

¹ DateLife: leveraging databases and analytical tools to reveal the dated Tree of Life

² Luna L. Sánchez Reyes^{1,2}, Emily Jane McTavish¹, & Brian O'Meara²

³ ¹ University of California, Merced

⁴ ² University of Tennessee, Knoxville

⁵ Author Note

6 School of Natural Sciences, University of California, Merced, Science and Engineering
7 Building 1.

8 Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville,
9 425 Hesler Biology Building, Knoxville, TN 37996, USA.

10 The authors made the following contributions. Luna L. Sánchez Reyes: Data curation,
11 Investigation, Software, Visualization, Validation, Writing - Original Draft Preparation,
12 Writing - Review & Editing; Emily Jane McTavish: Resources, Software, Writing - Review &
13 Editing; Brian O'Meara: Conceptualization, Funding acquisition, Methodology, Resources,
14 Software, Supervision, Writing - Review & Editing.

15 Correspondence concerning this article should be addressed to Luna L. Sánchez Reyes, .
16 E-mail: sanchez.reyes.luna@gmail.com

17

Abstract

18 Date estimates for times of evolutionary divergences are key data for research in the natural
19 sciences. These estimates also provide valuable information for education, science
20 communication and policy decisions. Although achieving a high-quality reconstruction of a
21 phylogenetic tree with branch lengths proportional to absolute time (chronogram), is a
22 difficult and time-consuming task, the increased availability of fossil and molecular data, and
23 time-efficient analytical techniques has resulted in many recent publications of large
24 chronograms for a large number and wide diversity of organisms. When these estimates are
25 shared in public, open databases this wealth of expertly-curated and peer-reviewed data on
26 time of evolutionary origin is exposed in a programmatic and reusable way. Intensive and
27 localized efforts have improved data sharing practices, as well as incentivized open science
28 in biology. Here we present DateLife, a service implemented as an R package and an Rshiny
29 website application available at www.datelife.org/query/, that provides functionalities for
30 efficient and easy finding, summary, reuse, and reanalysis of expert, peer-reviewed, public
31 data on time of evolutionary origin. The main DateLife workflow constructs a chronogram
32 for any given combination of taxon names, by searching a local chronogram database
33 constructed and curated from the Open Tree of Life Phylesystem phylogenetic database,
34 which incorporates phylogenetic data from TreeBASE database as well. We implement and
35 test methods for summarizing time data from multiple source chronograms using supertree
36 and congruification algorithms, and using age data extracted from source chronograms as
37 secondary calibration points to add branch lengths proportional to absolute time to a tree
38 topology. DateLife will be useful to increase awareness on the existing variation in expert
39 time of divergence data, and can foster exploration of the effect of alternative divergence
40 time hypothesis on the results of analyses, providing a framework for a more informed
41 interpretation of evolutionary results.

42

Keywords: Tree; Phylogeny; Scaling; Dating; Ages; Divergence times; Open Science;

⁴³ Congruification; Supertree; Calibrations; Secondary calibrations

⁴⁴ Word count: 4201

45 DateLife: leveraging databases and analytical tools to reveal the dated Tree of Life

46 **Introduction**

47 Chronograms –phylogenies with branch lengths proportional to time– provide key data
48 for the study of natural processes in many areas of biological research, such as developmental
49 biology (Delsuc et al., 2018; Laubichler & Maienschein, 2009), conservation biology
50 (Felsenstein, 1985; C. Webb, 2000), historical biogeography (Posadas, Crisci, & Katinas,
51 2006), and species diversification (Magallon & Sanderson, 2001; Morlon, 2014).

52 Building a chronogram is not an easy task. It requires obtaining and curating data to
53 construct a phylogeny; selecting and placing appropriate calibrations on the phylogeny using
54 independent age data points from the fossil record or other dated events, and inferring the
55 full dated tree. Estimating accurate chronograms generally requires specialized biological
56 training, taxonomic domain knowledge, and a non-negligible amount of research time,
57 computational resources and funding.

58 Here we present the DateLife software application, available as an R package and as an
59 online Rshiny interactive website at www.datelife.org/query/, which captures data from
60 published chronograms, and make these data readily accessible to users. DateLife features a
61 versioned, open and fully public chronogram database (McTavish et al., 2015) storing age
62 information in a computer readable format (Vos et al., 2012), an automated and
63 programmatic way of accessing the data (Stoltzfus et al., 2013) and methods to summarize
64 and compare age data.

65 **Description**

66 The DateLife algorithm is fully implemented using the R language. The latest stable
67 version of the R package **datelife** is available from the CRAN repository (v0.6.2;
68 Sanchez-Reyes et al. (2022)), and relies on functionalities from various biological R packages:
69 **ape** (Paradis, Claude, & Strimmer, 2004), **bold** (Chamberlain et al., 2019), **geiger** (Harmon,

70 Weir, Brock, Glor, & Challenger, 2008), paleotree (Bapst, 2012), phyloch (Heibl, 2008),
71 phylocomr (Ooms & Chamberlain, 2018), phytools (Revell, 2012), rotl (Michonneau, Brown,
72 & Winter, 2016), and taxize (Chamberlain & Szöcs, 2013; Chamberlain et al., 2019). Figure
73 1 provides a graphical summary of the three main steps of the DateLife algorithm: providing
74 an input, searching a chronogram database, and summarizing results from the search.

75 **Providing an input**

76 DateLife starts with an input query consisting of at least two taxon names, which can
77 be provided as a comma separated character string, or as tip labels on a tree. If the input is
78 a tree, it can be provided as a classic newick character string (Archie et al., 1986), or as a
79 “phylo” R object (Paradis et al., 2004). The input tree is not required to have branch lengths,
80 and its topology is used in the summary steps described below.

81 DateLife accepts scientific names as input. These names can belong to any inclusive
82 taxonomic group (e.g., genus, family, tribe, etc.) or binomial specific. Subspecies and
83 variants are ignored. If an input taxon name belongs to an inclusive taxonomic group the
84 algorithm has two alternative behaviors defined by the “get species from taxon” flag. If the
85 flag is active, the DateLife algorithm retrieves all species names within the inclusive
86 taxonomic group and adds them to the input. If the flag is inactive, DateLife ignores the
87 inclusive taxon names from the input.

88 Input scientific names are processed using a Taxonomic Name Resolution Service
89 (TNRS), which increases the probability of correctly finding the queried taxon names in the
90 chronogram database. TNRS detects, corrects and standardizes name misspellings and typos,
91 variant spellings and authorities, and nomenclatural synonyms to a single taxonomic
92 standard. DateLife implements TNRS using OpenTree’s taxonomy as standard (Open Tree
93 Of Life et al., 2016; Rees & Cranston, 2017).

94 The processed input taxon names are saved as an R object of a newly defined class

95 `datelifeQuery` that is used in the following steps. This object contains the processed
96 names, the corresponding OpenTree taxonomic id numbers, and the topology of the input
97 tree if any was provided.

98 **Searching the database**

99 A DateLife search consists of matching processed taxon names to tip labels in a
100 chronogram database. Chronograms with at least two matching tip labels are identified and
101 pruned down to preserve only the matched tips.

102 Matching pruned chronograms are stored as individual patristic distance matrices
103 (Figure 1 subfigure X). This matrix consists of . . . ???? the pairwise distance between pairs
104 of query taxa which are in that input tree, in units of millions of years.

105 This format speeds up extraction of pairwise taxon ages of the queried taxa, as opposed
106 to searching the ancestor node of a pair of taxa in a “phylo” object or newick string. The
107 patristic matrices are also associated to the study citation where the original chronogram
108 was published, and stored as an R object of the newly defined class `datelifeResult`.

109 DateLife’s chronogram database latest version consist of 253 chronograms published in
110 187 different studies. It is constructed from OpenTree’s phylogenetic database, the
111 Phylesystem, which constitutes an open source of expert phylogenetic knowledge with rich
112 metadata (McTavish et al., 2015) that allows automatic and reproducible construction of a
113 chronogram database. New chronograms can be added to Phylesystem by any user and are
114 immediately publicly available, and the DateLife database can be updated to include those
115 new data within a run.

116 **Summarizing search results**

117 At this point, summary information is extracted from the `datelifeResult` object to
118 inform decisions for the subsequent steps in the user workflow. Age data from the matching

119 pruned chronograms is summarized and used to generate a single summary chronogram.

120 Other basic summary information available to the user is:

- 121 1. The matching pruned chronograms as newick strings or “phylo” objects.
- 122 2. The ages of the root of all matching pruned chronograms. This can correspond to the
123 age of the most recent common ancestor (mrca) of your group of interest if the pruned
124 chronograms have all taxa belonging to the group. If not, the root corresponds to the
125 mrca of a subgroup within your group of interest.
- 126 3. Study citations where original chronograms were published.
- 127 4. A report of input taxon names matches across pruned chronograms.
- 128 5. The single matching pruned chronogram with the most input taxon names.

129 ***Identifying groves.***— To generate a single summary chronogram, the DateLife
130 algorithm starts by identifying the matching pruned chronograms that form a grove, roughly,
131 a sufficiently overlapping set of taxa between trees, by implementing definition 2.8 for
132 n-overlap from Ané et al. (2009). In rare cases, a group of trees can have multiple groves. By
133 default, DateLife chooses the grove with the most taxa, however, the “criterion = trees” flag
134 allows the user to choose the grove with the most trees instead.

135 ***Choosing a topology.***— DateLife requires a tree topology to summarize age data
136 upon. Users can provide one as input from the literature, or one of their own making. If no
137 topology is provided, DateLife automatically subsets one from the OpenTree synthetic tree
138 (Open Tree Of Life et al., 2019).

139 DateLife can also reconstruct branch lengths proportional to substitution rates on a
140 fixed tree topology using available genetic data from BOLD.

141 ***Congruifying nodes.***— DateLife then implements the congruification method
142 (Eastman, Harmon, & Tank, 2013) to find nodes belonging to the same clade across
143 matching pruned chronograms. Congruified node ages stored as a

¹⁴⁴ congruifiedCalibrations object are then matched to nodes in the chosen tree topology
¹⁴⁵ and stored as a matchedCalibrations object.

