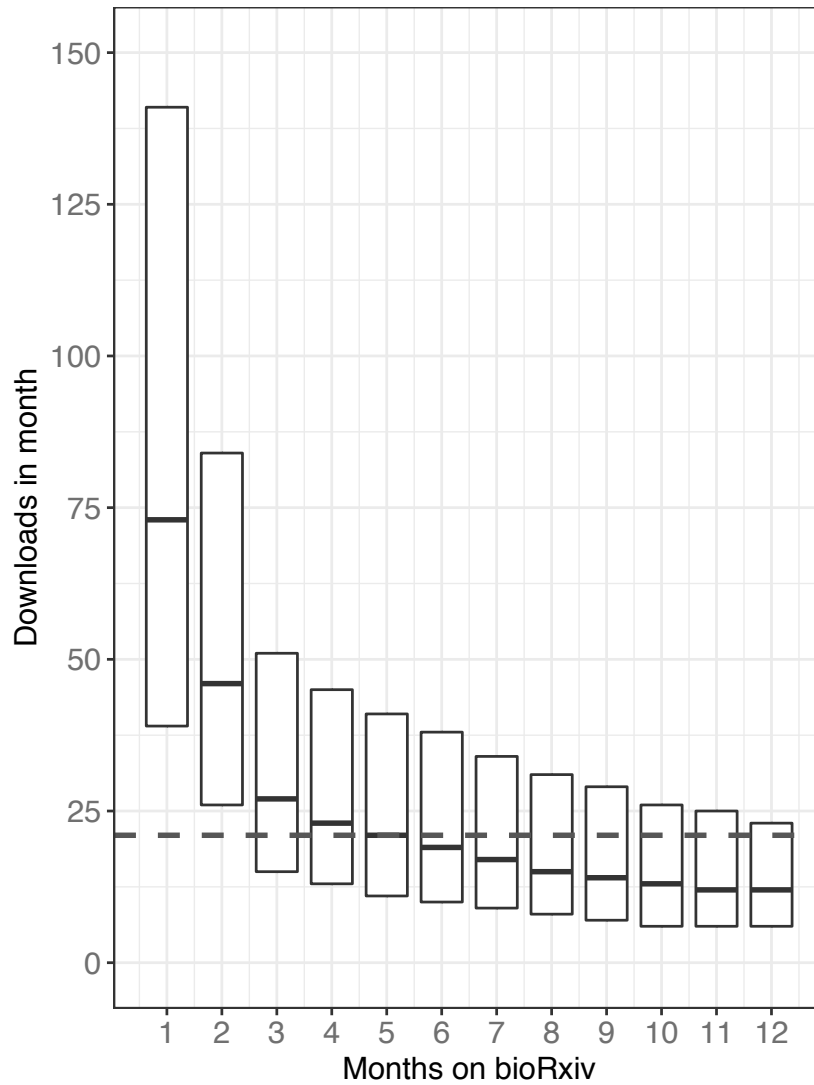


1 Supplements

2 Supplementary figures



3

4

Figure S1. The distribution of downloads that preprints accrue in their first months on bioRxiv. For example, the box at "1" on the x axis indicates the downloads that all preprints have received in their first month online.

