# Assessing Value of Biomedical Digital Repositories[*]

Chun-Nan Hsu[1], Anita Bandrowski[2], Jeffrey S. Grethe[3], Maryann E. Martone[3]

[1]Division of Biomedical Informatics, Department of Medicine, University of California, San Diego, La Jolla, CA, USA.

[2]SciCrunch Inc. San Diego, CA, USA.

[3]Center for Research in Biological Structure, University of California, San Diego, La Jolla, CA, USA.

Corresponding Author:

Chun-Nan Hsu

Email address: chunnan@ucsd.edu

---

[*] This manuscript was originally prepared as a response to *NOT-OD-16-133 Request for Information (RFI): Metrics to Assess Value of Biomedical Digital Repositories* by the National Institute of Health. C.-N. Hsu presented the contents orally in the 2016 workshop on collaborative data projects, Taipei, Taiwan, December 8, 2016. http://odw.tw/2016/

## Abstract

18   Digital repositories bring direct impact and influence on the research community and society but
19   measuring their value using formal metrics remains challenging.   their value. It is challenging to
20   define a single perfect metric that covers all quality aspects. Here, we distinguish here between
21   impact and influence and discuss measures and mentions as the basis of quality metrics of a
22   digital repository. We argue that these challenges may potentially be overcome through the
23   introduction of standard resource identification and data citation practices. We briefly summarize
24   our research and experience in the Neuroscience Information Framework, the BD2K BioCaddie
25   project on data citation, and the Resource Identification Initiative. Full implementation of these
26   standards will depend on cooperation from all stakeholders --- digital repositories, authors,
27   publishers, and funding agencies, but both resource and data citation have been gaining support
28   with researchers and publishers.

## Impact vs. Influence

30   Assessing the value of digital repositories shares many similar challenges to assessing the value
31   of any scholarly work. One challenge is whether to distinguish between direct impact and broad
32   influence. By direct impact we refer to actual changes that the work brings to the field in terms
33   of outcomes, practices, and methodologies. In biomedical sciences, these include, for example,
34   new drugs or treatments for disease, new models of molecular interactive pathways, new
35   experimental methods, etc.   By influence, we refer to how widely the work has been
36   disseminated and viewed across a broad community so that a work can influence other work,
37   either by inspiring new research ideas or preliminary testing of hypotheses.   Impact and
38   influence may be correlated but that is not always the case. A highly influential work may have a
39   low direct impact and vice versa. A digital repository may have a high influence in that it is
40   viewed many times, but low impact in that there is no evidence that the actual products are used
41   to advance science. However, the products may be very useful for educational purposes or in
42   planning research. The converse is also true; a digital repository may not be well known across a
43   wide swath of the community, but its products may be highly impactful in a smaller community.
44   Understanding where each resource fits and therefore how to evaluate their success and perhaps

45 improve both dimensions requires that it be possible to measure these in some objective and
46 preferably automated or semi-automated way.

# Measure vs. Mention

48 While traditional metrics of a scientific work are based on citations -- whether the work is
49 *mentioned* in scientific publications, digital repositories allow *measures* through the count of
50 access in different ways, URL connections, data transferring, etc. One may argue that measures
51 of access more accurately reflect the value of a digital repository for without access, a digital
52 repository is not used and cannot create value. However, as discussed above, the value of a work
53 may present as impact or influence. Usually, mention-based metrics, such as citations, reflect
54 influence better, for a work can be mentioned only after it is known. However, citations can also
55 reflect actual use of the resource within a published study.  Currently, both are hard to track; this
56 makes proper citation of resources and their data products in the literature extremely important.
57 Measure and mention are not always correlated for a digital repository (Huang et al. 2015;
58 Huang 2016; Rose & Hsu 2016). Moreover, different measure-based metrics, for example, URL
59 connection count, and FTP download count, size of data transferring, are not always correlated.
60 The lack of correlation applies not only when comparing digital repositories but also when
61 comparing content units within a digital repository. Results in (Huang 2016; Rose & Hsu 2016)
62 show that ranking protein structures in RCSB PDB (Protein Data Bank), a data repository of
63 protein structure data, by different measures of access give uncorrelated results. In the study, we
64 ranked protein structures according to their frequencies of Web access (http views) and FTP
65 access (file downloads). We found that *no* protein structures were shared in the top 20 of the two
66 resulting ranked lists. Moreover, the two frequencies are not correlated, in the sense that a
67 protein structure that is highly accessed by Web browsers is not necessarily highly accessed by
68 FTP, and vice versa.
69 Meanwhile, in addition to citations in publications, mention-based metrics may include citations
70 in press reports, blogs, social media, and other forms of publications, currently measured by
71 services such as Altmetrics (2016).  These may not be correlated either, and may better reflect
72 the influence of a work than its impact.  Citation analysis is currently hampered by a lack of
73 standard format for such references.  Citations may be in different forms, including directly
74 mentioning various names of a digital repository, citing the publications describing a digital

