# FlywheelTools:

# Data Curation and Manipulation on the Flywheel Platform

1   **Tinashe M. Tapera[1], Matthew Cieslak, PhD[1], Max Bertolero, PhD[1], Azeez Adebimpe, PhD[1],**
2   **Geoffrey K. Aguirre, MD, PhD[2], Ellyn R. Butler[1], Philip A. Cook[3], Diego Davila[1], Mark A.**
3   **Elliott, PhD[3], Sophia Linguiti[1], Kristin Murtha[1], William Tackett[2], John A. Detre, MD[2],**
4   **Theodore D. Satterthwaite, MD[1*]**

5   [1]Penn Lifespan Informatics & Neuroimaging Center, Department of Psychiatry, University of
6   Pennsylvania, Philadelphia, PA, USA

7   [2]Department of Neurology, Perelman School of Medicine, University of Pennsylvania, Philadelphia,
8   PA, USA.

9   [3]Department of Radiology, University of Pennsylvania, Philadelphia, PA, USA

10  **\* Address correspondence to:** sattertt@pennmedicine.upenn.edu

11 **ABSTRACT**

12 The recent and growing focus on reproducibility in neuroimaging studies has led many major
13 academic centers to use cloud-based imaging databases for storing, analyzing, and sharing complex
14 imaging data. Flywheel is one such database platform that offers easily accessible, large-scale data
15 management, along with a framework for reproducible analyses through containerized pipelines. The
16 Brain Imaging Data Structure (BIDS) is the de facto standard for neuroimaging data, but curating
17 neuroimaging data into BIDS can be a challenging and time-consuming task. In particular, standard
18 solutions for BIDS curation are limited on Flywheel. To address these challenges, we developed
19 "FlywheelTools", a software toolbox for reproducible data curation and manipulation on Flywheel.
20 FlywheelTools includes two elements: *fw-heudiconv*, for heuristic-driven curation of data into BIDS,
21 and *flaudit*, which audits and inventories projects on Flywheel. Together, these tools accelerate
22 reproducible neuroscience research on the widely used Flywheel platform.

23 **KEYWORDS:** neuroimaging, neuroinformatics, BIDS, curation, Python

24

## 1    INTRODUCTION

Many fields in science are grappling with failures of scientific reproducibility (Botvinik-Nezer et al. 2020). Given the high dimensionality of the data, the need for complex image processing, and a plethora of analytic techniques, this crisis is particularly acute for neuroimaging research. As such, major academic centers and large consortia have increasingly adopted platforms that leverage database technologies that have become standard in other fields. In addition to providing functionality for searching and categorizing complex source data, imaging databases enhance reproducible research by providing a clear audit trail of image processing applied to the data and its results, including both derived images and other data. Widely used imaging databases include Collaborative Informatics Neuroimaging Suite (COINS) (Landis et al. 2016), eXtensible Neuroimaging Archive Toolkit (XNAT) (Herrick et al. 2016), Longitudinal Online Research and Imaging System (LORIS) (Vaccarino et al. 2018), and others (Book et al. 2016; Helmer et al. 2011; Rogovin et al. 2020; Helmer et al. 2011; Poldrack and Gorgolewski 2017; Sherif et al. 2014). More recently, the commercial platform Flywheel has become widely used due to its modern technology, ease of use, and scalability.

Many neuroimaging databases now leverage the Brain Imaging Data Structure (BIDS) (Gorgolewski et al. 2016). BIDS is an open-source standard for neuroimaging data organization that specifies how files should be named, how directories should be organized, and how metadata should be structured. As such, BIDS provides users with a well-documented structure to understand both imaging data and metadata. Importantly, as BIDS provides transparent format for recording imaging parameters and key aspects of the experimental design, it enhances both data accessibility and data sharing. BIDS has rapidly evolved to become the standard in the neuroimaging community for data organization. Importantly, BIDS is supported by a large community that contributes to its development and adoption. Further, proposals for BIDS schema pass through a rigorous testing process before being adopted.

Notably, BIDS allows users to leverage BIDS-apps — image processing pipelines, (e.g., fMRIPrep, C-PAC, and QSIprep) that read the metadata defined by BIDS (Esteban et al. 2019; Craddock et al. 2013; Cieslak et al. 2020). As BIDS-apps can auto-configure to ensure that analytic parameters are appropriate for the input data provided, they dramatically reduce barriers to implementing best practices in image processing. Importantly, containerized BIDS-apps encompass all software dependencies, further enhancing reproducibility.