¹⁴⁶ ***Summarizing node ages.***— DateLife summarizes matched calibrations into a single
¹⁴⁷ patristic distance matrix using different methods. Summarizing options implemented include
¹⁴⁸ Super Distance Matrix method (SDM, Criscuolo, Berry, Douzery, & Gascuel, 2006) and
¹⁴⁹ summary statistics such as median, minimum and maximum ages.

¹⁵⁰ ***Dating the tree topology.***— Summarized calibrations can be applied as secondary
¹⁵¹ calibrations with different dating methods currently supported within DateLife: MrBayes
¹⁵² (Huelsenbeck & Ronquist, 2001; Ronquist & Huelsenbeck, 2003), PATHd8 (Britton,
¹⁵³ Anderson, Jacquet, Lundqvist, & Bremer, 2007), BLADJ (Campbell O. Webb, Ackerly, &
¹⁵⁴ Kembel, 2008; Campbell O. Webb & Donoghue, 2005), and treePL (Stephen A. Smith &
¹⁵⁵ O'Meara, 2012).

¹⁵⁶ By default, DateLife implements the Branch Length Adjuster (BLADJ) algorithm that
¹⁵⁷ assigns ages to nodes with no data evenly between nodes with age data, which minimizes age
¹⁵⁸ variance in the resulting chronogram (Campbell O. Webb et al., 2008). When there is
¹⁵⁹ conflict in ages across node with age data, the algorithm ignores ages that are older than
¹⁶⁰ parent nodes and/or younger than descendant nodes.

¹⁶¹ If there is no information on the age of the root in the chronogram database, users can
¹⁶² provide an estimate from the literature. If none is provided, DateLife assigns an arbitrary
¹⁶³ age to the root as 10% older than the oldest age available within the group.

¹⁶⁴ ***Visualizing results.***— Finally, users can save all source and summary chronograms in
¹⁶⁵ formats that permit reuse and reanalyses (newick and R “phylo” format), as well as view
¹⁶⁶ and compare results graphically, or construct their own graphs using datelife’s chronogram
¹⁶⁷ plot generation functions.

168

Benchmark

169 `datelife`'s code speed was tested on an Apple iMac with one 3.4 GHz Intel Core i5
170 processor. We registered variation in computing time of query processing and search through
171 the database relative to number of queried taxon names. Query processing time increases
172 roughly linearly with number of input taxon names, and increases considerably if TNRS is
173 activated. Up to ten thousand names can be processed and searched in less than 30 minutes
174 with the most time consuming settings. Once names have been processed as described in
175 methods, a name search through the chronogram database can be performed in less than a
176 minute, even with a very large number of taxon names (Fig. 2). `datelife`'s code
177 performance was evaluated with a set of unit tests designed and implemented with the R
178 package testthat (R Core Team, 2018) that were run both locally with the devtools package
179 (R Core Team, 2018), and on a public server –via GitHub, using the continuous integration
180 tool Travis CI (<https://travis-ci.org>). At present, unit tests cover more than 40% of
181 `datelife`'s code (<https://codecov.io/gh/phylotastic/datelife>).

182

Case study

183 We illustrate the DateLife algorithm using a group within the Passeriform birds
184 encompassing the family of true finches, Fringillidae and allies as case study. The first
185 example analyses 6 bird species and shows all steps of the algorithm. The second example is
186 a real life application

187 Small example

188 We chose 6 bird species associated to true finches at random. The sample includes two
189 species of cardinals: the black-thighed grosbeak – *Pheucticus tibialis* and the crimson-collared
190 grosbeak – *Rhodothraupis celaeno*; three species of buntings: the yellowhammer – *Emberiza*
191 *citrinella*, the pine bunting – *Emberiza leucocephalos* and the yellow-throated bunting –
192 *Emberiza elegans*; and one species of tanager, the vegetarian finch – *Platyspiza crassirostris*.

Processing input names found that *Emberiza elegans* is synonym for *Schoeniclus elegans* in the reference taxonomy. DateLife used the processed input names to search the local chronogram database and found 9 matching chronograms in 6 different studies. Three studies matched five input names (Barker, Burns, Klicka, Lanyon, & Lovette, 2015; Hedges, Marin, Suleski, Paymer, & Kumar, 2015; Jetz, Thomas, Joy, Hartmann, & Mooers, 2012), one study matched four input names (Hooper & Price, 2017) and two studies matched two input names (Barker, Burns, Klicka, Lanyon, & Lovette, 2013; Burns et al., 2014). No studies matched all input names. Together, matching chronograms have 28 unique age data points. All nodes have age data. As fixed tree topology, DateLife used OpenTree's synthetic tree as default and mapped age data to nodes in the tree. As expected, more inclusive nodes (e.g., node "n1") have more age data than less inclusive nodes (e.g., node "n5"). The processing step allowed discovering five data points for node "n4" that would not have had any data otherwise. Age summary statistics per node were calculated and tested as secondary calibrations to date the tree topology using the BLADJ algorithm. Age data for node "n2" was excluded as final calibration because it is older than age data of a more inclusive node.

Real life application

A college educator wishes to obtain state-of-the-art data on time of evolutionary origin of species belonging to the true finches for their class. They decide to use `datelife` because they are teaching best practices for reproducibility. Students have the option to go to the website at www.datelife.org and perform an interactive run. However, the educator also wants the students to practice their R skills. The first step is to run a `datelife` query using the "get species from taxon" flag. This will get all recognised species names within their chosen inclusive taxon. The Fringillidae has 289 species, according to the Open Tree of Life taxonomy. Once with a curated set of species taxon names, the next step is to run a `datelife` search that will find all chronograms that contain at least two species names. The algorithm proceeds to prune the trees to keep matching species names on tips only, and

transform the pruned trees to pairwise distance matrices. There are 13 chronograms containing at least two Fringillidae species, published in 9 different studies (Barker et al., 2013, 2015; Burns et al., 2014; Claramunt & Cracraft, 2015; Gibb et al., 2015; Hedges et al., 2015; Hooper & Price, 2017; Jetz et al., 2012; Price et al., 2014). The final step is to summarize the available information using two alternative types of summary chronograms, median and SDM. As explained in the “Description” section, data from source chronograms is first summarised into a single distance matrix and then the available node ages are used as fixed node calibrations over a consensus tree topology, to obtain a fully dated tree with the program BLADJ (Fig. 5). Median summary chronograms are older and have wider variation in maximum ages than chronograms obtained with SDM. ?????????????????? Say some things about the results!

Cross-validation test

Data from source chronograms can be used to date tree topologies with no branch lengths, as well as trees with branch lengths as relative substitution rates (Figs. 6 to 14). As a form of cross validation, we took tree topologies from each input study and calibrated them using time of lineage divergence data from all other source chronograms.

In the absence of branch lengths, the ages of internal nodes were recovered with a high precision in almost all cases (except for studies 3, and 5; Fig. ??). Maximum tree ages were only recovered in one case (study 2; Fig. 7). We also demonstrate the usage of PATHd8 (Britton et al., 2007) as an dating method alternative to BLADJ. For this, we run a **datelife** branch length reconstruction that searches for DNA sequence data from the Barcode of Life Data System [BOLD; Ratnasingham and Hebert (2007)] to generate branch lengths. We were able to successfully generate a tree with BOLD branch lengths for all of the Fringillidae source chronograms. However, dating with PATHd8 using congruified calibrations, was only successful in three cases (studies 3, 5, and 9, shown in Fig. ??). From these, two trees have a different sampling than the original source chronogram, mainly

245 because DNA BOLD data for some species is absent from the database. Maximum ages are
246 quite different from source chronograms, but this might be explained also by the differences
247 in sampling between source chronograms and BOLD trees. More examples and code used to
248 generate these trees were developed on an open repository that is available for consultation
249 and reuse at https://github.com/LunaSare/datelife_examples.

250

Discussion

251 The main goal of `datelife` is to make state-of-the-art information on time of lineage
252 divergence easily accessible for comparison, reuse, and reanalysis, to researchers in all areas
253 of science and with all levels of expertise in the matter. It is an open service that does not
254 require any expert biological knowledge from users –besides the names of the organisms they
255 want to work with, for any of its functionality.

256 At the time of writing of this manuscript (Apr 04, 2022), `datelife`'s database has 253
257 chronograms, pulled entirely from OpenTree's database, the Phylesystem (McTavish et al.,
258 2015). A unique feature of OpenTree's Phylesystem is that the community can add new
259 state-of-the-art chronograms any time. As chronograms are added to Phylesystem, they are
260 incorporated into an updated `datelife`'s database that is assigned a new version number,
261 followed by a package release on CRAN. `datelife`'s chronogram database is updated as new
262 chronogram data is added to Phylesystem, at a minimum of once a month and a maximum
263 of every 6 months. Users can also upload new chronograms to OpenTree themselves, and
264 trigger an update of their local `datelife` database to incorporate the new chronograms, to
265 have them immediately available for analysis.

266 Incorporation of more chronograms into `datelife`'s database is crucial to improve its
267 services. One option to increase chronogram number in the database is the Dryad data
268 repository. Methods to automatically mine chronograms from Dryad could be designed and
269 implemented. However, Dryad's metadata system has no information to automatically detect

270 branch length units, and those would still need to be determined manually by a curator.

271 The largest, and taxonomically broadest, summary chronogram currently available
272 from OpenTree was constructed using age data from 2,274 published chronograms (Hedges et
273 al., 2015). However the source chronograms used as input data for this tree are not available
274 in computer readable format for reuse or reanalysis. As this tree is part of datelife's
275 database, the amount of lineages that can be queried using **datelife** (99474 unique
276 terminal taxa) is substantial. Access to the input chronograms used to generate the Hedges
277 et al. summary tree would improve measures of uncertainty in DateLife, but they are
278 available only as image files and not as usable data (timetree.org). We would like to
279 emphasize on the importance of sharing chronogram data for the benefit of the scientific
280 community as a whole, into repositories that require expert input and manual curation, such
281 as OpenTree's Phylesystem (McTavish et al., 2015).

