

1 **Quality Matters: Biocuration Experts on the Impact of Duplication and Other**  
2 **Data Quality Issues in Biological Databases**

3  
4 Qingyu Chen<sup>1,\*a</sup>, Ramona Britto<sup>2,b</sup>, Ivan Erill<sup>3,c</sup>, Constance J. Jeffery<sup>4,d</sup>, Arthur Liberzon<sup>5</sup>,  
5 Michele Magrane<sup>2,e</sup>, Jun-ichi Onami<sup>6,7,f</sup>, Marc Robinson-Rechavi<sup>8,9,g</sup>, Jana Sponarova<sup>10,h</sup>, Justin  
6 Zobel<sup>1,\*i</sup>, Karin Verspoor<sup>1,\*j</sup>

7  
8 <sup>1</sup> *School of Computing and Information Systems, The University of Melbourne, Melbourne, VIC*  
9 *3010, Australia. Current affiliation: National Center for Biotechnology Information, National*  
10 *Library of Medicine, National Institutes of Health, MD 20892, USA*

11 <sup>2</sup> *European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI),*  
12 *Wellcome Trust Genome Campus, Cambridge CB10 1SD, UK*

13 <sup>3</sup> *Department of Biological Sciences, University of Maryland Baltimore County, MD 21250, USA*

14 <sup>4</sup> *Department of Biological Sciences, The University of Illinois at Chicago, IL 60607, USA*

15 <sup>5</sup> *The Broad Institute of MIT and Harvard University, MA 02142, USA*

16 <sup>6</sup> *Japan Science and Technology agency, National Bioscience Database Center, Tokyo 102-8666,*  
17 *Japan*

18 <sup>7</sup> *National Institute of Health Sciences, Tokyo 158-8501, Japan*

19 <sup>8</sup> *SIB Swiss Institute of Bioinformatics, CH-1015 Lausanne, Switzerland*

20 <sup>9</sup> *Department of Ecology and Evolution, University of Lausanne, CH-1015 Lausanne,*  
21 *Switzerland*

22 <sup>10</sup> *Nebion AG, Hohlstrasse 515, 8048 Zurich, Switzerland*

23 <sup>11</sup> *National Center for Biotechnology Information, National Library of Medicine, National*  
24 *Institutes of Health, Bethesda, MD 20894, USA*

25

26 \*Corresponding authors.

27 E-mails: [qingyu.chen@nih.gov](mailto:qingyu.chen@nih.gov) (Chen Q), [jzobel@unimelb.edu.au](mailto:jzobel@unimelb.edu.au) (Zobel J),

28 [karin.verspoor@unimelb.edu.au](mailto:karin.verspoor@unimelb.edu.au) (Verspoor K).

29

30 **Running title: *Chen et al / Expert Opinions on Duplication and Other Quality Issues***

31

32 <sup>a</sup> ORCID: 0000-0002-6036-1516.

33 <sup>b</sup> ORCID: 0000-0003-1011-5410.

34 <sup>c</sup> ORCID: 0000-0002-7280-7191.

35 <sup>d</sup> ORCID: 0000-0002-2147-3638.

36 <sup>e</sup> ORCID: 0000-0003-3544-996X.

37 <sup>f</sup> ORCID: 0000-0003-0790-8313.

38 <sup>g</sup> ORCID: 0000-0002-3437-3329.

39 <sup>h</sup> ORCID: 0000-0002-6345-6879.

40 <sup>i</sup> ORCID: 0000-0001-6622-032X.

41 <sup>j</sup> ORCID: 0000-0002-8661-1544.

42

43 **Total counts of words: 6186**

44 **Figures: 4**

45 **Tables: 1**

46 **Supplementary materials: survey questions, results on general data quality issues, and**  
47 **interview details are provided.**

48

## 49 **Abstract**

50 The volume of biological database records is growing rapidly, populated by complex records  
51 drawn from heterogeneous sources. A specific challenge is duplication, that is, the presence of  
52 redundancy (records with high similarity) or inconsistency (dissimilar records that correspond to  
53 the same entity). The characteristics (which records are duplicates), impact (why duplicates are  
54 significant), and solutions (how to address duplication), are not well understood. Studies on the  
55 topic are neither recent nor comprehensive. In addition, other data quality issues, such as  
56 inconsistencies and inaccuracies, are also of concern in the context of biological databases. A  
57 primary focus of this paper is to present and consolidate the opinions of over 20 experts and  
58 practitioners on the topic of duplication in biological sequence databases. The results reveal that  
59 survey participants believe that duplicate records are diverse; that the negative impacts of  
60 duplicates are severe, while positive impacts depend on correct identification of duplicates; and  
61 that duplicate detection methods need to be more precise, scalable, and robust. A secondary  
62 focus is to consider other quality issues. We observe that biocuration is the key mechanism used

63 to ensure the quality of this data, and explore the issues through a case study of curation in  
64 UniProtKB/Swiss-Prot as well as an interview with an experienced biocurator. While biocuration  
65 is a vital solution for handling of data quality issues, a broader community effort is needed to  
66 provide adequate support for thorough biocuration in the face of widespread quality concerns.

67

68 **KEYWORDS:** Duplication; Redundancy; Data quality; Biocuration; Biological databases

## 69 **Introduction**

70 The major biological databases represent an extraordinary collective volume of work. Diligently  
71 built up over decades and comprised of many millions of contributions from the biomedical  
72 research community, biological databases provide worldwide access to a massive number of  
73 records (also known as *entries*) [1]. Starting from individual laboratories, genomes are  
74 sequenced, assembled, annotated, and ultimately submitted to primary nucleotide databases such  
75 as GenBank [2], ENA [3], and DDBJ [4] (collectively known as INSDC). Translations of those  
76 nucleotide records, protein records, are deposited into central protein databases such as the  
77 UniProt KnowledgeBase (UniProtKB) [5] and the Protein Data Bank [6]. Sequence records are  
78 further accumulated into different databases for more specialised purposes: RFam [7] and PFam  
79 [8] for RNA and protein families respectively, such as DictyBase [9] and PomBase [10] for  
80 model organisms, ArrayExpress [11] and GEO [12] for gene expression profiles. These  
81 databases are selected as examples; the list is not intended to be exhaustive. However, they are  
82 representative of biological databases that have been named in the “golden set” of the 24th  
83 Nucleic Acids Research database issue. The introduction of that issue highlights the databases  
84 that “consistently served as authoritative, comprehensive, and convenient data resources widely  
85 used by the entire community and offer some lessons on what makes a successful database” [13].  
86 The associated information about sequences is also propagated into non-sequence databases,  
87 such as PubMed (<https://www.ncbi.nlm.nih.gov/pubmed/>) for the scientific literature, or GO [14]  
88 for function annotations. Those databases in turn benefit individual studies, many of which use  
89 these public available records as the basis for their own research.

90 Inevitably, given the scale of these databases, some submitted records are redundant [15],  
91 inconsistent [16], inaccurate [17], incomplete [18], or outdated [19]. Such quality issues can be  
92 addressed by manual curation, with the support of automatic tools, and by processes such as  
93 reporting of the issues by contributors detecting mistakes. Biocuration plays a vital role in  
94 biological database curation [20]. It de-duplicates database records [21], resolves inconsistencies  
95 [22], fixes errors [17], and resolves incomplete and outdated annotations [23]. Such curated  
96 records are typically of high quality and represent the latest scientific and medical knowledge.  
97 However, the volume of data prohibits exhaustive curation, and some records with those quality  
98 issues remain undetected.

99 In other work, we (Chen, Verspoor, and Zobel) have explored a particular form of quality  
100 issue, which we have characterized as *duplication* [24,25]. As described in that work, duplicates  
101 are characterized in different ways in different contexts, but they can be broadly categorized as  
102 *redundancies* or *inconsistencies*. The perception of a pair of records as duplicates depends on the  
103 task. As we wrote in previous work, “*a pragmatic definition for duplication is that a pair of*  
104 *records A and B are duplicates if the presence of A means that B is not required, that is, B is*  
105 *redundant in the context of a specific task or is superseded by A.*” [24]. Many such duplicates  
106 have been found through curation, but the prevalence of undetected duplicates is unknown, as is  
107 the accuracy and sensitivity of automated tools for duplicate or redundancy detection. Other  
108 work has explored the detection of duplicates, but often under assumptions that limit the impact.  
109 For example, some researchers have assumed that similarity of genetic sequence is the sole  
110 indicator of redundancy, whereas in practice some highly similar sequences may represent  
111 distinct information and some rather different sequences may in fact represent duplicates [26].  
112 We detail the notion and impacts of duplication in the next section.