On a filesystem, conversion of raw DICOM images to NIfTIs that conform to BIDS can be accomplished with a variety of tools including *HeuDiConv*, *dcm2bids*, and others (Halchenko et al. 2018). However, this crucial step, a process typically called "BIDS curation," is incompletely implemented on Flywheel. While Flywheel provides automated BIDS curation, flexibility is limited. As BIDS curation is one of the very first steps performed on the data, flexibility in curation is essential. Here we introduce FlywheelTools: software that provides flexible and reproducible methods for BIDS curation on the Flywheel platform. Documentation and code can be found online at: https://fw-heudiconv.readthedocs.io/en/latest/.

## 2    METHODS

The FlywheelTools toolkit allows users to follow a reproducible workflow for BIDS curation and auditing of their data. This workflow typically includes inspection of sequences collected during a study, design of a curation schema, implementation of that curation schema, and auditing the curated data.

69 **2.1 Programming Languages & Technologies**

70 FlywheelTools is built primarily in Python 3.6 (Van Rossum and Drake 2009) to leverage Flywheel's
71 highly accessible Software Development Kit (SDK). Additionally, R 3.4.1 (R Core Team 2019) is
72 used for HTML report generation. For reproducibility and workflow management, the modules of
73 FlywheelTools are packaged as Docker container images (Merkel 2014). It should be noted that
74 FlywheelTools relies on users adopting BIDS as their data standard.

75 **2.2 Flywheel**

76 Flywheel is a data management and analysis platform that is tailored for neuroimaging research. The
77 platform focuses heavily on collaborative and reproducible science. User-facing components of the
78 platform itself are the web User Interface (UI), the Command Line Interface (CLI), the Flywheel
79 Software Development Kit (SDK), and the Application Programming Interface (API).

80 **2.3 Flywheel Web UI**

81 The web UI is accessible through any modern web browser. Through this point-and-click interface,
82 users are able to upload, view, download, and analyze data with ease. However, accomplishing tasks
83 with many repetitive steps or over a large number of participants/sessions can be tiresome and error-
84 prone. In addition to interactions via web GUI, many users also make use of the API and SDK to
85 manipulate and analyze data programmatically.

86 **2.4 Flywheel API & SDK**

87 Flywheel's database uses MongoDB for data storage and access, meaning that all Flywheel data are
88 represented by hierarchical relationships between document objects. This allows users to create and
89 store complex structures with ease, and query data rapidly (Banker 2011). To access these data,
90 Flywheel uses a RESTful Application Programming Interface (REpresentational State Transfer)
91 (Biehl 2016), making each document or data object accessible through a specific URL that a web
92 browser or SDK can access by requesting the data and waiting for a response from the server. The
93 Flywheel Python SDK[1] provides a powerful interface for inspecting and manipulating data through
94 this API. By standardising this underlying data model into Pythonic objects, the Flywheel SDK is
95 effectively an object relationship mapper, similar to the popular SQLAlchemy software.

96 **2.5 Flywheel Data Model**

97 Objects in Flywheel's data model follow a specific hierarchical structure — at the top level is a
98 Flywheel instance, a process that serves the API to an organization (for example, a neuroimaging
99 center). Within the Flywheel instance, there are multiple groups, which are typically labs or research
100 units that collaborate on one or more projects. Each project object can have one or many subjects (i.e.
101 participants), and each subject can have one or many sessions (i.e. scanning visits). Within a session,
102 there may be one or many acquisition objects which represent the scanning sequences collected
103 during a particular scan or examination (e.g., sMRI, rs-fMRI, dMRI). Finally, the data files
104 associated with the sequence (e.g., NIfTIs or DICOMs) are attached to each acquisition. Note that a
105 file can additionally be attached to any object type, and each object can have metadata associated
106 with it. Hence, a "subject" object may have metadata associated with that participant (such as

---

[1] Flywheel also provides a MATLAB SDK, however we use the term SDK in this work to refer to the
Python SDK, which we use exclusively in FlywheelTools.

107 demographic information) and may also have a text file attached to it (such as clinical data). A
108 notable exception to this hierarchical structure is the analysis object, which behaves in much the
109 same way as others but can be a child object of any other object, allowing researchers to create
110 analyses of entire projects, for example, each with their own associated metadata and files.
111
112 Abstracting this data model in Python results in simple hierarchical objects, each with methods for
113 handling metadata and files, and methods for accomplishing object-specific tasks like traversing the
114 hierarchical structure or running analyses. The modules of FlywheelTools make use of this data
115 model to accomplish a wide range of tasks.