282 By default, **datelife** currently summarizes all source chronograms that overlap with
283 at least two species names. Users can exclude source chronograms if they have reasons to do
284 so. Strictly speaking, the best chronogram should reflect the real time of lineage divergence
285 accurately and precisely. To our knowledge, there are no good measures to determine
286 independently if a chronogram is better than another. Some measures that have been
287 proposed are the proportion of lineage sampling and the number of calibrations used
288 Magallón, Gómez-Acevedo, Sánchez-Reyes, & Hernández-Hernández (2015). Several
289 characteristics of the data used for dating analyses as well as from the output chronogram
290 itself, could be used to score quality of source chronograms. Some characteristics that are
291 often cited in published studies as a measure of improved age estimates as compared to
292 previously published estimates are: quality of alignment (missing data, GC content), lineage
293 sampling (strategy and proportion), phylogenetic and dating inference method, number of
294 fossils used as calibrations, support for nodes and ages, and magnitude of confidence
295 intervals. DateLife provides an opportunity to capture concordance and conflict among date

296 estimates, which can also be used as a metric for chronogram reliability.

297 Scientists usually also favor chronograms constructed using primary calibrations (ages
298 obtained from the fossil or geological record) to ones constructed with secondary calibrations
299 (ages coming from other chronograms)(Schenk, 2016). It has been observed with simulations
300 that divergence times inferred with secondary calibrations are significantly younger than
301 those inferred with primary calibrations in analyses performed with Bayesian inference
302 methods when priors are implemented in similar ways in both analyses (Schenk, 2016).
303 However, secondary calibrations can be applied using other dating methods that do not
304 require setting priors, such as penalized likelihood (Sanderson, 2003), or as fixed ages,
305 potentially mitigating the bias reported with Bayesian methods. Certainly, further studies
306 are required to fully understand the effect of using secondary calibrations on time estimates
307 and downstream analyses.

308 Furthermore, even chronograms obtained with primary fossil data can vary
309 substantially in time estimates between lineages, as observed from the comparison of source
310 chronograms in the Fringillidae example. This observation is often encountered in the
311 literature (see, for example, the ongoing debate about crown group age of angiosperms
312 (Barba-Montoya, Reis, Schneider, Donoghue, & Yang, 2018; Magallón et al., 2015; Ramshaw
313 et al., 1972; Sanderson & Doyle, 2001; Sauquet, Ramírez-Barahona, & Magallón, 2021). For
314 some studies, especially ones based on branch lengths (e.g., studies of species diversification,
315 timing of evolutionary events, phenotypic trait evolution), using a different chronogram may
316 return different results (Title & Rabosky, 2016). Stitching together these chronograms can
317 create a larger tree that uses information from multiple studies, but the effect of
318 uncertainties and errors at this level on downstream analyses is still largely unknown.

319 Summarizing chronograms might also imply summarizing fundamentally distinct
320 evolutionary hypotheses. For example, two different researchers working on the same clade
321 both carefully select and argument their choices of fossil calibrations. Still, if one researcher

322 decides a fossil will calibrate the ingroup of a clade, while another researcher uses the same
323 one to calibrate outside the clade, the resulting age estimates will often differ substantially,
324 as the placement of calibrations as stem or crown group is proved to deeply affect estimated
325 times of lineage divergence (Sauquet, 2013). Trying to summarize the resulting chronograms
326 into a single one using simple summary statistics can erase many types of relevant
327 information from the source chronograms. Accordingly, the prevailing view is that we should
328 favor time of lineage divergence estimates obtained from a single analysis, using fossil data as
329 primary sources of calibrations, and using fossils that have been widely discussed and
330 curated as calibrations to date other trees, making sure that all data used in the analysis
331 reflect a coherent evolutionary history (Antonelli et al., 2017). However, the exercise of
332 summarizing different chronograms has the potential to help getting a single global
333 evolutionary history for a lineage by putting together evidence from different hypothesis.
334 Choosing the elements of the chronograms that we are going to keep and the ones that we
335 are going to discard is key, since we are potentially loosing important parts of the
336 evolutionary history of a lineage that might only be reflected in source chronograms and not
337 on the summary chronogram (Sauquet et al., 2021).

338 Nonetheless, in ecology and conservation biology, incorporating at least some data on
339 lineage divergence times represents a relevant improvement for testing alternative hypothesis
340 using phylogenetic distance (Campbell O. Webb et al., 2008). Hence, we integrated into
341 datelife's workflow different ways of estimating node ages in the absence of calibrations and
342 branch length information for taxa lacking this information. "Making up" branch lengths is
343 an accepted practice in scientific publications: Jetz et al. (2012), created a time-calibrated
344 tree of all 9,993 bird species, where 67% had molecular data and the rest was simulated;
345 Rabosky et al. (2018) created a time-calibrated tree of 31,536 ray-finned fishes, of which only
346 37% had molecular data; Stephen A. Smith and Brown (2018) constructed a tree of 353,185
347 seed plants where only 23% had molecular data. Obviously, there are risks in this practice!
348 Taken to the extreme, one could make a fully resolved, calibrated tree of all modern and

349 extinct taxa using a single taxonomy and a single calibration with the polytomy resolution
350 and branch estimation methods. There has yet to be a thorough analysis of what can go
351 wrong when one extends inferences beyond the data in this way, so we urge caution; we also
352 urge readers to follow the example of many of the large tree papers cited above and make
353 carefully consider the statistical assumptions being made, and assess the consistency of the
354 results with prior work.

355 **Conclusions**

356 Divergence time information is key to many areas of evolutionary studies: trait
357 evolution, diversification, biogeography, macroecology and more. It is also crucial for science
358 communication and education, but generating chronograms is difficult, especially for those
359 who want to use phylogenies but who are not systematists, or do not have the time to
360 acquire and develop the necessary knowledge and data curation skills. Moreover, years of
361 primarily public funded research have resulted in vast amounts of chronograms that are
362 already available on scientific publications, but hidden to the public and scientific community
363 for reuse.

364 The `datelife` R package allows easy and fast summarization of publicly available
365 information on time of lineage divergence. This provides a straightforward way to get an
366 informed idea on the state of knowledge of the time frame of evolution of different regions of
367 the tree of life, and allows identification of regions that require more research or that have
368 conflicting information. It is available as an R package, or a web-based R shiny app at
369 dates.opentreeloflife.org/datelife. Both summary and newly generated trees are useful to
370 evaluate evolutionary hypotheses in different areas of research. The DateLife project helps
371 with awareness of the existing variation in expert time of divergence data, and will foster
372 exploration of the effect of alternative divergence time hypothesis on the results of analyses,
373 nurturing a culture of more cautious interpretation of evolutionary results.

Availability

374 **Availability**

375 **datelife** is free and open source and it can be used through its current website
376 <http://www.datelife.org/query/>, through its R package, and through Phylotastic's project
377 web portal <http://phylo.cs.nmsu.edu:3000/>. **datelife**'s website is maintained using
378 RStudio's shiny server and the shiny package open infrastructure, as well as Docker.
379 **datelife**'s R package stable version is available for installation from the CRAN repository
380 (<https://cran.r-project.org/package=datelife>) using the command `install.packages(pkgs`
381 `= "datelife"`) from within R. Development versions are available from the GitHub
382 repository (<https://github.com/phylotastic/datelife>) and can be installed using the
383 command `devtools::install_github("phylotastic/datelife")`.

Supplementary Material

384 **Supplementary Material**

385 Code used to generate all versions of this manuscript, the biological examples, as well
386 as the benchmark of functionalities are available at datelifeMS1, datelife_examples, and
387 datelife_benchmark repositories in LLSR's GitHub account.

Funding

388 **Funding**

389 Funding was provided by the US National Science Foundation (NSF) grants
390 ABI-1458603 to Datelife project and DBI-0905606 to the National Evolutionary Synthesis
391 Center (NESCent), and the Phylotastic project Grant ABI-1458572, and the OpenTree grant
392 ABI-1759846.

Acknowledgements

393 **Acknowledgements**

394 The DateLife project was born as a prototype tool aiming to provide these services,
395 and was developed over a series of hackathons at the National Evolutionary Synthesis
396 Center, NC, USA (Stoltzfus et al., 2013). We thank colleagues from the O'Meara Lab at the
397 University of Tennessee Knoxville for suggestions, discussions and software testing. The late

398 National Evolutionary Synthesis Center (NESCent), which sponsored hackathons that led to
399 initial work on this project. The team that assembled `datelife`'s first proof of concept:
400 Tracy Heath, Jonathan Eastman, Peter Midford, Joseph Brown, Matt Pennell, Mike Alfaro,
401 and Luke Harmon. The Open Tree of Life project that provides the open, metadata rich
402 repository of trees used for `datelife`. The many scientists who publish their chronograms in
403 an open, reusable form, and the scientists who curate them for deposition in the Open Tree
404 of Life repository. The NSF for funding nearly all the above, in addition to the ABI grant
405 that funded this project itself.

References

- Ané, C., Eulenstein, O., Piaggio-Talice, R., & Sanderson, M. J. (2009). Groves of phylogenetic trees. *Annals of Combinatorics*, 13(2), 139–167.
- Antonelli, A., Hettling, H., Condamine, F. L., Vos, K., Nilsson, R. H., Sanderson, M. J., ... Vos, R. A. (2017). Toward a self-updating platform for estimating rates of speciation and migration, ages, and relationships of Taxa. *Systematic Biology*, 66(2), 153–166. <https://doi.org/10.1093/sysbio/syw066>
- Archie, J., Day, W. H., Felsenstein, J., Maddison, W., Meacham, C., Rohlf, F. J., & Swofford, D. (1986). The Newick tree format. Retrieved from %7B<https://evolution.genetics.washington.edu/phylip/newicktree.html>%7D
- Bapst, D. W. (2012). Paleotree: An R package for paleontological and phylogenetic analyses of evolution. *Methods in Ecology and Evolution*, 3(5), 803–807. <https://doi.org/10.1111/j.2041-210X.2012.00223.x>
- Barba-Montoya, J., Reis, M. dos, Schneider, H., Donoghue, P. C., & Yang, Z. (2018). Constraining uncertainty in the timescale of angiosperm evolution and the veracity of a cretaceous terrestrial revolution. *New Phytologist*, 218(2), 819–834.
- Barker, F. K., Burns, K. J., Klicka, J., Lanyon, S. M., & Lovette, I. J. (2013). Going to extremes: Contrasting rates of diversification in a recent radiation of new world passerine birds. *Systematic Biology*, 62(2), 298–320.
- Barker, F. K., Burns, K. J., Klicka, J., Lanyon, S. M., & Lovette, I. J. (2015). New insights into new world biogeography: An integrated view from the phylogeny of blackbirds, cardinals, sparrows, tanagers, warblers, and allies. *The Auk: Ornithological Advances*, 132(2), 333–348.
- Britton, T., Anderson, C. L., Jacquet, D., Lundqvist, S., & Bremer, K. (2007). Estimating Divergence Times in Large Phylogenetic Trees. *Systematic Biology*, 56(788777878), 741–752. <https://doi.org/10.1080/10635150701613783>
- Burns, K. J., Shultz, A. J., Title, P. O., Mason, N. A., Barker, F. K., Klicka, J., ...