75    repository or mentioning the URL links to a digital repository. For example, an author may cite

76    RCSB PDB by its various publications, URL links to its portal Web page (with different versions

77    throughout the years after it went online), PDB IDs or URL links of protein structures.

78    Authors not only cite RCSB PDB in different forms, the annual growth rates of the counts of

79    these different citations forms are not correlated for either the data repository as a whole (Huang

80    et al. 2015), or for protein structures (Huang 2016; Rose & Hsu 2016). Authors most frequently

81    chose to cite publications, because usually that is how repositories instruct authors to do in a

82    "how to cite us" page. However, URL link mentions are growing rapidly. Though the PDB ID is

83    designed as a unique ID to mention specifically to a protein structure in PDB, the ID itself is not

84    globally unique without a prefix, and may coincide with a wide variety of entities (Rose & Hsu

85    2016). PDB IDs are always 4 characters in length. The first character is a numeral in the range 1-

86    9, while the last three characters can be either numerals (in the range 0-9) or letters. Examples of

87    other IDs and/or entities matching this format include "1USD" as currency, "2NO3" as a

88    chemical compound, and "1E10" as a floating-point number; while "1USD", "2NO3" and

89    "1E10" are all legitimate PDB IDs.

90

91                                    (Table 1)

92

93    Table 1 shows all the issued PDB IDs presented in full-text format articles. The statistics were

94    obtained from publications containing mentions of PDB ID from the PubMed Central (PMC),

95    where we obtained 1,015,179 articles in NXML format, and 1,093,980 articles in plain text

96    format as of August 2015. Removing duplicate PMC IDs yielded a total of 1,015,233 articles.

97

98                                    (Table 2)

99

100   Table 2 compares the top 10 PDB protein structures by the frequency of PDB ID mentions and

101   the top 10 ordered by the frequency that their original publications were cited in the references

102   by subsequent articles in the PubMed. The two lists share only two PDB protein structures

103   (2RH1 and 2A79), suggesting that high PDB ID mentions and high publication citations are not

104 necessarily correlated (Huang 2016). Therefore, relying on either frequency as a sole metric will
105 lead to different assessments of the influence of protein structures.


# Standardization of Mentions and Use

107 Currently, one of the most difficult problems facing assessments of digital repositories is the lack
108 of formal systems of citation that allow measures of influence and direct impact to be calculated
109 using modern information technology. As documented by (Huang et al. 2015), the current
110 means of referencing a digital repository or its content in the literature or any other work involve
111 a range of styles including URLs, reference to a publication describing the resource, accession
112 numbers and free text. Because of this, a very simple question like: how many people have
113 documented use of this resource cannot be answered without resorting to extensive manual labor
114 or advanced natural language processing (NLP) (Rose & Hsu 2016; Ozyurt et al. 2016).

115 Through the Neuroscience Information Framework, the NIDDK Information Network, the
116 Resource Identification Initiative, and the Data Citation Working groups at FORCE11, we've
117 successfully worked to change this by developing and promoting standards for both resource use
118 and data citation, with a focus on the literature.