### 2.6 Flywheel Gears

117 Flywheel encourages the use of pre-packaged computational workflows, called "gears". Gears are
118 run by virtual machines using Docker; as such, they are version-controlled and can be executed on
119 any platform supporting Docker. Gears can accomplish tasks such as data manipulation, pre-
120 processing, and analysis. In addition to the existing gears available on the platform, users are able to
121 package their own software in a gear and use it for running analysis workflows on their Flywheel
122 data via the web UI or SDK. The complexity and frequency of the task help to guide if a task should
123 be accomplished using the web UI, programmatically using the SDK, or by wrapping a workflow
124 into a gear. Gear developers are able to construct configurable options and necessary inputs for their
125 gear in a standardized manifest file, written as a JSON. These configurations appear as clickable
126 options in the GUI, can be set programmatically using the SDK, or specified using the command line
127 interface. Once a gear is launched, the gear queries Flywheel for the specified input data (such as
128 images or file attachments) and runs the pipeline. Once a workflow has completed running, Flywheel
129 collects any files remaining in the pre-defined output directory of the container and attaches them to a
130 resulting analysis object. The output of a gear (such as an HTML report or tabulated data) can be
131 viewed on the Flywheel UI, downloaded to disk for further sharing or analysis, or used as input to a
132 subsequent gear.

### 3 RESULTS

135 FlywheelTools is implemented using the Flywheel SDK to enable easy inspection, curation,
136 validation, and audit of Flywheel data through a handful of user-friendly gears and command-line
137 interfaces. The first module of the package is *fw-heudiconv*, a toolbox for reproducible curation of
138 neuroimaging data into BIDS on Flywheel. The second module, *flaudit*, is a tool for auditing a
139 Flywheel project, giving users an overview of the key elements of their dataset.

### 3.1 FW-HEUDICONV

141 The first tool, fw-heudiconv[2], is a multi-purpose command-line interface and Flywheel gear designed
142 for BIDS curation on Flywheel (**Figure 1**). It is designed to be intuitive, flexible, and reproducible.
143 Users of the Flywheel gear have used it to successfully run 27,251 jobs at the time of this writing.

### 3.1.1 Architecture & Design

---

2 https://github.com/PennLINC/fw-heudiconv

145 To curate data into BIDS format, *fw-heudiconv* first considers DICOM data to be the "ground truth"
146 and builds its curation approach using data in the DICOM headers. DICOMs are added to the
147 Flywheel database, either through manual upload or automatically from a linked scanner. The
148 DICOMs are then automatically converted into NIfTI files by Flywheel's automated gears. The result
149 is an acquisition object with both DICOMs and one or more NIfTIs. Ultimately, *fw-heudiconv* only
150 has permission to manipulate metadata associated with a NIfTI file. By not manipulating DICOMs or
151 their associated metadata, BIDS curation can safely be reproduced from ground truth data.

152 *fw-heudiconv* can be downloaded as a command-line tool from the Python Package Index using pip
153 or can be run from the Flywheel GUI as a gear. Running fw-heudiconv as a gear has the added
154 advantage of containerization, allowing gear configuration and all changes to the data to be tracked.
155 There are a number of commands available in *fw-heudiconv*, and each of them starts by querying data
156 from Flywheel. Users can filter their queries to operate on an entire Flywheel project, a subset of
157 subjects, or a subset of sessions. Notably, with the `--dry-run` option, each command has the
158 ability to test and evaluate its effects without actually manipulating metadata in the Flywheel
159 database or writing data to disk. Below, we consider each of the five available commands.

### 160 3.1.1.1 fw-heudiconv-tabulate

161 The tabulate tool is used to parse and extract DICOM header information in a project (or within a
162 filtered subset of that project) and compile these data into a table for the user to examine. By
163 collecting DICOM header information into a tabular format, the tabulate tool gives users a
164 comprehensive overview of the different scanning sequences that have been collected in the query,
165 including the sequence parameters. Additionally, users have the option to limit the tabulation to a
166 unique combination of common DICOM header fields, which significantly decreases the complexity
167 of the table. When used at the command line, the table produced by this command is written to a
168 local disk. When used as a gear on Flywheel, once the workflow has completed running, Flywheel
169 collects any files remaining in the pre-defined output directory of the container and attaches them to
170 the resulting analysis object.

