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

2          Luna L. Sánchez Reyes[1,2], Emily Jane McTavish[1], & Brian O'Meara[2]

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.

15        Correspondence concerning this article should be addressed to Luna L. Sánchez Reyes, .

16   E-mail: sanchez.reyes.luna@gmail.com

<sub>17</sub>                                        Abstract

<sub>18</sub>   Date estimates for times of evolutionary divergences are key data for research in the natural

<sub>19</sub>   sciences. These estimates also provide valuable information for for education, science

<sub>20</sub>   communication and policy decisions. Although achieving a high-quality reconstruction of a

<sub>21</sub>   phylogenetic tree with branch lengths proportional to absolute time (chronogram), is a

<sub>22</sub>   difficult and time-consuming task, the increased availability of fossil and molecular data, and

<sub>23</sub>   time-efficient analytical techniques has resulted in many recent publications of large

<sub>24</sub>   chronograms for a large number and wide diversity of organisms. When these estimates are

<sub>25</sub>   shared in public, open databases this wealth of expertly-curated and peer-reviewed data on

<sub>26</sub>   time of evolutionary origin is exposed in a programatic and reusable way. Intensive and

<sub>27</sub>   localized efforts have improved data sharing practices, as well as incentivizited open science

<sub>28</sub>   in biology. Here we present DateLife, a service implemented as an R package and an Rshiny

<sub>29</sub>   website application available at www.datelife.org/query/, that provides functionalities for

<sub>30</sub>   efficient and easy finding, summary, reuse, and reanalysis of expert, peer-reviewed, public

<sub>31</sub>   data on time of evolutionary origin. The main DateLife workflow constructs a chronogram

<sub>32</sub>   for any given combination of taxon names, by searching a local chronogram database

<sub>33</sub>   constructed and curated from the Open Tree of Life Phylesystem phylogenetic database,

<sub>34</sub>   which incorporates phylogenetic data from TreeBASE database as well. We implement and

<sub>35</sub>   test methods for summarizing time data from multiple source chronograms using supertree

<sub>36</sub>   and congruification algorithms, and using age data extracted from source chronograms as

<sub>37</sub>   secondary calibration points to add branch lengths proportional to absolute time to a tree

<sub>38</sub>   topology. DateLife will be useful to increase awereness on the existing variation in expert

<sub>39</sub>   time of divergence data, and can foster exploration of the effect of alternative divergence

<sub>40</sub>   time hypothesis on the results of analyses, providing a framework for a more informed

<sub>41</sub>   interpretation of evolutionary results.

<sub>42</sub>      *Keywords:* Tree; Phylogeny; Scaling; Dating; Ages; Divergence times; Open Science;

43 Congruification; Supertree; Calibrations; Secondary calibrations

44       Word count: 4204

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

## Introduction

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

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

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

## Description

The DateLife algorithm is fully implemented using the R language. The latest stable version of the R package `datelife` is available from the CRAN repository (v0.6.2; Sanchez-Reyes et al. (2022)), and relies on functionalities from various biological R packages: 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.

## 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 withing 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.

<sup>168</sup> **Benchmark**

<sup>169</sup>        `datelife`'s code speed was tested on an Apple iMac with one 3.4 GHz Intel Core i5

<sup>170</sup> processor. We registered variation in computing time of query processing and search through

<sup>171</sup> the database relative to number of queried taxon names. Query processing time increases

<sup>172</sup> roughly linearly with number of input taxon names, and increases considerably if TNRS is

<sup>173</sup> activated. Up to ten thousand names can be processed and searched in less than 30 minutes

<sup>174</sup> with the most time consuming settings. Once names have been processed as described in

<sup>175</sup> methods, a name search through the chronogram database can be performed in less than a

<sup>176</sup> minute, even with a very large number of taxon names (Fig. **??**). `datelife`'s code

<sup>177</sup> performance was evaluated with a set of unit tests designed and implemented with the R

<sup>178</sup> package testthat (R Core Team, 2018) that were run both locally with the devtools package

<sup>179</sup> (R Core Team, 2018), and on a public server –via GitHub, using the continuous integration

<sup>180</sup> tool Travis CI (https://travis-ci.org). At present, unit tests cover more than 40% of

<sup>181</sup> `datelife`'s code (https://codecov.io/gh/phylotastic/datelife).

<sup>182</sup> **Case study**

<sup>183</sup>        We illustrate the DateLife algorithm using a group within the Passeriform birds

<sup>184</sup> encompassing the family of true finches, Fringillidae and allies as case study. The first

<sup>185</sup> example analyses 6 bird species and shows all steps of the algorithm. The second example is

<sup>186</sup> a real life application

<sup>187</sup> **Small example**

<sup>188</sup>        We chose 6 bird species associated to true finches at random. The sample includes two

<sup>189</sup> species of cardinals: the black-thighed grosbeak – *Pheucticus tibialis* and the crimson-collared

<sup>190</sup> grosbeak – *Rhodothraupis celaeno*; three species of buntings: the yellowhammer – *Emberiza*

<sup>191</sup> *citrinella*, the pine bunting – *Emberiza leucocephalos* and the yellow-throated bunting –

<sup>192</sup> *Emberiza elegans*; and one species of tanager, the vegetarian finch – *Platyspiza crassirostris*.