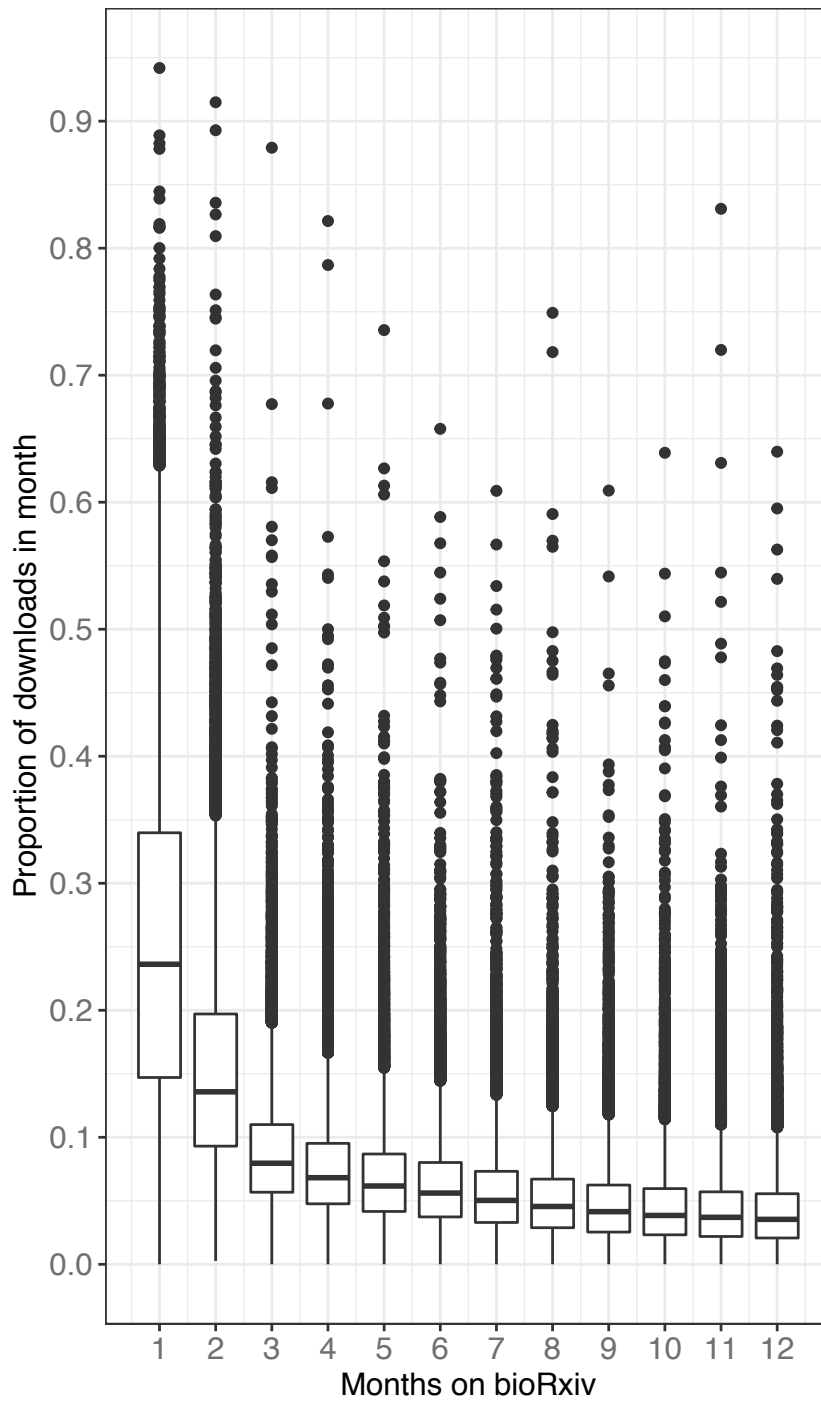
5

6

7

Supplementary file: *downloads_by_months.csv*

8

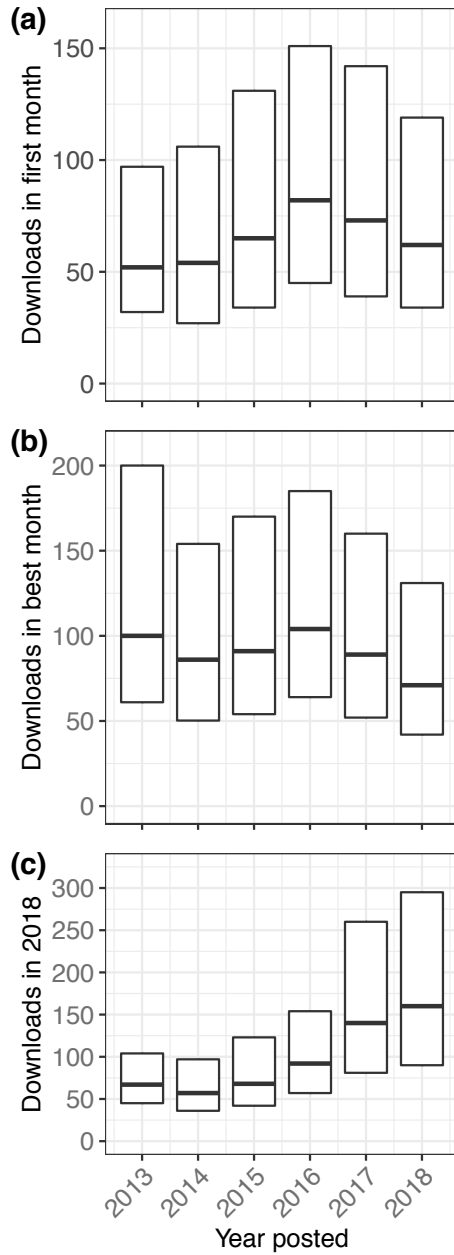


9

10

11

Figure S2. The proportion of downloads that preprints accrue in their first months on bioRxiv.