119 **Perspectives from the Neuroscience Information Framework.** The Neuroscience Information
120 Framework (NIF; Gupta et al. 2008; Gardner et al. 2008) and its sister project, the NIDDK
121 Information Network (dkNET; Whetzel et al., 2015) has been cataloging and tracking the digital
122 research resource landscape for over 8 years. We maintain a large database that tracks how a
123 resource has evolved over the years, including whether it is no longer in service. Currently, a
124 relatively small number of resources (229, List from NIF 2016) are completely out of service;
125 many more, however, grow stale over time. Over time, we have developed some criteria for
126 determining whether a resource is vibrant and growing or moribund: 1) when was the last time a
127 web page was updated on the site; 2) when was the last time data were added; 3) do the data
128 represent a significant fraction of data available in a community or a very limited amount? 4)
129 When a resource is down, does anyone complain? We call the latter the "squawk factor".
130 As the charge for NIF, established by the NIH Blueprint Consortium, was to determine what
131 resources had been created through NIH-funding and to make them available to the research
132 community, NIF early on worked to develop NLP pipelines to identify resources within the

133     biomedical literature, as most researchers creating resources will publish a paper about a
134     resource like a database or genetically modified organism, or will mention use of specific
135     resources within the materials and methods section of the paper.

136     The lack of formal or machine-readable standards for referencing these resources within the
137     literature uncovered significant problems in the way that researchers were recording resource
138     usage. These poor reporting standards directly led to the inability of funders or resource
139     providers to track usage or for researchers to identify research resources or find other papers that
140     used them. To address this, NIF worked through FORCE11 to launch the Resource
141     Identification Initiative.

142     **The Resource Identification (#RRID) Initiative** (Bandrowski et al. 2016) is designed to help
143     researchers sufficiently cite the key resources used to produce the scientific findings reported in
144     the biomedical literature. A diverse group of collaborators are involved in the project, including
145     the Neuroscience Information Framework which launched and has been leading the initiative, the
146     Oregon Health & Science University Library which contributed significantly to the early pilot
147     project, with the support of the National Institutes of Health and the International
148     Neuroinformatics Coordinating Facility (2016). Resources (e.g. antibodies, model organisms,
149     cell lines and digital tools) reported in the biomedical literature often lack sufficient detail to
150     enable reproducibility or reuse (Vasilevsky et al., 2013). For example, databases are cited by a
151     URL that is no longer available leading to 404 errors, and the version numbers for software
152     programs used for data analysis are often omitted as is the access date of digital repositories.

153     The Resource Identification Initiative aims to enable resource transparency within the
154     biomedical literature through promoting the use of unique Research Resource Identifiers
155     (RRIDs). In addition to being unique, RRID's meet three key criteria, they are:

156             1. Machine readable and search friendly.

157             2. Free to generate and access.

158             3. Consistent across publishers and journals.

159     RRID's depend on comprehensive resource registries which provide an authoritative source for
160     each resource type. Each is covered by a different database, e.g., the Antibody Registry, the
161     SciCrunch (NIF) Resource Registry. These databases were aggregated and made available
162     through the Resource Identification Portal (2017), supporting NIH's new guidelines for Rigor

163    and Transparency in biomedical publications. The portal aims to promote research resource

164    identification, discovery, and reuse and offers a central location for obtaining and exploring

165    RRIDs.   The current number of digital tools, including databases, software projects as well as

166    commercial tools, listed in the Registry is over 14K (Bandrowski et al. 2016). The number of

167    antibodies is $> 2.5M$, model animals are in the hundreds of thousands and cell lines over 60K.

168    The project has been running since 2014.  Currently, over 2,500 papers have appeared with

169    RRID's from over 200 biomedical journals.  Cell Press has just adopted the standard (Marcus et

170    al, 2016; Cell 2016) and eLife and the Endocrine Society announced that they will be strongly

171    encouraging authors to use RRID's in their journals.

