

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: 4362

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; it also generally requires specialized biological training, taxonomic domain
56 knowledge, and a non-negligible amount of research time, computational resources and
57 funding.

58 Here we present the DateLife project and its core software application, available as an
59 R package (Sanchez-Reyes et al., 2022), and as an online Rshiny interactive website at
60 www.datelife.org/query/, which captures data from published chronograms, and make these
61 data readily accessible to users for reuse and reanalysis. The software features key elements
62 for scientific reproducibility, such as a versioned, open and fully public chronogram database
63 (McTavish et al., 2015), age data stored in a computer readable format (Vos et al., 2012),
64 automated and programmatic ways of accessing the data (Stoltzfus et al., 2013) and
65 methods to summarize and compare age data.

66 **Description**

67 DateLife’s core software application consists of the R package `datelife`. Its latest
68 stable version – v0.6.2, is available from the CRAN repository (Sanchez-Reyes et al., 2022),
69 and relies on functionalities from various biological R packages: `ape` (Paradis, Claude, &

70 Strimmer, 2004), bold (Chamberlain et al., 2019), geiger (Harmon, Weir, Brock, Glor, &
71 Challenger, 2008), paleotree (Bapst, 2012), phyloch (Heibl, 2008), phylocomr (Ooms &
72 Chamberlain, 2018), phytools (Revell, 2012), rotl (Michonneau, Brown, & Winter, 2016),
73 and taxize (Chamberlain & Szöcs, 2013; Chamberlain et al., 2019). Figure 1 provides a
74 graphical summary of the three main steps of the DateLife algorithm: providing an input,
75 searching a chronogram database, and summarizing results from the search.

76 Processing an input

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

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

89 DateLife processes input scientific names using a Taxonomic Name Resolution Service
90 (TNRS), which increases the probability of correctly finding the queried taxon names in the
91 chronogram database. TNRS detects, corrects and standardizes name misspellings and typos,
92 variant spellings and authorities, and nomenclatural synonyms to a single taxonomic
93 standard (**boyle2013?**). DateLife implements TNRS using OpenTree’s unified taxonomy as
94 standard (Open Tree Of Life et al., 2016; Rees & Cranston, 2017), storing OpenTree’s

95 Taxonomy identification numbers for further processing.

96 The processed input taxon names are saved as an R object of a newly defined class
97 `datelifeQuery` that is used in the following steps. This object contains the standardized
98 names, the corresponding taxonomic id numbers, and the topology of the input tree if any
99 was provided.

100 **Searching the database**

101 DateLife's chronogram database latest version consist of 253 chronograms published in
102 187 different studies. It is curated from OpenTree's phylogenetic database, the Phylesystem,
103 which constitutes an open source of expert and peer-reviewed phylogenetic knowledge with
104 rich metadata (McTavish et al., 2015), which allows automatic and reproducible assembly of
105 our chronogram database. Datelife's chronogram database is navigable as an R data object
106 within the `datelife` R package. Published chronograms can be added to Phylesystem by
107 any user, at any time, and are immediately publicly available
108 (<https://tree.opentreeoflife.org/curator>). This facilitates an immediate update of DateLife's
109 chronogram database to include new chronogram data on a following search.

110 A DateLife search is implemented by matching processed taxon names provided by the
111 user, to tip labels in the chronogram database. Chronograms with at least two matching
112 taxon names on their tip labels are identified and pruned down to preserve only the matched
113 taxa. These matching pruned chronograms are referred to as source chronograms. Total
114 distance (in units of millions of years) between taxon pairs within each source chronogram
115 are stored as a patristic distance matrix (Figure 1). The matrix format speeds up extraction
116 of pairwise taxon ages of any queried taxa, as opposed to searching the ancestor node of a
117 pair of taxa in a “phylo” object or newick string. Finally, the patristic matrices are
118 associated to the study citation where the original chronogram was published, and stored as
119 an R object of the newly defined class `datelifeResult`.

120 **Summarizing search results**

121 Summary information is extracted from the `datelifeResult` object to inform
122 decisions for subsequent steps in the analysis workflow. Basic summary information available
123 to the user is:

- 124 1. The matching pruned chronograms as newick strings or “phylo” objects.
- 125 2. The ages of the root of all source chronograms. These ages can correspond to the age
126 of the most recent common ancestor (mrca) of the user’s group of interest if the source
127 chronograms have all taxa belonging to the group. If not, the root corresponds to the
128 mrca of a subgroup within the group of interest.
- 129 3. Study citations where original chronograms were published.
- 130 4. A report of input taxon names matches across source chronograms.
- 131 5. The source chronogram(s) with the most input taxon names.
- 132 6. Last but not least, single summary chronograms resulting from summarizing age data
133 from available source chronograms, as follows:

134 ***Choosing a topology.***— DateLife requires a tree topology to summarize age data
135 upon. We recommend that users provide a tree topology as input from the literature, or one
136 of their own making. If no topology is provided, DateLife automatically subsets one from the
137 OpenTree synthetic tree (Open Tree Of Life et al., 2019). Alternatively, DateLife can
138 reconstruct a tree with branch lengths proportional to substitution rates from a starting tree
139 topology using genetic data from the Barcode of Life Data System, BOLD (Ratnasingham &
140 Hebert, 2007), or combine topologies from source chronograms using a supertree approach.

141 ***Reconstructing branch lengths.***— DateLife starts by mining the BOLD database to
142 obtain genetic markers for the input taxa, and aligning them with MUSCLE (Edgar, 2004;
143 or MAFFT, Katoh, Asimenos, & Toh, 2009). Currently, branch length reconstruction is
144 performed with parsimony and the likelihood of the phylogenetic tree given a sequence

145 alignment is computed (Schliep, 2011).

146 **Combining source chronograms.**— To combine topologies from source
147 chronograms into a single summary (or supertree) topology, the DateLife algorithm starts by
148 identifying the source chronograms that form a grove, roughly, a sufficiently overlapping set
149 of taxa between trees, by implementing definition 2.8 for n-overlap from Ané et al. (2009).
150 In rare cases, a group of trees can have multiple groves. By default, DateLife chooses the
151 grove with the most taxa, however, the “criterion = trees” flag allows the user to choose the
152 grove with the most trees instead. If source chronograms do not form a grove, the supertree
153 reconstruction will fail.

154 **Congruifying nodes.**— Once with a chosen topology, DateLife applies the
155 congruification method (Eastman, Harmon, & Tank, 2013) to find nodes belonging to the
156 same clade across source chronograms, and extract the corresponding node ages from the
157 patristic distance matrices stored as `datelifeResult`. By definition, the matrices store total
158 distance (time from tip to tip), hence, node ages correspond to half the values stored in the
159 patristic distance matrices. A table of congruified node ages that can be used as secondary
160 calibrations is stored as a `congruifiedCalibrations` object.

161 **Summarizing congruified ages.**— For each congruent node, the pairwise distance
162 that traverse that node are summarized into a single summary matrix using classic statistics
163 (mean, median, minimum and maximum ages). These single summary ages per taxon pair
164 are then used as secondary calibrations to date the tree topology.

165 **Dating the tree topology.**— By default, DateLife implements the Branch Length
166 Adjuster (BLADJ) algorithm to obtain a fully dated topology. BLADJ fixes node ages that
167 have calibration data, and distributes time between nodes with no data evenly between
168 nodes with calibration data. This minimizes age variance in the resulting chronogram
169 (Campbell O. Webb, Ackerly, & Kembel, 2008). When there is conflict in ages between nodes
170 with calibration data, the algorithm ignores ages that are older than ages of parent nodes

171 and/or younger than ages from descendant nodes.

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

175 Summarized calibrations can be applied as secondary calibrations with different dating
176 methods currently supported within DateLife: MrBayes (Huelsenbeck & Ronquist, 2001;
177 Ronquist & Huelsenbeck, 2003), PATHd8 (Britton, Anderson, Jacquet, Lundqvist, &
178 Bremer, 2007), BLADJ (Campbell O. Webb et al., 2008; Campbell O. Webb & Donoghue,
179 2005), and treePL (Stephen A. Smith & O'Meara, 2012).

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

184 Benchmark

185 `datelife`’s code speed was tested on an Apple iMac with one 3.4 GHz Intel Core i5
186 processor. We registered variation in computing time of query processing and search through
187 the database relative to number of queried taxon names. Query processing time increases
188 roughly linearly with number of input taxon names, and increases considerably if Taxonomic
189 Name Resolution Service (TNRS) is activated. Up to ten thousand names can be processed
190 and searched in less than 30 minutes with the most time consuming settings. Once names
191 have been processed as described in methods, a name search through the chronogram
192 database can be performed in less than a minute, even with a very large number of taxon
193 names (Fig. 2). `datelife`’s code performance was evaluated with a set of unit tests designed
194 and implemented with the R package `testthat` (R Core Team, 2018) that were run both
195 locally with the `devtools` package (R Core Team, 2018), and on a public server –via GitHub,

196 using the continuous integration tool Travis CI (<https://travis-ci.org>). At present, unit tests
197 cover more than 40% of datelife's code (<https://codecov.io/gh/phylotastic/datelife>).