12

13

14

15

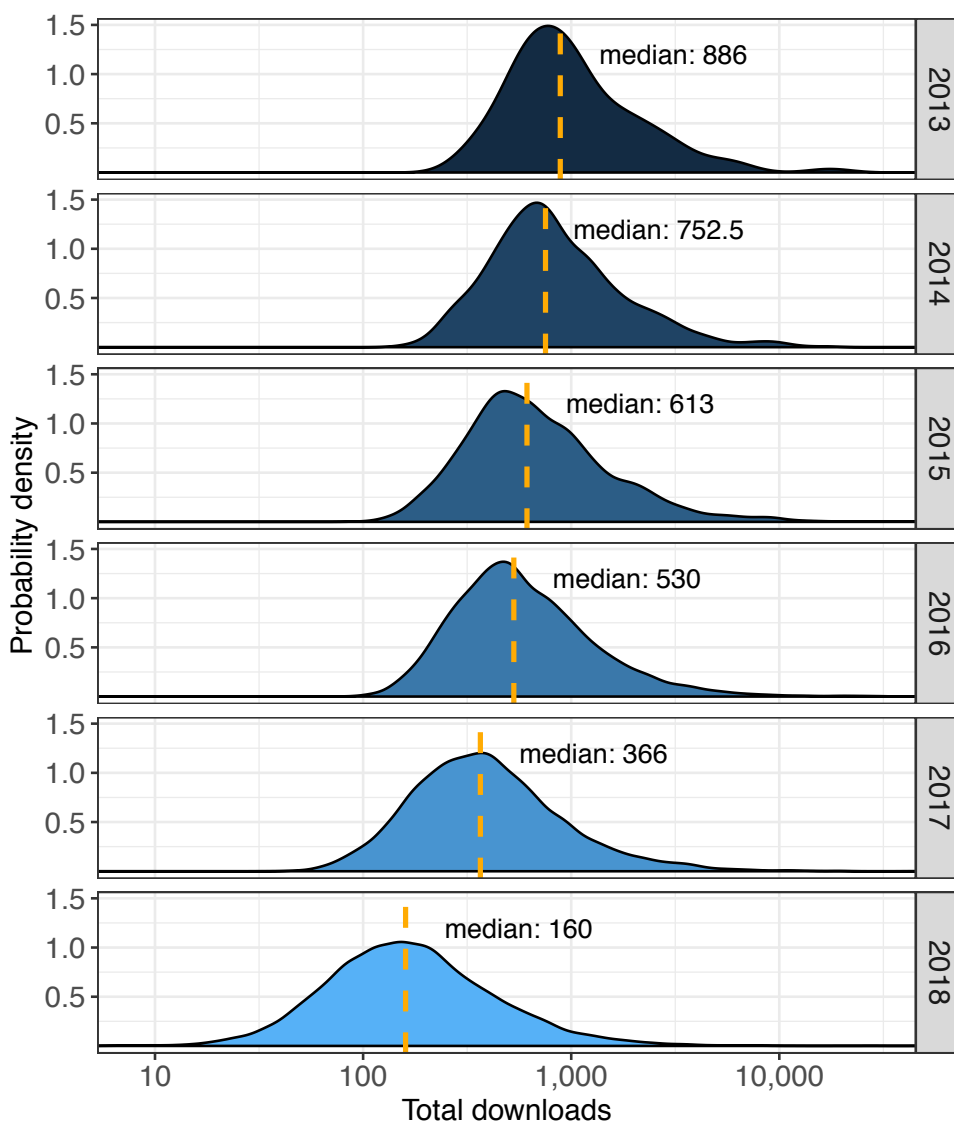
16

17

Figure S3. Multiple perspectives on per-preprint download statistics. Preprints in all three plots are categorized by the year in which they were first posted. **(a)** The median number of downloads for each preprint's first calendar month on bioRxiv. (There was no correction done for preprints posted toward the end of their first month.)**(b)** The median number of

18 downloads in each preprint's *best* month, in any year. **(c)** The median
19 downloads per preprint in 2018, for preprints posted in any year.

20 **Supplementary files:** *downloads_by_first_month.csv*,
21 *downloads_max_by_year_posted.csv*, *2018_downloads_per_year.csv*



22
23 **Figure S4.** Total downloads per preprint, segmented by the year in which
24 each preprint was posted. The Y axis indicates a probability density rather
25 than overall count.

26 **Supplementary file:** *downloads_per_year.csv*

Source	Articles
<i>Cell</i> vol. 174(6)	17
<i>Cell</i> vol. 175(1)	18
<i>Genetics</i> vol. 210(1)	23
<i>Jour of Biochem</i> vol. 164(3)	8
<i>PLoS Biology</i> vol. 16(9)	19
bioRxiv, 1–3 Sep 2018	100

27 **Table S1.** Full-length articles published by a selection of journals in September
 28 2018. The *Cell* count is limited to the "Articles" and "Resources" categories;
 29 *Genetics* is limited to their "Investigations" category, and *PLoS Biology* to
