

1

2

3

## 4 **Measuring the temporal quality of a biodiversity database**

5

6 Les G Underhill

7

8 Department of Biological Sciences, University of Cape Town, Rondebosch, South Africa

9

### 10 **Introduction**

11

12 Sutherland et al. (2015) developed a set of 10 priorities for biological recording, with a focus  
13 on monitoring biodiversity through citizen science. This paper develops an 11th priority  
14 which, though not made explicit by Sutherland et al. (2015), was probably in mind. In the  
15 context of climate change, it is necessary to articulate the conscious and continuous need to  
16 keep biodiversity databases up-to-date; a corollary to this priority is the need for a metric to  
17 assess the concept of up-to-dateness. The objective of this paper is to meet this need to  
18 quantify up-to-dateness, and thus provide one measure of the temporal quality of a biological  
19 database.

20

21 The context of this paper is the ongoing collection of “atlas”-type data for a taxon, used to  
22 generate on-demand, up-to-date distribution map for any species in that taxon. Such maps are  
23 of fundamental value in establishing conservation priorities (Underhill & Gibbons 2002,  
24 Harrison et al. 2008). For example, the Second Southern African Bird Atlas Project began in

25 2007 to serve as a five-year “snapshot” of bird distributions, as had been the objective of the  
26 initial bird atlas. Instead, the project remains ongoing in 2021, having morphed from the aim  
27 of capturing a “snapshot” to recording a “video” of changing distributions (Harrison et al.  
28 2008, Underhill 2016, Underhill et al. 2017). The key idea is that up-to-date distribution maps  
29 should be built on recent data. Representing the occurrence of a species at a location with an  
30 “old” record is unsatisfactory, and even misleading.

31

32 Therefore, the focus of this paper is to devise an algorithm which, in some way, measures the  
33 decay in temporal quality of biodiversity data. The algorithm attempts to answer the question:  
34 “how up-to-date is this collection of biodiversity records?”. The algorithm is illustrated using  
35 the South African component of the Biodiversity and Development Institute’s Virtual  
36 Museum database, which includes South Africa, Lesotho and eSwatini; however, the method  
37 is readily adapted to other contexts.

38

## 39 **Methods**

40

### 41 *Algorithm*

42

43 A biodiversity record is comprised of three components: species, place and date. The record  
44 provides evidence that a particular species was recorded at the place on the date. The  
45 algorithm proposed here requires each data point to be allocated to a “grid cell”. For the  
46 Virtual Museum, a 15-minute geographical grid is used, generating what are popularly known  
47 as quarter-degree grid cells (even though there are 16 of them in a one-degree grid cell). For  
48 the African Bird Atlas Project, a five-minute grid is in use, generating grid cells known as  
49 pentads (Underhill 2016). North of c. 40°N, geographical grid cells are no longer viable,

50 because the east-west distance is far smaller than the north-south distance. The convention  
51 then is to use grid cells measured in kilometres, as is done in the bird atlas projects of Europe  
52 (e.g. Hagemeijer & Blair 1997). In Britain, early projects made use of “counties” and “vice-  
53 counties”; the proposed method works with any spatial units, provided they are standardized  
54 through time.

55

56 For each record, the algorithm needs the three components: the date, the species, and the grid  
57 cell to which it belongs. The algorithm itself has four steps:

58

- 59 1. For each grid cell, find the most recent date for each species;
- 60 2. Find the median of these dates, one per species (this date provides a measure of the  
61 temporal quality of the data for the grid cell);
- 62 3. Find the median dates for all grid cells in the region under consideration;
- 63 4. Calculate the median of the grid cell medians (this date, the median of the medians, is  
64 defined to be the temporal quality of the data in the region, i.e. up-to-dateness).

65

66 The algorithm is simple and transparent. At Step 2, the median date calculated for a grid cell  
67 is the halfway point for the most recent dates for the species recorded in that grid cell. At  
68 Step 4, the halfway point for all the grid cells in the region is calculated. The median dates of  
69 half of the grid cells fall before this date, and half fall later.