- 433 Lovette, I. J. (2014). Phylogenetics and diversification of tanagers (passeriformes:
434 Thraupidae), the largest radiation of neotropical songbirds. *Molecular*
435 *Phylogenetics and Evolution*, 75, 41–77.
- 436 Chamberlain, S. A., & Szöcs, E. (2013). taxize : taxonomic search and retrieval in R
437 [version 2; referees: 3 approved]. *F1000Research*, 2(191), 1–29.
438 <https://doi.org/10.12688/f1000research.2-191.v2>
- 439 Chamberlain, S. A., Szöcs, E., Foster, Z., Arendsee, Z., Boettiger, C., Ram, K., ...
440 Li, G. (2019). *taxize: Taxonomic information from around the web*. Retrieved
441 from <https://github.com/ropensci/taxize>
- 442 Claramunt, S., & Cracraft, J. (2015). A new time tree reveals earth history's imprint
443 on the evolution of modern birds. *Science Advances*, 1(11), e1501005.
- 444 Criscuolo, A., Berry, V., Douzery, E. J. P., & Gascuel, O. (2006). SDM: A fast
445 distance-based approach for (super)tree building in phylogenomics. *Systematic*
446 *Biology*, 55(5), 740–755. <https://doi.org/10.1080/10635150600969872>
- 447 Delsuc, F., Philippe, H., Tsagkogeorga, G., Simion, P., Tilak, M.-K., Turon, X., ...
448 Douzery, E. J. (2018). A phylogenomic framework and timescale for comparative
449 studies of tunicates. *BMC Biology*, 16(1), 1–14.
- 450 Eastman, J. M., Harmon, L. J., & Tank, D. C. (2013). Congruification: Support for
451 time scaling large phylogenetic trees. *Methods in Ecology and Evolution*, 4(7),
452 688–691. <https://doi.org/10.1111/2041-210X.12051>
- 453 Felsenstein, J. (1985). Phylogenies and the Comparative Method. *The American*
454 *Naturalist*, 125(1), 1–15. Retrieved from <http://www.jstor.org/stable/2461605>
- 455 Gibb, G. C., England, R., Hartig, G., McLenachan, P. A., Taylor Smith, B. L.,
456 McComish, B. J., ... Penny, D. (2015). New zealand passerines help clarify the
457 diversification of major songbird lineages during the oligocene. *Genome Biology*
458 and *Evolution*, 7(11), 2983–2995.
- 459 Harmon, L., Weir, J., Brock, C., Glor, R., & Challenger, W. (2008). GEIGER:

- 460 investigating evolutionary radiations. *Bioinformatics*, 24, 129–131.
- 461 Hedges, S. B., Marin, J., Suleski, M., Paymer, M., & Kumar, S. (2015). Tree of life
462 reveals clock-like speciation and diversification. *Molecular Biology and Evolution*,
463 32(4), 835–845. <https://doi.org/10.1093/molbev/msv037>
- 464 Heibl, C. (2008). *PHYLOCH: R language tree plotting tools and interfaces to diverse*
465 *phylogenetic software packages*. Retrieved from
466 <http://www.christophheibl.de/Rpackages.html>
- 467 Hooper, D. M., & Price, T. D. (2017). Chromosomal inversion differences correlate
468 with range overlap in passerine birds. *Nature Ecology & Evolution*, 1(10), 1526.
- 469 Huelsenbeck, J. P., & Ronquist, F. (2001). MRBAYES: Bayesian inference of
470 phylogenetic trees. *Bioinformatics*, 17(8), 754–755.
471 <https://doi.org/10.1093/bioinformatics/17.8.754>
- 472 Jetz, W., Thomas, G., Joy, J. J. B., Hartmann, K., & Mooers, A. (2012). The global
473 diversity of birds in space and time. *Nature*, 491(7424), 444–448.
474 <https://doi.org/10.1038/nature11631>
- 475 Laubichler, M. D., & Maienschein, J. (2009). *Form and function in developmental*
476 *evolution*. Cambridge University Press.
- 477 Magallon, S., & Sanderson, M. J. (2001). Absolute diversification rates in angiosperm
478 clades. *Evolution*, 55(9), 1762–1780.
- 479 Magallón, S. (2010). Using fossils to break long branches in molecular dating: A
480 comparison of relaxed clocks applied to the origin of angiosperms. *Systematic*
481 *Biology*, 59(4), 384–399.
- 482 Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L. L., & Hernández-Hernández, T.
483 (2015). A metacalibrated time-tree documents the early rise of flowering plant
484 phylogenetic diversity. *New Phytologist*, 207(2), 437–453.
- 485 McTavish, E. J., Hinchliff, C. E., Allman, J. F., Brown, J. W., Cranston, K. A.,
486 Holder, M. T., ... Smith, S. A. (2015). Phylesystem: A git-based data store for

- 487 community-curated phylogenetic estimates. *Bioinformatics*, 31(17), 2794–2800.
- 488 Michonneau, F., Brown, J. W., & Winter, D. J. (2016). rotl: an R package to interact
489 with the Open Tree of Life data. *Methods in Ecology and Evolution*, 7(12),
490 1476–1481. <https://doi.org/10.1111/2041-210X.12593>
- 491 Morlon, H. (2014). Phylogenetic approaches for studying diversification. *Ecology*
492 *Letters*, 17(4), 508–525. <https://doi.org/10.1111/ele.12251>
- 493 Ooms, J., & Chamberlain, S. (2018). *Phylocomr: Interface to 'phylocom'*. Retrieved
494 from <https://CRAN.R-project.org/package=phylocomr>
- 495 Open Tree Of Life, Redelings, B., Cranston, K. A., Allman, J., Holder, M. T., &
496 McTavish, E. J. (2016). Open Tree of Life APIs v3.0. *Open Tree of Life Project*,
497 (Online Resources). Retrieved from
498 <https://github.com/OpenTreeOfLife/germinator/wiki/Open-Tree-of-Life->
499 Web-APIs%7D
- 500 Open Tree Of Life, Redelings, B., Sánchez Reyes, L. L., Cranston, K. A., Allman, J.,
501 Holder, M. T., & McTavish, E. J. (2019). Open tree of life synthetic tree v12.3.
502 *Zenodo*. Retrieved from <https://doi.org/10.5281/zenodo.3937742>
- 503 Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and
504 evolution in R language. *Bioinformatics*, 20(2), 289–290.
- 505 Posadas, P., Crisci, J. V., & Katinas, L. (2006). Historical biogeography: A review of
506 its basic concepts and critical issues. *Journal of Arid Environments*, 66(3),
507 389–403.
- 508 Price, T. D., Hooper, D. M., Buchanan, C. D., Johansson, U. S., Tietze, D. T.,
509 Alström, P., ... others. (2014). Niche filling slows the diversification of himalayan
510 songbirds. *Nature*, 509(7499), 222.
- 511 R Core Team. (2018). *R: a language and environment for statistical computing*.
512 Vienna, Austria: R Foundation for Statistical Computing.
- 513 Rabosky, D. L., Chang, J., Title, P. O., Cowman, P. F., Sallan, L., Friedman, M., ...

- 514 others. (2018). An inverse latitudinal gradient in speciation rate for marine fishes.
515 *Nature*, 559(7714), 392.
- 516 Ramshaw, J., Richardson, D., Mealyard, B., Brown, R., Richardson, M., Thompson,
517 E., & Boulter, D. (1972). The time of origin of the flowering plants determined by
518 using amino acid sequence data of cytochrome c. *New Phytologist*, 71(5), 773–779.
- 519 Ratnasingham, S., & Hebert, P. D. (2007). BOLD: The barcode of life data system
520 (<http://www.Barcodinglife.org>). *Molecular Ecology Notes*, 7(3), 355–364.
- 521 Rees, J. A., & Cranston, K. (2017). Automated assembly of a reference taxonomy for
522 phylogenetic data synthesis. *Biodiversity Data Journal*, (5).
- 523 Revell, L. J. (2012). Phytools: An r package for phylogenetic comparative biology
524 (and other things). *Methods in Ecology and Evolution*, 3, 217–223.
- 525 Ronquist, F., & Huelsenbeck, J. P. (2003). MrBayes 3: Bayesian phylogenetic
526 inference under mixed models. *Bioinformatics*, 19(12), 1572–1574.
527 <https://doi.org/10.1093/bioinformatics/btg180>
- 528 Sanchez-Reyes, L. L., O'Meara, B., Eastman, J., Heath, T., Wright, A., Schliep, K.,
529 ... Alfaro, M. (2022). datelife: Scientific Data on Time of Lineage Divergence for
530 Your Taxa. *R Package Version 0.6.2*. Retrieved from
531 <https://doi.org/10.5281/zenodo.593938>
- 532 Sanderson, M. J. (2003). r8s: Inferring absolute rates of molecular evolution and
533 divergence times in the absence of a molecular clock. *Bioinformatics*, 19(2),
534 301–302.
- 535 Sanderson, M. J., & Doyle, J. A. (2001). Sources of error and confidence intervals in
536 estimating the age of angiosperms from rbcL and 18S rDNA data. *American
537 Journal of Botany*, 88(8), 1499–1516.
- 538 Sauquet, H. (2013). A practical guide to molecular dating. *Comptes Rendus Palevol*,
539 12(6), 355–367.
- 540 Sauquet, H., Ramírez-Barahona, S., & Magallón, S. (2021). The age of flowering