198

Case studies

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

203 **Small example**

204 We randomly chose 6 bird species related to the family Fringillidae of true finches. The
205 sample includes two species of cardinals: the black-thighed grosbeak – *Pheucticus tibialis*
206 and the crimson-collared grosbeak – *Rhodothraupis celaeno*; three species of buntings: the
207 yellowhammer – *Emberiza citrinella*, the pine bunting – *Emberiza leucocephalos* and the
208 yellow-throated bunting – *Emberiza elegans*; and one species of tanager, the vegetarian finch –
209 *Platyspiza crassirostris*.

210 Processing input names found that *Emberiza elegans* is synonym for *Schoeniclus*
211 *elegans* in the default reference taxonomy [Open Tree of Life Taxonomy v3.3, June 1, 2021].
212 For a detailed discussion on the state of the synonym refer to Avibase (Avibase, 2022;
213 Lepage, 2004; Lepage, Vaidya, & Guralnick, 2014). DateLife used the processed input names
214 to search the local chronogram database and found 9 matching chronograms in 6 different
215 studies. Three studies matched five input names (Barker, Burns, Klicka, Lanyon, & Lovette,
216 2015; Hedges, Marin, Suleski, Paymer, & Kumar, 2015; Jetz, Thomas, Joy, Hartmann, &
217 Mooers, 2012), one study matched four input names (Hooper & Price, 2017) and two studies
218 matched two input names (Barker, Burns, Klicka, Lanyon, & Lovette, 2013; Burns et al.,
219 2014). No studies matched all input names. Together, matching chronograms have 28 unique
220 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

247 program BLADJ (Fig. 5). Median summary chronograms are older and have wider variation
248 in maximum ages than chronograms obtained with SDM.

249 **Cross-validation test**

250 To perform a cross validation analysis of the DateLife workflow, we used resulting data
251 from the previous section (Casestudy: Real Life Application). We took individual tree
252 topologies from each of the 19 source chronograms found (Supplementary data XX). Then
253 we congruified age data of source chronograms from studies, and used this ages to date the
254 tree topology with the program BLADJ.

255 We found that node ages from original study, and ages estimated using all other age
256 data available are generally correlated (Supplementary Fig. 6). In 5 studies, more inclusive
257 nodes have older original ages, and less inclusive nodes have younger original ages than their
258 cross-validated age estimates. Accordingly, root ages are generally older in the original study
259 than estimated using cross-validated ages. Root ages were similar in original and cross
260 validated ages in three vases (Supplementary Fig. 7). Notably, chronograms have different
261 species sampling, hence roots are not comparable across studies. Yet, chronograms with a
262 higher sampling number can potentially inform the age of the root of chronograms with less
263 sampling.

264 **Discussion**

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

270 At the time of writing of this manuscript (May 05, 2022), `datelife`'s database has 253
271 chronograms, pulled entirely from OpenTree's database, the Phylesystem (McTavish et al.,

272 2015). A unique feature of OpenTree's Phylesystem is that the community can add new
273 state-of-the-art chronograms any time. As chronograms are added to Phylesystem, they are
274 incorporated into an updated `datelife`'s database that is assigned a new version number,
275 followed by a package release on CRAN. `datelife`'s chronogram database is updated as new
276 chronogram data is added to Phylesystem, at a minimum of once a month and a maximum
277 of every 6 months. Users can also upload new chronograms to OpenTree themselves, and
278 trigger an update of their local `datelife` database to incorporate the new chronograms, to
279 have them immediately available for analysis.

280 Incorporation of more chronograms into `datelife`'s database is crucial to improve its
281 services. One option to increase chronogram number in the database is the Dryad data
282 repository. Methods to automatically mine chronograms from Dryad could be designed and
283 implemented. However, Dryad's metadata system has no information to automatically detect
284 branch length units, and those would still need to be determined manually by a curator.

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

296 By default, `datelife` currently summarizes all source chronograms that overlap with
297 at least two species names. Users can exclude source chronograms if they have reasons to do

so. Strictly speaking, the best chronogram should reflect the real time of lineage divergence accurately and precisely. To our knowledge, there are no good measures to determine independently if a chronogram is better than another. Some measures that have been proposed are the proportion of lineage sampling and the number of calibrations used Magallón, Gómez-Acevedo, Sánchez-Reyes, & Hernández-Hernández (2015). Several characteristics of the data used for dating analyses as well as from the output chronogram itself, could be used to score quality of source chronograms. Some characteristics that are often cited in published studies as a measure of improved age estimates as compared to previously published estimates are: quality of alignment (missing data, GC content), lineage sampling (strategy and proportion), phylogenetic and dating inference method, number of fossils used as calibrations, support for nodes and ages, and magnitude of confidence intervals. DateLife provides an opportunity to capture concordance and conflict among date estimates, which can also be used as a metric for chronogram reliability.

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

Furthermore, chronograms can be obtained with primary fossil data or with molecular substitution rates obtained experimentally, which can deepen the already substantial

324 variation in time estimates between lineages, as observed from the comparison of source
325 chronograms in the Fringillidae example. This observation is often encountered in the
326 literature (see, for example, the ongoing debate about crown group age of angiosperms
327 (Barba-Montoya, Reis, Schneider, Donoghue, & Yang, 2018; Magallón et al., 2015; Ramshaw
328 et al., 1972; Sanderson & Doyle, 2001; Sauquet, Ramírez-Barahona, & Magallón, 2021). For
329 some studies, especially ones based on branch lengths (e.g., studies of species diversification,
330 timing of evolutionary events, phenotypic trait evolution), using a different chronogram may
331 return different results (Title & Rabosky, 2016). Stitching together these chronograms can
332 create a larger tree that uses information from multiple studies, but the effect of
333 uncertainties and errors at this level on downstream analyses is still largely unknown.

334 Summarizing chronograms might also imply summarizing fundamentally distinct
335 evolutionary hypotheses. For example, two different researchers working on the same clade
336 both carefully select and argument their choices of fossil calibrations. Still, if one researcher
337 decides a fossil will calibrate the ingroup of a clade, while another researcher uses the same
338 one to calibrate outside the clade, the resulting age estimates will often differ substantially,
339 as the placement of calibrations as stem or crown group is proved to deeply affect estimated
340 times of lineage divergence (Sauquet, 2013). Trying to summarize the resulting chronograms
341 into a single one using simple summary statistics can erase many types of relevant
342 information from the source chronograms. Accordingly, the prevailing view is that we should
343 favor time of lineage divergence estimates obtained from a single analysis, using fossil data as
344 primary sources of calibrations, and using fossils that have been widely discussed and
345 curated as calibrations to date other trees, making sure that all data used in the analysis
346 reflect a coherent evolutionary history (Antonelli et al., 2017). However, the exercise of
347 summarizing different chronograms has the potential to help getting a single global
348 evolutionary history for a lineage by putting together evidence from different hypothesis.
349 Choosing the elements of the chronograms that we are going to keep and the ones that we
350 are going to discard is key, since we are potentially loosing important parts of the

351 evolutionary history of a lineage that might only be reflected in source chronograms and not
352 on the summary chronogram (Sauquet et al., 2021).

353 Nonetheless, in ecology and conservation biology, incorporating at least some data on
354 lineage divergence times represents a relevant improvement for testing alternative hypothesis
355 using phylogenetic distance (Campbell O. Webb et al., 2008). Hence, we integrated into
356 datelife's workflow different ways of estimating node ages in the absence of calibrations and
357 branch length information for taxa lacking this information. "Making up" branch lengths is
358 an accepted practice in scientific publications: Jetz et al. (2012), created a time-calibrated
359 tree of all 9,993 bird species, where 67% had molecular data and the rest was simulated;
360 Rabosky et al. (2018) created a time-calibrated tree of 31,536 ray-finned fishes, of which only
361 37% had molecular data; Stephen A. Smith and Brown (2018) constructed a tree of 353,185
362 seed plants where only 23% had molecular data. Obviously, there are risks in this practice!
363 Taken to the extreme, one could make a fully resolved, calibrated tree of all modern and
364 extinct taxa using a single taxonomy and a single calibration with the polytomy resolution
365 and branch estimation methods. There has yet to be a thorough analysis of what can go
366 wrong when one extends inferences beyond the data in this way, so we urge caution; we also
367 urge readers to follow the example of many of the large tree papers cited above and make
368 carefully consider the statistical assumptions being made, and assess the consistency of the
369 results with prior work.

370

Conclusions

371 Divergence time information is key to many areas of evolutionary studies: trait
372 evolution, diversification, biogeography, macroecology and more. It is also crucial for science
373 communication and education, but generating chronograms is difficult, especially for those
374 who want to use phylogenies but who are not systematists, or do not have the time to
375 acquire and develop the necessary knowledge and data curation skills. Moreover, years of
376 primarily public funded research have resulted in vast amounts of chronograms that are

377 already available on scientific publications, but hidden to the public and scientific community
378 for reuse.