70

71 The algorithm can be applied not only to the current database, but also to historic databases  
72 created by only considering records submitted prior to a chosen date. The crucial statistic  
73 then becomes the time interval, expressed in months or years, between the median of the  
74 medians and the chosen date. If this interval is shortening through time, then the temporal

75 quality of the database is improving, and *vice versa*. The trend in this statistic through time is  
76 a powerful communication tool.

77

78

## 79 *Application*

80

81 The algorithm was applied to the databases of six of the projects within the Virtual Museum.  
82 The Virtual Museum is a citizen science initiative which assembles photographic records for  
83 a series of taxa. It was originally developed for citizen scientists to contribute to reptile and  
84 butterfly atlases in South Africa (Bates et al. 2014, Mecenero et al. 2013), and gradually  
85 expanded to cover a selection of other taxa. The databases for reptiles, butterflies, dragonflies  
86 and damselflies (Underhill et al. 2016), lacewings (Underhill et al. 2019), birds and frogs  
87 were evaluated. The algorithm used all records up to 22 June 2021. The “popular” names for  
88 the six databases are included in Table 1.

89

90 To illustrate the interpretation of the pattern in trends through time, the algorithm was also  
91 applied to dragonflies and damselflies records (OdonataMAP) on an annual basis over the  
92 decade since the project commenced in September (Underhill et al. 2016). To assess whether  
93 the rate of renewal of old records resulted in distribution maps becoming more or less up-to-  
94 date through time, the algorithm was applied to observations uploaded by the ends of the six  
95 calendar years from 2015 to 2020.

96

97

## 98 Results

99

100 Of the six Virtual Museum projects for which up-to-dateness was evaluated in June 2021, it is  
101 clear that the distribution maps produced using the ReptileMAP data cannot be considered  
102 “up-to-date”. For half of the grid cells, the median date of the most recent record was 41  
103 years prior to the date of the analysis. Similarly, the LacewingMAP database was 25 years  
104 out of date, the FrogMAP database 21 years out of date, the LepiMAP database 12 years out  
105 of date, and the OdonataMAP database 4.3 years out of date. The BirdPix database was  
106 classified as 2.4 years, out of date (Table 1).

107

108 Table 1. Summary statistics, including up-to-dateness, for six projects within the Virtual Museum, in  
109 the 2,027 quarter degree grid cells in South Africa, Lesotho and eSwatini. Calculations were  
110 performed on 25 June 2021, and the up-to-dateness is calculated as the difference between the median  
111 date and 25 June 2021.

112

Project	Grid cells	Species	Grid cell-species	Records	Median date	Up-to-dateness (years)
BirdPix	1,479	843	54,935	143,207	05 Jan 2019	2.4
FrogMAP	1,763	132	15,253	50,545	17 Mar 2000	21
LacewingMAP	837	408	5,918	11,259	20 Feb 1996	25
LepiMAP (butterflies only)	1,698	886	61,694	469,852	08 Jan 2009	12
OdonataMAP	1,126	162	17,437	108,842	09 Mar 2017	4.3
ReptileMAP	1,946	478	35,926	149,924	21 Sep 1980	41

113

114 For OdonataMAP, the pattern of up-to-dateness reflects the typical behaviour anticipated as a  
 115 project evolves (Table 2). Over the first years to the end of 2013, coverage expanded rapidly,  
 116 and up-to-dateness remained constant. Thereafter, without an incentive to “refresh” records  
 117 of a species in a grid cell the up-to-dateness steadily moves backwards. Since 2016, a large  
 118 number of historical museum specimen records were included in the analysis, causing a 24-  
 119 month retrogression in up-to-dateness. In the subsequent four years from the end of 2016 to  
 120 the end of 2020, the up-to-dateness of the OdonataMAP database had slipped by four months,  
 121 from 46 months at the end of 2016 to 50 months at the end of 2020 (Table 2). At the end of  
 122 2020, for half of the 1098 grid cells visited, the median date of the most recent observation of  
 123 a species was after 10 November 2016, and before this date for the other half.

124

125 Table 2. Summary statistics, including up-to-dateness, for the OdonataMAP project of the Virtual  
 126 Museum in the 2,027 quarter degree grid cells of South Africa, Lesotho and eSwatini. These statistics  
 127 are calculated from the database using the records that had been uploaded by the end of each of the six  
 128 calendar years from 2011 to 2020.