- 541 plants is unknown.
- 542 Schenk, J. J. (2016). Consequences of secondary calibrations on divergence time
543 estimates. *PLoS ONE*, 11(1). <https://doi.org/10.1371/journal.pone.0148228>
- 544 Smith, Stephen A., & Brown, J. W. (2018). Constructing a broadly inclusive seed
545 plant phylogeny. *American Journal of Botany*, 105(3), 302–314.
- 546 Smith, Stephen A., & O'Meara, B. C. (2012). TreePL: Divergence time estimation
547 using penalized likelihood for large phylogenies. *Bioinformatics*, 28(20),
548 2689–2690. <https://doi.org/10.1093/bioinformatics/bts492>
- 549 Stoltzfus, A., Lapp, H., Matasci, N., Deus, H., Sidlauskas, B., Zmasek, C. M., ...
550 Jordan, G. (2013). Phylotastic! Making tree-of-life knowledge accessible, reusable
551 and convenient. *BMC Bioinformatics*, 14.
552 <https://doi.org/10.1186/1471-2105-14-158>
- 553 Title, P. O., & Rabosky, D. L. (2016). Do Macrophylogenies Yield Stable
554 Macroevolutionary Inferences? An Example from Squamate Reptiles. *Systematic
555 Biology*, syw102. <https://doi.org/10.1093/sysbio/syw102>
- 556 Vos, R. A., Balhoff, J. P., Caravas, J. A., Holder, M. T., Lapp, H., Maddison, W. P.,
557 ... others. (2012). NeXML: Rich, extensible, and verifiable representation of
558 comparative data and metadata. *Systematic Biology*, 61(4), 675–689.
- 559 Webb, C. (2000). Exploring the Phylogenetic Structure of Ecological Communities :
560 An Example for Rain Forest Trees. *The American Naturalist*, 156(2), 145–155.
- 561 Webb, Campbell O., Ackerly, D. D., & Kembel, S. W. (2008). Phylocom: Software for
562 the analysis of phylogenetic community structure and trait evolution.
563 *Bioinformatics*, 24(18), 2098–2100.
564 <https://doi.org/10.1093/bioinformatics/btn358>
- 565 Webb, Campbell O., & Donoghue, M. J. (2005). Phylomatic: Tree assembly for
566 applied phylogenetics. *Molecular Ecology Notes*, 5(1), 181–183.

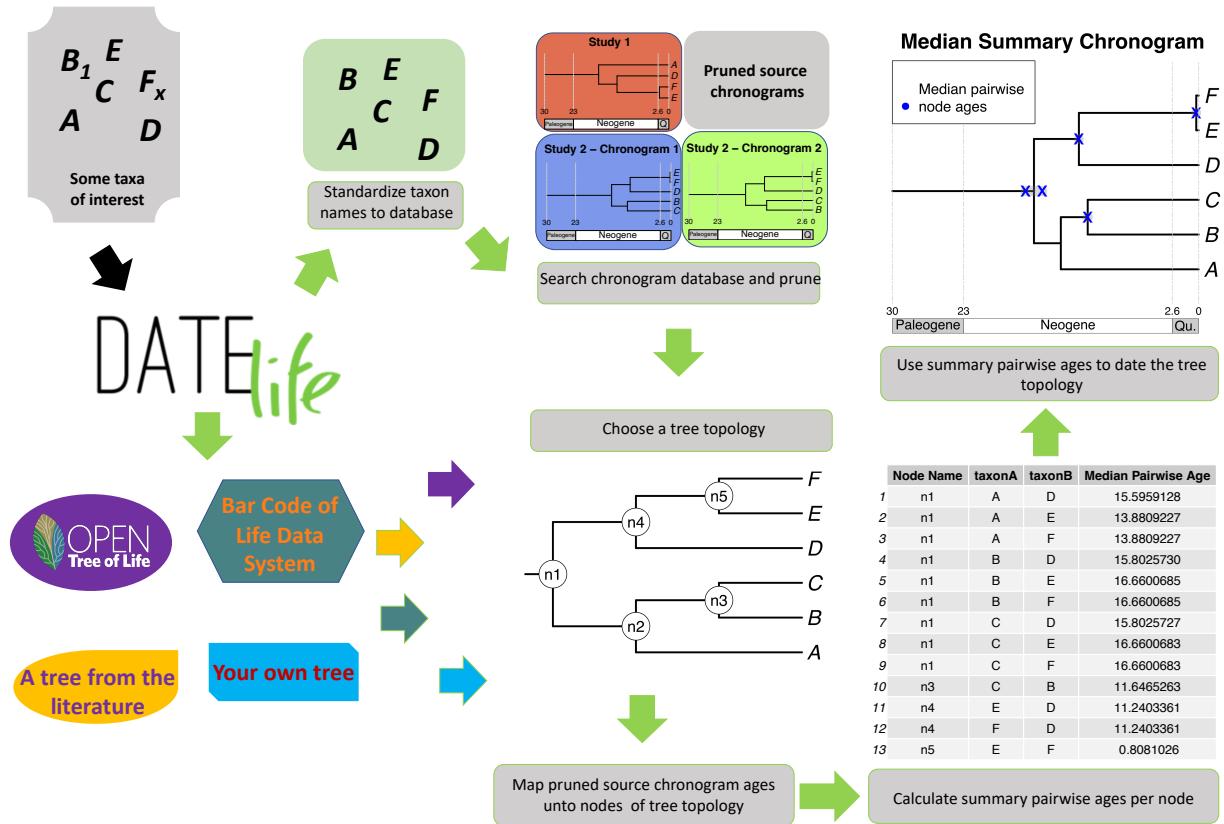


FIGURE 1. Stylized DateLife workflow. This shows the general workflows and analyses that can be performed with `datelife`, via the R package or through the website at <http://www.datelife.org/>. Details on the functions involved on each workflow are shown in `datelife`'s R package vignette.

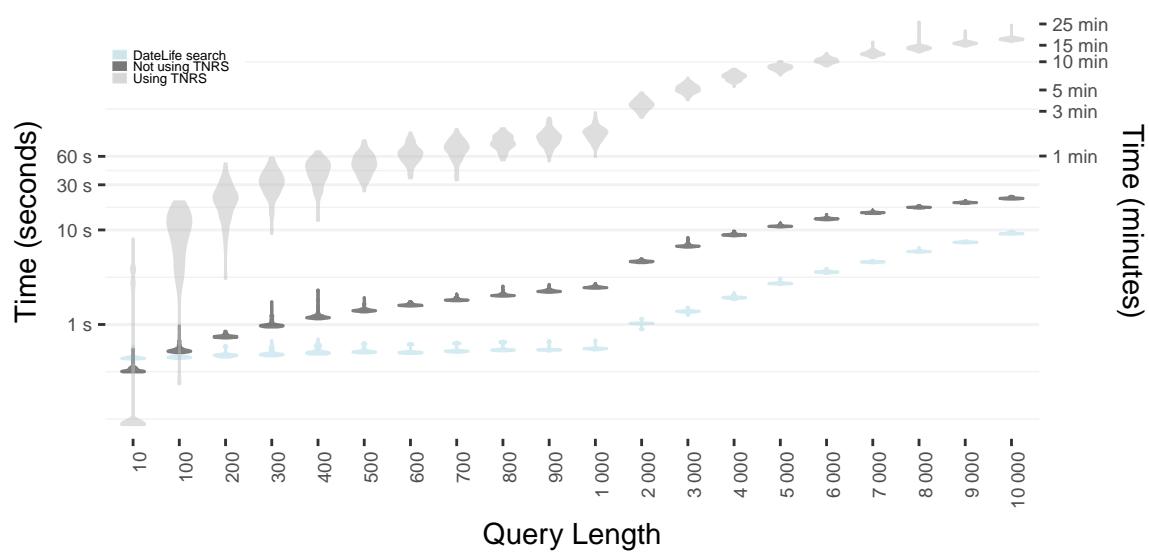


FIGURE 2. Computation time of query processing and search across **datelife**'s chronogram database relative to number of input taxon names. We sampled N names from the class Aves for each cohort 100 times and then performed a search with query processing not using the Taxon Names Resolution Service (TNRS; dark gray), and using TNRS (light gray). We also performed a search using the already processed query for comparison (light blue).

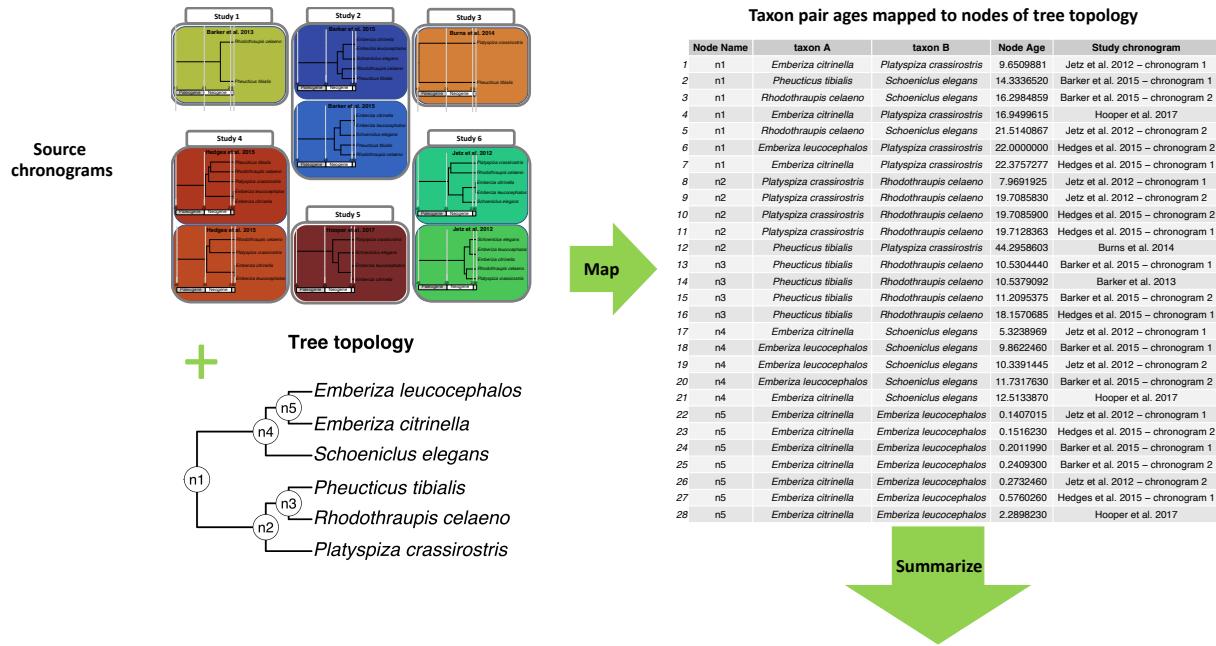


FIGURE 3. Age data results of a DateLife search of a small sample of 6 bird species within the Passeriformes. Input names were found across 9 chronograms within 6 independent studies (Barker et al. (2012), Barker et al. (2015), Burns et al. (2014), Hedges et al. (2015), Hooper and Price (2017), Jetz et al. (2012).) This revealed 28 age data points for the queried species names.

