

## FlywheelTools:

### Data Curation and Manipulation on the Flywheel Platform

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

5 <sup>1</sup>Penn Lifespan Informatics & Neuroimaging Center, Department of Psychiatry, University of  
6 Pennsylvania, Philadelphia, PA, USA

7 <sup>2</sup>Department of Neurology, Perelman School of Medicine, University of Pennsylvania, Philadelphia,  
8 PA, USA.

9 <sup>3</sup>Department of Radiology, University of Pennsylvania, Philadelphia, PA, USA

10 \* **Address correspondence to:** [sattertt@pennmedicine.upenn.edu](mailto: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 a data storage specification for neuroimaging data, but  
17 curating neuroimaging data into BIDS can be a challenging and time-consuming task. In particular,  
18 standard solutions for BIDS curation are not designed for use on cloud-based systems such as  
19 Flywheel. To address these challenges, we developed “FlywheelTools”, a software toolbox for  
20 reproducible data curation and manipulation on Flywheel. FlywheelTools includes two elements: *fw-*  
21 *heudiconv*, for heuristic-driven curation of data into BIDS, and *flaudit*, which audits and inventories  
22 projects on Flywheel. Together, these tools accelerate reproducible neuroscience research on the  
23 widely used Flywheel platform.

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

25

## 26 1 INTRODUCTION

27 Many fields in science are grappling with failures of scientific reproducibility (Botvinik-Nezer et al.  
28 2020). Given the high dimensionality of the data, the need for complex image processing, and a  
29 plethora of analytic techniques, this crisis is particularly acute for neuroimaging research.  
30 Increasingly, the field has recognized that unstructured filesystems are a major source of non-  
31 reproducible practices. As such, major academic centers and large consortia have increasingly  
32 adopted platforms that leverage database technologies that have become standard in other fields. In  
33 addition to providing functionality for searching and categorizing complex source data, imaging  
34 databases enhance reproducible research by providing a clear audit trail of image processing applied  
35 to the data and its results, including both derived images and other data. Widely used imaging  
36 databases include Collaborative Informatics Neuroimaging Suite (COINS) (Landis et al. 2016),  
37 eXtensible Neuroimaging Archive Toolkit (XNAT) (Herrick et al. 2016), Longitudinal Online  
38 Research and Imaging System (LORIS) (Vaccarino et al. 2018), and others (Book et al. 2016;  
39 Helmer et al. 2011; Rogovin et al. 2020; Helmer et al. 2011; Poldrack and Gorgolewski 2017). More  
40 recently, the commercial platform Flywheel has seen growing adoption by major academic imaging  
41 centers (e.g., Stanford, Columbia, Duke, and Penn) due to its modern technology, ease of use, and  
42 scalability.

43 Most neuroimaging databases now leverage the standards defined by the Brain Imaging Data  
44 Structure (BIDS) (Gorgolewski et al. 2016). BIDS is an open-source standard for neuroimaging data  
45 organization that specifies how files should be named, how directories should be organized, and how  
46 metadata should be structured. As such, BIDS provides users with a well-documented format to  
47 understand both imaging data and metadata. Importantly, as BIDS provides a transparent format for  
48 recording imaging parameters and key aspects of the experimental design, BIDS enhances both data  
49 accessibility and data sharing.

50 Furthermore, BIDS allows users to leverage BIDS-apps — image processing pipelines, such as  
51 fMRIPrep, c-Pac, and QSIprep, that read the metadata defined by BIDS (Esteban et al. 2019;  
52 Craddock et al. 2013; Cieslak et al. 2020). As BIDS-apps can auto-configure to ensure that analytic  
53 parameters are appropriate for the input data provided, they dramatically reduce barriers to  
54 implementing best practices in image processing. Furthermore, containerized BIDS-apps encompass  
55 all software dependencies, further enhancing reproducibility.

56 On a filesystem, conversion of raw DICOM images to NIFTIs that conform to BIDS can be  
57 accomplished with a variety of tools including *HeuDiConv*, *dcm2bids*, and others. However, this  
58 crucial step, a process typically called “BIDS curation,” is often difficult to implement within the  
59 database environment, as it requires a mechanism for converting data in hierarchical databases into  
60 BIDS. As BIDS curation is one of the very first steps performed on the data, lack of transparency in  
61 curation can impact the reproducibility of all subsequent steps. Here we introduce FlywheelTools, a  
62 software suite that provides flexible and reproducible methods for BIDS curation on the Flywheel  
63 platform. Documentation and code presented in this work can be found online at: [https://fw-  
64 heudiconv.readthedocs.io/en/latest/](https://fw-heudiconv.readthedocs.io/en/latest/).

## 65 2 METHODS

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

## 70 **2.1 Programming Languages & Technologies**

71 FlywheelTools is built primarily in Python 3.6 (Van Rossum and Drake 2009) to leverage Flywheel's  
72 highly accessible Software Development Kit (SDK). Additionally, R 3.4.1 (R Core Team 2019) is  
73 used for HTML report generation. For reproducibility and workflow management, the modules of  
74 FlywheelTools are packaged in version-controlled software containers built and managed in Docker  
75 (Merkel 2014). Lastly, the FlywheelTools package relies on users adopting BIDS to curate their data.  
76 BIDS has rapidly evolved to become the directory standard in the neuroimaging community for  
77 reproducible data organization. Importantly, BIDS is supported by a large community that  
78 contributes to its development and adoption. Further, proposals for BIDS schema pass through a  
79 rigorous testing process before being adopted. Software developers leverage this ubiquity by creating  
80 BIDS-apps: analysis and processing pipelines that operate directly on BIDS datasets as inputs,  
81 enhancing reproducibility and interoperability.