113

#### 114 **Authors’ contributions**

115 In this work, a main focus is to explore the characteristics, impacts, and solutions to duplication  
116 in biological databases; a secondary focus is to further investigate other quality issues. We  
117 present and consolidate the opinions of over 20 experts and practitioners on the topic of  
118 duplication and other data quality issues, via a questionnaire-based survey. To address different  
119 quality issues, we introduce biocuration as a key mechanism for ensuring the quality of  
120 biological databases. To our knowledge, there is no one-size-fits-all solution even to a single  
121 quality issue [27]. We thus explain the complete UniProtKB/Swiss-Prot curation process, via a  
122 descriptive report and an interview with its curation team leader, which provides a reference  
123 solution to different quality issues. Overall, the observations on duplication and other data  
124 quality issues highlight the significance of biocuration in data resources, but a broader  
125 community effort is needed to provide adequate support to facilitate thorough biocuration.

126

## 127 **The notion and impact of duplication**

128 Our focus is on database records – that is, entries in structured databases – not on biological  
129 processes such as gene duplication. Superficially, the question of what constitutes an *exact*  
130 *duplicate* in this context can seem obvious: two records that are exactly identical in both data  
131 (*e.g.*, sequence) and annotation (*e.g.*, meta-data including species and strain of origin) are  
132 duplicates. However, the notion of duplication varies. We demonstrate a generic biological data  
133 analysis pipeline involving biological databases and illustrate different notions of duplication.

134 **Figure 1** shows the pipeline; we explain the three stages of the pipeline using the databases  
135 managed by the UniProt Consortium (<http://www.uniprot.org/>) as examples.

136 **“pre-database” stage:** records from various sources are submitted to databases. For instance,  
137 UniProt protein records come from translations of primary INSDC nucleotide records (directly  
138 submitted by researchers), direct protein sequencing, gene prediction and other sources  
139 ([http://www.uniprot.org/help/sequence\\_origin](http://www.uniprot.org/help/sequence_origin)).

140 **“within database” stage:** database curation, search, and visualisation. Records are annotated  
141 in this stage, automatically (UniProtKB/TrEMBL) or through curation (UniProtKB/Swiss-Prot).  
142 Biocuration plays a vital role at this stage. For instance, UniProt manual curation not only  
143 merges records and documents discrepancies, it also annotates the records with biological  
144 knowledge drawn from the literature [28]. Also, the databases need to manage the records for  
145 search and visualisation purposes [29]. During this stage, UniProt undertakes extensive cross-  
146 referencing by linking hundreds of databases to provide centralized knowledge and resolve  
147 ambiguities [30]. **“post-database” stage:** record download, analysis, and inference. Records  
148 are downloaded and analysed for different purposes. For instance, both UniProtKB records and  
149 services have been extensively used in the research areas of biochemistry and molecular biology,  
150 biotechnology and computational biology, according to citation patterns [31]. The findings of  
151 studies may in turn contribute to new sources.

152 Duplication occurs in all of these stages, but its relevance varies. Continuing with the UniProt  
153 example, the first stage primarily concerns *entity duplicates* (often referred to as *true duplicates*):  
154 records that correspond to the same biological entities regardless of whether there are differences  
155 in the content of the database records. Merging those records into a single entry is the first step in

156 UniProtKB/Swiss-Prot manual curation [28]. The second stage primarily concerns *near-identical*  
157 *duplicates* (often referred to as *redundant records*): the records may not refer to the same  
158 entities, but nevertheless have high similarity. UniProt has found those records lead to  
159 uninformative BLAST search results ([http://www.uniprot.org/help/proteome\\_redundancy](http://www.uniprot.org/help/proteome_redundancy)). The  
160 third stage primarily concerns *study-dependent duplicates*: studies may further de-duplicate sets  
161 of records for their own purposes. For instance, studies on secondary protein structure prediction  
162 may further remove protein sequences at a 75% sequence similarity threshold [32]. This clearly  
163 shows that the notion of duplication varies and in general has two characteristics: *redundancy*  
164 and *inconsistency*. Thus it is critical to understand their characteristics, impacts, and solutions.

165 We have found numerous discussions of duplicates in the previous literature. As early as in  
166 1996, Korning et al. [33] observed duplicates from the GenBank *Arabidopsis thaliana* dataset  
167 when curating those records. The duplicates were of two main types: the same genes that were  
168 submitted twice (either by the same or different submitters), and different genes from the same  
169 gene family that were similar enough that only one was retained. Similar cases were also  
170 reported by different groups [21, 34–36]. Recently, the most significant case was the duplication  
171 in UniProtKB/TrEMBL [15]: in 2016, UniProt removed 46.9 million records corresponding to  
172 duplicate proteomes (for example, over 5.9 million of these records belong to 1,692 strains of  
173 *Mycobacterium tuberculosis*). They identified duplicate proteome records based on three criteria:  
174 belonging to the same organisms; sequence identity of over 90%; and the proteome ranks  
175 designed by biocurators (such as whether they are Reference proteome and the annotation level).

176 As this history shows, investigation of duplication has persisted for at least 20 years.  
177 Considering the type of duplicates, as the above discussion illustrates, duplication appears to be  
178 richer and more diverse than was originally described (we again note the definition of  
179 ‘duplication’ we are following in this paper, which includes the concept of redundancy). This  
180 motivates continued investigation of duplication.

181 An underlying question is: does duplication have positive or negative impact? There has been  
182 relatively little investigation of the impact of duplication, but there are some observations in the  
183 literature: (1) “The problem of duplicates is also existent in genome data, but duplicates are less  
184 interfering than in other application domains. Duplicates are often accepted and used for  
185 validation of data correctness. In conclusion, existing data cleansing techniques do not and

186 cannot consider the intricacies and semantics of genome data, or they address the wrong  
187 problem, namely duplicate elimination.” [38]; (2) “Biological data duplicates provide hints of the  
188 redundancy in biological datasets ... but rigorous elimination of data may result in loss of critical  
189 information.” [34]; (3) “The bioinformatics data is characterized by enormous diversity matched  
190 by high redundancy, across both individual and multiple databases. Enabling interoperability of  
191 the data from different sources requires resolution of data disparity and transformation in the  
192 common form (data integration), and the removal of redundant data, errors, and discrepancies  
193 (data cleaning).” [39]. Thus the answers to questions on the impact of duplicates are not clear.  
194 The above views are inconsistent, are opinions rather than conclusions drawn from studies, and  
195 are not supported by extensive examples. Moreover, they are not recent, and may not represent  
196 the current environment. Answering the question of the impact of duplications requires a more  
197 comprehensive and rigorous investigation.

198

## 199 **From duplication to other data quality issues**

200 Biological sources suffer from data quality issues other than duplication. We summarise diverse  
201 biological data quality issues reported in the literature: inconsistencies (such as conflicting  
202 results reported in the literature) [22], inaccuracies (such as erroneous sequence records and  
203 wrong gene annotations) [40–42], incompleteness (such as missing exons and incomplete  
204 annotations) [38, 40] and outdatedness (such as out-dated sequence records and annotations)  
205 [41]. This shows that while duplication is a primary data quality issue, other quality issues are  
206 also of concern. Collectively, there are five primary data quality issues: duplication,  
207 inconsistency, inaccuracy, incompleteness and outdatedness identified in general domains [43].  
208 It is thus also critical to understand what quality issues have been observed and how they impact  
209 database stakeholders under the context of biological databases.

210

## 211 **Practitioner viewpoint: survey questions**

212 Studies on data quality broadly take one of three approaches: domain expertise, theoretical or  
213 empirical. The first is opinion-based: accumulating views from (typically a small group of)  
214 domain experts [44–46]. For example, one book summarises opinions from domain experts on



215 elements of spatial data quality [44]. The second is theory-based: inference of potential data  
216 quality issues from a generic process of data generation, submission, and usage [47–49]. For  
217 example, a data quality framework was developed by inferring the data flow of a system (such as  
218 input and output for each process) and estimating the possible related quality issues [47]. The  
219 third is empirically based: analysis of data quality issues in a quantitative manner [50–52]. For  
220 example, an empirical investigation on what data quality means to stakeholders was performed  
221 via a questionnaire [50]. Each approach has its own strengths and weaknesses; for example,  
222 opinion-based studies represent high domain expertise, but may be narrow due to the small group  
223 size. Quantitative surveys in contrast have a larger number of participants, but the level of  
224 expertise may be relatively lower.

225 Our approach integrates opinion-based and empirical-based approaches: the study presents  
226 opinions from domain experts; but the data was gathered via a questionnaire; the survey  
227 questions are provided in the Supplementary Material File S1. We surveyed 23 practitioners on  
228 the questions of duplicates and other general data quality issues. These practitioners are from  
229 diverse backgrounds (including experimental biology, bioinformatics, and computer science),  
230 with a range of affiliation types (such as service providers, universities, or research institutes) but  
231 all have domain expertise. These practitioners include senior database staff, project and lab  
232 leaders, and biocurators. The publications of the participants are directly relevant to databases,  
233 data quality and curation; as illustrated by some instances [10, 15, 28, 53–69]. They were  
234 selected by personal approach at conferences and in a small number of cases by email; most of  
235 the practitioners were not known to the originating authors (Chen, Verspoor, Zobel) before this  
236 study.