 30 "Research Articles," "Methods and Resources" and "Meta-Research Articles."
 31 Links to each issue's table of contents are included in supplementary file
 32 *figures.md*.

Author name	bioRxiv preprints	Primary field	Email addresses
George Davey Smith	97	Epidemiology	4
Ian J. Deary	61	Genetics	4
Andrew M. McIntosh	57	Genetics	1
Mark J. Daly	47	Genetics	3
Richard M. Murray	45	Synthetic biology	2

George M. Church	43	Synthetic biology	4
Wei Wang	39	Bioinformatics	23
Benjamin M. Neale	39	Genetics	2
Alkes L. Price	36	Genetics	1
Jian Yang	36	Genetics	6
Po-Ru Loh	35	Genetics	2
Caroline Hayward	35	Genetics	1
Aarno Palotie	34	Genetics	4
Jay Shendure	34	Genomics	3
Ole A. Andreassen	34	Genetics	2

33 **Table S2.** The top 15 authors with the most preprints on bioRxiv. Names are
34 listed as they appear on biorxiv.org, after making corrections outlined in the
35 Methods section. An author's "Primary field" is the bioRxiv collection to which
36 they have submitted the most preprints. Preprint count is does not account for
37 duplicates: For example, Ian J. Deary and Andrew M. McIntosh are both high on
38 the list, but their counts include multiple preprints that they co-authored together.
39 The "Email addresses" field lists the number of email addresses observed in that
40 author's preprints that is attributed to them, and is used to approximate the risk
41 that the author is actually a conglomeration of multiple researchers with the same
42 name.