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

389 Availability

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

399 Supplementary Material

400 Code used to generate all versions of this manuscript, the biological examples, as well
401 as the benchmark of functionalities are available at datelifeMS1, datelife_examples, and

402 datelife_benchmark repositories in LLSR's GitHub account.

403 **Funding**

404 Funding was provided by the US National Science Foundation (NSF) grants
405 ABI-1458603 to Datelife project; DBI-0905606 to the National Evolutionary Synthesis Center
406 (NESCent), ABI-1458572 to the Phylotastic project, and ABI-1759846 to the Open Tree of
407 Life project.

408 **Acknowledgements**

409 The DateLife project was born as a prototype tool aiming to provide these services,
410 and was developed over a series of hackathons at the National Evolutionary Synthesis
411 Center, NC, USA (Stoltzfus et al., 2013). We thank colleagues from the O'Meara Lab at the
412 University of Tennessee Knoxville for suggestions, discussions and software testing. The late
413 National Evolutionary Synthesis Center (NESCent), which sponsored hackathons that led to
414 initial work on this project. The team that assembled **datelife**'s first proof of concept:
415 Tracy Heath, Jonathan Eastman, Peter Midford, Joseph Brown, Matt Pennell, Mike Alfaro,
416 and Luke Harmon. The Open Tree of Life project that provides the open, metadata rich
417 repository of trees used for **datelife**. The many scientists who publish their chronograms in
418 an open, reusable form, and the scientists who curate them for deposition in the Open Tree
419 of Life repository. The NSF for funding nearly all the above, in addition to the ABI grant
420 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
- Avibase. (2022). Yellow-throated Bunting. *Avibase - The World Bird Database*, (Online Resource). Retrieved from %7B<https://avibase.bsc-eoc.org/species.jsp?lang=EN&avibaseid=82D1EE0049D8D927%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).

- 448 Estimating Divergence Times in Large Phylogenetic Trees. *Systematic Biology*,
449 56(788777878), 741–752. <https://doi.org/10.1080/10635150701613783>
- 450 Burns, K. J., Shultz, A. J., Title, P. O., Mason, N. A., Barker, F. K., Klicka, J., ...
451 Lovette, I. J. (2014). Phylogenetics and diversification of tanagers (passeriformes:
452 Thraupidae), the largest radiation of neotropical songbirds. *Molecular
453 Phylogenetics and Evolution*, 75, 41–77.
- 454 Chamberlain, S. A., & Szöcs, E. (2013). taxize : taxonomic search and retrieval in R
455 [version 2; referees: 3 approved]. *F1000Research*, 2(191), 1–29.
456 <https://doi.org/10.12688/f1000research.2-191.v2>
- 457 Chamberlain, S. A., Szöcs, E., Foster, Z., Arendsee, Z., Boettiger, C., Ram, K., ...
458 Li, G. (2019). *taxize: Taxonomic information from around the web*. Retrieved
459 from <https://github.com/ropensci/taxize>
- 460 Claramunt, S., & Cracraft, J. (2015). A new time tree reveals earth history's imprint
461 on the evolution of modern birds. *Science Advances*, 1(11), e1501005.
- 462 Delsuc, F., Philippe, H., Tsagkogeorga, G., Simion, P., Tilak, M.-K., Turon, X., ...
463 Douzery, E. J. (2018). A phylogenomic framework and timescale for comparative
464 studies of tunicates. *BMC Biology*, 16(1), 1–14.
- 465 Eastman, J. M., Harmon, L. J., & Tank, D. C. (2013). Congruification: Support for
466 time scaling large phylogenetic trees. *Methods in Ecology and Evolution*, 4(7),
467 688–691. <https://doi.org/10.1111/2041-210X.12051>
- 468 Edgar, R. C. (2004). MUSCLE: Multiple sequence alignment with high accuracy and
469 high throughput. *Nucleic Acids Research*, 32(5), 1792–1797.
- 470 Felsenstein, J. (1985). Phylogenies and the Comparative Method. *The American
471 Naturalist*, 125(1), 1–15. Retrieved from <http://www.jstor.org/stable/2461605>
- 472 Gibb, G. C., England, R., Hartig, G., McLenachan, P. A., Taylor Smith, B. L.,
473 McComish, B. J., ... Penny, D. (2015). New zealand passerines help clarify the
474 diversification of major songbird lineages during the oligocene. *Genome Biology*

- 475 and *Evolution*, 7(11), 2983–2995.
- 476 Harmon, L., Weir, J., Brock, C., Glor, R., & Challenger, W. (2008). GEIGER:
477 investigating evolutionary radiations. *Bioinformatics*, 24, 129–131.
- 478 Hedges, S. B., Marin, J., Suleski, M., Paymer, M., & Kumar, S. (2015). Tree of life
479 reveals clock-like speciation and diversification. *Molecular Biology and Evolution*,
480 32(4), 835–845. <https://doi.org/10.1093/molbev/msv037>
- 481 Heibl, C. (2008). *PHYLOCH: R language tree plotting tools and interfaces to diverse*
482 *phylogenetic software packages*. Retrieved from
483 <http://www.christophheibl.de/Rpackages.html>
- 484 Hooper, D. M., & Price, T. D. (2017). Chromosomal inversion differences correlate
485 with range overlap in passerine birds. *Nature Ecology & Evolution*, 1(10), 1526.
- 486 Huelsenbeck, J. P., & Ronquist, F. (2001). MRBAYES: Bayesian inference of
487 phylogenetic trees. *Bioinformatics*, 17(8), 754–755.
488 <https://doi.org/10.1093/bioinformatics/17.8.754>
- 489 Jetz, W., Thomas, G., Joy, J. J. B., Hartmann, K., & Mooers, A. (2012). The global
490 diversity of birds in space and time. *Nature*, 491(7424), 444–448.
491 <https://doi.org/10.1038/nature11631>
- 492 Katoh, K., Asimenos, G., & Toh, H. (2009). Multiple alignment of DNA sequences
493 with MAFFT. In *Bioinformatics for DNA sequence analysis* (pp. 39–64).
494 Springer.
- 495 Laubichler, M. D., & Maienschein, J. (2009). *Form and function in developmental*
496 *evolution*. Cambridge University Press.
- 497 Lepage, D. (2004). *Avibase: The world bird database*. Bird Studies Canada.
- 498 Lepage, D., Vaidya, G., & Guralnick, R. (2014). Avibase—a database system for
499 managing and organizing taxonomic concepts. *ZooKeys*, (420), 117.
- 500 Magallon, S., & Sanderson, M. J. (2001). Absolute diversification rates in angiosperm
501 clades. *Evolution*, 55(9), 1762–1780.

- 502 Magallón, S. (2010). Using fossils to break long branches in molecular dating: A
503 comparison of relaxed clocks applied to the origin of angiosperms. *Systematic*
504 *Biology*, 59(4), 384–399.
- 505 Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L. L., & Hernández-Hernández, T.
506 (2015). A metacalibrated time-tree documents the early rise of flowering plant
507 phylogenetic diversity. *New Phytologist*, 207(2), 437–453.
- 508 McTavish, E. J., Hinchliff, C. E., Allman, J. F., Brown, J. W., Cranston, K. A.,
509 Holder, M. T., . . . Smith, S. A. (2015). Phylesystem: A git-based data store for
510 community-curated phylogenetic estimates. *Bioinformatics*, 31(17), 2794–2800.
- 511 Michonneau, F., Brown, J. W., & Winter, D. J. (2016). rotl: an R package to interact
512 with the Open Tree of Life data. *Methods in Ecology and Evolution*, 7(12),
513 1476–1481. <https://doi.org/10.1111/2041-210X.12593>
- 514 Morlon, H. (2014). Phylogenetic approaches for studying diversification. *Ecology*
515 *Letters*, 17(4), 508–525. <https://doi.org/10.1111/ele.12251>
- 516 Ooms, J., & Chamberlain, S. (2018). *Phylocomr: Interface to 'phylocom'*. Retrieved
517 from <https://CRAN.R-project.org/package=phylocomr>
- 518 Open Tree Of Life, Redelings, B., Cranston, K. A., Allman, J., Holder, M. T., &
519 McTavish, E. J. (2016). Open Tree of Life APIs v3.0. *Open Tree of Life Project*,
520 (Online Resources). Retrieved from
521 <https://github.com/OpenTreeOfLife/germinator/wiki/Open-Tree-of-Life->
522 [Web-APIs](#)
- 523 Open Tree Of Life, Redelings, B., Sánchez Reyes, L. L., Cranston, K. A., Allman, J.,
524 Holder, M. T., & McTavish, E. J. (2019). Open tree of life synthetic tree v12.3.
525 *Zenodo*. Retrieved from <https://doi.org/10.5281/zenodo.3937742>
- 526 Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and
527 evolution in R language. *Bioinformatics*, 20(2), 289–290.
- 528 Posadas, P., Crisci, J. V., & Katinas, L. (2006). Historical biogeography: A review of