## 82 **2.2 Flywheel**

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

## 87 **2.3 Flywheel Web UI**

88 The web UI is accessible through any modern web browser. Through this point-and-click interface,  
89 users are able to upload, view, download, and analyze data with ease. However, accomplishing tasks  
90 with many repetitive steps or over a large number of participants/sessions can be tiresome and error-  
91 prone. Alongside knowledge of navigating the web UI, many users also make use of the API and  
92 SDK to manipulate and analyze data programmatically.

## 93 **2.4 Flywheel API & SDK**

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

## 103 **2.5 Flywheel Data Model**

104 Objects in Flywheel's data model follow a specific hierarchical structure — at the top level is a  
105 Flywheel instance, a process that serves the API to an organization (for example, a neuroimaging  
106 center). Within the Flywheel instance, there are multiple groups, which are typically labs or research  
107 units that collaborate on one or more projects. Each project object can have one or many subjects (i.e.

---

<sup>1</sup> 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.

108 participants), and each subject can have one or many sessions (i.e. scanning visits). Within a session,  
109 there may be one or many acquisition objects which represent the scanning sequences collected  
110 during a particular scan or examination (e.g., sMRI, rs-fMRI, dMRI). Finally, the data files  
111 associated with the sequence (e.g., NIfTIs or DICOMs) are attached to each acquisition. Note that a  
112 file can additionally be attached to any object type, and each object can have metadata associated  
113 with it. Hence, a “subject” object may have metadata associated with that participant (such as  
114 demographic information) and may also have a text file attached to it (such as clinical data). A  
115 notable exception to this hierarchical structure is the analysis object, which behaves in much the  
116 same way as others but can be a child object of any other object, allowing researchers to create  
117 analyses of entire projects, for example, each with their own associated metadata and files.

118  
119 Abstracting this data model in Python results in simple hierarchical objects, each with methods for  
120 handling metadata and files, and methods for accomplishing object-specific tasks like traversing the  
121 hierarchical structure or running analyses. The modules of FlywheelTools make use of this data  
122 model to accomplish a wide range of tasks.

## 123 **2.6 Flywheel Gears**

124 Flywheel encourages the use of pre-packaged computational workflows, called “gears”. Gears are  
125 run by virtual machines/containers using Docker and hence are version-controlled and  
126 software/platform agnostic. Gears can accomplish tasks such as data manipulation, pre-processing,  
127 and analysis. In addition to the existing gears available on the platform, users are able to package  
128 their own software in a gear and use it for running analysis workflows on their Flywheel data via the  
129 web UI or SDK. The complexity and frequency of the task guides if a task should be accomplished  
130 using the web UI, programmatically using the SDK, or by wrapping it as a workflow into a gear.  
131 Gears ingest existing Flywheel data (such as images or file attachments) as inputs to the workflow  
132 and can be created with clickable configuration options. Once a workflow has completed running,  
133 Flywheel collects any files remaining in the pre-defined output directory of the container and attaches  
134 them to a resulting analysis object. The output of a gear (such as an HTML report or tabulated data)  
135 can be viewed on the Flywheel UI, downloaded to disk for further sharing or analysis, or used as  
136 input to a subsequent gear.

137

## 138 **3 RESULTS**

139 FlywheelTools is implemented using the Flywheel SDK to enable easy inspection, curation,  
140 validation, and audit of Flywheel data through a handful of user-friendly gears and command-line  
141 interfaces. The first module of the package is called *fw-heudiconv* and is largely inspired by the  
142 popular Heuristic DICOM Converter package (Halchenko et al. 2018). *fw-heudiconv* is a multi-part  
143 toolbox for reproducible curation of neuroimaging data into BIDS on Flywheel. The second module,  
144 *flaudit*, is a tool for auditing a Flywheel project, giving users an overview of the key elements of their  
145 data set.

### 146 **3.1 FW-HEUDICONV**

147 The first tool, *fw-heudiconv*, is a multi-purpose command-line interface and Flywheel gear designed  
148 for BIDS curation on Flywheel (**Figure 1**). It is designed to be intuitive, flexible, and reproducible.

#### 149 **3.1.1 Architecture & Design**

150 *fw-heudiconv* is inspired in large part by the Heuristic DICOM Converter (*HeuDiConv*) package, and  
151 shares much of its design practices. To curate data into BIDS format, *fw-heudiconv* first considers  
152 DICOM data to be the “ground truth” and builds its curation approach using data in the DICOM  
153 headers. On Flywheel, when DICOMs are added to the database, either through manual upload or  
154 automatically from a linked scanner. The DICOMs are then automatically converted into NIfTI files  
155 by Flywheel’s automated gears. The result is an acquisition object with both DICOMs and one or  
156 more NIFTIs. Ultimately, *fw-heudiconv* only has permission to manipulate metadata associated with a  
157 NIFTI file. By not manipulating DICOMs or their associated metadata, BIDS curation can safely be  
158 reproduced from ground truth data.