43 **Supplementary file:** *papers_per_author.csv*

Institution	Authors	Preprints
Stanford University	1,473	1,045
University of Oxford	1,192	902
University of Cambridge	1,109	842
University of Washington	924	609
University College London	801	644
University of Pennsylvania	764	544
University of Michigan	763	484
University of California, San Francisco	750	511
University of California, San Diego	725	495
Imperial College London	703	472
University of Edinburgh	646	487
University of California, Berkeley	620	528
Yale University	555	392
Duke University	554	323
Harvard University	532	557
Harvard Medical School	529	453

Columbia University	520	422
Cornell University	486	365
University of Toronto	462	334
Johns Hopkins University	461	407
University of California, Davis	461	291
Icahn School of Medicine at Mount Sinai	448	281
University of Chicago	444	353
University of British Columbia	431	281
University of Minnesota	430	310

44 **Table S3.** Top 25 institutions with the most authors listing them as their
45 affiliation, and how many papers have been published by those authors.
46 Each institution's count of total preprints is based on the number of papers
47 posted by authors currently listed with those affiliations, but preprints
48 attributed to authors from multiple institutions count toward the total for all
49 institutions mentioned. A paper with multiple authors from the same
50 institution is counted only once for that institution.

51 **Supplementary file:** *authors_per_institution.csv*

52 Supplementary information

53 Database schema

54 Rxivist stores bioRxiv data in 16 tables in a PostgreSQL database. It is possible
55 these fields may change in the future to add more data and solidify the schema (more
56 specific constraints, established foreign key relationships, improved primary keys, etc.);
57 changes will be published on rxivist.org and as addendums in the data dumps. Postgres
58 also uses organizational tools called "schemas," namespace-like entities that hold
59 tables. There are two schemas currently available: "prod" contains the most up-to-date
60 information pulled by the Rxivist web crawler. The "paper" schema stores a static copy
61 of the data used for this manuscript, and contains data on all papers posted through the
62 end of November 2018. Currently, the table definitions for both the "prod" and "paper"
63 schemas are identical, except for one table: The "publication_dates" table contains
64 information that was pulled for this manuscript but is not continuously updated, so is
65 omitted from the "prod" schema.

66 There are several tables that contain data that do not contain the most current
67 data: Any tables in the database with the "_working" suffix have identical fields to the
68 version of the table without that suffix, but contain the previous version of that table.
69 Temporary tables are used to ensure new rankings are loaded properly before
70 "activating" them by renaming the temporary table to the one used by the API, and if an
71 error occurs, the previous rankings can be restored. All tables are listed below in
72 alphabetical order except for "articles" and "authors," which contain the entities central
73 to the data model. See supplementary file *figures.md* for examples of queries that may
74 be helpful.

75

76 **Table: articles**

77 Each entry represents an individual preprint.

- 78 ● id (serial, primary key) – An arbitrary integer assigned to each preprint when it is
79 first recorded by Rxivist. These values start at 386.
- 80 ● url (text, unique) – The web address of the preprint's main page on biorxiv.org.
81 Revisions of a paper are given new URLs; this field stores the most recently
82 observed version.
- 83 ● title (text) – The title of the preprint. This field is updated for revisions.
- 84 ● abstract (text) – The abstract for the preprint, as specified on bioRxiv. This field is
85 updated for revisions, but abstracts are recorded in a separate step from the one
86 in which preprints and revisions are recorded, so this field may occasionally be
87 null.
- 88 ● doi (text, unique) – The Digital Object Identifier assigned to the preprint. BioRxiv
89 does *not* issue new DOI numbers for revisions.
- 90 ● collection (text) – The bioRxiv "collection" in which this preprint has been posted.
91 Because a preprint's collection can only be determined by its appearance in
92 bioRxiv's chronological listing of papers for each category, there is sometimes a
93 delay in which this field will be null.
- 94 ● title_vector (tsvector) – A sorted list of distinct lexemes found in the preprint's
95 title, as determined by Postgres. Used to facilitate text search. Updated for
96 revisions.