Period	Grid cells	Species	Grid cell-species	Cumulative records	Median date	Up-to-dateness (months)
Up to 2011	99	81	466	534	04 Dec 2010	13
Up to 2012	242	118	1,344	2,332	12 Dec 2011	13
Up to 2013	395	133	2,863	6,001	06 Jan 2013	12
Up to 2014	529	145	4,438	9,591	05 May 2013	20
Up to 2015	679	147	5,954	14,994	22 Feb 2014	22
Up to 2016	852	159	11,620	36,900	17 Mar 2013	46
Up to 2017	925	160	13,227	48,441	09 Feb 2014	47
Up to 2018	1,008	160	14,656	65,759	06 Feb 2015	47

Up to 2019	1,062	162	15,739	82,497	14 Dec 2015	49
Up to 2020	1,098	162	16,556	98,505	10 Nov 2016	50

129

130

## 131 **Discussion**

132

### 133 *Multiple measures of quality*

134

135 When applied to a biodiversity database, the term “data quality” is frequently used  
136 qualitatively rather than quantitatively. For example, Wetzel et al. (2018) repeatedly used the  
137 word “quality” in describing the needs of biodiversity data in Europe, but never clearly  
138 defined the term. There is a need to quantify the concept of “quality.” In this context,  
139 “quality” is a multi-dimensional concept; there is no single measure of quality.

140

141 Up-to-dateness, as presented here, is only one such measure of quality. It must be supported  
142 by supplementary information which captures other dimensions of quality, i.e. the total  
143 number of records within the region, the number of species represented, the number of grid  
144 cells with records, and the number of grid cell-species records (the sum of the number of  
145 species recorded in each grid cell). A further family of data quality refer to taxonomic issues:  
146 the up-to-dateness of the taxonomy used for the database, and the percentage of records not  
147 correctly identified to species level.

148

149 There are several qualitative dimensions implicit in the concept of biodiversity data quality.  
150 The most frequently employed relates to gaps in coverage, i.e. spatial quality. Data are of  
151 poor quality if they display gaps. These are either geographical gaps, areas for which little or

152 no biodiversity data exists for an entire taxon; or gaps in the range map for an individual  
153 species, i.e. places where the species probably occurs, but has not yet been recorded (false  
154 negatives). A second dimension is yet another measure of temporal quality, which is  
155 generally interpreted to mean that there are data for an extended time period, usually  
156 measured in years or decades. In this context, biodiversity data for a region are said to have  
157 poor temporal quality if they lack historical records, making it difficult to examine trends  
158 through time.

159  
160 For example, with data from 1,946 of 2,027 grid cells, ReptileMAP has the best spatial  
161 coverage (96%) of the six projects (Table 1). More nuanced measures of spatial quality  
162 would need to 1) account of whether gaps in coverage tend to be in adjacent or scattered grid  
163 cells, and 2) estimate the percentage of false negatives in the database. For 2), determining  
164 the number which goes into the numerator when estimating this percentage is  
165 straightforward; it is the sum of the numbers of species in each grid cell, 35,926 in the case of  
166 ReptileMAP (Table 1). The denominator for the percentage requires an estimate of the  
167 number of grid cells in the range of each species, a quantity which the atlas project aims to  
168 determine. The number 35,926 is the total of the total number of grid cells in which the 478  
169 reptile species have been recorded, and hence shaded in the distribution maps (Table 1).

170  
171 Table 2 details improvements of other measures of quality in OdonataMAP data through  
172 time. A set of distribution maps for 147 species made at the end of 2015 would have been  
173 based on 14,994 records from 679 grid cells. In these maps, the total number of grid cells  
174 which would have shown a species as present would have been 5,954 (Table 2). One year  
175 later, after the inclusion of the historical data, maps for 159 species would have been based  
176 on 36,900 records from 852 grid cells, and the total number of shaded grid cells would have



177 almost doubled to 11,620 (Table 2). The improved spatial coverage was obtained by  
178 sacrificing up-to-dateness (Table 2). By June 2021, the set of 162 maps including the  
179 additional records collected from 2017 onwards would have had a total of 17,437 grid cells  
180 showing a species as present (Table 1). This is an increase of 50% since the end of 2016,  
181 providing another useful measure of the improvement of the spatial quality of the  
182 OdonataMAP database in South Africa, Lesotho and eSwatini in 4.5 years.