159 *fw-heudiconv* can be downloaded as a Python command-line interface from the Python Package  
160 Index using *pip*. It is also available as a point-and-click gear that is managed by Docker  
161 containerization. This containerization allows for reproducible version control, as every computation  
162 can be recorded, including information about the software, OS, and algorithm versions. There are a  
163 number of commands available in *fw-heudiconv*, and each of them starts by querying data from  
164 Flywheel. Users can filter their queries to operate on an entire Flywheel project, a subset of subjects,  
165 or a subset of sessions. Notably, with the `--dry-run` option, each command has the ability to test  
166 and evaluate its effects without actually manipulating metadata in the Flywheel database or writing  
167 data to disk. Below, we consider each of the five available commands.

### 168 **3.1.1.1 fw-heudiconv-tabulate**

169 The tabulate tool is used to parse and extract DICOM header information in a project (or within a  
170 filtered subset of that project) and compile these data into a table for the user to examine. By  
171 collecting DICOM header information into a tabular format, the tabulate tool gives users a  
172 comprehensive overview of the different scanning sequences that have been collected in the query,  
173 including the sequence parameters. Additionally, users have the option to limit the tabulation to a  
174 unique combination of common DICOM header fields, which significantly decreases the complexity  
175 of the table. When used at the command line, the table produced by this command is written to a  
176 local disk. As a gear on Flywheel, the table is automatically saved in the output section of the gear.

### 177 **3.1.1.2 fw-heudiconv-curate**

178 The curate tool is used to curate a dataset on Flywheel into BIDS format. Much like *HeuDiConv*,  
179 curation is accomplished through the use of a heuristic: a Python file that programmatically defines  
180 the templates for a range of BIDS-valid filenames, and defines the boolean logic that would assign a  
181 given scanning sequence to each template. This boolean logic is usually based on the sequence  
182 information users find in the tabulation of sequences, but all fields available in the DICOM header  
183 can be used to determine which template a particular file can be assigned to. Additionally, the curate  
184 tool can be used to manipulate BIDS metadata that may need to be added to the dataset. The process  
185 of curation only manipulates the BIDS metadata of NIFTI files, and hence can be repeated or updated  
186 at any time at the user’s discretion.

### 187 **3.1.1.3 fw-heudiconv-export**

188 The export tool is used to export a BIDS dataset on Flywheel to disk. It can also be used by other  
189 gears and scripts to quickly and easily extract their BIDS data into the workspace of their analysis  
190 pipeline.

### 191 **3.1.1.4 fw-heudiconv-validate**

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

### 199 3.1.1.5 fw-heudiconv-clear

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

### 203 3.1.2 The Heuristic File

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

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

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

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

```
223 IntendedFor = {  
224   fieldmap1: ['sub-{{SubjectLabel}}_ses-{{SessionLabel}}_task-rest_bold.nii.gz']  
225 }
```

226 By reserving select keywords for functions and metadata, heuristic files become versatile tools for  
227 defining and manipulating a wide array of metadata in Flywheel BIDS curation.

228 Importantly, because this heuristic file is plain text Python code, users are able to version control  
229 their files using Git and share these files via Github. Finally, when run on Flywheel as a gear, the

230 heuristic file is automatically attached as an input to the analysis object created by *fw-heudiconv-*  
231 *curate*, allowing users to easily access the version history of their curation.

### 232 3.1.3 Curation Workflow

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

242 If *fw-heudiconv* is run from the Flywheel GUI, each of the commands is available as a Flywheel gear.  
243 This option is beneficial for data provenance, as all of a gear's commands and inputs, as well as  
244 outputs and logs, are stored and attached to each gear run.

## 245 3.2 FLAUDIT

246 The second module of FlywheelTools is a Flywheel project auditor, named *flaudit*. The module is  
247 intended to give Flywheel users a broad understanding of their entire Flywheel project, by  
248 summarizing the available data and illustrating analysis workflows. The output of this module, a  
249 portable HTML report, presents this information using a number of visualizations built in R  
250 Markdown using HTML, Javascript, and ggplot2, in two main sections: project overview and project  
251 completeness.

### 252 3.3 Architecture & Design

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

### 258 3.4 Flaudit: Project Overview

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

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

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



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

### 272 **3.5 Flaudit: Project Completeness**

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

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

## 291 **4 DISCUSSION**

292 FlywheelTools provides new capabilities for the popular and powerful Flywheel platform, allowing  
293 researchers to maximize reproducibility and enhance scalability. Specifically, *fw-heudiconv* provides  
294 users a flexible and reproducible way of curating data into BIDS on Flywheel. Complementary  
295 function is provided by *flaudit*, which provides intuitive visual reports of raw, curated, and processed  
296 data.