Summary of mapped taxon pair age data

Node Name	taxon A	taxon B	Pairwise Median Age	Node Median Age
1	<i>Pheucticus tibialis</i>	<i>Emberiza citrinella</i>	16.298486	
2	<i>Pheucticus tibialis</i>	<i>Emberiza leucocephalos</i>	16.298486	
3	<i>Platyspiza crassirostris</i>	<i>Emberiza citrinella</i>	21.514085	
4	<i>Platyspiza crassirostris</i>	<i>Emberiza leucocephalos</i>	21.514085	
5 n1	<i>Rhodothraupis celaeno</i>	<i>Emberiza citrinella</i>	20.408031	19.301977
6	<i>Rhodothraupis celaeno</i>	<i>Emberiza leucocephalos</i>	20.408031	
7	<i>Schoeniclus elegans</i>	<i>Pheucticus tibialis</i>	15.316069	
8	<i>Schoeniclus elegans</i>	<i>Platyspiza crassirostris</i>	19.301977	
9	<i>Schoeniclus elegans</i>	<i>Rhodothraupis celaeno</i>	17.800231	
10 n2	<i>Platyspiza crassirostris</i>	<i>Pheucticus tibialis</i>	32.004348	25.856467327225
11	<i>Rhodothraupis celaeno</i>	<i>Platyspiza crassirostris</i>	19.708587	
12 n3	<i>Rhodothraupis celaeno</i>	<i>Pheucticus tibialis</i>	10.873723	10.87372335475
13 n4	<i>Schoeniclus elegans</i>	<i>Emberiza citrinella</i>	10.647794	10.6477935
14	<i>Schoeniclus elegans</i>	<i>Emberiza leucocephalos</i>	10.647794	
15 n5	<i>Emberiza leucocephalos</i>	<i>Emberiza citrinella</i>	0.273246	0.273246

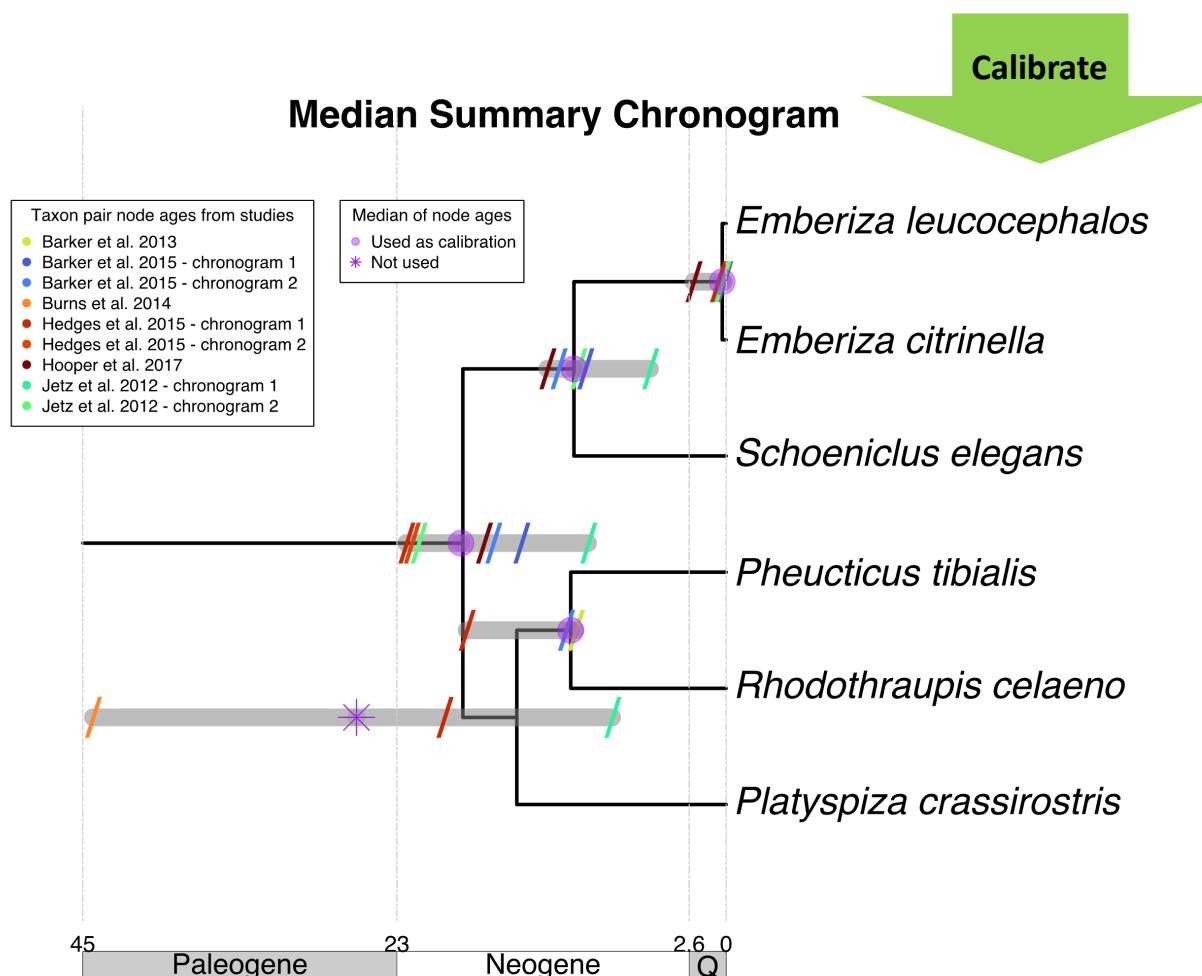


FIGURE 4. Summarized age data is used as secondary calibrations to date a tree topology as a summary chronogram.

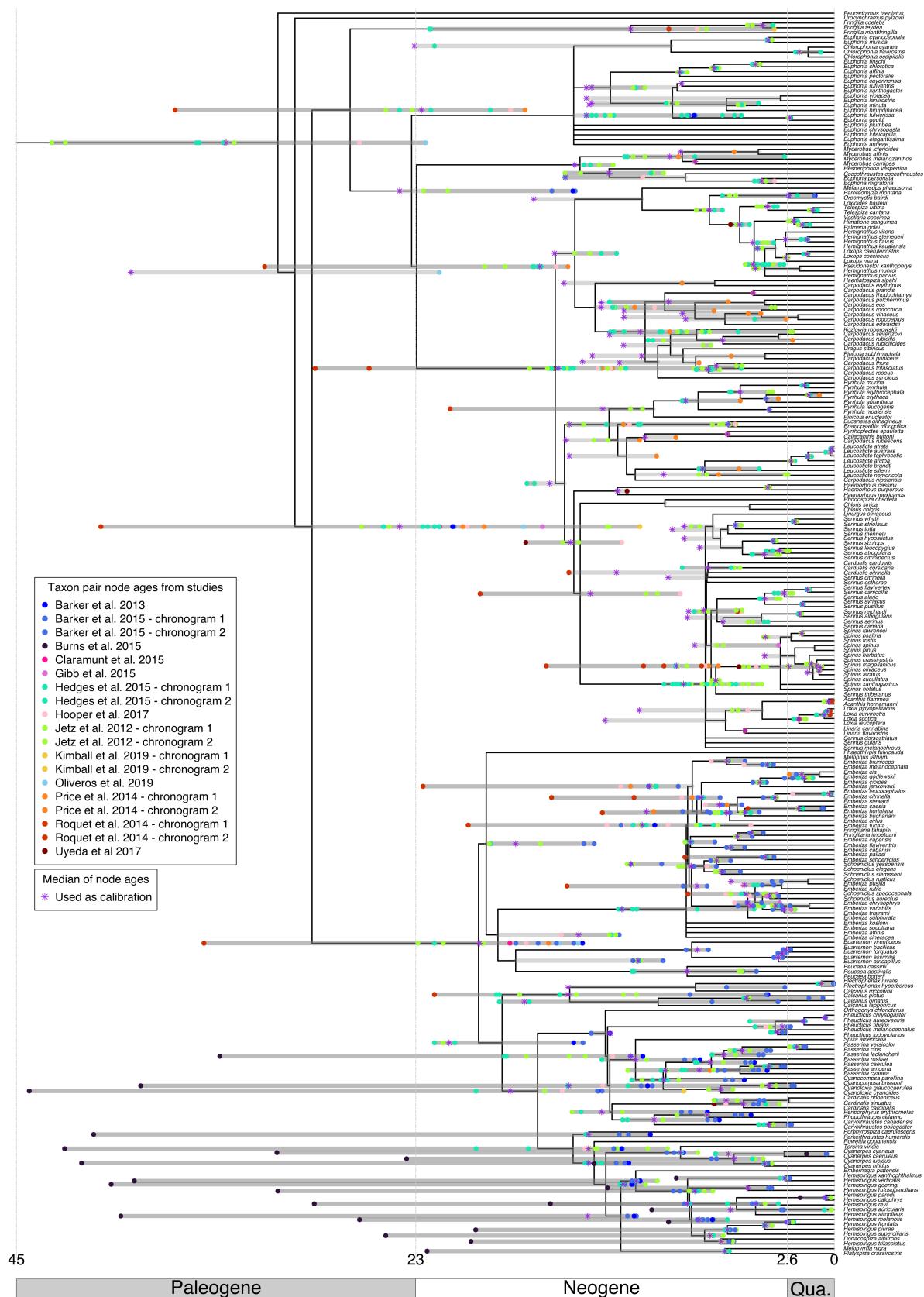


FIGURE 5. Fringillidae median summary chronogram generated with DateLife. It has 256 tips and 233 nodes.