237 A limitation is that the small participant size may mean that we have collected unrepresentative  
238 opinions. However, the community of biocuration is small and the experience represented by  
239 these 23 is highly relevant. A 2012 survey conducted by the International Society of Biocuration  
240 (ISB) had 257 participants [67]. Of those 257 participants, 57% of them were employed in short-  
241 term contracts and only 9% were principal investigators. A similar study initiated by the  
242 BioCreative team had only 30 participants, including all the attendees of the BioCreative  
243 conference in that year [68]. Therefore, the number of participants of this study reflects the size

244 of the biocuration community; moreover, the relatively high expertise ensures the validity of the  
245 opinions.

246 The survey asked three primary questions about duplication: (1) *What* are duplicates? We  
247 asked practitioners what records they think should be regarded as duplicated; (2) *Why* care about  
248 duplicates? We asked practitioners what impact duplicates have; (3) *How* to manage duplicates?  
249 We asked practitioners whether, and how, duplicates should be resolved.

250 In detail, the questions and their possible responses were as follows:

251 *Defining duplicate records (The ‘what’ question).* We provided five options for experts to  
252 select: (1) Exact duplicate records: two or more records are exactly identical; (2) Near identical  
253 duplicates: two or more records are not identical but similar; (3) Partial or fragmentary records:  
254 one record is a fragment of another; (4) Duplicate records with low similarity: records have  
255 relatively low similarity but belong to the same entity; (5) Other types: if practitioners also  
256 consider other cases as duplicates.

257 Respondents were asked to comment on their choices. We also requested examples to support  
258 the choice of options 4 or 5, given that in our review of the literature we observed that the first  
259 three options were prevalent [70, 71]. Option 1 refers to exact duplicates, option 2 refers to  
260 (highly) similar or redundant records or to some quantitative extent, records share X% similarity,  
261 option 3 refers to partial or incomplete records, and option 4 refers to entity duplicates that are  
262 inconsistent. The “Other types” option provides capture of remaining types of duplicates.

263 *Quantifying the impacts of duplication (The ‘Why’ question).* We asked in two steps: first,  
264 whether respondents believed that duplicates have impact. The second question was presented  
265 only if the answer to the first was yes. It is used to comment on positive and negative impacts  
266 respectively. We also asked respondents to explain their opinion or give examples.

267 *Addressing duplication (The ‘How’ question).* We offered three subquestions: (1) Do you  
268 believe that duplicate detection is useful/needed? (2) Do you believe that current duplicate  
269 detection methods/software are sufficient to satisfy your requirements? (We also asked  
270 respondents to explain what they expected if they selected ‘no’.) (3) How would you prefer that  
271 duplicate records be handled? These were the suggested options: label and remove duplicates,  
272 label and make duplicates obsolete, label but leave duplicates active, and other solutions.

## 273 **Practitioner viewpoints: summary**

274 In this section, we present the survey results on duplication and other quality issues.

### 275 **Duplication: practitioners' opinions**

276 The responses are summarized below, in the same order as the three primary questions. For each  
277 question, we detail the response statistics, summarise the common patterns, augmented by  
278 detailed responses, and draw conclusions.

279 The views on *what are duplicates* are summarised in **Figure 2**. Out of 23 practitioners, 21 have  
280 made a choice by selecting at least one option. While the other two did not select any options,  
281 they have considered that duplicates have impacts for later questions. We therefore do not regard  
282 the empty responses as an opinion that duplication does not exist; rather simply do not track the  
283 response in this case.

284 The results show (1) all types of duplicates have been observed by some of practitioners, but  
285 none is universal. The commonest type is *similar record*, which was selected by over half of the  
286 respondents; but the other types (*exact duplicates*, *partial records*, and *low similarity duplicates*)  
287 were also selected by at least a third of the respondents. Three of them considered *other*  
288 *duplicate* types, and (2) more than 80% of respondents indicated that they have observed at least  
289 two types.

290 Also recall that existing literature rarely covers the fourth type of duplication – that is,  
291 relatively different records that should in fact be considered as duplicates. However, close to  
292 40% of respondents acknowledge having seen such cases and further point out that identifying  
293 them requires significant manual effort. The following summarises several cases (each identified  
294 by respondent ID, tabulated at the end of this paper).

295 *Low similarity duplicates within a single database*. Representative comments are “We have  
296 such records in ClinVar [64]. We receive independent submissions from groups that define  
297 variants with great precision, and groups that define the same variant in the same paper, but  
298 describe it imprecisely. Curators have to review the content to determine identity.” [R19] and  
299 “Genomes or proteomes of the same species can often be different enough even they are  
300 redundant.” [R24]

301 *Low similarity duplicates in databases having cross-references.* Representative comments are  
302 “Protein-Protein Interaction databases: the same publication may be in BioGRID [72] annotated  
303 at the gene level and in one of the IMEx databases (<http://www.imexconsortium.org/>) annotated  
304 at the protein level.” [R20] and “Also secondary databases import data (*e.g.* STRING sticking to  
305 the PPI example) but will only import a part of what is available.” [R20].

306 *Low similarity duplicates in databases having the same kinds of contents.* For instance,  
307 “Pathway databases (KEGG<sup>29</sup>, Reactome<sup>30</sup>, EcoCyc<sup>31</sup> etc) tend to look at same pathways but are  
308 open to curator interpretation and may differ.” [R20]

309 The results of the “why care about duplicates” question are shown in **Figure 3**. All  
310 practitioners made a choice. Most (21 out of 23) believe that duplication does matter. Moreover,  
311 19 out of 21 experts weighted on potential impact of duplicates: only one believed that the  
312 impact is purely positive, compared to 8 viewing it solely negative; the remaining 10 thought the  
313 impact has both positive and negative sides. We assembled all responses on impacts of  
314 duplicates as follows below.

315 *Impact on database storage, search and mapping.* Representative comments are (1) “When  
316 duplicates (sequence only) are in big proportion they will have an impact on sequence search  
317 tool like BLAST, when pre-computing the database to search against. Then it’ll affect the  
318 statistics on the E-value returned.” [R10], (2) “Duplicates in one resource make exact mappings  
319 between 2 resources difficult.” [R21], “Highly redundant records can result in: Increasing bias in  
320 statistical analyses; Repetitive hits in BLAST searches.” [R24], and (3) “Querying datasets with  
321 duplicate records impacts the diversity of hits and increase overall noise; we have discussed this  
322 in our paper on hallmark signatures” [56]. [R8]

323 *Impact on meta-analysis in biological studies.* Representative comments are (1) “Duplicate  
324 transcriptome records can impact the statistics of meta-analysis.” [R1], (2) “Authors often state a  
325 fact is correct because it has been observed in multiple resources. If the resources are re-using, or  
326 recycling the same piece of information, this statement (or statistical measure), is incorrect.”  
327 [R20] (Note that it has been previously observed that cascading errors may arise due to this type  
328 of propagation of information [73].) and (3) “Duplicates affect enrichments if duplicate records  
329 used in background sets.” [R21]

330 *Impact on time and resources.* Representative comments are (1) “Archiving and storing  
331 duplicated data may just be a waste of resources.” [R12], (2) “Result in time wasted by the  
332 researcher.” [R19], and (3) “As a professional curation service; our company suffers from the  
333 effects of data duplication daily. Unfortunately there is no pre-screening of data done by  
334 Biological DBs and thus it is up to us to create methods to identify data duplication before we  
335 commit time to curate samples. Unfortunately, with the onset of next generation data, it has  
336 become hard to detect duplicate data where the submitter has intentionally re-arranged the reads  
337 without already committing substantial computational resources in advance”. [R9]

338 *Impact on users.* Representative comments are (1) “Duplicate records can result in confusion  
339 by the novice user. If the duplication is of the ‘low similarity’ type, information may be  
340 misleading.” [R19], “Duplicate gene records may be misinterpreted as species paralogs.” [R21],  
341 (2) “When training students, they can get very confused when a protein in a database has  
342 multiple entries -which one should they use, for example. Then I would need to compare the  
343 different entries and select one for them to use. It would be better if the information in the  
344 duplicate entries was combined into one correct and more complete entry.” [R23], and (3) “Near  
345 identical duplicate records: two or more records are not strictly identical but very similar and can  
346 be considered duplicates; because users don't realise they are the same thing or don't understand  
347 the difference between them.” [R25].

348 In contrast, practitioners also pointed out two primary positive impacts: (1) identified  
349 duplicates enrich the information about an entity; for example, “When you try to look sequence  
350 homology across species, it is good to keep duplicates as it allows to build orthologous trees.”  
351 [R10] and “When they are isoforms of each other - so while they are for the same entity, they  
352 have distinct biological significance.” [R25], and (2) identified duplicates verify the correctness  
353 as replications; for example, “On the other hand, if you have many instances of the same data, or  
354 near identical data, one could feel more confident on that data point.” [R12] (Note that  
355 confidence information ontology can be used to capture “confidence statement from multiple  
356 evidence lines of same type” [74].), and “If it is a duplicate record that has arisen from different  
357 types of evidence, this could strengthen the claim.” [R13]

358 The cases outlined above detail the impact of duplication.. Clearly duplication does matter. The  
359 negative impacts are broad. They range from databases to studies, from research to training, and

360 from curators to students. The potential impacts are severe: valuable search results may be  
361 missed, statistical results may be biased, and study interpretations may be misled. Management  
362 of duplication during is a significant amount of labour.