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

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

308 Once data are curated with *fw-heudiconv*, *flaudit* allows users to audit the data in a Flywheel project.  
309 Specifically, *flaudit* provides intuitive summaries at each stage of a typical workflow, concisely  
310 visualizing raw data, curated data in BIDS, and data processed by containerized analytic gears. These

311 accessible reports allow users to rapidly assess the overall organizational state of a project, while  
312 interactive tables allow for more granular inspection of data. This approach facilitates understanding  
313 the diverse data types typically collected in multi-modal imaging studies.

314 There are of course limitations of FlywheelTools. First, it should be acknowledged that  
315 FlywheelTools is built for the Flywheel platform, and as such does not generalize to other imaging  
316 databases that are in use. However, given the rapid adoption of this platform by the imaging  
317 community, we anticipate that this toolkit will fill an important need for the many large research  
318 institutions that rely upon Flywheel. Second, some understanding of Python is necessary to build the  
319 heuristic for *fw-heudiconv*. We attempt to minimize this issue by providing both extensive  
320 documentation and heuristic templates for various uses of *fw-heudiconv*<sup>2</sup>, but usage is ultimately a  
321 programming task.

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

325

## 326 **5 CONFLICT OF INTEREST**

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

## 329 **6 AUTHOR CONTRIBUTIONS**

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

## 335 **7 FUNDING**

336 Support was provided by R01MH120482, R01MH112847, R01MH113550, and RF1MH116920.

337

## 338 **8 ACKNOWLEDGMENTS**

339 None.

340

---

<sup>2</sup><https://fw-heudiconv.readthedocs.io/en/latest/index.html>

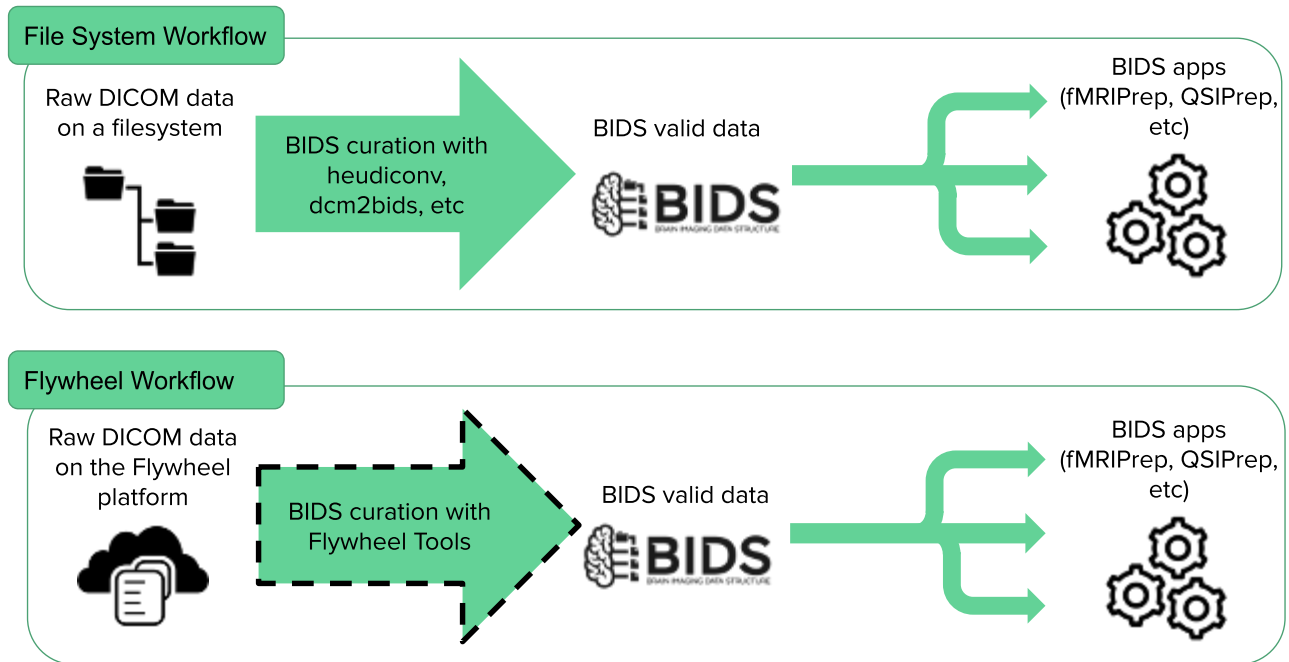
341 **9 REFERENCES**

- 342 Banker, Kyle. 2011. MongoDB in Action. USA: Manning Publications Co.
- 343 Biehl, Matthias. 2016. RESTful Api Design. Vol. 3. API-University Press.
- 344 Book, Gregory A, Michael C Stevens, Michal Assaf, David C Glahn, and Godfrey D Pearlson. 2016.  
345 “Neuroimaging Data Sharing on the Neuroinformatics Database Platform.” *Neuroimage* 124: 1089–  
346 92.
- 347 Botvinik-Nezer, Rotem, Felix Holzmeister, Colin F Camerer, Anna Dreber, Juergen Huber, Magnus  
348 Johannesson, Michael Kirchler, et al. 2020. “Variability in the Analysis of a Single Neuroimaging  
349 Dataset by Many Teams.” *Nature*, 1–7.
- 350 Cieslak, Matthew, Philip A Cook, Xiaosong He, Fang-Cheng Yeh, Thijs Dhollander, Azeez  
351 Adebimpe, Geoffrey K Aguirre, et al. 2020. “QSIprep: An Integrative Platform for Preprocessing  