172    RRID's provide the means for users to unambiguously reference the resources used within a

173    study in their publication.  Authors are asked to insert RRID's for resources *used* in their studies

174    after the first reference to the resource in the materials and methods.  To ensure that RRID's are

175    easily identified and extractable from the literature, authors are asked to prepend the namespace

176    RRID: before using the database accession number.  Thus, RRID's specifically target the use of

177    resource resources as opposed to mentions in an introduction or discussion.  A simple search

178    through Google Scholar for an RRID will return papers that have used a resource, e.g., 6 articles

179    have appeared to date that used the PDB (Google Scholar PDB 2016).

180    RRID's also provide a convenient means for authors to access which digital resources used in

181    papers. Research resource providers can update the registry in the portal when there is a need to

182    transfer the data and software to another repository, but the RRID will remain the same to ensure

183    that readers can always locate the data and software through a centralized registry. This new

184    approach solves data access, sharing, archiving, and preservation at the same time. In addition, it

185    provides a standard citation format that can be easily extracted to show what resources were used

186    in a published study - allowing for measurement of impact. For example, consider FSL, a

187    widely-used software library for functional MRI. Querying Google Scholar with query string

188    "RRID:SCR_002823" will precisely match 26 recent articles reporting research results using

189    FSL. In contrast, querying Google Scholar with keyword "FSL" will overwhelm a user with tens

190    of thousands of hits. In many cases, "FSL" matches an author's initial while others are not

191    related to the software library in question. Since maintaining a correct reference of the RRID

192    increases visibility and thus influence of a research resource, and will bring direct impact

J Preprints

193　eventually, providers of research sources will be highly motivated to maintain its correctness,

194　closing a healthy positive feedback loop to sustain the whole system.

195　Early adoption of RRIDs already allows us to perform a preliminary study of digital research

196　resources including software, data sets, and Web services in neuroscience. We deployed a

197　software robot called "Scibot" to automatically annotate RRID mentions in research articles on

198　the Hypothes.is (2017) Web annotation platform. A team of curators of SciCrunch then

199　continued to manually curate each automatically annotated RRID mentions with Hypothes.is.  In

200　this way, we could rapidly collect RRID mentions. By August 1, 2016, we have collected 2493

201　curated mentions from 757 articles. Figure 1 shows a histogram of the number of unique RRID

202　mentions identified in an article. The histogram shows that most of articles contain more than

203　one unique RRID mentions, providing an opportunity to analyze correlation between research

204　resources based on their co-mention partners in publications.

205

206　　　　　　　　　　　　　　　　　　　(Figure 1)

207

208　Table 3 shows the top 30 highly mentioned RRIDs from this data set. Though the data set is

209　biased toward early adopters of RRIDs, the list shows an interesting mix of general-purpose

210　image analysis/statistical tools and neuroscience specialized resources that are widely used by

211　the neuroscience research community. We have also developed fully automatic text mining

212　methods (Ozyurt et al. 2016) to complement the curation approach to perform comprehensive

213　analysis of the impact and influence of research resources in life sciences.

214

215　　　　　　　　　　　　　　　　　　　(Table 3)

216

217　**Data Citation Implementation Pilot Project (2017).**  RRID's address the citation of digital

218　repositories and associated tools at a high level; however, we also need a system to cite

219　individual data sets that may include only a subset of data in a repository or be assembled from

220　multiple data sources. Precisely referring to which subset of data is retrieved and used can be a

221 computationally intractable problem, which leads to some pessimistic views regarding data
222 citation (Buneman et al. 2016).

223 We argue that the ultimate purpose of data citation is not only to identify precisely a data subset
224 for facilitating reproducibility, but also to ensure that both the individuals contributing data and
225 the repositories housing them receive proper credit and attribution, as specified in the Joint
226 Declaration of Data Citation Principles (JDDCP, Data citation 2014). The JDDCP has been
227 endorsed by 253 individual scientists and 114 organizations, representing different sectors of
228 stakeholders, including data centers/data repositories, educational institutions, funding
229 agencies/organizations, libraries, publishers, registries/social networks/research networks,
230 societies/associations/consortiums, and technology providers.