363 Our survey respondents identified duplicates as having two main positive impacts: enriching  
364 the information and verifying the correctness. This has an implicit yet important prerequisite: **the**  
365 **duplicates need to be detected and labelled beforehand**. For instance, in order to achieve  
366 information richness, duplicate records must first be accurately identified and cross-references  
367 should be explicitly made. Similarly, for confirmation of results, the duplicate records need to be  
368 labelled beforehand. Researchers then can seek labelled duplicates to find additional interesting  
369 observations made by other researchers on the same entities, that is, to find out whether their  
370 records are consistent with others.

371 The views on *how to manage duplicates* are summarised in **Figure 4**. None of the practitioners  
372 regards duplicate detection as unnecessary; 10 practitioners further believe that current duplicate  
373 detection methods are not sufficient. We propose the following suggestions accordingly.

374 *Precision matters.* The methods need to find duplicates accurately: “It should correctly remove  
375 duplicate records, while leaving legitimate similar entries in the database.” [R15] and “Duplicate  
376 detection method need to be invariant to small changes (at the file level, or biological sample  
377 level); otherwise we would miss the vast majority of these.” [R9]

378 *Automation matters.* In some fields few duplicate detection methods exist: “We re-use GEO  
379 public data sets, to our knowledge there is no systematic duplicate detection.” [R7], “Not aware  
380 of any software.” [R3] and “I do not use any duplicate detection methods, they are often difficult  
381 to spot are usually based on a knowledge of the known size of the gene set.” [R21]

382 *Characterisation matters.* The methods should analyse the characteristics of duplicates: “A  
383 measure of how redundant the database records are would be useful.” [R24]

384 *Robustness and generalisation matter.* “All formats of data need to be handled cross-wise; it  
385 does not help trying to find duplicates only within a single file format for a technology.” [R9]

386 To our knowledge, there is no universal approach to managing duplication. Similar databases  
387 may use different de-duplication techniques. For instance, as sequencing databases, ENCODE  
388 uses standardized metadata organisation, multiple validation identifiers, and its own merging



389 mechanism for the detection and management of duplicate sequencing reads; the Sequence Read  
390 Archive (SRA) uses hash functions whereas GEO uses manual curation in addition to hash  
391 functions [27]. Likewise, different databases may choose different parameters even when using  
392 the same de-duplication approach. For instance, protein databases often use clustering methods  
393 to handle redundant records. However, the values of chosen similarity thresholds for clustering  
394 range from 30% to 100% in different databases [75]. Thus, it is impossible to provide a uniform  
395 solution to handling of duplication (as well as other quality issues). We introduce sample  
396 solutions used in UniProtKB/Swiss-Prot that demonstrate how quality issues are handled in a  
397 single database. The approaches or software used in the UniProtKB/Swiss-Prot curation pipeline  
398 may also provide insights into others.

399

#### 400 **Beyond duplication: other data quality issues**

401 We also extend the investigation to general quality issues other than duplication, to complement  
402 the key insights. We asked the respondents for their opinions on general data quality issues. The  
403 two primary questions asked were: *what* data quality issues have been observed in biological  
404 databases? and *why* care about data quality? The style is the same as the above questions on  
405 duplication. The detailed results are summarized in Supplementary Material File S2. Overall it  
406 shows the quality issues can be widespread; for example, each data quality issue has been  
407 observed by at least 80% of the respondents.

408

#### 409 **Limitations**

410 It is worth noting that while we have carefully phrased the questions in the survey, it may still be  
411 the case that different respondents may have different internal definitions of duplicates in mind  
412 when responding. For example, some respondents may only consider records with minor  
413 differences as redundant records whereas others may also include records with larger differences,  
414 even though they selected the same option. We acknowledge that this diversity of interpretation  
415 is inevitable – data is multifaceted; hence so is data quality and the associated perspectives on it.  
416 The internal definitions of duplicate records depend on more specific context and there is indeed  
417 no universal agreement [24]. However, we argue that this does not detract from the results of the  
418 survey; respondents provided clear examples to support their choices and those examples

419 demonstrate that the duplicate types do impact biological studies, regardless of internal variation  
420 in specific definitions. Such internal differences are also observed in other data quality studies,  
421 such as reviews on general data quality [76] and detection of duplicate videos [77].

422 It is also noteworthy that some databases primarily serve an archival purpose, such as INSDC  
423 and GEO. The records in these databases are directly coordinated by record submitters;  
424 therefore, the databases have had relatively little curation compared to databases like  
425 UniProtKB/Swiss-Prot. Arguably, data quality issues are not major concerns from an archival  
426 perspective. We do not examine the quality issues in archival databases; rather, we suggest  
427 labelling duplicate records or records with other quality issues (without withdrawing or removing  
428 the records) could potentially facilitate database usage. The archival purpose does not limit other  
429 uses; for example, studies including BLAST searches against GenBank for sequence  
430 characterization [78–80]. In such cases, the sequences and annotations would impact the related  
431 analyses.

432 However, quality issues may be important in archival databases. Indeed, in some instances the  
433 database managers have been aware of data quality issues and are working on solutions. A recent  
434 work proposed by the ENCODE database team concerns the quality issues, in particular  
435 duplication in sequencing repositories such as ENCODE, GEO and SRA [27]. They  
436 acknowledge that, while archival databases are responsible for data preservation, duplication  
437 affects data storage and could mislead users. As a result, they propose three guidelines to prevent  
438 duplication in ENCODE and summarise other de-duplication approaches in GEO and SRA;  
439 furthermore, the ENCODE work encourages making a community effort (such as archival  
440 databases, publishers, and submitters) to handle quality issues.

441

## 442 **Biocuration: a solution to data quality issues in biological databases**

443 In this section, we introduce solutions to data quality issues in biological databases. Biocuration  
444 is a general term that refers to addressing data quality issues in biological databases. We provide  
445 a concrete case study on the UniProtKB/Swiss-Prot curation pipeline – consisting of a detailed  
446 description on the curation procedure and an interview with the curation team leader. It provides  
447 an example of a solution to different quality issues.

448



#### 449 **The curation pipeline of UniProtKB/Swiss-Prot**

450 UniProtKB has two data sections: UniProtKB/Swiss-Prot and UniProtKB/TrEMBL. Sequence  
451 records are first deposited in UniProtKB/TrEMBL and then selected records are transferred into  
452 UniProt/Swiss-Prot. Curation in UniProtKB has two stages: (1) automatic curation in  
453 UniProt/TrEMBL, where records are curated by software automatically without manual review,  
454 and (2) expert (or manual) curation in UniProtKB/Swiss-Prot on selected records from  
455 UniProtKB/TrEMBL. A major task in automatic curation is to annotate records using annotation  
456 systems; for example, UniRules, which contains rules created by biocurators, and external rules  
457 from other annotation systems, such as RuleBase [81] and HAMAP [82], are used in this task.  
458 Rule UR000031345 is an example of UniRules (<http://www.uniprot.org/unirule/UR000031345>);  
459 Record B1YYB is also a sequence record example that was annotated using the rules during  
460 automatic curation. For expert curation, biocurators run a comprehensive set of software, search  
461 supporting information from range of databases, manually review the results and interpret the  
462 evidence level [31]. **Table 1** describes representative software and databases used in expert  
463 curation [14, 83–98]. This expert curation in UniProtKB/Swiss-Prot has 6 dedicated steps, shown  
464 in Table 1 and explained below.

465 *Sequence curation.* This step focuses on de-duplication. It has two components: (1) Detect and  
466 merge duplicate records. (2) Analyse and document the inconsistencies caused by duplication. In  
467 this specific case ‘duplicates’ are records belonging to the same genes: an example of entity  
468 duplicates. Biocurators perform BLAST searches and also search other database resources to  
469 confirm whether two records are the same genes, and merge them if they are. The merged  
470 records are explicitly documented in the record’s *Cross-reference* section. Sometimes the  
471 merged records do not have the same sequences, mostly due to errors. Biocurators have to  
472 analyse the causes of those differences and document the errors.

473 *Sequence analysis.* Biocurators analyse sequence features after addressing duplication and  
474 inconsistencies. They run standard prediction tools, review and interpret the results, and annotate  
475 the records. The complete annotations for sequence features cover 39 annotation fields under 7  
476 categories: Molecule processing, Regions, Sites, Amino acid modifications, Natural variations,  
477 Experimental info, and Secondary structure ([http://www.uniprot.org/help/sequence\\_annotation](http://www.uniprot.org/help/sequence_annotation)).

478 As such, it involves a comprehensive range of software and databases to facilitate sequence  
479 analysis, some of which are shown in Table 1.

