they are interested in. Brain imaging genetics is at the intersection of two disciplines with rich and complex data for which dictionaries and ontologies have been developed to facilitate data search and analysis. Here we present the genetics brain imaging data structure extension, a data descriptor of the genomic and transcriptomic data associated, possibly in different databases, with human brain imaging data. This extension will facilitate identifying micro-scale molecular features that are linked to macro-scale imaging repositories facilitating data aggregation across studies.
<b>Corresponding Author:</b>	Cyril Pernet, PhD The University of Edinburgh Edinburgh, Scotland UNITED KINGDOM
<b>Corresponding Author Secondary Information:</b>	
<b>Corresponding Author's Institution:</b>	The University of Edinburgh
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Cyril Pernet, PhD
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Cyril Pernet, PhD Clara Moreau, PhD Martineau Jean-Louis Ross Blair Christopher Markiewicz, PhD Jessica Turner, PhD Vince Calhoun, PhD Thomas Nichols, PhD
<b>Order of Authors Secondary Information:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	No
Full details of the experimental design and	

<p>statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "<b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p> <p>"</p>	<p>commentary paper on metadata descriptor</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>No</p>
<p>If not, please give reasons for any</p>	<p>not relevant</p>

<p>omissions below.</p> <p>as follow-up to "<b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p> <p>"</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

## **The genetics-BIDS extension: Easing the search for genetic data associated with human brain imaging**

Clara Moreau<sup>1</sup>, Martineau Jean-Louis<sup>1</sup>, Ross Blair<sup>2</sup>, Chris Markiewicz<sup>2</sup>, Jessica Turner<sup>3</sup>, Vince Calhoun<sup>4</sup>, Thomas E. Nichols<sup>5</sup> & Cyril R Pernet<sup>6</sup>

<sup>1</sup> Sainte Justine Research Center, University of Montreal, Montreal QC, CA

<sup>2</sup> Centre for Reproducible Neuroscience, Stanford University, USA

<sup>3</sup> Imaging Genetics and Informatics Lab, Georgia State University, USA

<sup>4</sup> Center for Translational Research in Neuroimaging and Data Science, Georgia State University, USA

<sup>5</sup> Oxford Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, Nuffield Department of Population Health, University of Oxford, UK

<sup>6</sup> Centre for Clinical Brain Sciences & Edinburgh Imaging, University of Edinburgh, UK

### **Abstract (100 words)**

Metadata are key in our ability to search databases. Without them, researchers would spend hours examining datasets in the hope to find data with features they are interested in. Brain imaging genetics is at the intersection of two disciplines with rich and complex data for which dictionaries and ontologies have been developed to facilitate data search and analysis. Here we present the genetics brain imaging data structure extension, a data descriptor of the genomic and transcriptomic data associated, possibly in different databases, with human brain imaging data. This extension will facilitate identifying micro-scale molecular features that are linked to macro-scale imaging repositories facilitating data aggregation across studies.

### **Keywords**

Human brain imaging, genomics, transcriptomics, Brain Imaging Data Structure

### **Introduction**

Brain imaging genetics aims at studying the association between brain structure or function and genetic variation [1]. Since gene expression influences cellular mechanisms, which in turn influences neural circuits underlying behaviour, studying associations at the brain level deepens our understanding of gene function at the system level. There is also evidence that using endophenotypes (i.e. brain phenotypes [2, 3]) is better suited to understand diseases, providing an intermediate description level between genes and clinical phenotypes.

Both brain imaging and genetics are fields in which researchers are used to share data in order to replicate findings and allow secondary usage. The way data are shared is however different with genetic data being sensitive personal information that must, therefore, be shared through secured and controlled access. Brain imaging, by contrast, is often shared via open databases, or with fewer restrictions on who can access. This has led to different approaches in data sharing for brain imaging genetics: fully secured and controlled for all data (e.g. UK biobank [4])

vs. splitting data with open access to brain images but secured access for genetic data (e.g. Human Connectome Project [5]). The former approach works for large homogenous projects requiring heavy data management while the latter approach is easier, especially for multiple individual smaller studies or multicentric studies with heterogeneous data collection. The Brain Imaging Data Structure (BIDS) describes a way of organising neuroimaging and behavioural data using dedicated names and dictionaries, documenting metadata [6]. Over time, extensions are being developed and integrated to address users' needs. Here we present the [BIDS genetics extension](#). The primary goal of the BIDS genetics extension is to link BIDS datasets to associated genetic data, especially those existing in separate repositories. The secondary goal is to provide a succinct description of the type of genetic data available, thus enabling searches through multiple imaging datasets.