231 Based on the eight principles given in JDDCP, FORCE 11 and other groups have been working
232 on developing practical standards to implement data citations. One of these is the Data Citation
233 Implementation Pilot Project (DCIP) as part of the NIH BD2K bioCADDIE (2017) project that
234 we have been working on. The primary goal is to provide basic coordination between publishers,
235 repositories and identifier / metadata services for early adopters of data citation according to the
236 JDDCP. To meet this goal, we will provide authoritative guidance and group consultation on
237 data citation implementation to help establish one or more benchmark implementations of data
238 citation based on the JDDCP and Starr et al (2015), its cross-domain implementation guidance.

239 The key ideas here include working with data repositories on best practices that repositories can
240 follow to support data citation with the support of community metadata standards, the use of
241 persistent identifiers (e.g., DOI's), and machine-readable landing pages, which provide essential
242 information on the content and accessibility of data within the data repository (Cousijn et al.
243 2017; Fenner et al. 2016). A landing page allows for an access point that is independent from any
244 multiple encodings of the data that may be available (Starr et al. 2015), and thus avoids the
245 complicated computational problem of citing arbitrary subsets of data precisely, as described in
246 (Buneman 2016). A landing page can also provide information on access controls required by
247 licensing or privacy considerations. In addition, user requested landing pages can be minted for
248 custom data aggregations as well.

249 We are often asked how RRID's differ from the referencing of a specific data sets as proposed
250 by the JDDCP. The issue is one of granularity. RRID's are meant to identify the parent entity

251     like the PDB, while additional identifiers may be used to identify the specific data set used. This

252     more granular data citation may comprise a subset of a data repository or a supra-set across

253     repositories. The RRID essentially functions as an ORCID to identify the organizational entities

254     involved, e.g., the data repository, while the DOI points to a specific and unique data set.

255     DataCite (2017) and Dryad (2017) are closely related to RRID. They assign persistent identifiers

256     (e.g., DOIs) for research data, especially data sets that do not fit well into thematic data

257     repositories that contain data sets organized to serve the research needs on common topics such

258     as PDB. They also provide permanent storage space for these data sets. DataCite and Dryad

259     complement the efforts of RRID, which covers a broader range of research resources including

260     thematic data repositories. GRID (2017) releases the IDs and metadata for research organizations

261     and data providers to use under the Creative Commons Public Domain 1.0 International licence.

262     GRID maintains well curated hierarchies of research organizations (e.g., a lab within a

263     department in a university) and is useful for accurate identification of research resources with its

264     uniquely distinguishable IDs of provider organizations.

## 265    Towards Reliable and Accurate Metrics

266     Though counting frequencies of standardized RRID mentions and data citations might not be the

267     single perfect metric of the value of a digital repository, widespread adoption of these standards

268     will lead to a more reliable and comparable metric than the status quo and open development of

269     more sophisticated metrics like the h-index (Hirsch 2005) and pagerank (Page 1999) derived

270     from raw frequencies of literature citations. However, unlike publications, authors may cease to

271     credit and mention highly used resources that become routine, such as PubMed. This is when

272     access statistics may provide a better metric to assess their value. Also, accurate mention

273     identification measures influence at best. Assessing impact will requires not only provenance of

274     research outcomes to their various digital and data repositories contributing to their development

275     but also the impact of the research outcome in question, for which a general acceptable metric is

276     not currently available, and usually the impact may take years or decades to reveal.

277     A potential remedy for these issues is to request authors to explicitly distinguish why they chose

278     to mention a digital repository -- whether they used the data or service to obtain their results, or

279     they are merely related. Even without explicit citation mechanisms, it may be possible to make

280   the distinction to some extent from the context where the mentions appear (e.g. in the methods

281   section it may suggest that the data was used), and therefore distinguishing whether the data or

282   service lead to direct impact (a mention indicates influence of the resource in some way already).