480 *Literature curation.* This step often contains two processes: retrieval of relevant literature and  
481 application of text mining tools to analysis of text data, such as recognising named entities [99]  
482 and identifying critical entity relationships [100]. The annotations are made using controlled  
483 vocabularies (the complete list is in the UniProt keyword documentation via  
484 <http://www.uniprot.org/docs/keywlist>) and are explicitly labelled as “*Manual assertion based on*  
485 *experiment in literature*”. Record Q24145 is an example that was annotated based on findings  
486 published in literature (<http://www.uniprot.org/uniprot/Q24145>).

487 *Family-based curation.* This step transitions curation from single-record level to family-level,  
488 finding relationships amongst records. Biocurators identify putative homologs using BLAST  
489 search results and phylogenetic resources and make annotations accordingly. The tools and  
490 databases are the same as those in the *Sequence curation* step.

491 *Evidence Attribution.* This step standardises the curations made in the previous steps. Curations  
492 are made manually or automatically from different types of sources, such as sequence similarity,  
493 animal model results and clinical study results. This step uses the Evidence and Conclusion  
494 Ontology (ECO) to describe evidence in a precise manner: it details the type of evidence and the  
495 assertion method (manual or automatic) used to support a curated statement [98]. As such,  
496 database users can know how the decision was made and on what basis. For example,  
497 ECO\_0000269 was used in the literature curation for Record Q24145.

498 *Quality assurance, integration and update.* The curation is complete at this point. This step  
499 finally checks everything and integrates curated records to the existing UniProtKB/Swiss-Prot  
500 knowledgebase. Those records will then be available in the new release. In turn, it helps further  
501 automatic curation within UniProtKB/Swiss-Prot. The newly made annotations will be used as  
502 the basis for creating automatic annotation rules.

503

#### 504 **The curation in UniProtKB/Swiss-Prot: an interview**

505 We interviewed UniProtKB/Swiss-Prot annotation team leader Sylvain Poux. The interview  
506 questions covered how UniProtKB/Swiss-Prot handles general data quality issues. Some of the

507 responses are also related to specific curation process in UniProtKB/Swiss-Prot which shows that  
508 the solutions are database-dependent as well. The detailed interview is summarized in the  
509 Supplementary Material File S3. We have edited the questions for clarity, and omitted answers  
510 where Poux did not offer a view.

511 The above case study demonstrates that biocuration is an effective solution to diverse quality  
512 issues. Indeed, since 2003, when the first regular meeting amongst biocurators was held [101],  
513 the importance of biocuration activities has widely been recognised [20, 102–104]. Yet, on the  
514 other hand, the biocuration community still lacks broader support. A survey of 257 former or  
515 current biocurators showed that biocurators suffered from a lack of secured funding for primary  
516 biological databases, exponential data growth, and underestimation of the importance of  
517 biocuration [69]; consistent results were also demonstrated in other studies [105, 106].  
518 According to recent reports, the funding for model-organism databases will be cut 30%-40% and  
519 the same threat applies to other databases [107–109].

520

## 521 **Conclusion**

522 In this study, we explored the perspectives of both database managers and database users on the  
523 issue of data duplication – one of several significant data quality issues. We also extended the  
524 investigation to other data quality issues to complement this primary focus. Our survey of  
525 individual practitioners showed that duplication in biological databases is of concern: its  
526 characteristics are diverse and complex, its impacts cover almost all stages of database creation  
527 and analysis, and methods for managing the problem of duplication, either manual or automatic,  
528 have significant limitations. The overall impacts of duplication are broadly negative, and the  
529 positive impacts such as enriched entity information and validation of correctness rely on the  
530 duplicate records being correctly labelled or cross-referenced. This suggests a need for further  
531 development of methods for precisely classifying duplicate records (accuracy), detecting  
532 different duplicate types (characterisation), and achieving scalable performance in different data  
533 collections (generalisation). In some specific domains duplicate detection software (automation)  
534 is a critical need.

535 The responses relating to general data quality further show that data quality issues go well  
536 beyond duplication. As can be inferred from the survey we conducted, curation – dedicated  
537 efforts to ensure that biological databases represent accurate and up-to-date scientific knowledge  
538 – is an effective tool for addressing quality issues. We provide a concrete case study on the  
539 UniProtKB/Swiss-Prot curation pipeline as a sample solution to quality issues. However, manual  
540 curation alone is not sufficient to resolve all data quality problems due to rapidly growing data  
541 volumes in a context of limited resources. A broader community effort is required to manage  
542 data quality and to provide support to facilitate data quality and curation.

543

### 544 **Authors' contributions**

545 QC, JZ and KV initiated the survey, analysed the results and wrote the paper. RB, IE, CJ, AL,  
546 MM, JO, MR, JS, and RY contributed to presenting the views and revising the paper. All authors  
547 read and approved the final manuscript.

548

### 549 **Competing interests**

550 The authors have declared no competing interests.

551

### 552 **Acknowledgments**

553 The project receives funding from the Australian Research Council through a Discovery Project  
554 grant, DP150101550. We thank Sylvain Poux for contributions to the UniProtKB/Swiss-Prot  
555 curation case study. We acknowledge the participation of the following people in the survey:  
556 Cecilia Arighi (University of Delaware), Ruth C Lovering (University College London), Peter  
557 McQuilton (University of Oxford), and Valerie Wood (University of Cambridge).

558

### 559 **References**

- 560 [1] Baxevanis A, Bateman A. The importance of biological databases in biological discovery. *Curr Protoc*  
561 *Bioinformatics* 2015;50:1–8.
- 562 [2] Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, et al. GenBank. *Nucleic*  
563 *Acids Res* 2017;45:D37.

- 564 [3] Toribio AL, Alako B, Amid C, Cerdeño-Tarrága A, Clarke L, Cleland I, et al. European nucleotide  
565 archive in 2016. *Nucleic Acids Res* 2017;45:D32–6.
- 566 [4] Cochrane G, Karsch-Mizrachi I, Takagi T. The international nucleotide sequence database  
567 collaboration. *Nucleic Acids Res* 2017;44:D48–51.
- 568 [5] The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids Res*  
569 2017;45:D158–69.
- 570 [6] Rose PW, Prlić A, Altunkaya A, Bi C, Bradley AR, Christie CH, et al. The RCSB protein data bank:  
571 integrative view of protein, gene and 3D structural information. *Nucleic Acids Res* 2017;45:D271–81.
- 572 [7] Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, et al. Rfam 12.0: updates to  
573 the RNA families database. *Nucleic Acids Res* 2015;43:D130–7.
- 574 [8] Finn RD, Coghill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al. The Pfam protein families  
575 database: towards a more sustainable future. *Nucleic Acids Res* 2016;44:D279–85.
- 576 [9] Basu S, Fey P, Pandit Y, Dodson R, Kibbe WA, Chisholm RL. DictyBase 2013: integrating multiple  
577 Dictyostelid species. *Nucleic Acids Res* 2013;41:D676–83.
- 578 [10] McDowall MD, Harris MA, Lock A, Rutherford K, Staines DM, Bähler J, et al. PomBase 2015:  
579 updates to the fission yeast database. *Nucleic Acids Res* 2015;43:D656–61.
- 580 [11] Kolesnikov N, Hastings E, Keays M, Melnichuk O, Tang YA, Williams E, et al. ArrayExpress  
581 update—simplifying data submissions. *Nucleic Acids Res* 2014:gku1057.
- 582 [12] Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO: archive  
583 for functional genomics data sets—update. *Nucleic Acids Res* 2013;41:D991–5.
- 584 [13] Galperin MY, Fernández-Suárez XM, Rigden DJ. The 24th annual *Nucleic Acids Research* database  
585 issue: a look back and upcoming changes. *Nucleic Acids Res* 2017;45:D1–11.
- 586 [14] The Gene Ontology Consortium. Expansion of the Gene Ontology knowledgebase and resources.  
587 *Nucleic Acids Res* 2017;45:D331–8.
- 588 [15] Bursteinas B, Britto R, Bely B, Auchincloss A, Rivoire C, Redaschi N, et al. Minimizing proteome  
589 redundancy in the UniProt Knowledgebase. *Database (Oxford)* 2016;2016.
- 590 [16] Bouadjenek MR, Verspoor K, Zobel J. Literature consistency of bioinformatics sequence databases  
591 is effective for assessing record quality. *Database (Oxford)* 2017;2017.
- 592 [17] Poux S, Arighi CN, Magrane M, Bateman A, Wei CH, Lu Z, et al. On expert curation and  
593 sustainability: UniProtKB/Swiss-Prot as a case study. *Bioinformatics* 2017:3454–60.
- 594 [18] Nellore A, Jaffe AE, Fortin J-P, Alquicira-Hernández J, Collado-Torres L, Wang S, et al. Human  
595 splicing diversity and the extent of unannotated splice junctions across human RNA-seq samples on the  
596 Sequence Read Archive. *Genome Biol* 2016;17:266.