183

184 Thus quality, in this context, is a multidimensional concept. Temporal quality, as developed  
185 here, is only one dimension in describing the value of the database.

186

187 *An alternative approach to measure the temporal decay in quality*

188

189

190 A biodiversity record provides evidence that a particular species was recorded at the place on  
191 the date. However, biodiversity literature seems to pay little or no attention to the reality that  
192 the value of a record steadily decreases as its date of occurrence recedes into the past. Aging  
193 records slowly but steadily become less and less valuable as evidence that a species still  
194 occurs at a given site.

195

196 In this paper, the concept of temporal quality of a biodiversity database is introduced, rooted  
197 in the notion that data quality decays over time. However, the shape of the function which  
198 describes the decay in value through time remains unknown, and needs to be quantified.

199

200 Here is an experimental approach to achieve this. Consider a substantially-sized sample (say  
201  $n = 1000$ ) of correctly georeferenced biodiversity records with accurate dates. The sample

202 needs to be stratified geographically and by date, possibly from one year ago to 100 years  
203 ago. Revisit the sites at which the records were made on the calendar date of the original  
204 record, and search for the species. Three outcomes are possible: (1) the species was recorded,  
205 and still occurs at the site, (2) there is still suitable habitat for the species, but it was not  
206 found, and (3) the site has been transformed to such an extent that the species almost  
207 certainly no longer occurs. The data analysis should aim to develop a function which  
208 describes the average rate of decay in value of biodiversity records, i.e. to estimate the “half-  
209 life” of a record. It is likely that these decay functions vary between species. One objective,  
210 probably unattainable, would be to establish a gold standard definition for the age at which a  
211 record no longer provides suitable evidence that a species still occurs at a locality.  
212 Distribution maps for a species could then exclude records older than this date, or plot them  
213 differently to indicate regions in which the species has either actually gone locally extinct, or  
214 where further search effort is needed to refresh evidence of its presence. Alternatively,  
215 distribution maps could weight records by their age, so that older records have smaller  
216 weights than newer ones.

217

218 The decay function also opens up new possibilities of developing a more nuanced measure of  
219 up-to-dateness than the one developed here; for instance, using the last recorded date for each  
220 species in a grid cell, one could apply the decay function to estimate the remaining evidential  
221 value of the record, and calculate appropriate summary statistics.

222

### 223 *Use of up-to-dateness to motivate citizen scientists*

224

225 As demonstrated above for OdonataMAP (Table 2), the algorithm is not only applicable to  
226 the current database. It can also be applied to generate trends from historical databases,

227 created by only considering records submitted up to a certain point in time. When applied in  
228 this way, these trends may serve as motivational guidance for project leaders and citizen  
229 scientists.

230

231 For instance: “In this era of rapid development and climate change, any record of a species at  
232 a locality which is more than two years old needs to be refreshed with a new one.” The  
233 awareness of last recorded dates for the species in a grid cell is in itself a powerful  
234 motivational force for citizen scientists. There is a real challenge in seeking to detect the  
235 species that have the most-out-of-date last recorded dates in a grid cell.

236

237 The background to the six projects of Table 1 provides insights into the extent of their out-of-  
238 dateness. This information must be considered in communications with the citizen scientists  
239 participating in each project. Three of these projects (FrogMAP, LepiMAP and ReptileMAP)  
240 are continuations of citizen science projects which produced published atlases for frogs  
241 (Minter et al. 2004), butterflies (Mecenero et al. 2013) and reptiles (Bates et al. 2014)  
242 respectively. In each case, the foundational data for the project were drawn from historical  
243 databases in museums, other non-museum specimen collections, and literature. The explicit  
244 objective of the citizen science fieldwork, which took place over five-year periods, was to fill  
245 in distribution gaps with a goal of targeting false negatives, grid cells where the species likely  
246 occurred but had not yet been recorded. This mentality of “filling in the gaps,” has prevailed  
247 in the continuation of these projects within the Virtual Museum. Thus, many observers do not  
248 upload records to the Virtual Museum if the species has already been recorded in the grid  
249 cell. The consequence is that the up-to-dateness of the databases slowly deteriorates.  
250 Changing this outlook has proved challenging.