- 97 ● `abstract_vector` (tsvector) – A sorted list of distinct lexemes found in the preprint's
98 abstract, as determined by Postgres. Used to facilitate text search. Updated for
99 revisions.
- 100 ● `last_crawled` (date) – A date indicating the last time the web crawler checked
101 bioRxiv for updated traffic metrics for the preprint.
- 102 ● `posted` (date) – The date on which the preprint first appeared on bioRxiv. This
103 field is *not* updated for revisions.
- 104 ● `author_vector` (tsvector) – A sorted list of distinct lexemes found in the preprint's
105 list of authors, as determined by Postgres. Used to facilitate text search. Updated
106 for revisions.

107 **Table: authors**

108 Each entry represents an individual preprint author.

- 109 ● `id` (serial, primary key) – An arbitrary integer assigned to each author when they
110 are first recorded by Rxivist. These values start at 200,058.
- 111 ● `name` (text) – The full name of the author as it was first recorded by Rxivist.
112 Because the back catalog of preprints was not recorded in chronological order
113 and the author name string comparison process changed over time, this name is
114 not necessarily the one attached to an author's earliest preprint.
- 115 ● `institution` (text) – The institutional affiliation associated with the author on their
116 most recent preprint. This is updated for revisions.
- 117 ● `orcid` (text, unique) – The ORCID unique identifier associated with the author.
118 This field is populated only if the author has listed the ID on at least one of their
119 preprints.

- 120 • noperiodname (text) – The author's full name stripped of all full stops. This field is
121 used to more quickly search for authors without accounting for the punctuation
122 that is applied most inconsistently.

123 **Table: alltime_ranks**

124 Each entry stores download ranking information for a single preprint. This table is
125 emptied and replaced every time rankings are calculated.

- 126 • article (integer, primary key) – The ID of a preprint.
127 • rank (integer) – The ordinal position of the preprint in the list of all indexed
128 preprints, organized in descending order according to the "downloads" field.
129 (These rankings do *not* account for ties, so two preprints with the same download
130 count will receive sequential ranks.)
131 • downloads (integer) – The total overall downloads for the preprint, as of the last
132 ranking calculation. Rankings are not necessarily re-calculated every time traffic
133 data is updated, so this number may be lower than the preprint's total downloads
134 as recorded in the "article_traffic" table, from which this number is calculated.

135 **Table: article_authors**

136 Associative table recording the many-to-many relationship between preprints and
137 authors.

- 138 • id (serial, primary key) – An arbitrary integer assigned to each association.
139 • article (integer) – The ID of the preprint being associated.
140 • author (integer) – The ID of the author being associated.

141 **Table: article_publications**

142 Each entry contains the publication information about a single preprint; only preprints
143 that have been published elsewhere have entries.

- 144 • article (integer, primary key) – The ID of a preprint.
- 145 • doi (text) – The new Digital Object Identifier associated with the published
146 version of the preprint.
- 147 • publication (text) – The name of the journal in which the preprint was published,
148 as reported by bioRxiv.

149 **Table: article_traffic**

150 Each entry contains bioRxiv traffic data for a single preprint in a single month. Metrics
151 are as reported by bioRxiv, in each preprint's "Metrics" page, in the "Article Usage"
152 section.

- 153 • id (serial, primary key) – An arbitrary integer assigned to each entry.
- 154 • article (integer) – The ID of a preprint.
- 155 • month (integer) – The (1-indexed) month in which the metrics were tallied.
- 156 • year (integer) – The four-digit year in which the metrics were tallied.
- 157 • abstract (integer) – The number of times the preprint's abstract was viewed on
158 bioRxiv in the specified month.
- 159 • pdf (integer) – The number of times the preprint's full-text PDF was downloaded
160 in the specified month.