- 597 [19] Huntley RP, Sitnikov D, Orlic-Milacic M, Balakrishnan R, D'Eustachio P, Gillespie ME, et al.  
598 Guidelines for the functional annotation of microRNAs using the Gene Ontology. *RNA* 2016;22:667-76.
- 599 [20] Howe D, Costanzo M, Fey P, Gojobori T, Hannick L, Hide W, et al. Big data: the future of  
600 biocuration. *Nature* 2008;455:47–50.
- 601 [21] Rosikiewicz M, Comte A, Niknejad A, Robinson-Rechavi M, Bastian FB. Uncovering hidden  
602 duplicated content in public transcriptomics data. *Database (Oxford)* 2013;2013:bat010.
- 603 [22] Poux S, Magrane M, Arighi CN, Bridge A, O'Donovan C, Laiho K. Expert curation in UniProtKB: a  
604 case study on dealing with conflicting and erroneous data. *Database (Oxford)* 2014;2014.
- 605 [23] Pfeiffer F, Oesterhelt D. A manual curation strategy to improve genome annotation: application to a  
606 set of haloarchael genomes. *Life* 2015;5:1427–44.
- 607 [24] Chen Q, Zobel J, Verspoor K. Duplicates, redundancies and inconsistencies in the primary nucleotide  
608 databases: a descriptive study. *Database (Oxford)* 2017;2017.
- 609 [25] Chen Q, Zobel J, Verspoor K. Benchmarks for measurement of duplicate detection methods in  
610 nucleotide databases. *Database (Oxford)* 2017.
- 611 [26] Chen Q, Zobel J, Zhang X, Verspoor K. Supervised Learning for Detection of Duplicates in  
612 Genomic Sequence Databases. *PLoS One* 2016;11:e0159644.
- 613 [27] Gabdank I, Chan ET, Davidson JM, Hilton JA, Davis CA, Baymuradov UK, et al. Prevention of data  
614 duplication for high throughput sequencing repositories. *Database (Oxford)* 2018;2018.
- 615 [28] The UniProt Consortium. Activities at the universal protein resource (UniProt). *Nucleic Acids Res*  
616 2014;42:D191–D8.
- 617 [29] Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH. UniRef clusters: a comprehensive and  
618 scalable alternative for improving sequence similarity searches. *Bioinformatics* 2015;31:926–32.
- 619 [30] Gasteiger E, Jung E, Bairoch AM. SWISS-PROT: connecting biomolecular knowledge via a protein  
620 database. *Curr Issues Mol Biol* 2001;3:47–55.
- 621 [31] The UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res* 2014:gku989.
- 622 [32] Cole C, Barber JD, Barton GJ. The Jpred 3 secondary structure prediction server. *Nucleic Acids Res*  
623 2008;36:W197–W201.
- 624 [33] Korning PG, Hebsgaard SM, Rouzé P, Brunak S. Cleaning the GenBank *Arabidopsis thaliana* data  
625 set. *Nucleic Acids Res* 1996;24:316–20.
- 626 [34] Koh J, Lee ML, Khan AM, Tan P, Brusica V. Duplicate detection in biological data using association  
627 rule mining. *Proceedings of the Second European Workshop on Data Mining and Text Mining in*  
628 *Bioinformatics* 2004;501:S22388.



- 629 [35] Salgado H, Santos-Zavaleta A, Gama-Castro S, Peralta-Gil M, Peñaloza-Spínola MI, Martínez-  
630 Antonio A, et al. The comprehensive updated regulatory network of *Escherichia coli* K-12. *BMC*  
631 *Bioinformatics* 2006;7:5.
- 632 [36] Bouffard M, Phillips MS, Brown AM, Marsh S, Tardif J-C, van Rooij T. Damming the genomic data  
633 flood using a comprehensive analysis and storage data structure. *Database (Oxford)* 2010;2010:baq029.
- 634 [37] Bastian F, Parmentier G, Roux J, Moretti S, Laudet V, Robinson-Rechavi M. Bgee: integrating and  
635 comparing heterogeneous transcriptome data among species. *International Workshop on Data Integration*  
636 *in the Life Sciences* 2008:124–31.
- 637 [38] Müller H, Naumann F, Freytag J-C. Data quality in genome databases. *Proceedings of the*  
638 *Conference on Information Quality* 2003.
- 639 [39] Chellamuthu S, Punithavalli DM. Detecting redundancy in biological databases? an efficient  
640 approach. *Global Journal of Computer Science and Technology* 2009;9.
- 641 [40] Bork P, Bairoch A. Go hunting in sequence databases but watch out for the traps. *Trends Genet*  
642 *1996;12:425–7.*
- 643 [41] Pennisi E. Keeping genome databases clean and up to date. *Science* 1999;286:447–50.
- 644 [42] Schnoes AM, Brown SD, Dodevski I, Babbitt PC. Annotation error in public databases:  
645 misannotation of molecular function in enzyme superfamilies. *PLoS Comput Biol* 2009;5:e1000605.
- 646 [43] Fan W. Data quality: From theory to practice. *Proc ACM SIGMOD Int Conf Manag Data*  
647 *2015;44:7–18.*
- 648 [44] Guptill SC, Morrison JL. *Elements of spatial data quality*. Elsevier, 2013.
- 649 [45] Abiteboul S, Dong L, Etzioni O, Srivastava D, Weikum G, Stoyanovich J, et al. The elephant in the  
650 room: getting value from Big Data. *Proceedings of the 18th international workshop on web and databases*  
651 *2015:1–5.*
- 652 [46] Sadiq S, Papotti P. Big data quality-whose problem is it? *IEEE 32nd International Conference on*  
653 *Data Engineering (ICDE)* 2016:1446–7.
- 654 [47] Ballou DP, Pazer HL. Modeling data and process quality in multi-input, multi-output information  
655 systems. *Manage Sci* 1985;31:150–62.
- 656 [48] Wang RY, Storey VC, Firth CP. A framework for analysis of data quality research. *IEEE Trans*  
657 *Knowl Data Eng* 1995;7:623–40.
- 658 [49] Yeganeh NK, Sadiq S, Sharaf MA. A framework for data quality aware query systems. *Inf Syst*  
659 *2014;46:24–44.*
- 660 [50] Wang RY, Strong DM. Beyond accuracy: What data quality means to data consumers. *J Manag Inf*  
661 *Syst* 1996;12:5–33.

- 662 [51] Wixom BH, Watson HJ. An empirical investigation of the factors affecting data warehousing  
663 success. *MIS quarterly* 2001;17–41.
- 664 [52] Coussement K, Van den Bossche FA, De Bock KW. Data accuracy's impact on segmentation  
665 performance: Benchmarking RFM analysis, logistic regression, and decision trees. *J Bus Res*  
666 2014;67:2751–8.
- 667 [53] Bultet LA, Aguilar Rodriguez J, Ahrens CH, Ahrne EL, Ai N, Aimo L, et al. The SIB Swiss Institute  
668 of Bioinformatics" resources: focus on curated databases. *Nucleic Acids Res* 2016;44:D27–D37.
- 669 [54] Magrane M, The UniProt Consortium. UniProt Knowledgebase: a hub of integrated protein data.  
670 *Database (Oxford)* 2011;2011:bar009.
- 671 [55] Mani M, Chen C, Amblee V, Liu H, Mathur T, Zwicke G, et al. MoonProt: a database for proteins  
672 that are known to moonlight. *Nucleic Acids Res* 2014;43:D277–D82.
- 673 [56] Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The molecular  
674 signatures database hallmark gene set collection. *Cell Syst* 2015;1:417–25.
- 675 [57] Kılıç S, White ER, Sagitova DM, Cornish JP, Erill I. CollecTF: a database of experimentally  
676 validated transcription factor-binding sites in Bacteria. *Nucleic Acids Res* 2014;42:D156–D60.
- 677 [58] Kılıç S, Sagitova DM, Wolfish S, Bely B, Courtot M, Ciufo S, et al. From data repositories to  
678 submission portals: rethinking the role of domain-specific databases in CollecTF. *Database (Oxford)*  
679 2016;2016.
- 680 [59] Rutherford KM, Harris MA, Lock A, Oliver SG, Wood V. Canto: an online tool for community  
681 literature curation. *Bioinformatics* 2014;30:1791–2.
- 682 [60] Pundir S, Martin MJ, O'Donovan C. Protein Bioinformatics: From Protein Modifications and  
683 Networks to Proteomics. 2017.
- 684 [61] Poux S, Arighi CN, Magrane M, Bateman A, Wei C-H, Lu Z, et al. On expert curation and  
685 sustainability: UniProtKB/Swiss-Prot as a case study. *Bioinformatics* 2017.
- 686 [62] Gaudet P, Michel P-A, Zahn-Zabal M, Britan A, Cusin I, Domagalski M, et al. The neXtProt  
687 knowledgebase on human proteins: 2017 update. *Nucleic Acids Res* 2017;45:D177–D82.
- 688 [63] Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu S, Chitsaz F, Geer LY, et al. CDD: NCBI's  
689 conserved domain database. *Nucleic Acids Res* 2014;43:D222–D6.
- 690 [64] Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, et al. ClinVar: public archive of  
691 interpretations of clinically relevant variants. *Nucleic Acids Res* 2016;44:D862–D8.
- 692 [65] Orchard S. Data standardization and sharing—the work of the HUPO-PSI. *Biochim Biophys Acta*  
693 2014;1844:82–7.
- 694 [66] Poux S, Gaudet P (2017), 'Best practices in manual annotation with the gene ontology', *The Gene*  
695 *Ontology Handbook*, pp. 41–54.