- 529 its basic concepts and critical issues. *Journal of Arid Environments*, 66(3),
530 389–403.
- 531 Price, T. D., Hooper, D. M., Buchanan, C. D., Johansson, U. S., Tietze, D. T.,
532 Alström, P., ... others. (2014). Niche filling slows the diversification of himalayan
533 songbirds. *Nature*, 509(7499), 222.
- 534 R Core Team. (2018). *R: a language and environment for statistical computing*.
535 Vienna, Austria: R Foundation for Statistical Computing.
- 536 Rabosky, D. L., Chang, J., Title, P. O., Cowman, P. F., Sallan, L., Friedman, M., ...
537 others. (2018). An inverse latitudinal gradient in speciation rate for marine fishes.
538 *Nature*, 559(7714), 392.
- 539 Ramshaw, J., Richardson, D., Mealyard, B., Brown, R., Richardson, M., Thompson,
540 E., & Boulter, D. (1972). The time of origin of the flowering plants determined by
541 using amino acid sequence data of cytochrome c. *New Phytologist*, 71(5), 773–779.
- 542 Ratnasingham, S., & Hebert, P. D. (2007). BOLD: The barcode of life data system
543 (<http://www.Barcodinglife.org>). *Molecular Ecology Notes*, 7(3), 355–364.
- 544 Rees, J. A., & Cranston, K. (2017). Automated assembly of a reference taxonomy for
545 phylogenetic data synthesis. *Biodiversity Data Journal*, (5).
- 546 Revell, L. J. (2012). Phytools: An r package for phylogenetic comparative biology
547 (and other things). *Methods in Ecology and Evolution*, 3, 217–223.
- 548 Ronquist, F., & Huelsenbeck, J. P. (2003). MrBayes 3: Bayesian phylogenetic
549 inference under mixed models. *Bioinformatics*, 19(12), 1572–1574.
550 <https://doi.org/10.1093/bioinformatics/btg180>
- 551 Sanchez-Reyes, L. L., O'Meara, B., Eastman, J., Heath, T., Wright, A., Schliep, K.,
552 ... Alfaro, M. (2022). datelife: Scientific Data on Time of Lineage Divergence for
553 Your Taxa. *R Package Version 0.6.2*. Retrieved from
554 <https://doi.org/10.5281/zenodo.593938>
- 555 Sanderson, M. J. (2003). r8s: Inferring absolute rates of molecular evolution and

- 556 divergence times in the absence of a molecular clock. *Bioinformatics*, 19(2),
557 301–302.
- 558 Sanderson, M. J., & Doyle, J. A. (2001). Sources of error and confidence intervals in
559 estimating the age of angiosperms from rbcL and 18S rDNA data. *American*
560 *Journal of Botany*, 88(8), 1499–1516.
- 561 Sauquet, H. (2013). A practical guide to molecular dating. *Comptes Rendus Palevol*,
562 12(6), 355–367.
- 563 Sauquet, H., Ramírez-Barahona, S., & Magallón, S. (2021). *The age of flowering*
564 *plants is unknown*.
- 565 Schenk, J. J. (2016). Consequences of secondary calibrations on divergence time
566 estimates. *PLoS ONE*, 11(1). <https://doi.org/10.1371/journal.pone.0148228>
- 567 Schliep, K. P. (2011). Phangorn: Phylogenetic analysis in r. *Bioinformatics*, 27(4),
568 592–593.
- 569 Smith, Stephen A., & Brown, J. W. (2018). Constructing a broadly inclusive seed
570 plant phylogeny. *American Journal of Botany*, 105(3), 302–314.
- 571 Smith, Stephen A., & O'Meara, B. C. (2012). TreePL: Divergence time estimation
572 using penalized likelihood for large phylogenies. *Bioinformatics*, 28(20),
573 2689–2690. <https://doi.org/10.1093/bioinformatics/bts492>
- 574 Stoltzfus, A., Lapp, H., Matasci, N., Deus, H., Sidlauskas, B., Zmasek, C. M., ...
575 Jordan, G. (2013). Phylotastic! Making tree-of-life knowledge accessible, reusable
576 and convenient. *BMC Bioinformatics*, 14.
577 <https://doi.org/10.1186/1471-2105-14-158>
- 578 Title, P. O., & Rabosky, D. L. (2016). Do Macrophylogenies Yield Stable
579 Macroevolutionary Inferences? An Example from Squamate Reptiles. *Systematic*
580 *Biology*, syw102. <https://doi.org/10.1093/sysbio/syw102>
- 581 Vos, R. A., Balhoff, J. P., Caravas, J. A., Holder, M. T., Lapp, H., Maddison, W. P.,
582 ... others. (2012). NeXML: Rich, extensible, and verifiable representation of

- 583 comparative data and metadata. *Systematic Biology*, 61(4), 675–689.
- 584 Webb, C. (2000). Exploring the Phylogenetic Structure of Ecological Communities :
585 An Example for Rain Forest Trees. *The American Naturalist*, 156(2), 145–155.
- 586 Webb, Campbell O., Ackerly, D. D., & Kembel, S. W. (2008). Phylocom: Software for
587 the analysis of phylogenetic community structure and trait evolution.
- 588 *Bioinformatics*, 24(18), 2098–2100.
589 <https://doi.org/10.1093/bioinformatics/btn358>
- 590 Webb, Campbell O., & Donoghue, M. J. (2005). Phylomatic: Tree assembly for
591 applied phylogenetics. *Molecular Ecology Notes*, 5(1), 181–183.

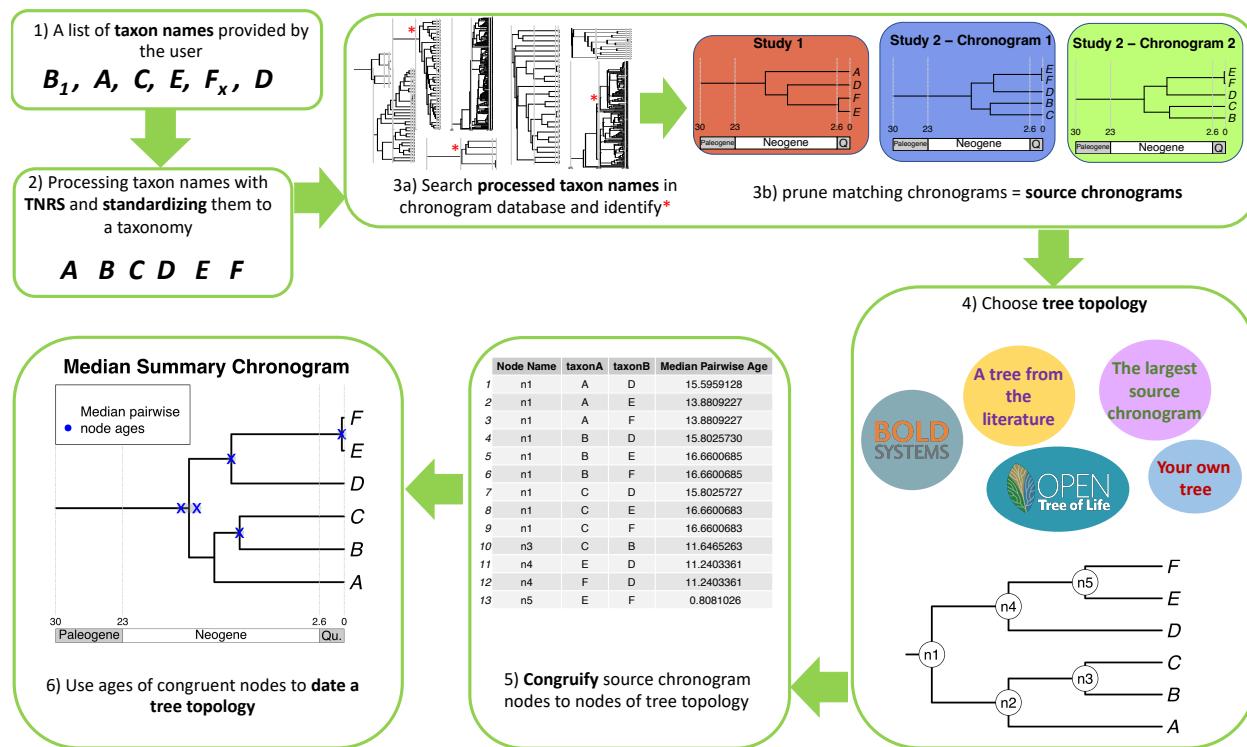


FIGURE 1. Stylized DateLife workflow. This shows the general workflows and analyses that can be performed with `datelife`, via the R package or through the website at <http://www.datelife.org/>. Details on the functions involved on each workflow are shown in `datelife`'s R package vignette.