251

252 OdonataMAP is a relatively new project (Underhill et al. 2016), and the cohort of citizen  
253 scientists involved are relatively strongly attuned to the importance of repeated submission of  
254 the same species from the same locality. To a large measure, this is because the project  
255 includes a focus on generating data from which the phenology of adult dragonfly and  
256 damselfly occurrence may be estimated. A by-product of this approach is that overall up-to-  
257 dateness of the database deteriorated by only four months over four years (Table 2).

258

259 The foundational data for LacewingMAP were developed by Mervyn Mansell over a career;  
260 they contain the global specimen database for the Neuroptera and Megaloptera, with records  
261 dating back to the nineteenth century (Underhill et al. 2019). The citizen science database is  
262 built on this platform. Contributions are opportunistic; no citizen scientist has a primary  
263 interest in this taxon (Underhill et al. 2019). The rate of record submission is thus relatively  
264 slow (but far faster than the rate of specimen acquisition in the decades when museum  
265 collections were growing fastest (Underhill et al. 2019)). Thus, the fact that the database is 25  
266 years out of date is not surprising (Table 2).

267

268 In contrast, although the BirdPix project of the Virtual Museum was started in 2012, active  
269 promotion of the project commenced in 2017. By then a large number of grid cells had  
270 received a small number of records each. Many citizen scientists have uploaded series of  
271 historical photos taken in the early days of digital photography, or have even scanned and  
272 submitted slides; this has effectively pushed the median of medians further into the past. For  
273 the other projects within the Virtual Museum discussed here, most historical image  
274 collections (for example of reptiles and butterflies) were uploaded during the formal atlas  
275 periods for each taxon. However, in spite of these challenges, the BirdPix database was 2.4  
276 years out of date in June 2021.

277

278 The continuous presentation of up-to-dateness, preferably in graphical form, provides an  
279 incentive to citizen scientists to keep the database up-to-date. This is an example of the  
280 application of gamification to motivate data collection by citizen scientists (Ainsley &  
281 Underhill 2017). Gamification is the development of built-in strategies to encourage project  
282 participation (in particular, it does not mean turning data-collection into a “game”) (Ainsley  
283 & Underhill 2017).

284

285 For a citizen scientist to contribute towards bringing the median date for a grid cell closer to  
286 the present, a straightforward strategy is to increase the number of species for that grid cell,  
287 especially if it is still relatively easy to add species; each added species then has a current  
288 date and moves the median for the grid cell forward. Similarly, to bring the median of  
289 medians closer to the present would require undertaking fieldwork in grid cells which have  
290 no data. This increases the number of grid cells with median dates; the new median dates are  
291 current, and move the median of medians forward. Both strategies not only improve up-to-  
292 dateness, they also positively impact other measures of completeness for the database (i.e.  
293 spatial coverage).

294

### 295 ***Local extinction***

296

297 Clearly, if a species goes extinct in a grid cell, its last recorded date can no longer be updated.  
298 At some point, decisions about local extinction need to be made, so that a species can be  
299 removed from the list over which the median for the grid cell is calculated. Mechanisms to do  
300 this need to be devised, but will almost certainly be a combination of quantitative and  
301 qualitative arguments. Species declared locally extinct need to remain on the species list, and

302 flagged appropriately. The awareness among citizen scientists of local extinctions has the  
303 potential to lead to civic awareness so that the impact of the project transcends science into  
304 grassroots action and policy making (Loos et al. 2015).

305

### 306 *Caveats to the measure of up-to-dateness*

307

308 The measure of up-to-dateness described here is conditional on two factors: first, it only takes  
309 account of grid cells for which data are available; and second, it only takes account of species  
310 which have already been recorded in the grid cell. An absolute measure of up-to-dateness  
311 would also account for the grid cells lacking data, and include an estimate of the number of  
312 species in each grid cell.