283   Similarly, it would be possible to distinguish whether the mention carries positive or negative

284   sentiment of the resource. The key is that the standards bring unambiguous and persistent

285   references to digital repositories.

## Funding Statement

# References

291   Altmetrics. What are altmetrics?https://www.altmetric.com/about-altmetrics/what-are-altmetrics/
292   (Accessed 27 December 2016).

293   Bandrowski A., et al., *The Resource Identification Initiative: A cultural shift in publishing.*
294   Journal of Comparative Neurology, 2016. **524**(1): p. 8-22.

295   BioCaddie. http://biocaddie.org (Accessed 9 May 2017).

296   Buneman P, Davidson S, Frew J. Why data citation is a computational problem.
297   Communications of the ACM. 2016 Aug 24;59(9):50-7.

298   Cell. The STAR Methods. http://www.cell.com/star-methods (Accessed 9 May 2017).

299   Cousijn H, Kenall A, Ganley E, Harrison M, Kernohan D, Murphy F, Polischuk P, Martone M,
300   and Clark, T. A Data Citation Roadmap for Scientific Publishers. bioRxiv. January 19, 2017 **doi:**
301   https://doi.org/10.1101/100784

302   Data Citation Implementation Pilot Project. https://www.force11.org/group/dcip (Accessed 9
303   May 2017).

304   Data Citation Synthesis Group: Joint Declaration of Data Citation Principles. Martone M. (ed.)
305   San Diego CA: FORCE11; 2014 https://www.force11.org/group/joint-declaration-data-citation-
306   principles-final (Accessed 14 November 2016).

307   DataCite. http://www.datacite.org (Accessed 9 May 2017).

308   Dryad. http://datadryad.org (Accessed 9 May 2017).

309   dkNET, The NIDDK Information Network. https://dknet.org/ (Accessed 26 December 2016).

Fenner M, Crosas M, Grethe J, Kennedy D, Hermjakob H, Rocca-Serra P, Berjon R, Karcher S, Martone M and Clark T. A Data Citation Roadmap for Scholarly Data Repositories. bioRxiv. December 28, 2016 **doi:** https://doi.org/10.1101/097196

Gardner D, Akil H, Ascoli GA, Bowden DM, Bug W, Donohue DE, Goldberg DH, Grafstein B, Grethe JS, Gupta A, Halavi M. The neuroscience information framework: a data and knowledge environment for neuroscience. Neuroinformatics. 2008 Sep 1;6(3):149-60.

Google Scholar listing of articles published with RRID for the Protein Data Bank. https://scholar.google.com/scholar?q=RRID%3ASCR_012820+OR+RRID%3Anif-0000-00135&hl=en&as_sdt=0%2C5&oq=RRID%3ASCR_012820+. (Accessed 17 October 2016).

Gupta A, Bug W, Marenco L, Qian X, Condit C, Rangarajan A, Müller HM, Miller PL, Sanders B, Grethe JS, Astakhov V. Federated access to heterogeneous information resources in the Neuroscience Information Framework (NIF). Neuroinformatics. 2008 Sep 1;6(3):205-17.

GRID. https://grid.ac (Accessed 9 May 2017).

Hirsch JE. An index to quantify an individual's scientific research output. Proceedings of the National academy of Sciences of the United States of America. 2005 Nov 15:16569-72.

Huang YH, Rose PW, Hsu CN. Citing a data repository: A case study of the protein data bank. PloS One. 2015 Aug 28;10(8):e0136631. http://dx.doi.org/10.1371/journal.pone.0136631

Huang YH, A study of data citation. PhD Dissertation. Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan. 2015. http://www.airitilibrary.com/Publication/alDetailedMesh1?DocID=U0001-2601201621083800 (Accessed 26 December 2016)

Hypothes.is. https://hypothes.is (Accessed 9 May 2017)

International Neuroinformatics Coordinating Facility. http://incf.org/ (Accessed 9 May 2017)

List from NIF, List of resources no longer in service according to the Neuroscience Information Framework. https://neuinfo.org/Resources/search?q=%2A&l=&facet[]=Availability:THIS%20RESOURCE%20IS%20NO%20LONGER%20IN%20SERVICE&sort=asc&column=resource_name&sort=asc (Accessed 17 October 2016).