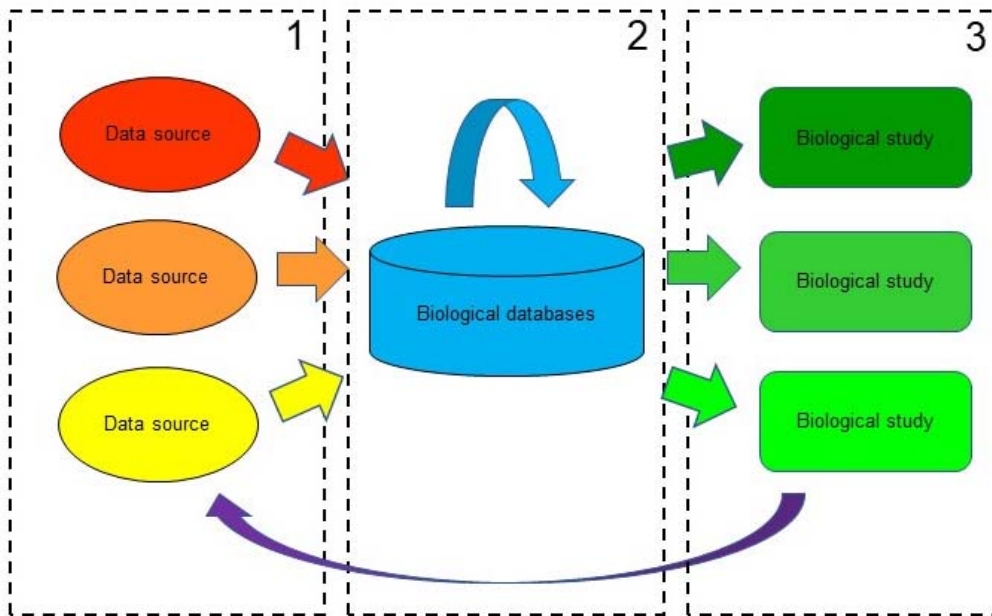


- 696 [67] Burge S, Attwood TK, Bateman A, Berardini TZ, Cherry M, O'Donovan C, et al. Biocurators and  
697 biocuration: surveying the 21st century challenges. *Database (Oxford)* 2012;2012:bar059.
- 698 [68] Hirschman L, Burns GAC, Krallinger M, Arighi C, Cohen KB, Valencia A, et al. Text mining for the  
699 biocuration workflow. *Database* 2012;2012:bas020.
- 700 [69] Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular  
701 signatures database (MSigDB) 3.0. *Bioinformatics* 2011;27:1739–40.
- 702 [70] Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation sequencing  
703 data. *Bioinformatics* 2012;28:3150–2.
- 704 [71] Song M, Rudniy A. Detecting duplicate biological entities using Markov random field-based edit  
705 distance. *Knowl Inf Syst* 2010;25:371–87.
- 706 [72] Chatr-aryamontri A, Oughtred R, Boucher L, Rust J, Chang C, Kolas NK, et al. The BioGRID  
707 interaction database: 2017 update. *Nucleic Acids Res* 2017;45:D369–D79.
- 708 [73] Gilks WR, Audit B, De Angelis D, Tsoka S, Ouzounis CA. Modeling the percolation of annotation  
709 errors in a database of protein sequences. *Bioinformatics* 2002;18:1641–9.
- 710 [74] Bastian FB, Chibucos MC, Gaudet P, Giglio M, Holliday GL, Huang H, et al. The Confidence  
711 Information Ontology: a step towards a standard for asserting confidence in annotations. *Database*  
712 (Oxford) 2015;2015:bav043.
- 713 [75] Chen Q, Wan Y, Zhang X, Lei Y, Zobel J, Verspoor K. Comparative Analysis of Sequence  
714 Clustering Methods for Deduplication of Biological Databases. *ACM J Data Inf Qual* 2018;9:17.
- 715 [76] Batini C, Scannapieco M. *Data and Information Quality: Dimensions, Principles and Techniques*.  
716 Springer, 2016.
- 717 [77] Liu J, Huang Z, Cai H, Shen HT, Ngo CW, Wang W. Near-duplicate video retrieval: Current  
718 research and future trends. *ACM Comput Surv* 2013;45:44.
- 719 [78] Chowdhary A, Kathuria S, Singh PK, Sharma B, Dolatabadi S, Hagen F, et al. Molecular  
720 characterization and in vitro antifungal susceptibility of 80 clinical isolates of mucormycetes in Delhi,  
721 India. *Mycoses* 2014;57:97–107.
- 722 [79] Qiao Y, Xu D, Yuan H, Wu B, Chen B, Tan Y, et al. Investigation on the Association of Soil  
723 Microbial Populations with Ecological and Environmental Factors in the Pearl River Estuary. *Journal of*  
724 *Geoscience and Environment Protection* 2018;6:8.
- 725 [80] Persson S, Al-Shuweli S, Yapici S, Jensen JN, Olsen KE. Identification of Clinical *Aeromonas*  
726 *Species by rpoB and gyrB Sequencing and Development of a Multiplex PCR Method for Detection of*  
727 *Aeromonas hydrophila, A. caviae, A. veronii, and A. media.* *J Clin Microbiol* 2015;53:653–6.
- 728 [81] Fleischmann W, Gateau A, Apweiler R. A novel method for automatic functional annotation of  
729 proteins. *Bioinformatics* 1999;15:228–33.

- 730 [82] Pedruzzi I, Rivoire C, Auchincloss AH, Coudert E, Keller G, De Castro E, et al. HAMAP in 2015:  
731 updates to the protein family classification and annotation system. *Nucleic acids research*  
732 2015;43:D1064–D70.
- 733 [83] Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-  
734 BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–402.
- 735 [84] Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, et al. Ensembl comparative  
736 genomics resources. *Database (Oxford)* 2016;2016.
- 737 [85] Notredame C, Higgins DG, Heringa J. T-coffee: a novel method for fast and accurate multiple  
738 sequence alignment1. *J Mol Biol* 2000;302:205–17.
- 739 [86] Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*  
740 *Acids Res* 2004;32:1792–7.
- 741 [87] Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity of progressive  
742 multiple sequence alignment through sequence weighting, position-specific gap penalties and weight  
743 matrix choice. *Nucleic Acids Res* 1994;22:4673–80.
- 744 [88] Emanuelsson O, Brunak S, Von Heijne G, Nielsen H. Locating proteins in the cell using TargetP,  
745 SignalP and related tools. *Nat Protoc* 2007;2:953.
- 746 [89] Krogh A, Larsson B, Von Heijne G, Sonnhammer EL. Predicting transmembrane protein topology  
747 with a hidden markov model: application to complete genomes1. *J Mol Biol* 2001;305:567–80.
- 748 [90] Julenius K, Mølgaard A, Gupta R, Brunak S. Prediction, conservation analysis, and structural  
749 characterization of mammalian mucin-type O-glycosylation sites. *Glycobiology* 2005;15:153–64.
- 750 [91] Monigatti F, Gasteiger E, Bairoch A, Jung E. The Sulfinator: predicting tyrosine sulfation sites in  
751 protein sequences. *Bioinformatics* 2002;18:769–70.
- 752 [92] Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, et al. InterPro in 2017—beyond  
753 protein family and domain annotations. *Nucleic Acids Res* 2016;45:D190–D9.
- 754 [93] Andrade MA, Ponting CP, Gibson TJ, Bork P. Homology-based method for identification of protein  
755 repeats using statistical significance estimates. *J Mol Biol* 2000;298:521–37.
- 756 [94] NCBI RC. Database resources of the National Center for Biotechnology Information. *Nucleic Acids*  
757 *Res* 2016;44:D7.
- 758 [95] Müller H-M, Kenny EE, Sternberg PW. Textpresso: an ontology-based information retrieval and  
759 extraction system for biological literature. *PLoS biology* 2004;2:e309.
- 760 [96] Veuthey A-L, Bridge A, Gobeil J, Ruch P, McEntyre JR, Bougueleret L, et al. Application of text-  
761 mining for updating protein post-translational modification annotation in UniProtKB. *BMC*  
762 *Bioinformatics* 2013;14:104.