193    Processing input names found that *Emberiza elegans* is synonym for *Schoeniclus*

194 *elegans* in the reference taxonomy. DateLife used the processed input names to search the

195 local chronogram database and found 9 matching chronograms in 6 different studies. Three

196 studies matched five input names (Barker, Burns, Klicka, Lanyon, & Lovette, 2015; Hedges,

197 Marin, Suleski, Paymer, & Kumar, 2015; Jetz, Thomas, Joy, Hartmann, & Mooers, 2012),

198 one study matched four input names (Hooper & Price, 2017) and two studies matched two

199 input names (Barker, Burns, Klicka, Lanyon, & Lovette, 2013; Burns et al., 2014). No

200 studies matched all input names. Together, matching chronograms have 28 unique age data

201 points. All nodes have age data. As fixed tree topology, DateLife used OpenTree's synthetic

202 tree as default and mapped age data to nodes in the tree. As expected, more inclusive nodes

203 (e.g., node "n1") have more age data than less inclusive nodes (e.g., node "n5"). The

204 processing step allowd discovering five data points for node "n4" that would not have had any

205 data otherwise. Age summary statistics per node were calculated and tested as secondary

206 calibrations to date the tree topology using the BLADJ algorithm. Age data for node "n2"

207 was excluded as final calibration because it is older than age data of a more inclusive node.


## Real life application

209    A college educator wishes to obtain state-of-the-art data on time of evolutionary origin

210 of species belonging to the true finches for their class. They decide to use `datelife` because

211 they are teaching best practices for reproducibility. Students have the option to go to the

212 website at www.datelife.org and perform an interactive run. However, the educator also

213 wants the students to practice their R skills. The first step is to run a `datelife` query using

214 the "get species from taxon" flag. This will get all recognised species names within their

215 chosen inclusive taxon. The Fringillidae has 289 species, according to the Open Tree of Life

216 taxonomy. Once with a curated set of species taxon names, the next step is to run a

217 `datelife` search that will find all chronograms that contain at least two species names. The

218 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. **??** and 6). 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. **??**). We also demonstrate the usage of PATHd8 (Britton et al., 2007) as an alternative method 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. 6). From these, two trees have a different sampling than the original source chronogram, mainly because DNA BOLD data for

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

## Discussion

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

At the time of writing of this manuscript (Mar 29, 2022), `datelife`'s database has 253 chronograms, pulled entirely from OpenTree's database, the Phylesystem (McTavish et al., 2015). A unique feature of OpenTree's Phylesystem is that the community can add new state-of-the-art chronograms any time. As chronograms are added to Phylesystem, they are incorporated into an updated `datelife`'s database that is assigned a new version number, followed by a package release on CRAN. `datelife`'s chronogram database is updated as new chronogram data is added to Phylesystem, at a minimum of once a month and a maximum of every 6 months. Users can also upload new chronograms to OpenTree themselves, and trigger an update of their local `datelife` database to incorporate the new chronograms, to have them immediately available for analysis.

Incorporation of more chronograms into `datelife`'s database is crucial to improve its services. One option to increase chronogram number in the database is the Dryad data repository. Methods to automatically mine chronograms from Dryad could be designed and 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.


## Conclusions

355

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.opentreeoflife.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

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

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

## Acknowledgements

³⁹³

## 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 %7Bhttps://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., . . .

Lovette, I. J. (2014). Phylogenetics and diversification of tanagers (passeriformes: Thraupidae), the largest radiation of neotropical songbirds. *Molecular Phylogenetics and Evolution*, *75*, 41–77.

Chamberlain, S. A., & Szöcs, E. (2013). taxize : taxonomic search and retrieval in R [version 2; referees: 3 approved]. *F1000Research*, *2*(191), 1–29. https://doi.org/10.12688/f1000research.2-191.v2

Chamberlain, S. A., Szöcs, E., Foster, Z., Arendsee, Z., Boettiger, C., Ram, K., . . . Li, G. (2019). *taxize: Taxonomic information from around the web.* Retrieved from https://github.com/ropensci/taxize

Claramunt, S., & Cracraft, J. (2015). A new time tree reveals earth history's imprint on the evolution of modern birds. *Science Advances*, *1*(11), e1501005.

Criscuolo, A., Berry, V., Douzery, E. J. P., & Gascuel, O. (2006). SDM: A fast distance-based approach for (super)tree building in phylogenomics. *Systematic Biology*, *55*(5), 740–755. https://doi.org/10.1080/10635150600969872

Delsuc, F., Philippe, H., Tsagkogeorga, G., Simion, P., Tilak, M.-K., Turon, X., . . . Douzery, E. J. (2018). A phylogenomic framework and timescale for comparative studies of tunicates. *BMC Biology*, *16*(1), 1–14.