352 and Reconstructing Diffusion Mri.” *bioRxiv*.
- 353 Craddock, Cameron, Sharad Sikka, Brian Cheung, Ranjeet Khanuja, Satrajit S Ghosh, Chaogan Yan,  
354 Qingyang Li, et al. 2013. “Towards Automated Analysis of Connectomes: The Configurable Pipeline  
355 for the Analysis of Connectomes (c-Pac).” *Front Neuroinform* 42.
- 356 Esteban, Oscar, Christopher J Markiewicz, Ross W Blair, Craig A Moodie, A Ilkay Isik, Asier  
357 Erramuzpe, James D Kent, et al. 2019. “fMRIPrep: A Robust Preprocessing Pipeline for Functional  
358 Mri.” *Nature Methods* 16 (1): 111–16.
- 359 Halchenko, Yaroslav, Mathias Goncalves, Matteo Visconti di Oleggio Castello, Satrajit Ghosh,  
360 Michael Hanke, Matthew Brett, Taylor Salo, et al. 2018. Nipy/Heudiconv: Heudiconv V0.5.1  
361 (version v0.5.1). Zenodo. <https://doi.org/10.5281/zenodo.1306159>.
- 362 Helmer, Karl G, Jose Luis Ambite, Joseph Ames, Rachana Ananthakrishnan, Gully Burns, Ann L  
363 Chervenak, Ian Foster, et al. 2011. “Enabling Collaborative Research Using the Biomedical  
364 Informatics Research Network (Birn).” *Journal of the American Medical Informatics Association* 18  
365 (4): 416–22.
- 366 Herrick, Rick, William Horton, Timothy Olsen, Michael McKay, Kevin A Archie, and Daniel S  
367 Marcus. 2016. “XNAT Central: Open Sourcing Imaging Research Data.” *NeuroImage* 124: 1093–6.
- 368 Landis, Drew, William Courtney, Christopher Dieringer, Ross Kelly, Margaret King, Brittny Miller,  
369 Runtang Wang, Dylan Wood, Jessica A Turner, and Vince D Calhoun. 2016. “COINS Data  
370 Exchange: An Open Platform for Compiling, Curating, and Disseminating Neuroimaging Data.”  
371 *NeuroImage* 124: 1084–8.
- 372 Merkel, Dirk. 2014. “Docker: Lightweight Linux Containers for Consistent Development and  
373 Deployment.” *Linux J.* 2014 (239).
- 374 Poldrack, Russell A, and Krzysztof J Gorgolewski. 2017. “OpenfMRI: Open Sharing of Task fMRI  
375 Data.” *Neuroimage* 144: 259–61.
- 376 R Core Team. 2019. R: A Language and Environment for Statistical Computing. Vienna, Austria: R  
377 Foundation for Statistical Computing. <https://www.R-project.org/>.

- 378 Rogovin, O, Y Zhao, S Chen, Z Wang, O Papaemmanouil, SD Van Hooser, and others. 2020. “NDI:  
379 A Platform-Independent Data Interface and Database for Neuroscience Physiology and Imaging  
380 Experiments.”
- 381 Vaccarino, Anthony L, Moyez Dharsee, Stephen Strother, Don Aldridge, Stephen R Arnott, Brendan  
382 Behan, Costas Dafnas, et al. 2018. “Brain-Code: A Secure Neuroinformatics Platform for  
383 Management, Federation, Sharing and Analysis of Multi-Dimensional Neuroscience Data.” *Frontiers  
384 in Neuroinformatics* 12: 28.
- 385 Van Rossum, Guido, and Fred L. Drake. 2009. *Python 3 Reference Manual*. Scotts Valley, CA:  
386 CreateSpace.