### 171 3.1.1.2 fw-heudiconv-curate

172 The curate tool is used to curate a dataset on Flywheel into BIDS format. Much like *HeuDiConv*,
173 curation is accomplished through the use of a heuristic: a Python file that programmatically defines
174 the templates for a range of BIDS-valid filenames, and defines the boolean logic that would assign a
175 given scanning sequence to each template. This logic is usually based on the sequence information
176 users find in the tabulation of sequences, but all fields available in the DICOM header can be used to
177 determine which template a particular file can be assigned to. Additionally, the curate tool can be
178 used to manipulate BIDS metadata that may need to be added to the dataset. The process of curation
179 only manipulates the BIDS metadata of NIfTI files, and hence can be repeated or updated at any time
180 at the user's discretion. In addition to hardcoding metadata, users can also use the heuristic to create
181 and upload additional BIDS files (e.g., a README, participants TSV, or events TSV). If the
182 required dataset description does not already exist in the project, it is auto-populated by fw-
183 heudiconv with the minimum required fields. The curate tool ensures that these files are tracked as
184 BIDS files by hardcoding their relative BIDS paths into the file's Flywheel metadata.

### 185 3.1.1.3 fw-heudiconv-export

186 The export tool is used to export a BIDS dataset on Flywheel to disk.  It can also be used by other
187 gears or scripts to easily extract their BIDS data into the workspace of an analysis pipeline.

188 Importantly, the export tool downloads all additional files created by the curate tool by ensuring they
189 have been created with BIDS-valid paths in their metadata.

190 **3.1.1.4 fw-heudiconv-validate**

191 The validate tool is a wrapper around the popular BIDS Validator package and is used to check if the
192 applied curation results in a BIDS-valid dataset (Gorgolewski et al., 2020). After exporting a dataset
193 with *fw-heudiconv-export*, the validate tool runs the BIDS Validator on the dataset and returns the
194 verbose output of the errors and warnings given by the BIDS Validator. Additionally, the results of
195 the validator can be tabulated for easy inspection. On the Flywheel GUI, *fw-heudiconv-validate* also
196 displays a green check mark in the analysis tab for a successful validation, and a red check mark
197 otherwise, allowing for quick visual inspection of BIDS curation status for each session.

198 **3.1.1.5 fw-heudiconv-clear**

199 The clear tool is used to clear BIDS information cleanly and safely from the project or subjects and
200 sessions queried. This can be useful when a user wants to rerun the curation. The previously created
201 persistent fields can be removed by running *fw-heudiconv-clear* before re-curating.

202 **3.1.2 The Heuristic File**

203 The heuristic file is a Python file used as input to the *fw-heudiconv-curate* command. The file
204 instructs *fw-heudiconv* on how to programmatically sort and parse through each acquisition object in
205 Flywheel and assign it to a valid BIDS naming template. This is done by checking the attributes of a
206 list of *seqInfo* objects — which are generated from each DICOM's header information — against
207 user-defined boolean rules. For example, if a T1-weighted image is present in a dataset, the user may
208 define a string with a BIDS-valid naming template for this type of file, such as:

209 
```
t1w = 'sub-{SubjectLabel}_ses-{SessionLabel}_T1w.nii.gz'
```

210 Where the SubjectLabel and SessionLabel portions are expected to be automatically generated for
211 each subject and session in the dataset. After the DICOM SeriesDescription field is added to the
212 SeriesDescription attribute of *seqInfo*, the user can create a simple boolean expression to check if the
213 string 'T1w' is in the SeriesDescription. If such a rule is met, this acquisition and its NIfTI file will
214 be assigned to the T1-weighted image naming template. The NIfTI file will ultimately have this
215 BIDS naming added to its metadata and be named correctly when exported to a filesystem. In more
216 complex naming scenarios, *fw-heudiconv* can flexibly use boolean expressions involving any number
217 of *seqInfo* attributes, which the user can access in the output of *fw-heudiconv-tabulate*.