Barker et al. 2013

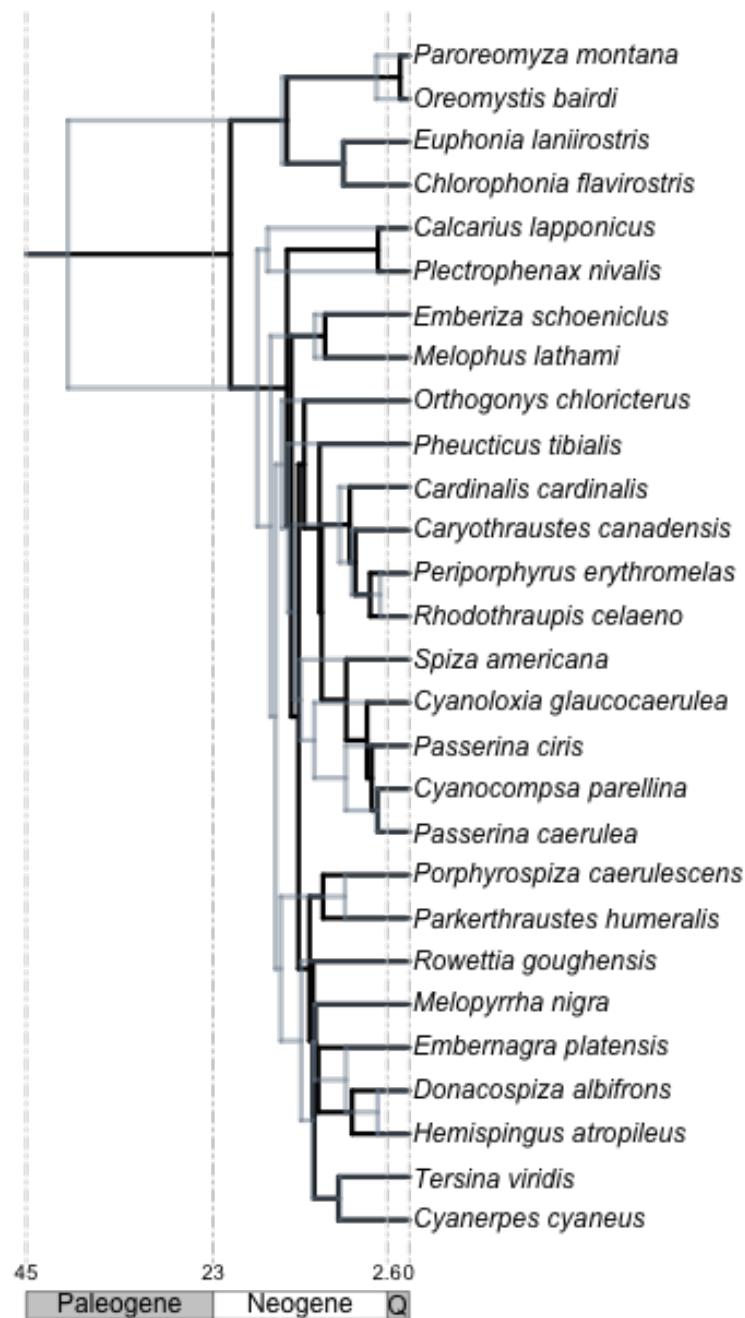


FIGURE 6. Cross validation of first source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ, i.e., under a fossil calibration.

Barker et al. 2015 - chronogram 1



FIGURE 7. Cross validation of second source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to

Barker et al. 2015 - chronogram 2

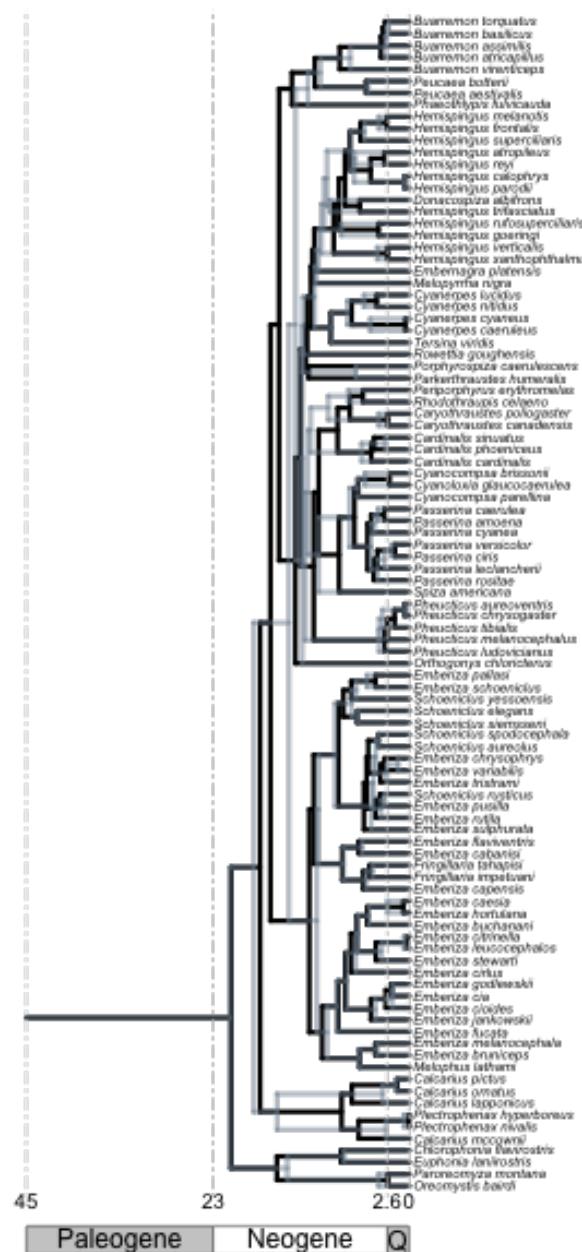


FIGURE 8. Cross validation of third source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ.

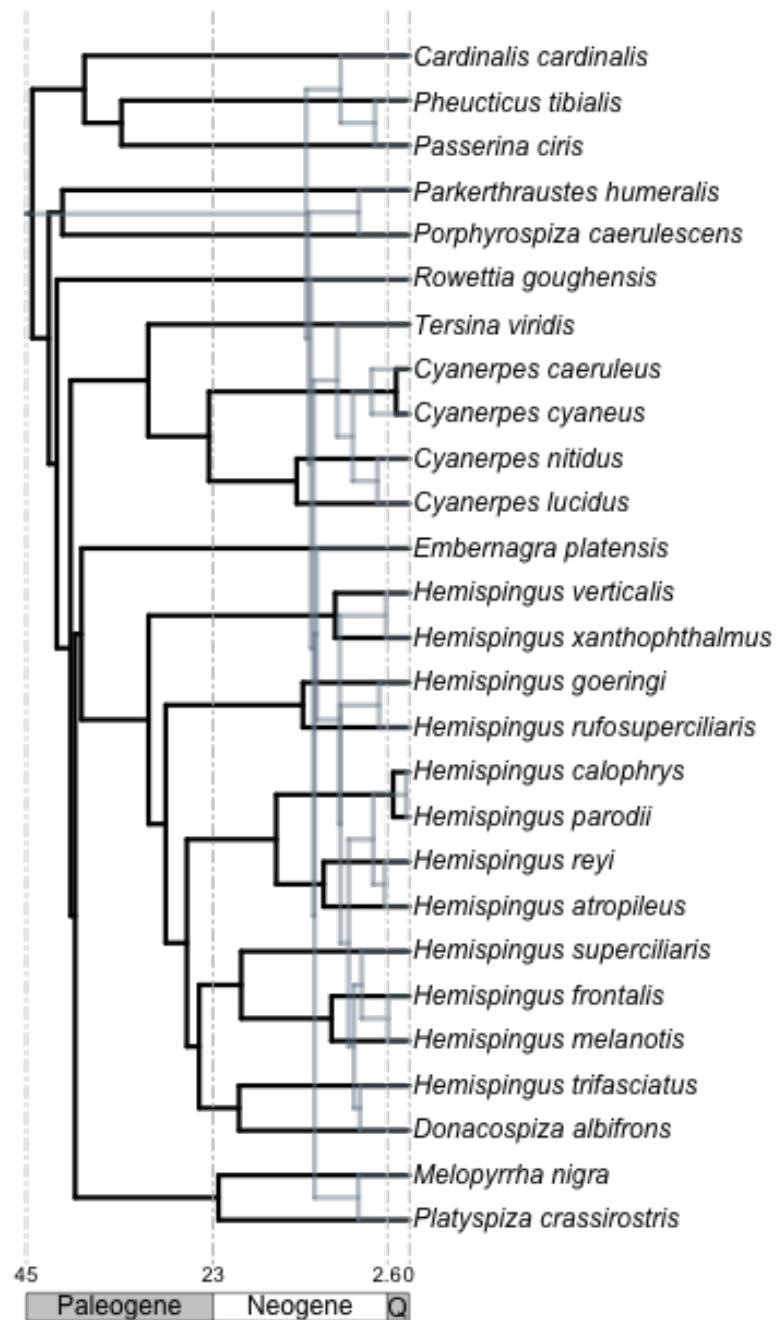
Burns et al. 2015

FIGURE 9. Cross validation of fourth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ, i.e., the mean of all the samples.

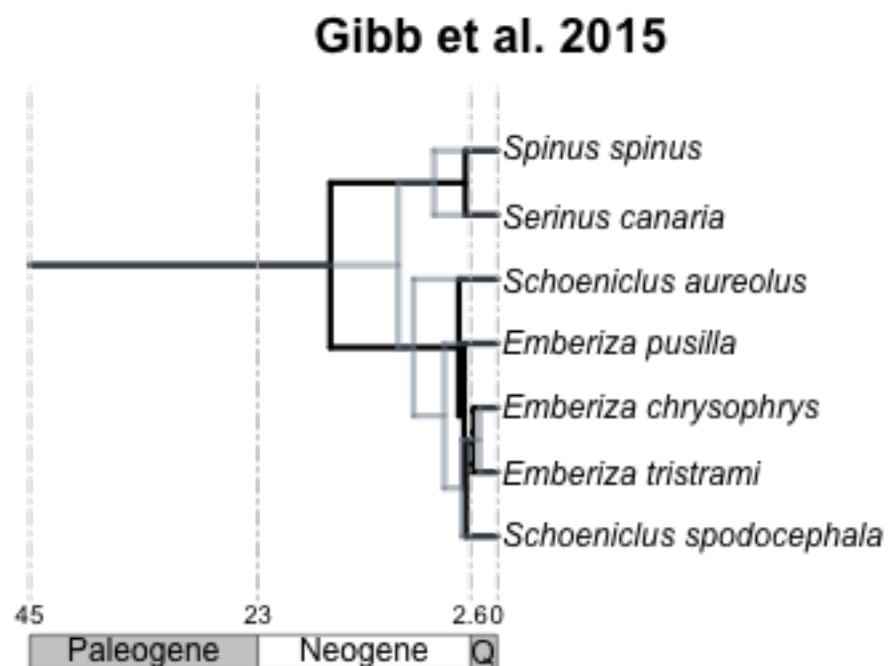


FIGURE 10. Cross validation of sixth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the same tree topology dated with BLADJ using node ages from all other source chronograms as secondary calibrations.