Marcus E and the whole Cell team. A STAR is born. Cell. 2016 August 25;166(5):1059-60.

Neuroscience Information Framework. http://neuinfo.org/ (Accessed 26 December 2016)

Ozyurt IB, Grethe JS, Martone ME, Bandrowski AE. Resource disambiguator for the web: extracting biomedical resources and their citations from the scientific literature. PLoS One. 2016 Jan 5;11(1):e0146300.

343  Page L, Brin S, Motwani R, Winograd T. The PageRank citation ranking: bringing order to the
344  web. Technical Report. Stanford InfoLab. 1999. http://ilpubs.stanford.edu:8090/422/  (Accessed
345  14 November 2016).

346  Resource Identification Portal. https://scicrunch.org/resources (Accessed 9 May 2017).

347  Rose PW and Hsu CN. bioCADDIE pilot project 3.2 Development of Citation and Data Access
348  Metrics   applied   to   RCSB   Protein   Data   Bank   and   related   Resources.   2015.
349  https://biocaddie.org/group/pilot-project/pilot-project-3-2-development-citation-and-data-access-
350  metrics-applied-rcsb (Accessed 14 November 2016).

351  Starr J, Castro E, Crosas M, Dumontier M, Downs RR, Duerr R, Haak LL, Haendel M, Herman
352  I, Hodson S, Hourclé J. Achieving human and machine accessibility of cited data in scholarly
353  publications. PeerJ Computer Science. 2015 May 27;1:e1.

354  Vasilevsky NA, Brush MH, Paddock H, Ponting L, Tripathy SJ, LaRocca GM, Haendel MA.
355  (2013) On the reproducibility of science: unique identification of research resources in the
356  biomedical literature. PeerJ 1:e148 https://doi.org/10.7717/peerj.148

357  Whetzel PL, Grethe JS, Banks DE, Martone ME. The NIDDK Information Network: A
358  Community Portal for Finding Data, Materials, and Tools for Researchers Studying Diabetes,
359  Digestive, and Kidney Diseases. PLoS ONE. 2015 Sep 22;10(9):e0136206. doi:
360  https://doi.org/10.1371/journal.pone.0136206. eCollection 2015. PMID: 26393351

361 # Tables.

362 Table 1. Different types of mentions of issued PDB IDs identified in PMC. The statistics of mentions may include
363 false positives due to errors by the text mining software for the last two types.

| Identifier | Example | Machine readable | Mentions(*) | % |
|---|---|---|---|---|
| PDB ID | PDB ID: **1STP** | yes | 14,888 | 4.8 |
| PDB DOI | http://dx.doi.org/10.2210/pdb**1stp**/pdb | yes | 155 | 0.05 |
| External link tag | <ext-link … ext-link-type=”pdb” xlink:href=”**1STP**”> | yes | 32,108 | 10 |
| PDB file name | **1stp**.pdb | yes | 895 | 0.03 |
| PDB URL | http://www.rcsb.org/… /structureId=**1stp** | yes, but URL may change | 657 | 0.2 |
| Non-standard PDB ID | PDB code: **1STP**, PDB reference **1STP**, PDB accession number **1STP**, Many variations... | yes/no | 22,081 | 7.1 |
| PDB in context | “We employed the following **PDB** coordinates: glycogen phosphorylase, **1gpy** …” | yes/no with text mining | 16,726 | 5.4 |
| Free text | “We first placed S2 bound to human PI3KC; (**3ene**) into the reference coordinates …” | yes/no with text mining | 221,287 | 72 |