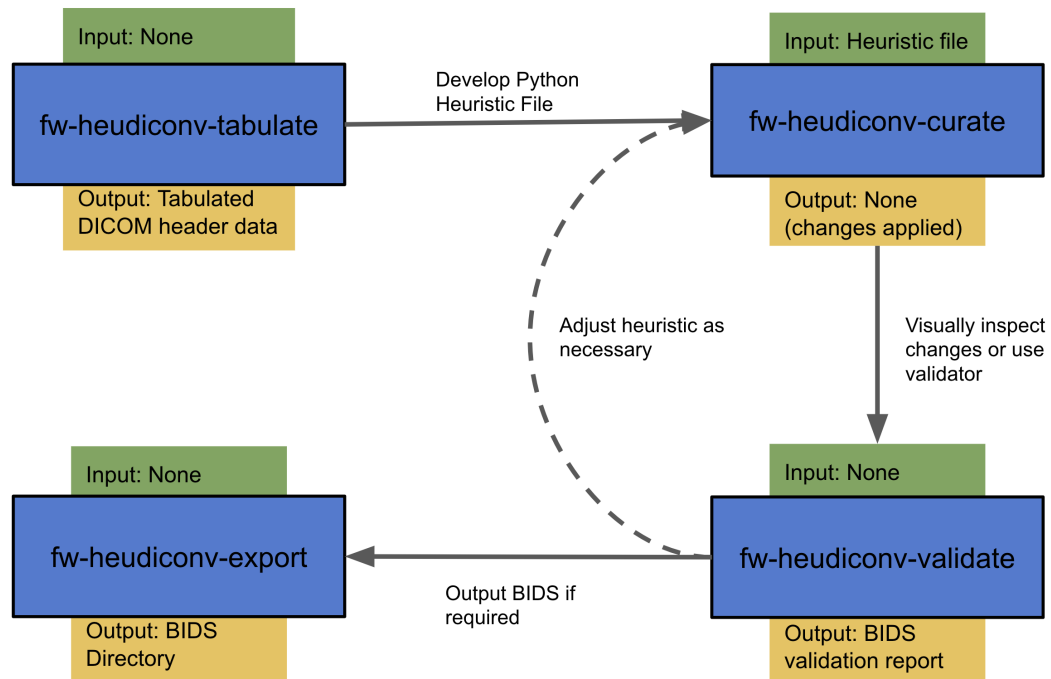
387

388 **10 FIGURES**



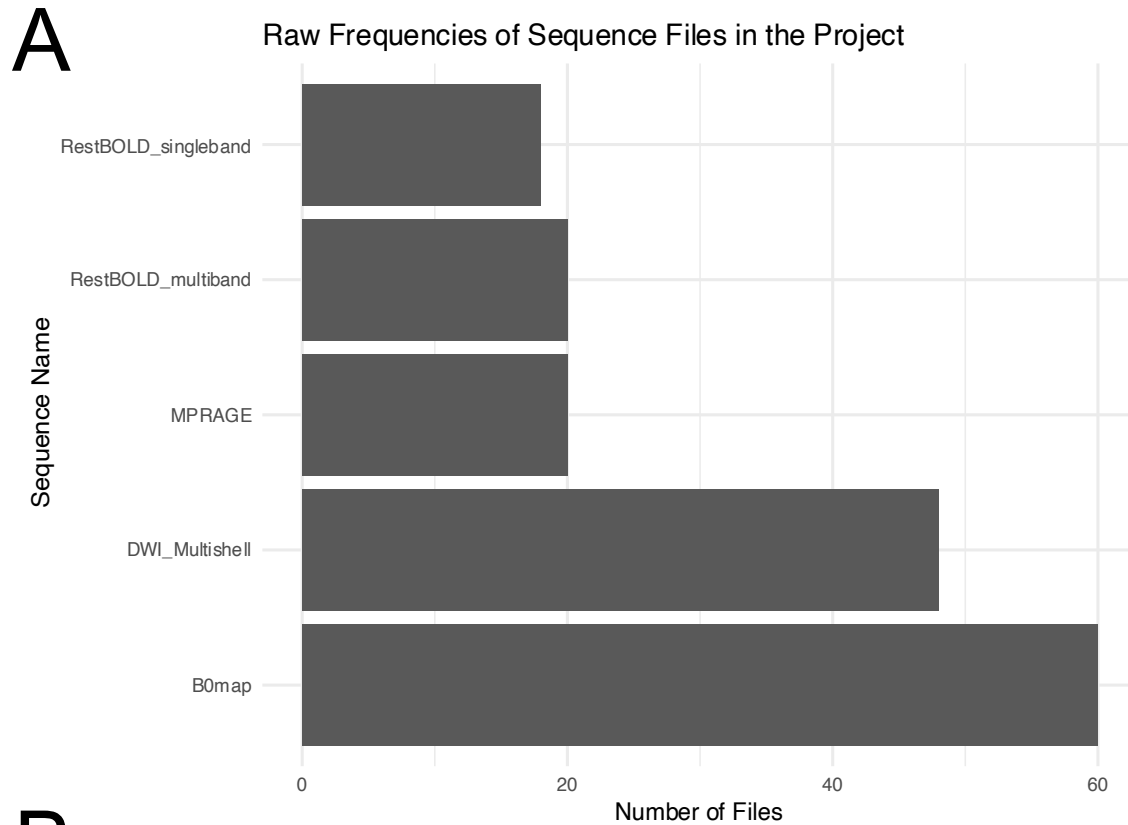
389

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



394

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



**B**

Show 10 ▾ entries

Raw Frequencies of Sequence Files in the Project

Search:

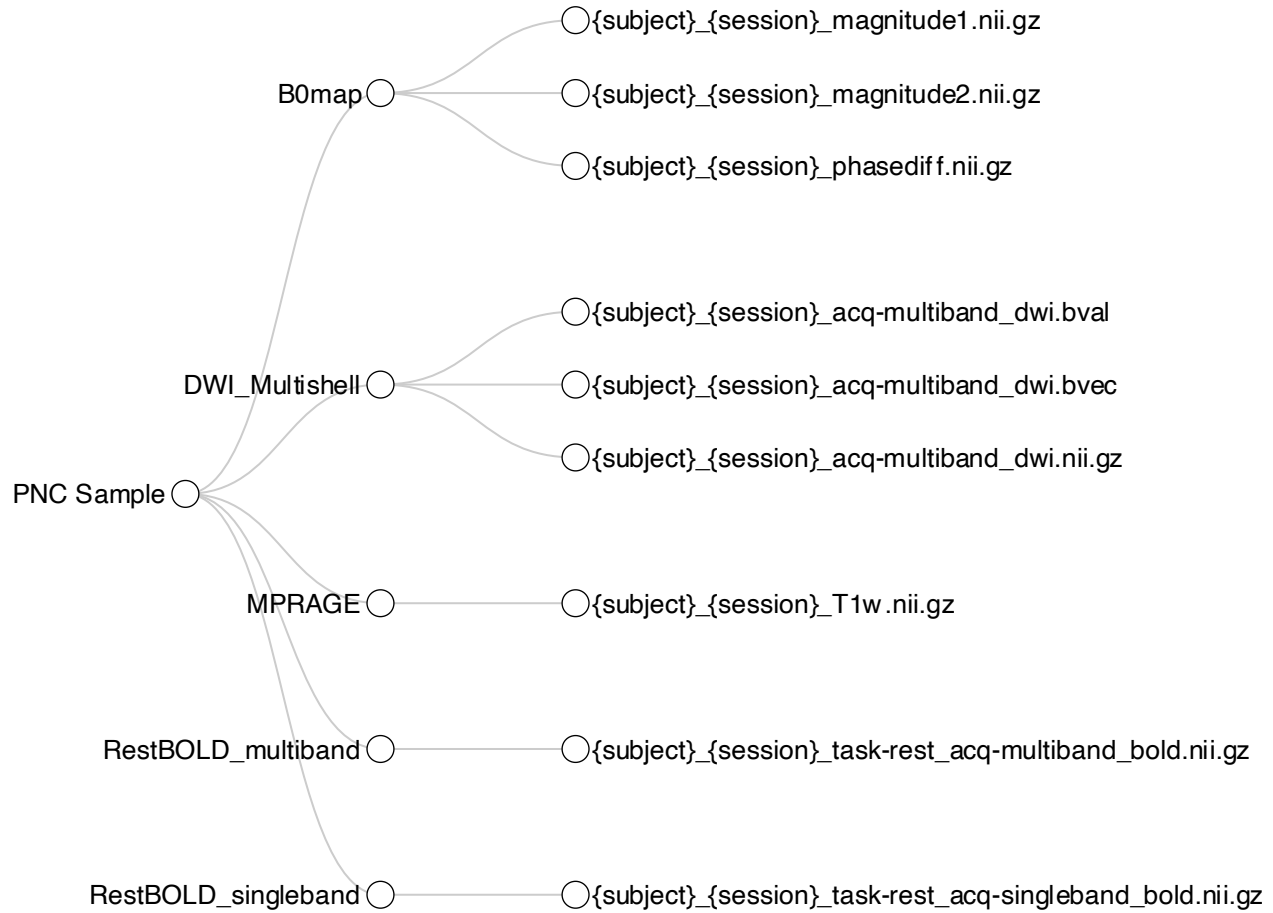
	Sequence Name	Freq
1	B0map	60
2	DWI_Multishell	48
3	MPRAGE	20
4	RestBOLD_multiband	20
5	RestBOLD_singleband	18