### **The brain imaging data structure genetic descriptor**

Data organized according to BIDS have a rigid folder structure and naming convention. Every dataset comes with a *dataset\_description.json* file that contains information relative to authors, funders, ethics, licence, etc. To refer to associated genetic data, this file must now include the URL pointing to the genetic data, and optionally the URL of the database, and other associated materials like dataset descriptor articles. This will allow searching quickly through BIDS compliant repositories for datasets with associated genetic data.

Another requirement of BIDS is that a *participants.tsv* file is included with, at minimum, the subject's identifier. This file is now also used to link the brain imaging and genetic datasets if different pseudo-identifiers are used, making it easy to associate pseudo-IDs without need to ever access personal information.

#### *Extension characteristics and imaging genetic information*

This extension of the BIDS project aims to help researchers to structure their molecular (multi-level) and imaging datasets side-by-side in order to improve data linkage and search performance. To facilitate metadata search, a *genetic\_info.json* file must be associated with a BIDS dataset describing which type of genetic information is available. Among the multiple available fields, it minimally requires the `GeneticLevel` at which genetic analyses were carried out: genetic, genomic, epigenomic, transcriptomic, metabolomic or proteomic [7] (figure 1), and the `SampleOrigin` data: blood, saliva, brain, csf, breast milk, bile, amniotic fluid, or other biospecimen. If the `SampleOrigin` value is brain, it is further recommended to specify the `TissueOrigin` (gray matter, white matter, csf, meninges, macrovascular or microvascular) as the genetic or the genomic information may be more specifically related to the available imaging data. This can further be refined by indicating the `CellType` analysis with values taken from the cell ontology [8] and, if the `TissueOrigin` is gray matter, white matter or CSF, by indicating the `BrainLocation` (either using MNI coordinates or labels from the Allen Brain Atlas [9]).

A last, recommended, field is the `AnalyticApproach`, that is the sampling methodology. While optional, this is of particular importance since it indicates in greater detail the type of genetic data available using values from the database of Genotypes and Phenotypes (dbGaP [10]). As an example, the single nucleotide polymorphisms (SNP) genotyping (Array) and whole genome sequencing approaches provide both a whole genome level of genetic information,

albeit with some critical differences. The SNP genotyping reports genomic data with lower density compared to the whole genome sequencing which cover over ~95% of the genomic DNA.

## Conclusion

BIDS is an openly developed, community-led standard to name, document and organize human brain imaging data, allowing FAIR data sharing and the automation of complex data preprocessing. In just four years of existence, it has revolutionized data sharing and analysis in neuroscience, from a wide-adoption and a reference for publications to supporting data repository architectures, and it is critical to many open-source analysis pipelines. Here, we present the genetic extension which is integrated into the BIDS specification providing a full documentation of the fields that may be provided, along with online examples and a Javascript validator to ensure datasets are compliant. By adding a genetic descriptor for imaging data, we hope to facilitate data mining to constitute large multi-scale heterogeneous analysis human datasets that reflect human variability, necessary to enhance our understanding of genetic influence on brain phenotypes.