- 763 [97] Wei C-H, Kao H-Y, Lu Z. PubTator: a web-based text mining tool for assisting biocuration. *Nucleic*  
764 *Acids Res* 2013;41:W518–W22.
- 765 [98] Chibucos MC, Mungall CJ, Balakrishnan R, Christie KR, Huntley RP, White O, et al. Standardized  
766 description of scientific evidence using the Evidence Ontology (ECO). *Database (Oxford)*  
767 2014;2014:bau075.
- 768 [99] Choi M, Liu H, Baumgartner W, Zobel J, Verspoor K. Coreference resolution improves extraction of  
769 Biological Expression Language statements from texts. *Database (Oxford)* 2016;2016:baw076.
- 770 [100] Peng Y, Wei C-H, Lu Z. Improving chemical disease relation extraction with rich features and  
771 weakly labeled data. *J Cheminform* 2016;8:53.
- 772 [101] Harding A. Rise of the Bio-librarian: the field of biocuration expands as the data grows. *Scientist*  
773 2006;20:82–4.
- 774 [102] Bourne PE, McEntyre J. Biocurators: contributors to the world of science. *PLoS Comput Biol*  
775 2006;2:e142.
- 776 [103] Bateman A. Curators of the world unite: the International Society of Biocuration. *Bioinformatics*  
777 2010.
- 778 [104] Mitchell CS, Cates A, Kim RB, Hollinger SK. Undergraduate biocuration: developing tomorrow's  
779 researchers while mining today's data. *J Undergrad Neurosci Educ* 2015;14:A56–A65.
- 780 [105] Reiser L, Berardini TZ, Li D, Muller R, Strait EM, Li Q, et al. Sustainable funding for biocuration:  
781 The Arabidopsis Information Resource (TAIR) as a case study of a subscription-based funding model.  
782 *Database (Oxford)* 2016;2016:baw018.
- 783 [106] Karp PD. How much does curation cost? *Database (Oxford)* 2016;2016:baw110.
- 784 [107] Hayden E. Funding for model-organism databases in trouble. *Nature* 2016.
- 785 [108] Kaiser J. Funding for key data resources in jeopardy. *Science* 2016;351:14.
- 786 [109] Bourne PE, Lorsch JR, Green ED. Perspective: Sustaining the big-data ecosystem. *Nature*  
787 2015;527:S16–S7.
- 788

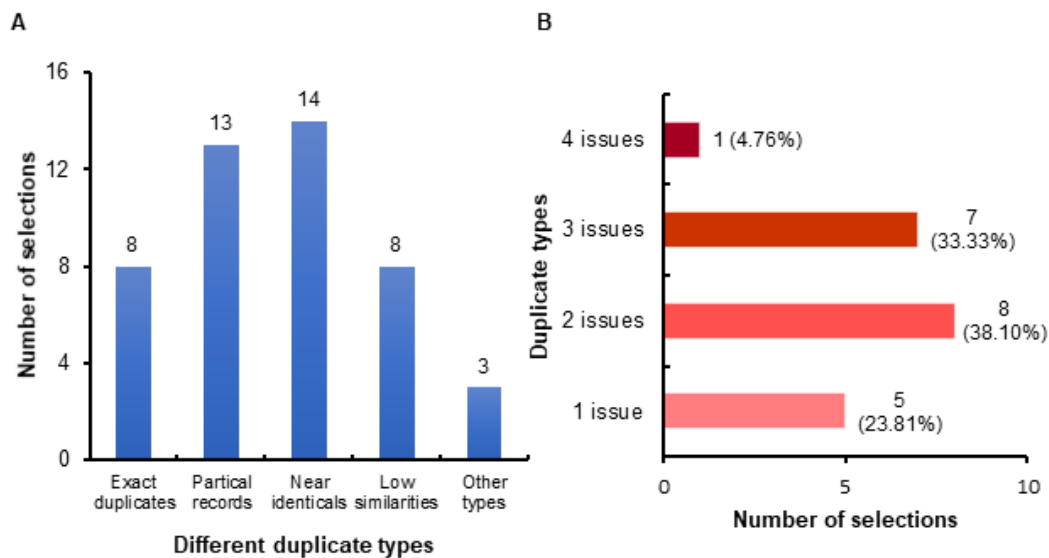
## 789 Figures and Tables



790

### 791 Figure 1 Biological analysis pipeline

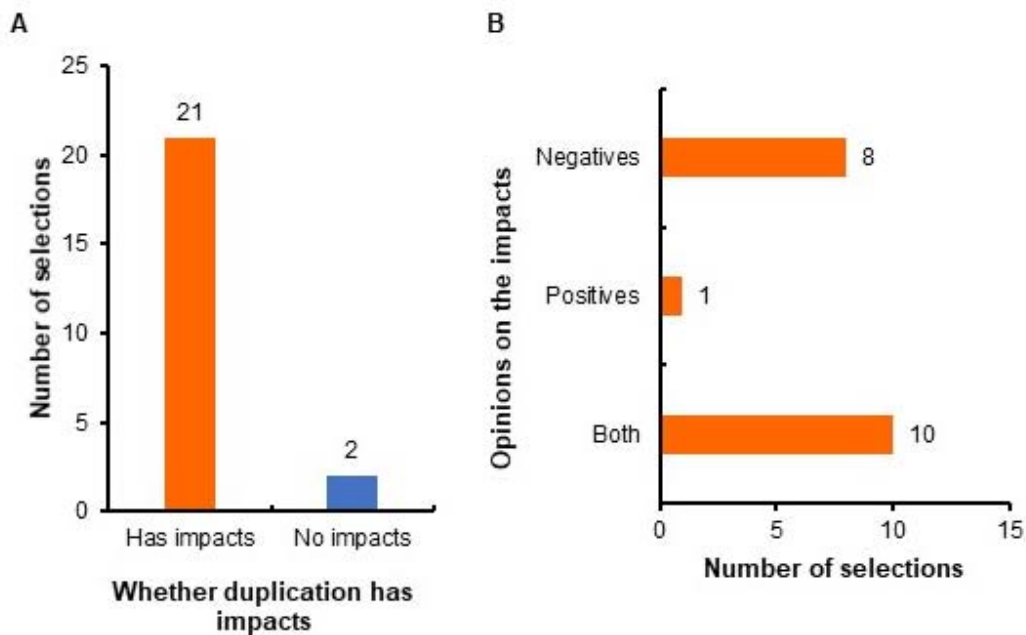
792 Three stages of a biological analysis pipeline, heavily involving biological databases, are  
793 presented.



794

### 795 Figure 2 Characteristics of duplicate records

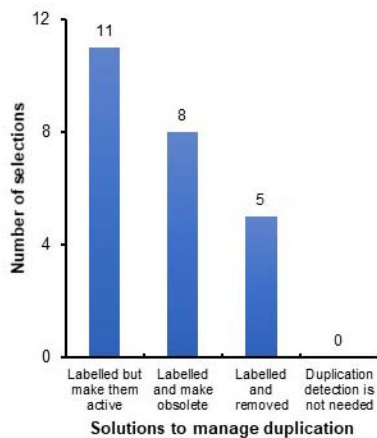
796 What are duplicates? The X-axis shows different duplicate types; the Y-axis shows the  
797 associated number of participants who selected that type.



798

799 **Figure 3 Impacts of duplicate records**

800 A. Do duplicates have impacts? The number of participants who believed whether duplication  
801 has impacts or not is shown. B. a more detailed breakdown by type of impact, for those who  
802 believed duplication has impacts, is illustrated.



803

804

805 **Figure 4 Solutions to duplicate records**

806 How to address duplication? The X-axis represents the options to address duplication; the Y-axis  
 807 represents the corresponding number of participants selected that option.

808

809 **Table 1 Representative software and resources used in expert curation**

<b>Curation steps</b>	<b>Software/ Databases</b>	<b>Purpose</b>	<b>Ref.</b>
<b>Sequence curation</b>			
Identify homologs	BLAST	Sequence alignment	[83]
	Ensembl	Phylogenetic resources	[84]
Document inconsistencies	T-Coffee	Analysis of causes of	[85]
	Muscle	inconsistencies due to	[86]
	ClustalW	duplication	[87]
<b>Sequence analysis</b>			
Predict topology	Signal P	Signal peptides prediction	[88]
	TMHMM	Transmembrane domain prediction	[89]
Post-translations	NetNGlyc	N-glycosylation sites prediction	[90]
	Sulfinator	Tyrosine sulfation sites prediction	[91]
Identify domains	InterPro	Retrievals of motif matches	[92]
	REPEAT	Identification of repeats	[93]
<b>Literature curation</b>			
Identify relevant literature	PubMed	Literature resources	[94]
	iHOP		[95]
Text mining	PTM	Information extraction	[96]
	PubTator		[97]
Assign GOs	GO	Gene ontology terms	[14]
<b>Family curation</b>	Same as identify homologs		
<b>Evidence attribution</b>	ECO	Evidence code ontology	[98]

810 *Note:* A complete set of the software, including the detailed versions of the software, can be  
811 found in UniProt manual curation standard operating procedure documentation  
812 ([www.uniprot.org/docs/sop\\_manual\\_curation.pdf](http://www.uniprot.org/docs/sop_manual_curation.pdf)).

813

814

## 815 **Supplementary material**

816 **File S1 Survey questions**

817 **File S2 Results and discussions on quality issues beyond duplication**

818 **File S3 UniProtKB/Swiss-Prot annotation team leader interview details**