Showing 1 to 5 of 5 entries

Previous  Next

399

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



403

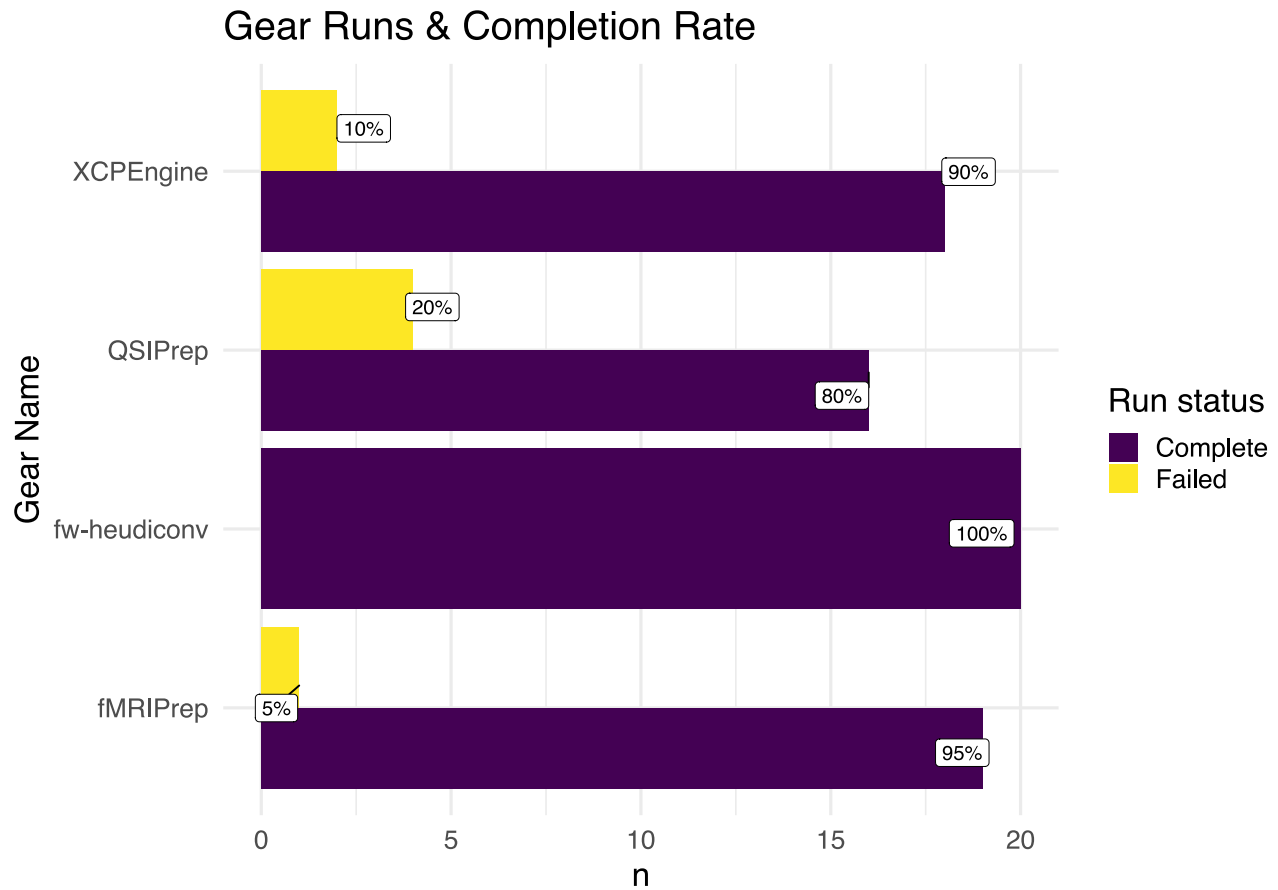
404

405

406

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





407

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

**A** Show 5 ▾ entries  
Sequence Completeness Compared to the Template Subject

Search:

	subject	session	B0map	DWI_Multishell	MPRAGE	RestBOLD_multiband	RestBOLD_singleband
	All	All	All	All	All	All	All
1	b9595315	ses-1	✓ complete	✓ complete	✓ complete	✓ complete	✓ complete
2	cec4ba54	ses-1	✓ complete	✗ incomplete	✓ complete	✓ complete	✓ complete
3	f53cd86f	ses-1	✓ complete	✗ incomplete	✓ complete	✓ complete	✓ complete
4	487753ab	ses-1	✓ complete	✓ complete	✓ complete	✓ complete	✓ complete
5	33c33e4a	ses-1	✓ complete	✓ complete	✓ complete	✓ complete	✓ complete

Showing 1 to 5 of 19 entries

Previous  2 3 4 Next

**B** Show 5 ▾ entries  
BIDS Curation Completeness Compared to the Template Subject

Search:

	subject	session	{subject}_{session}_acq-multiband_dwi.bval	{subject}_{session}_acq-multiband_dwi.bvec	{subject}_{session}_acq-multiband_dwi.nii.gz	{subject}_{session}
	All	All	All	All	All	All
1	b9595315	ses-1	✓ complete	✓ complete	✓ complete	✓ complete
2	cec4ba54	ses-1	✗ incomplete	✗ incomplete	✗ incomplete	✓ complete
3	f53cd86f	ses-1	✗ incomplete	✗ incomplete	✗ incomplete	✓ complete
4	487753ab	ses-1	✓ complete	✓ complete	✓ complete	✓ complete
5	33c33e4a	ses-1	✓ complete	✓ complete	✓ complete	✓ complete

Showing 1 to 5 of 19 entries

Previous  2 3 4 Next

**C** Show 5 ▾ entries  
Gear Analysis Completeness Compared to the Template Subject

Search:

	subject	session	fw-heudiconv	QSIPrep	fMRIprep	XCPengine
	All	All	All	All	All	All
1	b9595315	ses-1	✓ complete	✓ complete	✓ complete	✓ complete
2	cec4ba54	ses-1	✓ complete	✗ incomplete	✓ complete	✓ complete
3	f53cd86f	ses-1	✓ complete	✗ incomplete	✗ incomplete	✗ incomplete
4	487753ab	ses-1	✓ complete	✓ complete	✓ complete	✓ complete
5	33c33e4a	ses-1	✓ complete	✓ complete	✓ complete	✓ complete

Showing 1 to 5 of 19 entries

Previous  2 3 4 Next

411

412 **Figure 6: Project completeness tables compared to the template participant in a Flywheel**  
 413 **project.** Panel (A) compares the sequences available for each participant to the template subject and  
 414 identifies missing sequences. For example, this table illustrates that subjects cec4ba54 and f53cd86f  
 415 did not have DWI sequences collected. Panel (B) similarly shows completeness of BIDS curation.  
 416 As expected, the two participants who did not have DWI sequences (in A) did not have diffusion data  
 417 curated into BIDS. Panel (C) shows gear run completion; here, fludit reports that the same two  
 418 participants that lacked DWI data did not have a successful run of QSIPrep. Finally, the report notes  
 419 that participant f53cd86f did not yet complete fMRIprep or XCPengine successfully.