218 In addition to setting naming templates, the heuristic file can also be used to hard-code and assign
219 metadata in BIDS. These data are hard-coded into the metadata of the file object on Flywheel and are
220 assigned by using specially reserved functions and keywords in *fw-heudiconv*. For example, the
221 heuristic file can be used to point fieldmap scans to their intended sequences using a list:

222 
```
IntendedFor = {
```

223 
```
fieldmap1: ['sub-{SubjectLabel}_ses-{SessionLabel}_task-rest_bold.nii.gz']
```

224 
```
}
```

225  By reserving select keywords for functions and metadata, heuristic files become versatile tools for
226  defining and manipulating a wide array of metadata in Flywheel BIDS curation. For example, users
227  can make use of the keywords ReplaceSubject() or ReplaceSession() to create functions that
228  dynamically and programmatically manipulate subject or session labels shown in BIDS. The
229  MetadataExtras keyword is used for hardcoding metadata fields found in JSON sidecars, and the
230  AttachToProject() and AttachToSession() keywords can be used to dynamically attach BIDS files to
231  Flywheel objects.

232  Importantly, because this heuristic file is plain text Python code, users are able to version control
233  their files using Git and share these files via Github. Finally, when run on Flywheel as a gear, the
234  heuristic file is automatically attached as an input to the analysis object created by *fw-heudiconv-*
235  *curate*, allowing users to easily access the version history of their curation.

### 3.1.3 Curation Workflow

237  For most users, the curation workflow follows the sequence detailed above (**Figure 2**). After
238  DICOMs have been converted to NIfTIs, users can then begin by running *fw-heudiconv-tabulate* to
239  gather the information stored in the DICOM headers necessary for creating a heuristic. Once the
240  tabulation has been completed, the output file can be opened by any program that can read tabular
241  data. At this stage, users can begin creating a heuristic file and running *fw-heudiconv-curate*, using
242  the `--dry-run` flag to test the heuristic changes incrementally with informative logging. When
243  satisfied, users can simply remove the `--dry-run` flag to apply the changes. The user can then use
244  fw-heudiconv-validate to run the BIDS validator on the dataset or start over by removing all BIDS
245  metadata with fw-heudiconv-clear.

246  Importantly, when running at the command line, all *fw-heudiconv* tools run on the entire Flywheel
247  project by default, but all come with optional --subject and --session flags to allow the user to specify
248  operations. The recommended workflow at the command line is to first develop and test the heuristic
249  file on a single subject's session. Often, it is most useful to conduct testing on the session with the
250  most complete data.   When testing is complete and curation works as intended, this heuristic can be
251  applied to the full project. When using the Flywheel GUI, gears by default run at the session level.  In
252  this case, users can similarly develop and test their heuristic on a single session using the GUI. When
253  satisfied with that session's curation, the gear can be run from the project level to curate all the
254  available data. This option is beneficial for data provenance, as all of a gear's commands and inputs,
255  as well as outputs and logs, are stored and attached to each gear run.

## 3.2    FLAUDIT

257  The second module of FlywheelTools is a Flywheel project auditor, named *flaudit*[3]. The module is
258  intended to give Flywheel users a broad understanding of their entire Flywheel project, by
259  summarizing the available data and illustrating analysis workflows. The output of this module, a
260  portable HTML report, presents this information using a number of visualizations built in R
261  Markdown using HTML, Javascript, and ggplot2, in two main sections: project overview and project
262  completeness.

### 3.2.1 Architecture & Design

---

3 https://github.com/PennLINC/flaudit

264  Using internal machinery similar to *fw-heudiconv-tabulate*, *flaudit* loops over existing data in a
265  project and tabulates information about scanning sequences, BIDS metadata, and gear analyses that
266  have been run. These three tables are saved internally and then passed as input to an R markdown
267  script that generates an interactive HTML report. The data are also saved as output for the user to
268  further access and analyze in their software of choice.

269  ### 3.2.2 Flaudit: Project Overview

270  The overview section of the *flaudit* report provides a numerical overview of sequences, BIDS data,
271  gear runs, and gear runtimes.

272  The first visualization uses the sequence data input to create a bar chart visualizing the names of the
273  different sequences acquired across the entire Flywheel dataset. This chart is accompanied by an
274  interactive table that users can search to compare values (**Figure 3**).

275  Next, using the BIDS metadata input, the report provides an interactive tree viewer to examine BIDS
276  curation. In the tree, the nodes branch out from the project to show each sequence acquisition. For
277  each acquisition, if the data has been curated into BIDS, the node itself can also branch out to show a
278  BIDS name template, demonstrating what BIDS name that sequence has been given. Hovering over
279  the BIDS name will display the number of subjects whose data have been named as such (**Figure 4**).