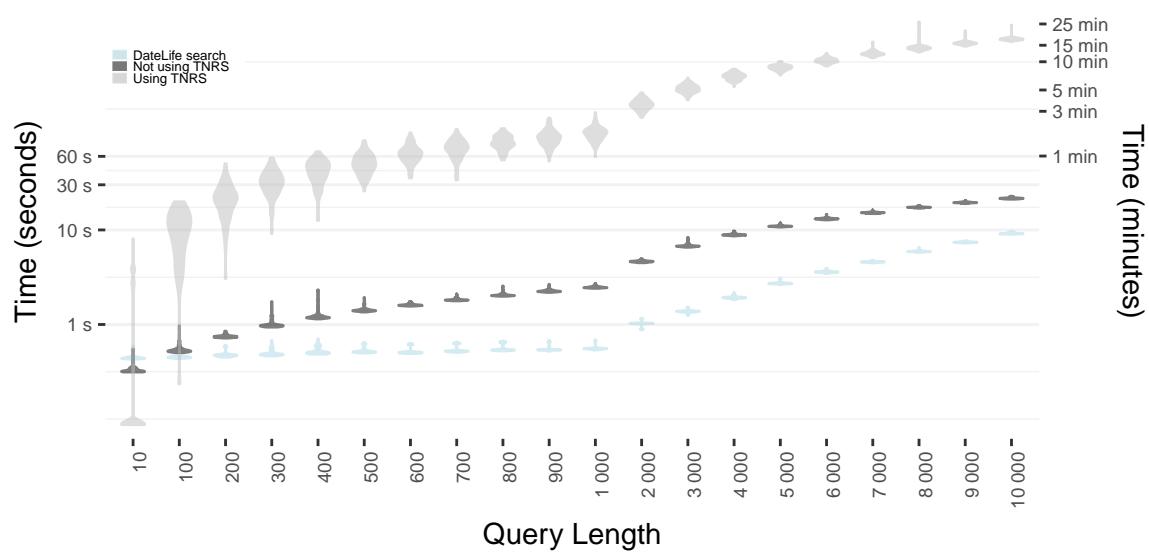


FIGURE 2. Computation time of query processing and search across **datelife**'s chronogram database relative to number of input taxon names. We sampled N names from the class Aves for each cohort 100 times and then performed a search with query processing not using the Taxon Names Resolution Service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

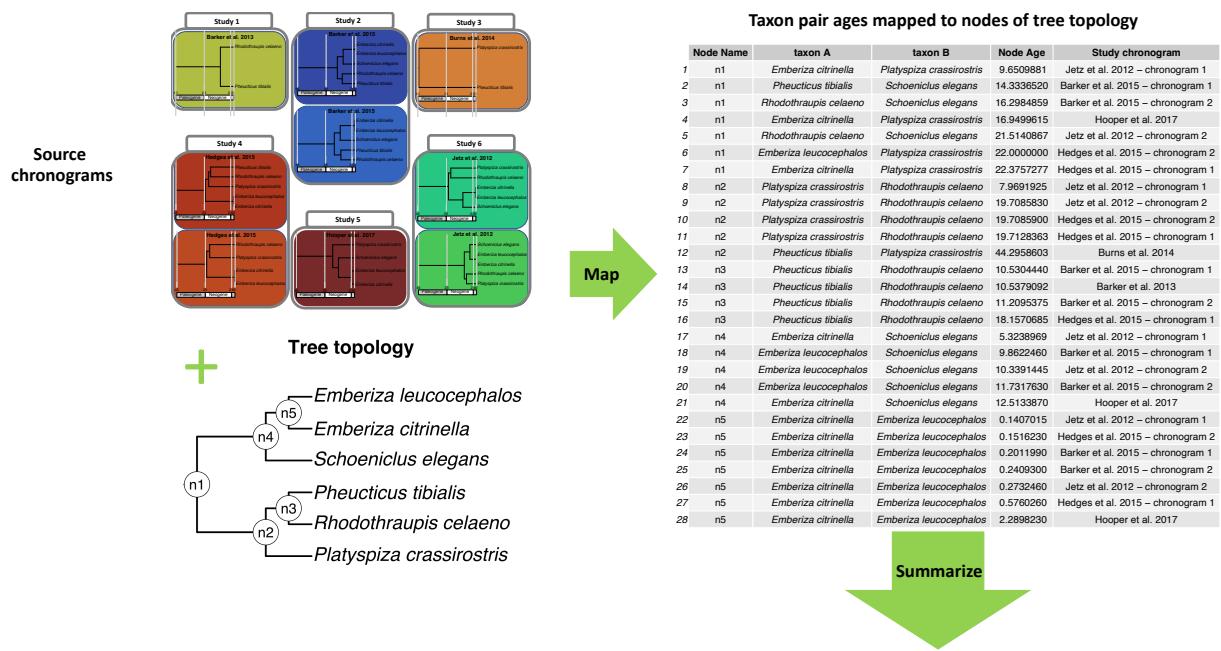


FIGURE 3. Age data results of a DateLife search of a small sample of 6 bird species within the Passeriformes. Input names were found across 9 chronograms within 6 independent studies (Barker et al. (2012), Barker et al. (2015), Burns et al. (2014), Hedges et al. (2015), Hooper and Price (2017), Jetz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

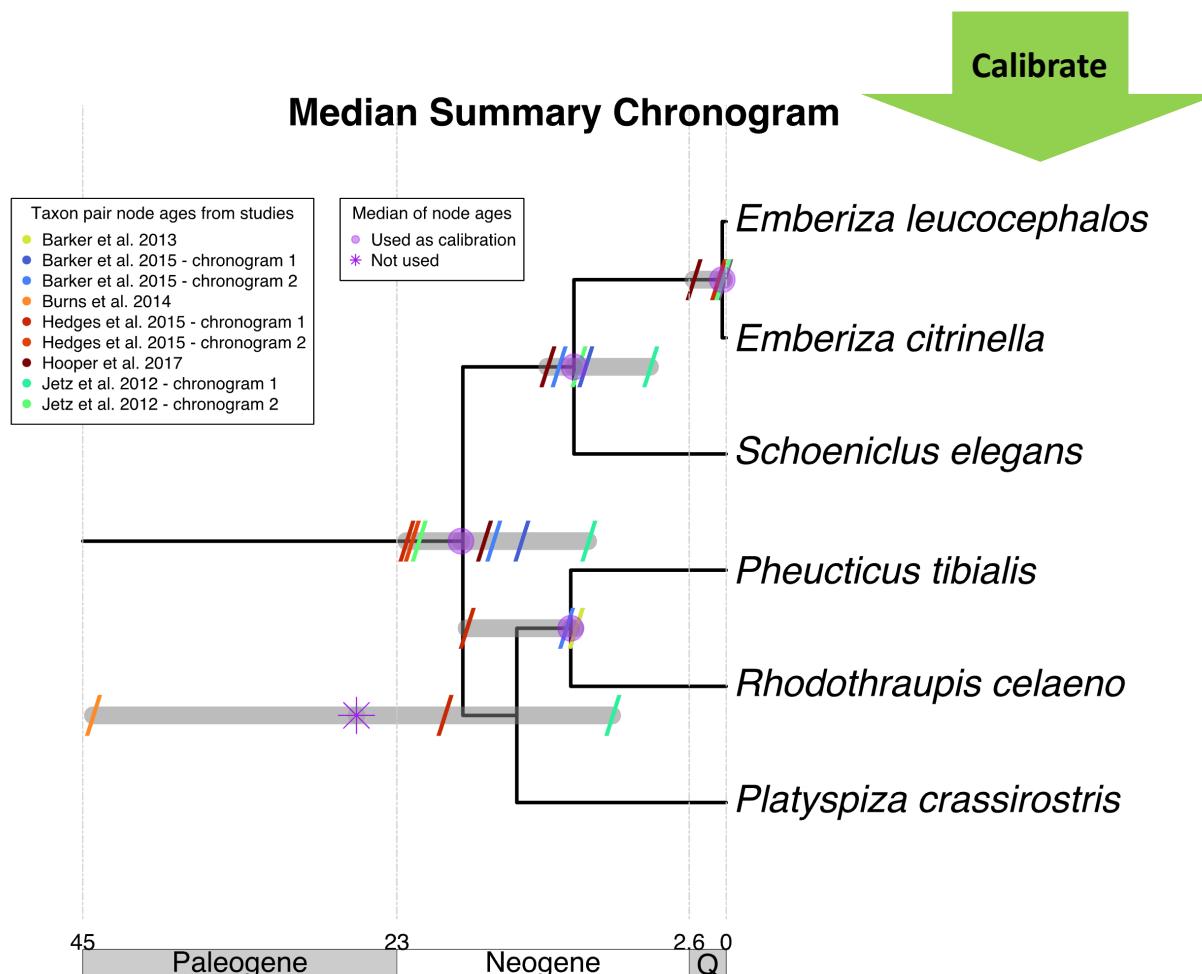


FIGURE 4. Summarized age data is used as secondary calibrations to date a tree topology as a summary chronogram.