313

314 The choice of the conditional measure was deliberate. Primarily, it is designed to measure the  
315 up-to-dateness of the data already collected; other aspects of database quality can be  
316 evaluated using other statistics, such as spatial coverage measures. Secondly, the choice of  
317 conditional measure facilitates the use of gamification, as described above.

318

### 319 **Acknowledgements**

320

321 Itxaso Quintana, Magda Remisiewicz, Karis Daniel, Johan van Rooyen, Greg Distiller and  
322 David Thomson made important comments on earlier drafts. The paper owes its existence to  
323 the thousands of citizen scientists who have participated in the Virtual Museum.

324

325

326

327 **References**

328

329 Ainsley J, Underhill LG. Gamification (persuasive design) in the Second Southern African  
330 Bird Atlas Project (SABAP2). *Vogelwelt* 2017;137:19-22.

331

332 Bates MF, Branch WR, Bauer AM, Burger M, Marais J, Alexander GJ, de Villiers MS,  
333 editors. 2014. Atlas and Red List of the reptiles of South Africa, Lesotho and Swaziland.  
334 Pretoria: South African National Biodiversity Institute; 2014.

335

336 Hagemeyer EJM, Blair MJ, editors. The EBCC atlas of European breeding birds: Their  
337 distribution and abundance. London: T & A D Poyser; 1997.

338

339 Harrison JA, Underhill LG, Barnard P. The seminal legacy of the Southern African Bird  
340 Atlas Project. *South African Journal of Science* 2008;102:82-84.

341

342 Loos J, Horcea-Milcu AI, Kirkland P, Hartel T, Osváth-Ferencz M, Fischer J. Challenges for  
343 biodiversity monitoring using citizen science in transitioning social-ecological systems.  
344 *Journal for Nature Conservation* 2015;26:45-48.

345

346 Mecenero S, Ball JD, Edge DA, Hamer ML, Henning GA, Krüger M, et al., editors.  
347 Conservation assessment of butterflies of South Africa, Lesotho and Swaziland: Red List  
348 and atlas. Johannesburg: Safronics and Cape Town: Animal Demography Unit; 2013.

349 Minter LR, Burger M, Harrison JA, Braack HH, Bishop PJ, Kloepfer D, editors. Atlas and  
350 Red Data book of the frogs of South Africa, Lesotho and Swaziland. SIMAB Series #9.

351 Washington, DC: Smithsonian Institution and Cape Town: Avian Demography Unit;  
352 2004.  
353  
354 Sutherland WJ, Roy DB, Amano T. An agenda for the future of biological recording for  
355 ecological monitoring and citizen science. *Biological Journal of the Linnean Society*  
356 2015;115:779-784.  
357  
358 Underhill LG. The fundamentals of the SABAP2 protocol. *Biodiversity Observations*  
359 2016;7.42:1-12.  
360  
361 Underhill LG, Brooks M, Loftie-Eaton M. The Second Southern African Bird Atlas Project:  
362 protocol, process, product. *Vogelwelt* 2017;137:64-70.  
363  
364 Underhill LG, Gibbons D. Mapping and monitoring bird populations: their conservation uses.  
365 In: Norris K, Pain DJ, editors. *Conserving bird biodiversity: general principles and their*  
366 *application*. Cambridge: Cambridge University Press; 2002. p. 34-60.  
367  
368 Underhill LG, Navarro R, Mansell M. LacewingMAP: Progress report on the atlas of African  
369 Neuroptera and Megaloptera, 2014–2019. *Biodiversity Observations* 2019;10.10:1-21.  
370  
371 Underhill LG, Navarro R, Manson AD, Labuschagne JP, Tarboton WR. OdonataMAP:  
372 progress report on the atlas of the dragonflies and damselflies of Africa, 2010-2016.  
373 *Biodiversity Observations* 2016;7.47:1-10.  
374



- 375 Wetzel FT, Bingham HC, Groom Q, Haase P, Kõljalg U, Kuhlmann M, et al. Unlocking  
376 biodiversity data: Prioritization and filling the gaps in biodiversity observation data in  
377 Europe. *Biological Conservation* 2018;221:78-85.