280  Finally, using the gear analysis data as input, the last section of the overview enumerates the gear
281  analyses that have been run successfully on any session within the project, and enumerates the
282  runtimes for these processes. The results are visualized in a bar chart (**Figure 5**).

283  ### 3.2.3 Flaudit: Project Completeness

284  As an optional input, *flaudit* allows users to specify a *template* subject — a subject from the Flywheel
285  project who serves as an exemplar for other subjects to be compared against. This subject should be
286  chosen based on the fact that they have both complete input data and analyses. This allows *flaudit* to
287  determine if other subjects in the project have equally complete data or are missing specific raw data
288  or analytic output. The project completion section of the *flaudit* report consists of three interactive
289  tables.

290  In the first table, it's assumed that the template subject has acquired a complete set of imaging
291  sequences. These sequences are listed as columns in the table. Each subsequent row is a subject in the
292  project, and each value in the table is a boolean (complete or incomplete) indicating if that subject
293  has each sequence. The table is searchable, meaning that users can simply filter each column for
294  "incomplete" to learn which subjects do not have the same data as the template (**Figure 6A**).
295  Likewise, the second table illustrates the completeness of BIDS data for other subjects in comparison
296  to the template (**Figure 6B**). In this case, rows indicate subjects while columns delineate the
297  specified BIDS naming template. Lastly, the third table illustrates completeness of analytic gear runs.
298  Researchers can use this table to compare the analytic output of all other subjects in the project to the
299  analysis pipelines run for the template subject. To ensure uniform versions of pipeline software, the
300  version of a pipeline that was used for each subject must match that of the template subject (**Figure
301  6C**).

302  ## 4    DISCUSSION

303  FlywheelTools provides new capabilities for the popular and powerful Flywheel platform, allowing
304  researchers to maximize reproducibility and enhance scalability. Specifically, *fw-heudiconv* provides

305     users a flexible and reproducible way of curating data into BIDS on Flywheel. Complementary
306     function is provided by *flaudit*, which provides intuitive visual reports of raw, curated, and processed
307     data.

308     Flywheel has been rapidly adopted by major imaging centers due to its ease of use, extensive
309     functionality, scalability, and emphasis on reproducible research. Despite these strengths, at present
310     there have been limited options for conversion of imaging data to BIDS format on Flywheel. This
311     step is absolutely critical, as BIDS provides a standardized format for important imaging metadata.
312     Notably, initial curation to BIDS has frequently been an important gap in workflows for reproducible
313     research.

314     Accordingly, *fw-heudiconv* provides critical functionality for reproducible and flexible BIDS curation
315     on Flywheel. The combination of containerized code and heuristics that are version controlled with
316     git maximizes reproducibility, ensuring that all curation steps have a clear audit trail. Furthermore,
317     the flexible architecture employed by *fw-heudiconv* allows workflows to be updated to accommodate
318     both new scanning protocols and the evolving specifications of the BIDS standard.

319     Once data are curated with *fw-heudiconv*, *flaudit* allows users to audit the data in a Flywheel project.
320     Specifically, *flaudit* provides intuitive summaries at each stage of a typical workflow, concisely
321     visualizing raw data, curated data in BIDS, and data processed by containerized analytic gears. These
322     accessible reports allow users to rapidly assess the overall organizational state of a project, while
323     interactive tables allow for more granular inspection of data. This approach facilitates understanding
324     the diverse data types typically collected in multi-modal imaging studies.

325     There are of course limitations of FlywheelTools. First, it should be acknowledged that
326     FlywheelTools is built for the Flywheel platform, and as such does not generalize to other imaging
327     databases that are in use. However, given the rapid adoption of this platform by the imaging
328     community, we anticipate that this toolkit will fill an important need for the many large research
329     institutions that rely upon Flywheel. Second, some understanding of Python is necessary to build the
330     heuristic for *fw-heudiconv*. We attempt to minimize this issue by providing both extensive
331     documentation and heuristic templates for various uses of *fw-heudiconv*[4], but usage is ultimately a
332     programming task.

333     A significant improvement to FlywheelTools would be automated curation of BIDS data similar to
334     that implemented by ReproNim's ReproIn workflow[5]. In the ReproIn workflow, researchers adhere
335     to specific naming conventions when scans are acquired, resulting in DICOMs that can be
336     automatically converted by HeuDiConv with the use of a turnkey heuristic that comes with the
337     software. At present, *fw-heudiconv* includes a prototype ReproIn heuristic that we plan to expand and
338     further evaluate moving forward.