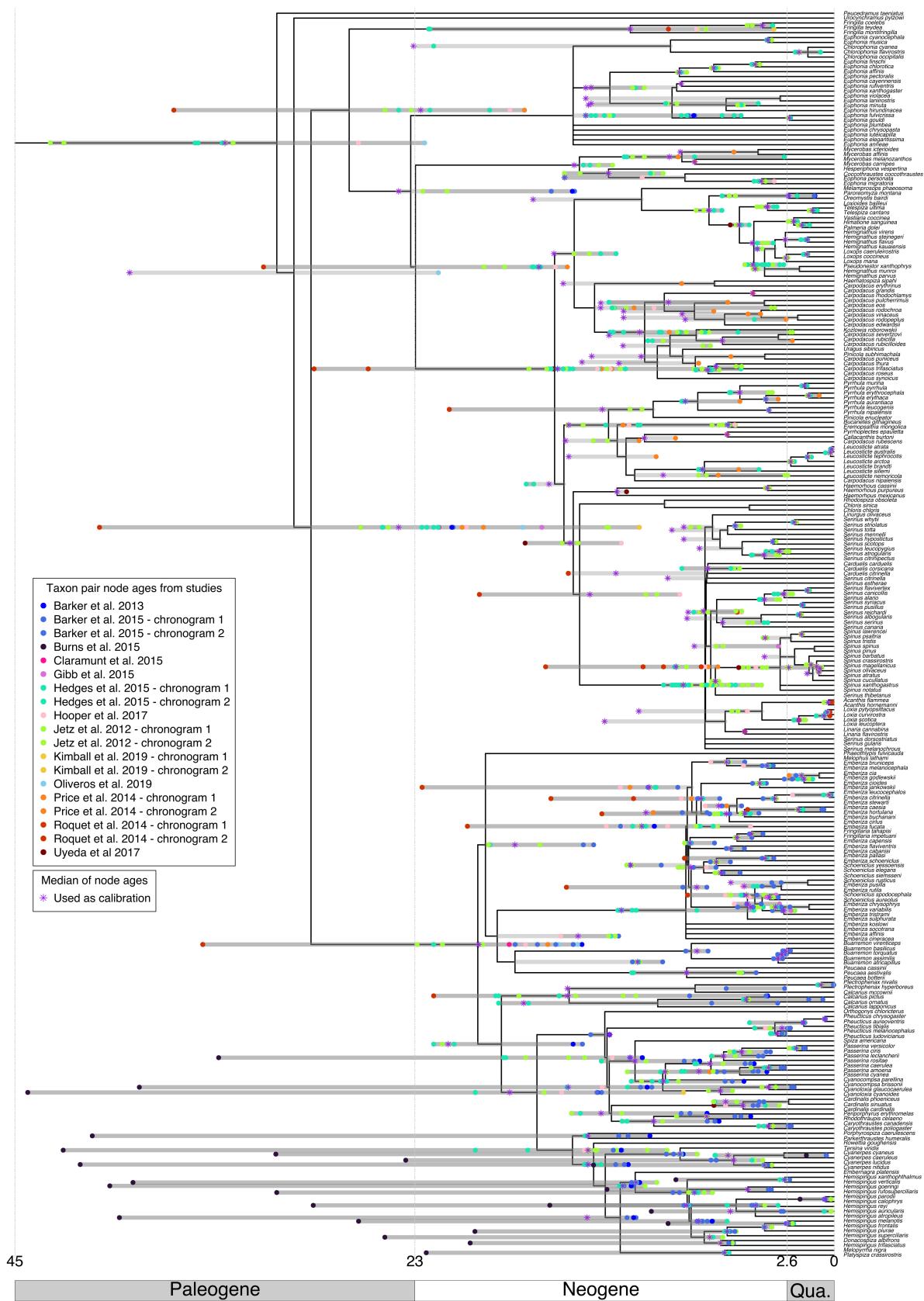


FIGURE 5. Fringillidae median summary chronogram generated with DateLife. It has 256 tips and 233 nodes.

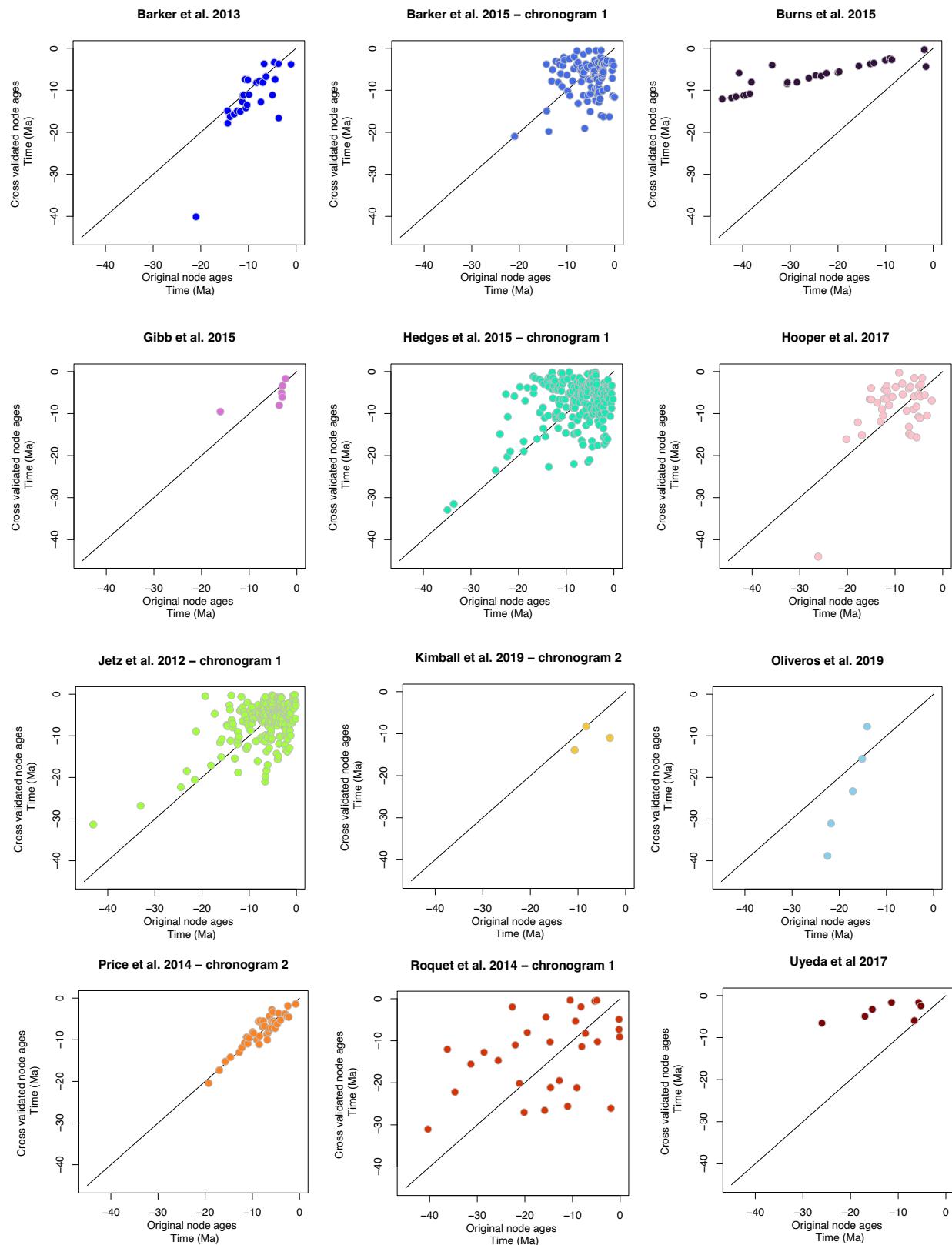


FIGURE 6. Results from cross validation analysis.