161 **Table: author_emails**

162 Each entry associates an author with an email address. Note there is no requirement
163 that the email be unique, so multiple authors may be associated with the same email

164 address, albeit in a denormalized format. For now, logic in the web crawler avoids
165 associating an author with the same email address multiple times.

- 166 • id (serial, primary key) – An arbitrary integer assigned to each entry.
- 167 • author (integer) – The ID of an author.
- 168 • email (text) – An email address that was associated with that author on the
169 bioRxiv page of one of their preprints.

170 **Table: author_ranks**

171 Each entry stores download ranking information for a single author. This table is
172 emptied and replaced every time rankings are calculated.

- 173 • author (integer, primary key) – The ID of an author.
- 174 • rank (integer) – The ordinal position of the author in the list of all authors,
175 organized in descending order according to the combined all-time downloads of
176 their preprints. (Unlike the alltime_ranks table, these rankings *do* account for ties,
177 so two authors with the same download count will receive identical values in this
178 field.)
- 179 • tie (boolean) – A flag indicating whether the author is tied with one or more other
180 authors at the same rank.
- 181 • downloads (integer) – The total overall downloads for all preprints associated
182 with the author, as of the last ranking calculation. Rankings are not necessarily
183 re-calculated every time traffic data is updated, so this number may be lower
184 than the total downloads as recorded in the "article_traffic" table, from which this
185 number is calculated.

186 **Table: author_ranks_category**

187 Each entry stores download ranking information for a single author in a single bioRxiv
188 category. Only authors with more than zero downloads in a category will have an entry
189 for that category. This table is emptied and replaced every time rankings are calculated.

- 190 ● id (serial, primary key) – An arbitrary integer assigned to each entry.
- 191 ● author (integer) – The ID of an author.
- 192 ● category (text) – Which of the 27 bioRxiv categories was used to limit the
193 download count.
- 194 ● rank (integer) – The ordinal position of the author in the list of all authors,
195 organized in descending order according to the combined all-time downloads of
196 their preprints in the specified category. (Unlike the alltime_ranks table, these
197 rankings *do* account for ties, so two authors with the same download count will
198 receive identical values in this field.)
- 199 ● tie (boolean) – A flag indicating whether the author is tied with one or more other
200 authors at the same rank.
- 201 ● downloads (integer) – The total overall downloads for all preprints in the specified
202 category that are associated with the author, as of the last ranking calculation.
203 Rankings are not necessarily re-calculated every time traffic data is updated, so
204 this number may be lower than the total downloads as recorded in the
205 "article_traffic" table, from which this number is calculated.

206 **Table: author_translations**

207 A deprecated table used to specify redirects to search engines that indexed outdated ID
208 numbers for authors that were later modified. Of no practical use elsewhere, but

209 included here because the table is included in the "paper" schema and the database
210 snapshot accompanying this manuscript.

- 211 • old (integer, primary key) – An author ID that may have been indexed by search
212 engines but has since changed.
- 213 • new (integer) – The new author ID associated with the same individual.

214 **Table: category_ranks**

215 Each entry stores download ranking information for a single preprint in the category to
216 which it was posted. (Since each preprint is posted only to a single category, the
217 category itself is not included here.) This table is emptied and replaced every time
218 rankings are calculated.

- 219 • article (integer, primary key) – The ID of a preprint.
- 220 • rank (integer) – The ordinal position of the preprint in the list of all preprints in the
221 same category, organized in descending order according to all-time download
222 count. (These rankings do *not* account for ties, so two preprints with the same
223 download count will receive sequential ranks.)

224 **Table: crossref_daily**

225 Each entry records Crossref social media information for a single preprint on a single
226 day. Note that there are many entries in this table that do *not* reference a preprint; to
227 avoid throwing away data for preprints that have not yet been indexed by Rxivist, the
228 web crawler saves tweet counts for all results with a DOI matching the bioRxiv DOI
229 prefix (10.1101). However, that prefix is shared by other organizations, so some entries
230 are irrelevant external references.