339     FlywheelTools provides essential functionality to the Flywheel platform. The flexible toolkit allows
340     for curation and description of complex imaging studies. Taken together, the toolkit is designed to
341     accelerate reproducible imaging research at scale.

---

[4]https://fw-heudiconv.readthedocs.io/en/latest/index.html

[5] https://www.repronim.org/do.html

342

## 5 CONFLICT OF INTEREST

344 The authors declare that the research was conducted in the absence of any commercial or financial
345 relationships that could be construed as a potential conflict of interest.

## 6 AUTHOR CONTRIBUTIONS

347 T.M.T. and M.C. contributed to the design and implementation of the code and software. T.D.S.,
348 M.B., and M.C. supervised writing of the manuscript. A.A., M.C., E.R.B., D.D. and W.T. provided
349 substantial software use-cases, software feature requests, test data, and critical bug reports. K.M. and
350 S.L. assisted with software documentation. M.A.E., G.K.A., P.AC., J.A.D., and T.D.S. were involved
351 in proposing, planning, and supervising all the work.

## 7 FUNDING

354

## 8 ACKNOWLEDGMENTS

356 None.

357

## 9    REFERENCES

Banker, Kyle. 2011. MongoDB in Action. USA: Manning Publications Co.

Biehl, Matthias. 2016. RESTful Api Design. Vol. 3. API-University Press.

Book, Gregory A, Michael C Stevens, Michal Assaf, David C Glahn, and Godfrey D Pearlson. 2016. "Neuroimaging Data Sharing on the Neuroinformatics Database Platform." Neuroimage 124: 1089–92.

Botvinik-Nezer, Rotem, Felix Holzmeister, Colin F Camerer, Anna Dreber, Juergen Huber, Magnus Johannesson, Michael Kirchler, et al. 2020. "Variability in the Analysis of a Single Neuroimaging Dataset by Many Teams." Nature, 1–7.

Cieslak, Matthew, Philip A Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez Adebimpe, Geoffrey K Aguirre, et al. 2020. "QSIPrep: An Integrative Platform for Preprocessing and Reconstructing Diffusion Mri." bioRxiv.

Craddock, Cameron, Sharad Sikka, Brian Cheung, Ranjeet Khanuja, Satrajit S Ghosh, Chaogan Yan, Qingyang Li, et al. 2013. "Towards Automated Analysis of Connectomes: The Configurable Pipeline for the Analysis of Connectomes (c-Pac)." Front Neuroinform 42.

Esteban, Oscar, Christopher J Markiewicz, Ross W Blair, Craig A Moodie, A Ilkay Isik, Asier Erramuzpe, James D Kent, et al. 2019. "fMRIPrep: A Robust Preprocessing Pipeline for Functional Mri." Nature Methods 16 (1): 111–16.

Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P., ... & Poldrack, R. A. (2016). The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. Scientific data, *3*(1), 1-9.

Gorgolewski, C., Hardcastle, N., Teal Hobson-Lowther, Nishikawa, D., Blair, R., Appelhoff, S., Suyash, Constellates, Mainak Jas, Holdgraf, C., Jones, A., Goyal, R., Oostenveld, R., Markiewicz, C., Noack, G., Zito, M., Durnez, J., Traut, N., Naveau, M., … Thomas, A. (2020). *Bids-standard/bids-validator: 1. 4. 3* (1.4.3) [Computer software]. Zenodo. https://doi.org/10.5281/ZENODO.3688707

Halchenko, Yaroslav, Mathias Goncalves, Matteo Visconti di Oleggio Castello, Satrajit Ghosh, Michael Hanke, Matthew Brett, Taylor Salo, et al. 2018. Nipy/Heudiconv: Heudiconv V0.5.1 (version v0.5.1). Zenodo. https://doi. org/10.5281/zenodo.1306159.

Helmer, Karl G, Jose Luis Ambite, Joseph Ames, Rachana Ananthakrishnan, Gully Burns, Ann L Chervenak, Ian Foster, et al. 2011. "Enabling Collaborative Research Using the Biomedical Informatics Research Network (Birn)." Journal of the American Medical Informatics Association 18 (4): 416–22.