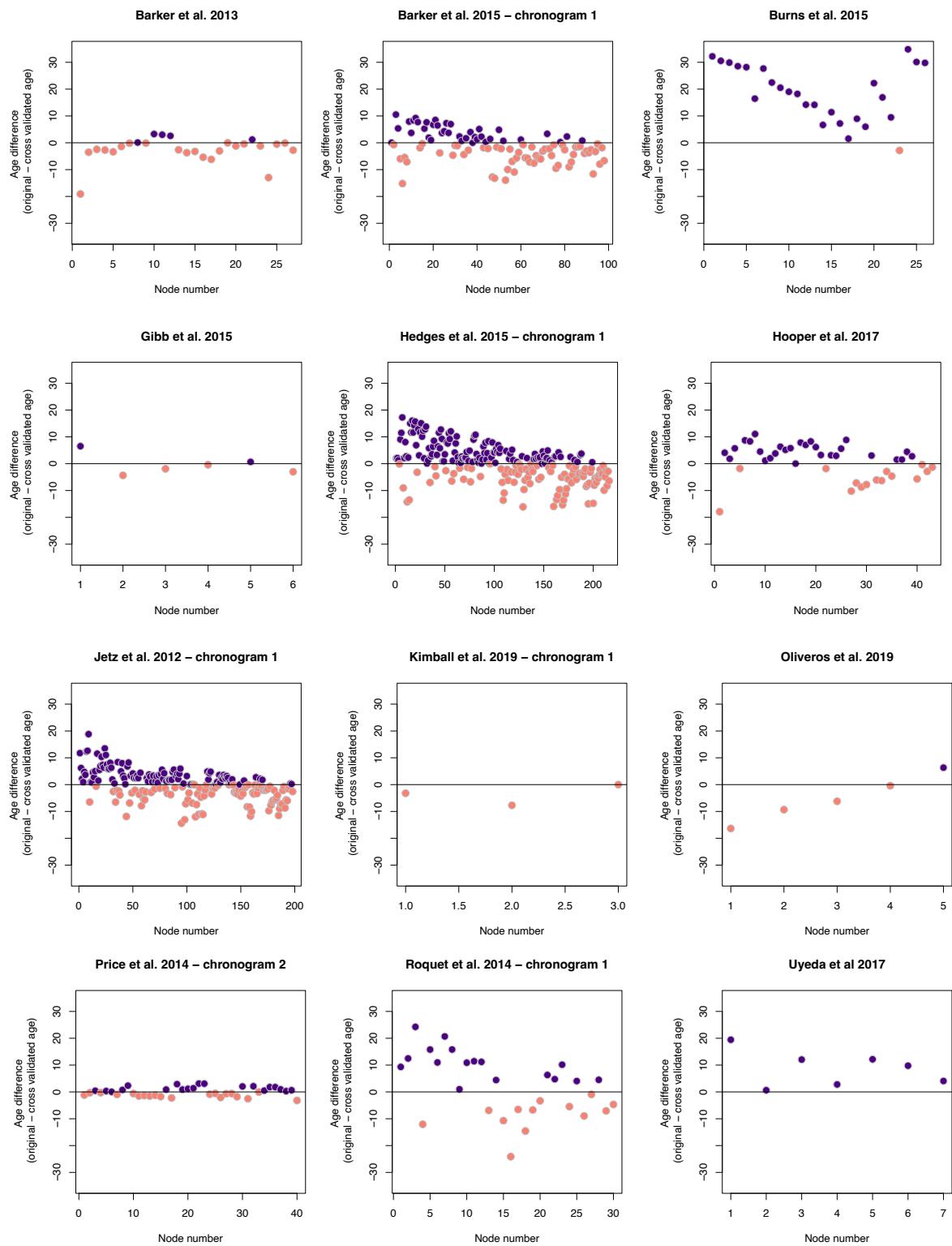


FIGURE 7. Results from cross validation analysis.

Barker et al. 2015 - chronogram 1



FIGURE 8. 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 9. 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, i.e., the same for all the nodes.

Burns et al. 2015

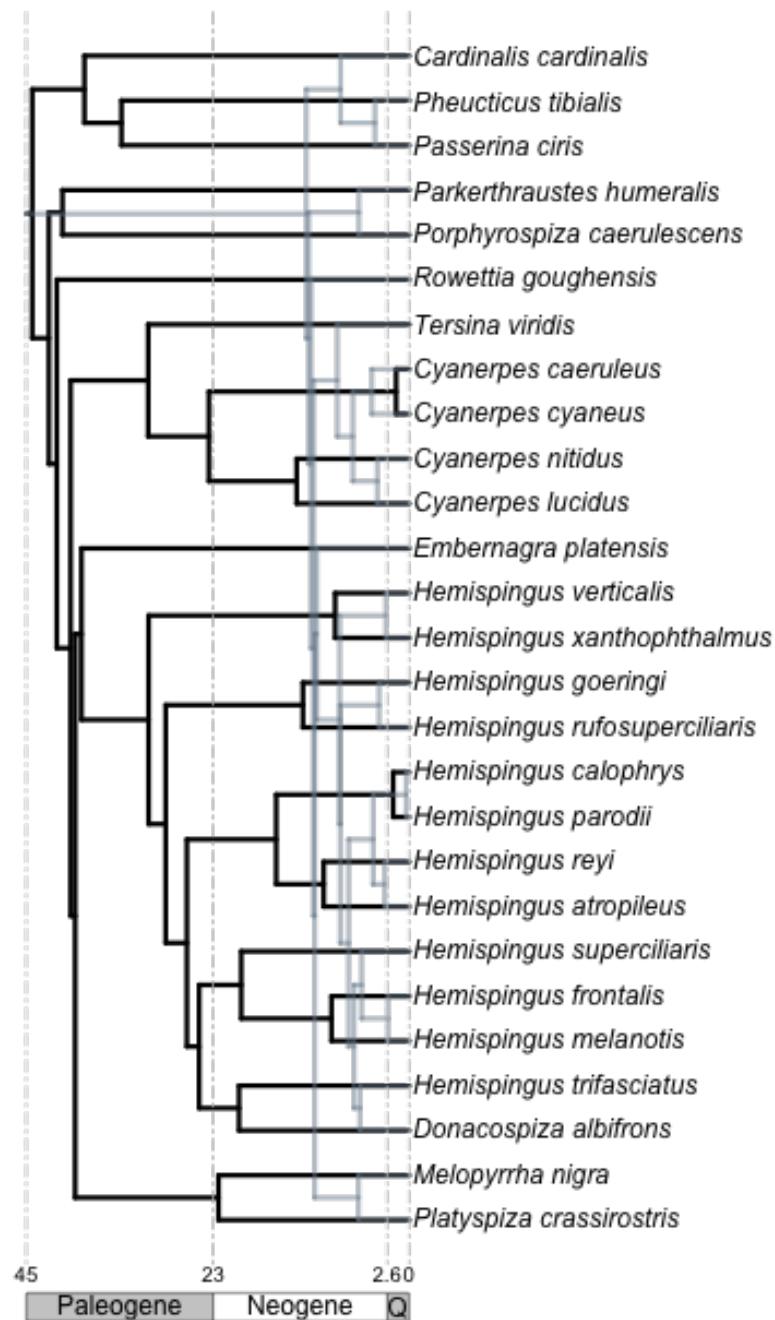


FIGURE 10. Cross validation of fourth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ, i.e., the same for all the nodes.

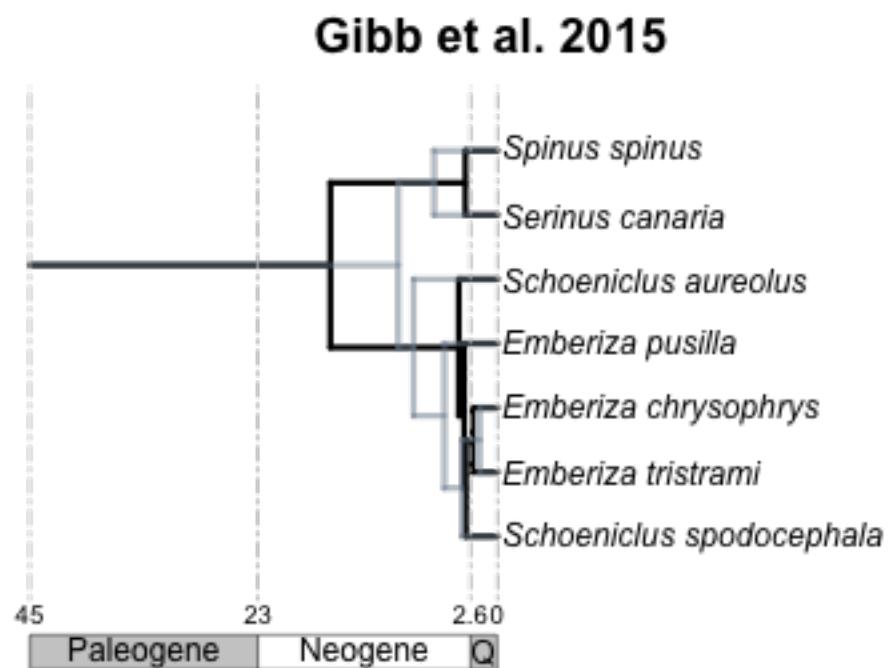


FIGURE 11. Cross validation of sixth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the same tree topology dated with BLADJ using node ages from all other source chronograms as secondary calibrations.