364

365 Table 2: Top 10 highly cited protein structures (top) and top 10 highly mentioned protein structures in PDB. "Year"
366 shows when the PDB ID was issued.

| Citation Rank | PDB ID | Year | # of Citations | # of Mentions | Mention Rank |
|---|---|---|---|---|---|
| 1 | 1AOI | 1997 | 1527 | 31 | 37 |
| 2 | 1BL8 | 1998 | 1234 | 35 | 24 |
| 3 | 1F88 | 2000 | 957 | 44 | 16 |
| 4 | 1GC1 | 1998 | 852 | 26 | 57 |
| 5 | 1RV1 | 2004 | 747 | 11 | 488 |
| 6 | 1FFK | 2000 | 746 | 31 | 34 |
| 7 | **2RH1** | 2007 | 682 | 124 | 1 |
| 8 | 1YSG | 2005 | 650 | 6 | 1984 |
| 9 | **2A79** | 2005 | 635 | 49 | 10 |
| 10 | 1AIK | 1997 | 561 | 12 | 403 |

| Mention Rank | PDB ID | Year | # of Mentions | # of Citations | Citation Rank |
|---|---|---|---|---|---|
| 1 | **2RH1** | 2007 | 124 | 682 | 7 |
| 2 | 1UBQ | 1987 | 96 | 222 | 142 |
| 3 | 1KX5 | 2002 | 69 | 272 | 87 |
| 4 | 2R9R | 2007 | 65 | 433 | 20 |
| 5 | 3EML | 2008 | 65 | 408 | 24 |
| 6 | 1U19 | 2004 | 64 | 227 | 134 |
| 7 | 1K4C | 2001 | 59 | 454 | 18 |
| 8 | 2VT4 | 2008 | 55 | 356 | 38 |
| 9 | 2B4C | 2005 | 55 | 289 | 71 |
| 10 | **2A79** | 2005 | 49 | 635 | 9 |

367

368      Table 3. Top 30 highly mentioned resources in neuroscience RRID early adopters.

| RRID | count | resource name |
|---|---|---|
| SCR_003070 | 260 | imageJ |
| SCR_002798 | 156 | Graphpad Prism |
| SCR_014199 | 138 | Adobe Photoshop CS5 CS6 |
| SCR_001622 | 106 | MATLAB |
| SCR_002865 | 64 | SPSS |
| SCR_001775 | 57 | Neurolucida |
| SCR_001905 | 50 | R software |
| SCR_011323 | 46 | pClamp |
| SCR_002677 | 35 | AxioVision |
| SCR_007037 | 34 | SPM |
| SCR_013672 | 32 | Zeiss Zen software |
| SCR_002078 | 31 | Adobe Photoshop CS2 |
| SCR_002285 | 29 | Fiji |
| SCR_014198 | 28 | Adobe illustrator |
| SCR_013726 | 23 | G*Power software |
| SCR_003238 | 22 | Open Science Framework |
| SCR_002823 | 21 | FSL |
| SCR_002526 | 20 | Stereo Investigator |
| SCR_008520 | 20 | FlowJo |
| SCR_001847 | 19 | FreeSurfer |
| SCR_000903 | 18 | Spike 2 software |

| SCR_007369 | 17 | Image-Pro Plus |
|---|---|---|
| SCR_002368 | 17 | MetaMorph Microscopy Automation and Image Analysis Software |
| SCR_000325 | 16 | IGOR Pro |
| SCR_003210 | 16 | Sigma Plot |
| SCR_007353 | 15 | Advanced 3D Visualization and Volume Modeling |
| SCR_013673 | 15 | Leica AS AF software |
| SCR_001818 | 12 | NeuroExplorer |
| SCR_008567 | 11 | Statistical Analysis System |
| SCR_002716 | 11 | Synapse Web Reconstruct |

369

# Figures.

371   Figure 1. 573 out of 757 articles contain more than one unique RRID mentions. One of them contains as many as 24
372   unique RRID mentions.



373