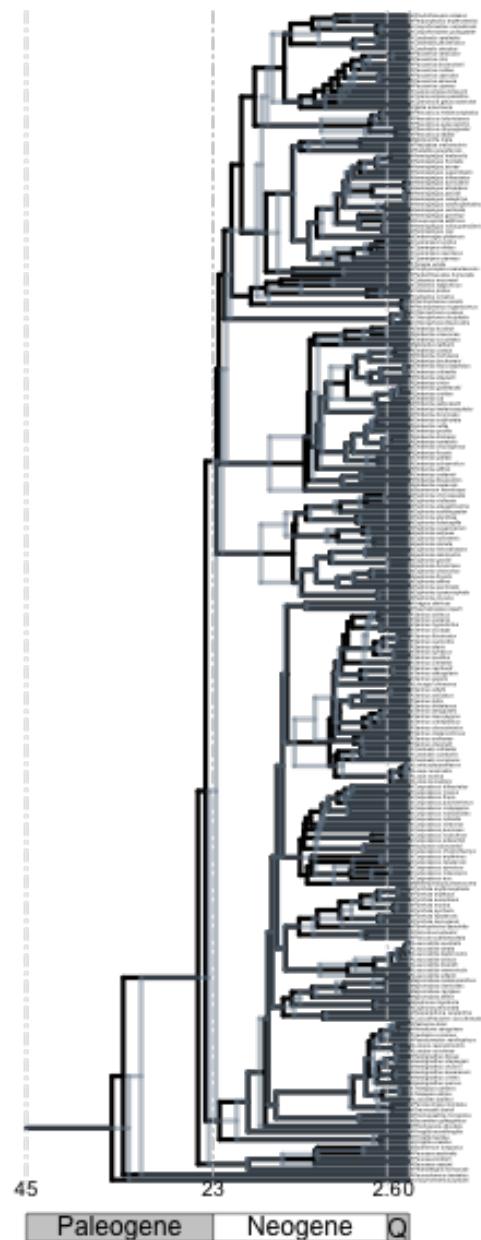
Hedges et al. 2015 - chronogram 1

FIGURE 11. Cross validation of seventh source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADe. In order to facilitate the comparison, the same color scheme was used.

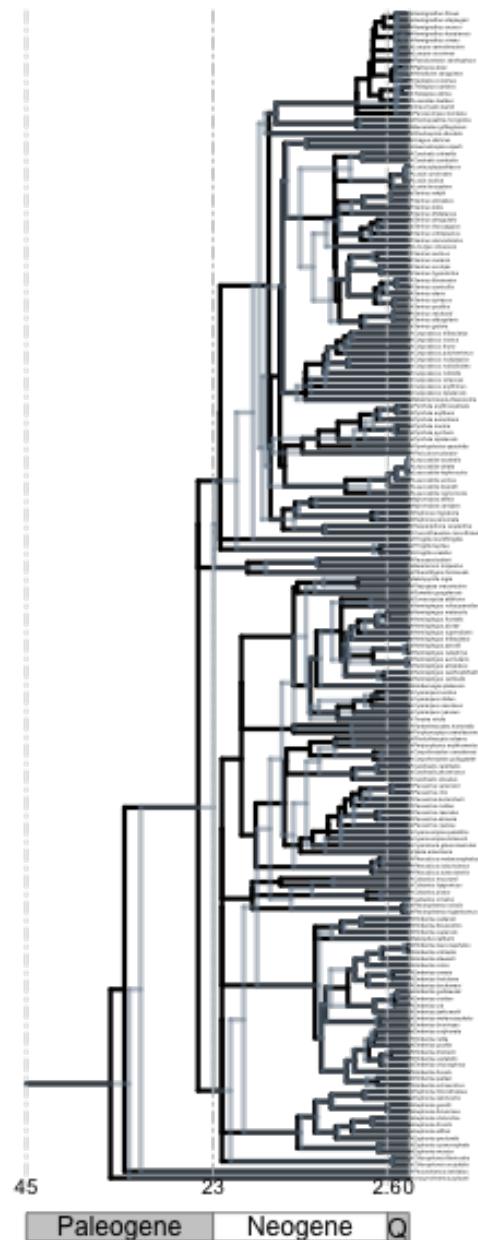
Hedges et al. 2015 - chronogram 2

FIGURE 12. Cross validation of eight source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ, i.e., the cross-validation procedure.

Hooper et al. 2017

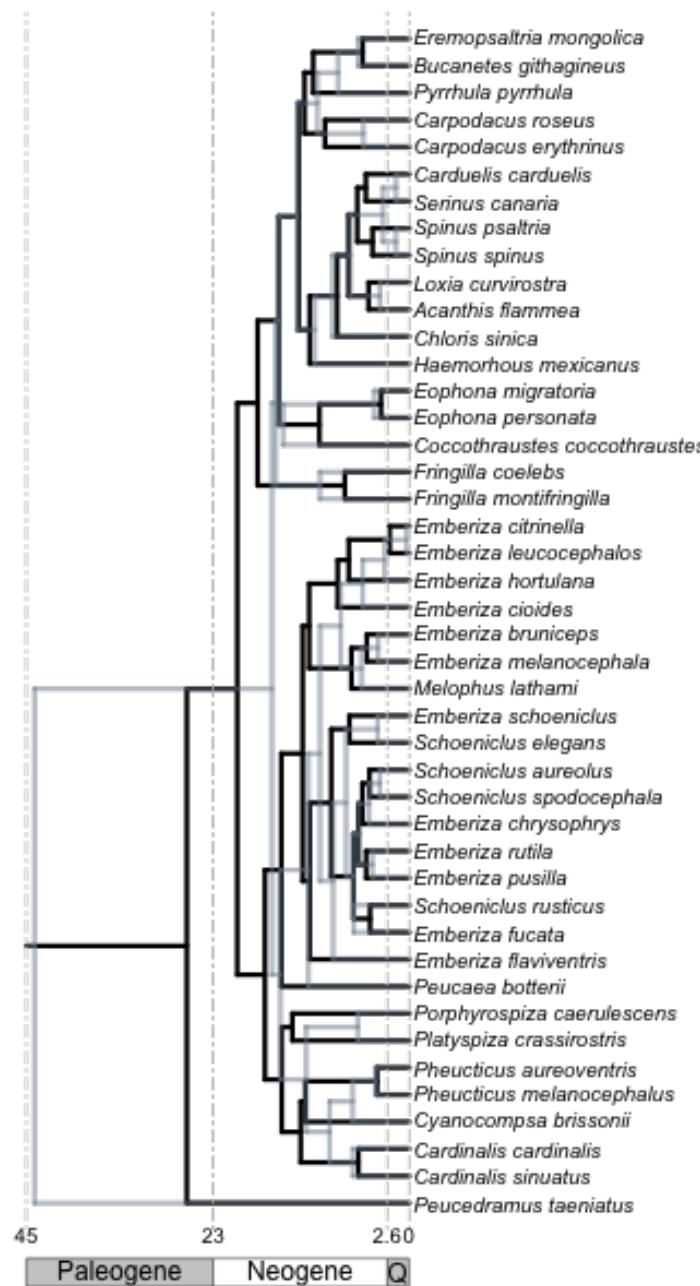


FIGURE 13. Cross validation of ninth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ.

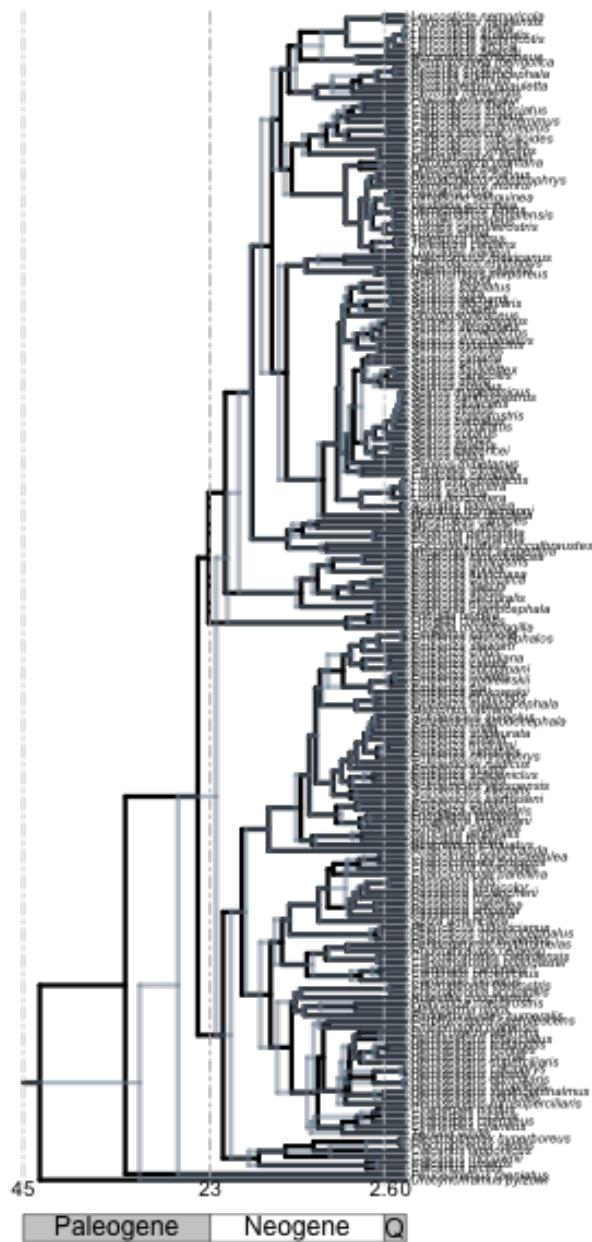
Jetz et al. 2012 - chronogram 1

FIGURE 14. Cross validation of tenth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADe. In each case, the tree is the same.