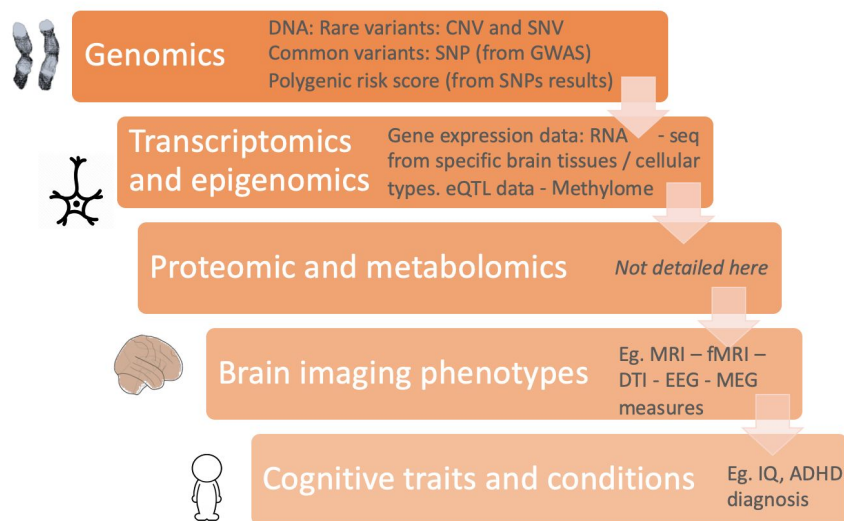


Figure 1: Linking micro- to macro-scale data. (CNV: copy number variation, SNV: single nucleotide variation, SNP: single nucleotide polymorphisms, PRS: polygenic risk score, eQTL: expression quantitative trait loci, IQ: intelligence quotient, ADHD: attention deficit hyperactivity disorder)

## Author contributions

CIM, JT, TN, VC and CP conceptualized the BIDS extension, CIM and CP wrote the manuscript draft; CIM, MJL, ChM and CP wrote the extension and example, RB wrote the javascript validator. All co-authors have contributed to the preparation of the manuscript, and/or read and approved the final version.

Authors declare no conflict of interest

### **Availability of supporting source code and requirements**

Project name: Brain Imaging Data Structure

Project home page: <https://bids.neuroimaging.io/>

Operating system(s): Platform independent

Programming language: Markdown and JavaScript

Other requirements: Node.js to run the validator locally or a web browser to run the validator online

RRID: SCR\_016124

License: CC-BY

### **Acknowledgements**

Thank you to Elizabeth Dupre and Stefan Appelhoff for helping to structure and correct the BIDS markdown extension.

### **References**

1. Poline, J-B., Breeze, J.L., Frouin, V. (2015) Imaging Genetics with fMRI. In: Uludag K, Ugurbil, K., Berliner, L. (eds) fMRI: From Nuclear Spins to Brain Functions. Springer, New-York, pp 699–738
2. John B, Lewis KR (1966) Chromosome Variability and Geographic Distribution in Insects. *Science* 152:711–721. <https://doi.org/10.1126/science.152.3723.711>
3. Gottesman II, Gould TD (2003) The Endophenotype Concept in Psychiatry: Etymology and Strategic Intentions. *AJP* 160:636–645. <https://doi.org/10.1176/appi.ajp.160.4.636>
4. UK Biobank. <https://www.ukbiobank.ac.uk/about-biobank-uk/>. Accessed 14 Jul 2020
5. Human Connectome Project | Mapping the human brain connectivity. <http://www.humanconnectomeproject.org/>. Accessed 14 Jul 2020
6. Gorgolewski KJ, Auer T, Calhoun VD, et al (2016) The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Scientific Data* 3:160044. <https://doi.org/10.1038/sdata.2016.44>
7. Hasin Y, Seldin M, Lusk A (2017) Multi-omics approaches to disease. *Genome Biology* 18:83. <https://doi.org/10.1186/s13059-017-1215-1>
8. Malladi VS, Erickson DT, Poddaturi NR, et al (2015) The Cell Ontology. *Database* 2015:. <https://doi.org/10.1093/database/bav010>
9. Hawrylycz MJ, Lein ES, Guillozet-Bongaarts AL, et al (2012) An anatomically comprehensive atlas of the adult human brain transcriptome. *Nature* 489:391–399. <https://doi.org/10.1038/nature11405>
10. Database of Genotypes and Phenotypes (dbGaP). <https://www.ncbi.nlm.nih.gov/gap/>. Accessed 14 Jul 2020