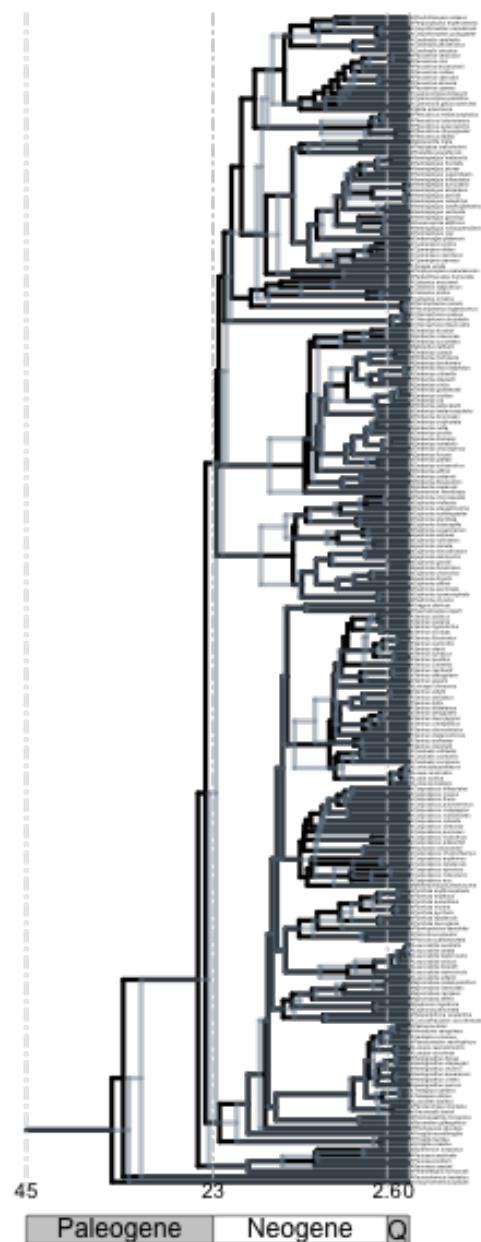
Hedges et al. 2015 - chronogram 1

FIGURE 12. 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

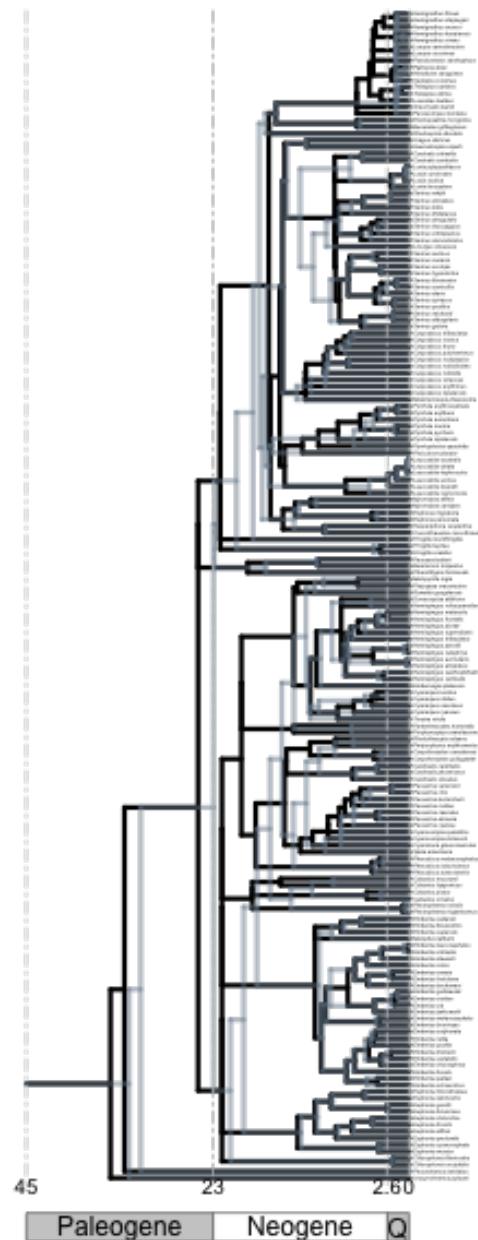
Hedges et al. 2015 - chronogram 2

FIGURE 13. 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

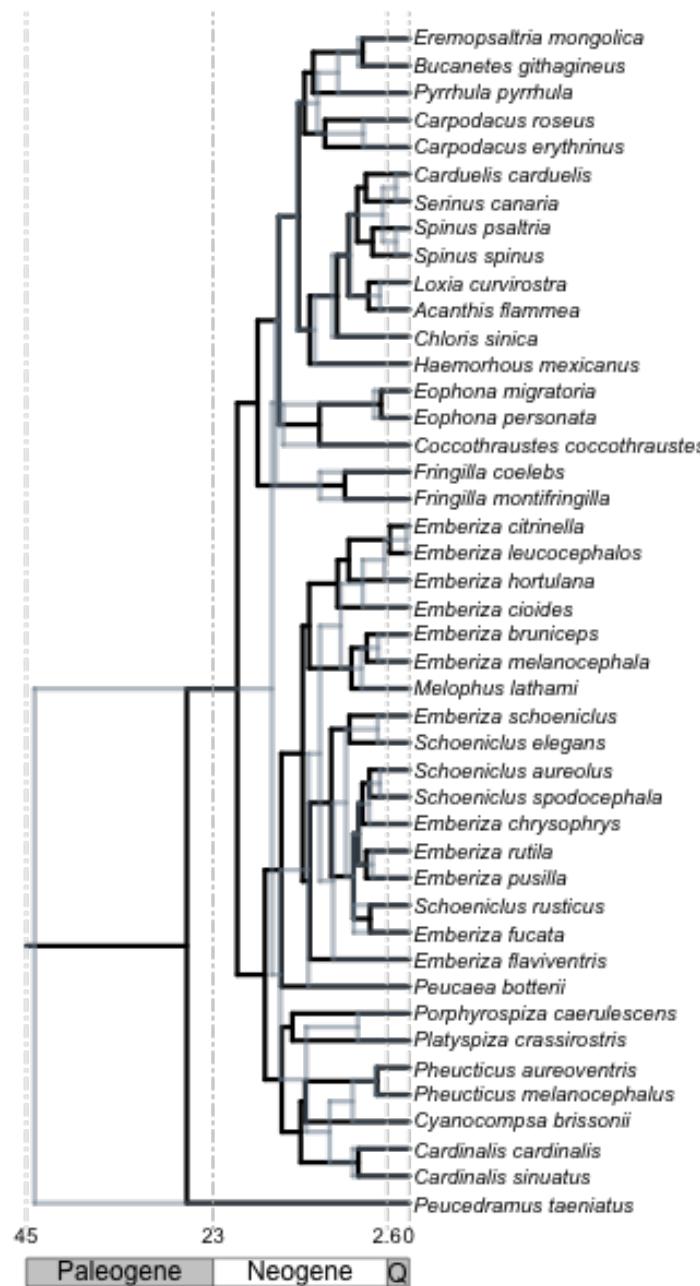


FIGURE 14. Cross validation of ninth source chronogram. The chronogram shown in black corresponds to the dates published in the original study. The gray chronogram corresponds to the dates calculated with BLADJ in our analysis. The tree is rooted on the left and branches to the right.

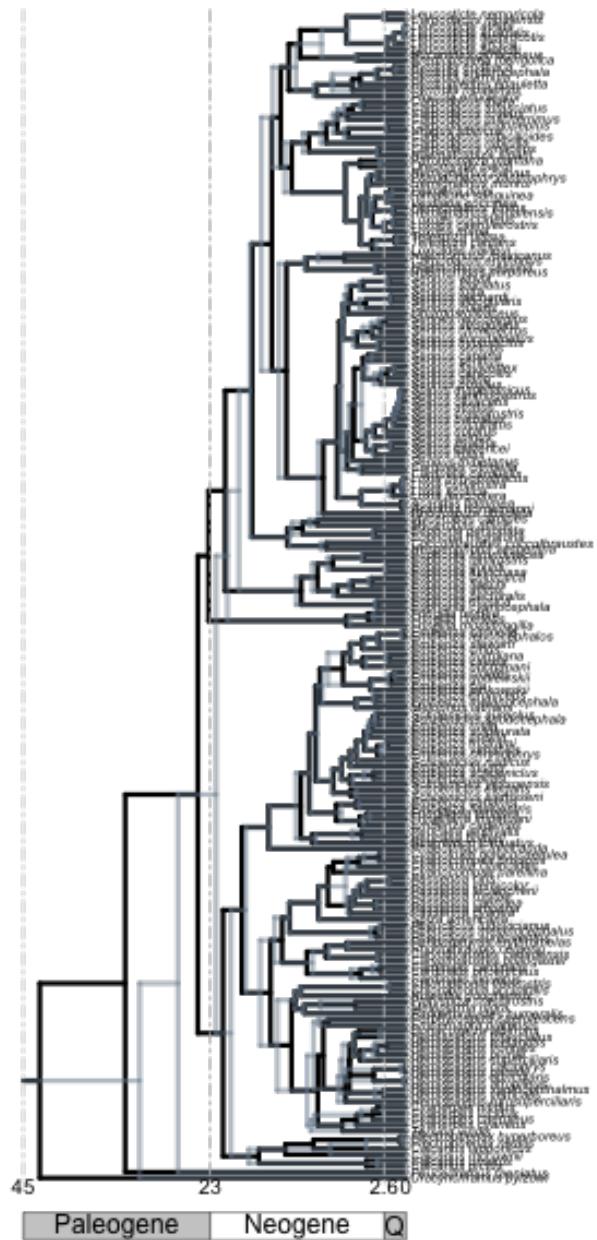
Jetz et al. 2012 - chronogram 1

FIGURE 15. 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, but the dates are different.