Herrick, Rick, William Horton, Timothy Olsen, Michael McKay, Kevin A Archie, and Daniel S Marcus. 2016. "XNAT Central: Open Sourcing Imaging Research Data." NeuroImage 124: 1093–6.

Landis, Drew, William Courtney, Christopher Dieringer, Ross Kelly, Margaret King, Brittny Miller, Runtang Wang, Dylan Wood, Jessica A Turner, and Vince D Calhoun. 2016. "COINS Data

395    Exchange: An Open Platform for Compiling, Curating, and Disseminating Neuroimaging Data."
396    NeuroImage 124: 1084–8.

397    Merkel, Dirk. 2014. "Docker: Lightweight Linux Containers for Consistent Development and
398    Deployment." Linux J. 2014 (239).

399    Poldrack, Russell A, and Krzysztof J Gorgolewski. 2017. "OpenfMRI: Open Sharing of Task fMRI
400    Data." Neuroimage 144: 259–61.

401    R Core Team. 2019. R: A Language and Environment for Statistical Computing. Vienna, Austria: R
402    Foundation for Statistical Computing. https://www.R- project.org/.

403    Rogovin, O, Y Zhao, S Chen, Z Wang, O Papaemmanouil, SD Van Hooser, and others. 2020. "NDI:
404    A Platform-Independent Data Interface and Database for Neuroscience Physiology and Imaging
405    Experiments."

406    Sherif, T., Rioux, P., Rousseau, M.E., Kassis, N., Beck, N., Adalat, R., Das, S., Glatard, T., & Evans,
407    A. (2014). CBRAIN: a web-based, distributed computing platform for collaborative neuroimaging
408    research. *Frontiers in neuroinformatics, 8, 54.*

409    Vaccarino, Anthony L, Moyez Dharsee, Stephen Strother, Don Aldridge, Stephen R Arnott, Brendan
410    Behan, Costas Dafnas, et al. 2018. "Brain-Code: A Secure Neuroinformatics Platform for
411    Management, Federation, Sharing and Analysis of Multi-Dimensional Neuroscience Data." Frontiers
412    in Neuroinformatics 12: 28.

413    Van Rossum, Guido, and Fred L. Drake. 2009. Python 3 Reference Manual. Scotts Valley, CA:
414    CreateSpace.

415

## 10    FIGURE LEGENDS

417

418    **Figure 1: BIDS-app workflow.** BIDS curation on file systems is a common task that can be
419    accomplished by existing tools (heudiconv, dcm2bids, etc) or manually, but mechanisms for BIDS
420    curation on many cloud databases have yet to be developed. FlywheelTools provides this
421    functionality for the Flywheel platform.

422    **Figure 2: FlywheelTools workflow.** Users first use the tabulate tool to extract sequence information
423    from their data, which they use to develop a heuristic that delineates how sequences are mapped into
424    BIDS. After this, they use the curate tool to convert their data into BIDS, and the validate tool to
425    assess their curation. The export tool can be used to export their BIDS data as necessary.

426    **Figure 3: Enumeration of available sequences in a flywheel project.** Panel (**A**) plots the count of
427    files in each collected sequence; in this example, there are 60 files collected for the B0map sequence,
428    as there are 20 subjects and 3 B0map sequences. Panel (**B**) shows the accompanying interactive table.

429    **Figure 4: Interactive tree diagram illustrating BIDS curation.** The tree shows how each sequence
430    has been curated into BIDS format; users can hover their mouse over each leaf to show how many
431    files have been curated into each BIDS filename template.

432    **Figure 5: Enumeration of gear runs in a Flywheel project.** The number of gear runs is shown for
433    various gears. For each gear, the percent of completed versus failed runs is shown. For example, 95
434    percent of the subjects (n=19) were successfully run through fMRI-prep.

435    **Figure 6: Project completeness tables compared to the template participant in a Flywheel**
436    **project.** Panel (**A**) compares the sequences available for each participant to the template subject and
437    identifies missing sequences. For example, this table illustrates that subjects cec4ba54 and f53cd86f
438    did not have DWI sequences collected. Panel (**B**) similarly shows completeness of BIDS curation.
439    As expected, the two participants who did not have DWI sequences (in **A**) did not have diffusion data
440    curated into BIDS. Panel (**C**) shows gear run completion; here, flaudit reports that the same two
441    participants that lacked DWI data did not have a successful run of QSIPrep. Finally, the report notes
442    that participant f53cd86f did not yet complete fMRIPrep or XCPEngine successfully.