- 231 • id (serial, primary key) – An arbitrary integer assigned to each entry.

- 232 • source_date (date) – The date for which the data was collected. (For example, a
233 date of "5 Dec 2017" would indicate events reported by Crossref on that date,
234 regardless of when Rxivist itself actually recorded it.)
- 235 • doi (text) – The DOI of an entity that may be a bioRxiv preprint.
- 236 • count (integer) – The total number of Twitter posts observed on the specified
237 date that referenced the specified DOI.
- 238 • crawled (date) – The date on which this information was retrieved from Crossref.

239 **Table: download_distribution**

240 Each entry indicates a bin in a histogram used to measure the distribution of all-time
241 downloads for either authors or preprints. Preprints are recorded in the category
242 "alltime"; authors are in the category "authors". Though it doesn't conform well to the
243 schema, this table is also used to record the mean and median download counts for
244 preprints for authors, preprints, and each of the 27 bioRxiv categories. The "category"
245 field for these entries is the category name plus either "_median" or "_mean" at the end.
246 The "bin" field for these entries is "0" and the "count" field is the value. These values are
247 used for visualizations on the Rxivist website.

- 248 • id (serial, primary key) – An arbitrary integer assigned to each entry.
- 249 • bucket (integer) – The maximum number of downloads for this bin.
- 250 • count (integer) – The number of entities with a total download count that falls
251 within the limits of this bin.
- 252 • category (text) – Which histogram this bin belongs to.

253 **Table: month_ranks**

254 Each entry stores download ranking information for a single preprint based on download
255 data going back to the beginning of the previous month. This table is emptied and
256 replaced every time rankings are calculated.

- 257 • article (integer, primary key) – The ID of a preprint.
- 258 • rank (integer) – The ordinal position of the preprint in the list of all indexed
259 preprints, organized in descending order according to the "downloads" field.
260 (These rankings do *not* account for ties, so two preprints with the same download
261 count will receive sequential ranks.)
- 262 • downloads (integer) – The total downloads for the preprint since the beginning of
263 the previous month, as of the last ranking calculation. Rankings are not
264 necessarily re-calculated every time traffic data is updated, so this number may
265 be lower than the total as recorded in the "article_traffic" table, from which this
266 number is calculated. However, ranks for all preprints are calculated at the same
267 time, so the timeframe covered by this field will be the same across the whole
268 table.

269 **Table: publication_dates**

270 Each entry stores publication information for a single preprint. Most (but *not* all)
271 preprints that are recorded in the "article_publications" table have a corresponding entry
272 in this table. This is the only table present in the "paper" schema but not the "prod"
273 schema and is not maintained.

- 274 • article (integer, primary key) – The ID of a preprint.
- 275 • date (date) – The date on which the preprint was published in a journal, based on
276 data from Crossref.

277 **Table: ytd_ranks**

278 Each entry stores download ranking information for a single preprint based on download
279 data going back to the beginning of the current year. This table is emptied and replaced
280 every time rankings are calculated.

- 281 ● article (integer, primary key) – The ID of a preprint.
- 282 ● rank (integer) – The ordinal position of the preprint in the list of all indexed
283 preprints, organized in descending order according to the "downloads" field.
284 (These rankings do *not* account for ties, so two preprints with the same download
285 count will receive sequential ranks.)
- 286 ● downloads (integer) – The total downloads for the preprint since the beginning of
287 the current year, as of the last ranking calculation. Rankings are not necessarily
288 re-calculated every time traffic data is updated, so this number may be lower
289 than the total as recorded in the "article_traffic" table, from which this number is
290 calculated. However, ranks for all preprints are calculated at the same time, so
291 the timeframe covered by this field will be the same across the whole table.