Eastman, J. M., Harmon, L. J., & Tank, D. C. (2013). Congruification: Support for time scaling large phylogenetic trees. *Methods in Ecology and Evolution*, *4*(7), 688–691. https://doi.org/10.1111/2041-210X.12051

Felsenstein, J. (1985). Phylogenies and the Comparative Method. *The American Naturalist*, *125*(1), 1–15. Retrieved from http://www.jstor.org/stable/2461605

Gibb, G. C., England, R., Hartig, G., McLenachan, P. A., Taylor Smith, B. L., McComish, B. J., . . . Penny, D. (2015). New zealand passerines help clarify the diversification of major songbird lineages during the oligocene. *Genome Biology and Evolution*, *7*(11), 2983–2995.

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

community-curated phylogenetic estimates. *Bioinformatics*, *31*(17), 2794–2800.

Michonneau, F., Brown, J. W., & Winter, D. J. (2016). rotl: an R package to interact with the Open Tree of Life data. *Methods in Ecology and Evolution*, *7*(12), 1476–1481. https://doi.org/10.1111/2041-210X.12593

Morlon, H. (2014). Phylogenetic approaches for studying diversification. *Ecology Letters*, *17*(4), 508–525. https://doi.org/10.1111/ele.12251

Ooms, J., & Chamberlain, S. (2018). *Phylocomr: Interface to 'phylocom'*. Retrieved from https://CRAN.R-project.org/package=phylocomr

Open Tree Of Life, Redelings, B., Cranston, K. A., Allman, J., Holder, M. T., & McTavish, E. J. (2016). Open Tree of Life APIs v3.0. *Open Tree of Life Project*, (Online Resources). Retrieved from %7Bhttps://github.com/OpenTreeOfLife/germinator/wiki/Open-Tree-of-Life-Web-APIs%7D

Open Tree Of Life, Redelings, B., Sánchez Reyes, L. L., Cranston, K. A., Allman, J., Holder, M. T., & McTavish, E. J. (2019). Open tree of life synthetic tree v12.3. *Zenodo*. Retrieved from https://doi.org/10.5281/zenodo.3937742

Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and evolution in R language. *Bioinformatics*, *20*(2), 289–290.

Posadas, P., Crisci, J. V., & Katinas, L. (2006). Historical biogeography: A review of its basic concepts and critical issues. *Journal of Arid Environments*, *66*(3), 389–403.

Price, T. D., Hooper, D. M., Buchanan, C. D., Johansson, U. S., Tietze, D. T., Alström, P., . . . others. (2014). Niche filling slows the diversification of himalayan songbirds. *Nature*, *509*(7499), 222.

R Core Team. (2018). *R: a language and environment for statistical computing.* Vienna, Austria: R Foundation for Statistical Computing.

Rabosky, D. L., Chang, J., Title, P. O., Cowman, P. F., Sallan, L., Friedman, M., . . .

others. (2018). An inverse latitudinal gradient in speciation rate for marine fishes. *Nature*, *559*(7714), 392.

Ramshaw, J., Richardson, D., Meatyard, B., Brown, R., Richardson, M., Thompson, E., & Boulter, D. (1972). The time of origin of the flowering plants determined by using amino acid sequence data of cytochrome c. *New Phytologist*, *71*(5), 773–779.

Ratnasingham, S., & Hebert, P. D. (2007). BOLD: The barcode of life data system (http://www. Barcodinglife. org). *Molecular Ecology Notes*, *7*(3), 355–364.

Rees, J. A., & Cranston, K. (2017). Automated assembly of a reference taxonomy for phylogenetic data synthesis. *Biodiversity Data Journal*, (5).

Revell, L. J. (2012). Phytools: An r package for phylogenetic comparative biology (and other things). *Methods in Ecology and Evolution*, *3*, 217–223.

Ronquist, F., & Huelsenbeck, J. P. (2003). MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics*, *19*(12), 1572–1574. https://doi.org/10.1093/bioinformatics/btg180

Sanchez-Reyes, L. L., O'Meara, B., Eastman, J., Heath, T., Wright, A., Schliep, K., . . . Alfaro, M. (2022). datelife: Scientific Data on Time of Lineage Divergence for Your Taxa. *R Package Version 0.6.2*. Retrieved from https://doi.org/10.5281/zenodo.593938

Sanderson, M. J. (2003). r8s: Inferring absolute rates of molecular evolution and divergence times in the absence of a molecular clock. *Bioinformatics*, *19*(2), 301–302.

Sanderson, M. J., & Doyle, J. A. (2001). Sources of error and confidence intervals in estimating the age of angiosperms from rbcL and 18S rDNA data. *American Journal of Botany*, *88*(8), 1499–1516.

Sauquet, H. (2013). A practical guide to molecular dating. *Comptes Rendus Palevol*, *12*(6), 355–367.

Sauquet, H., Ramírez-Barahona, S., & Magallón, S. (2021). The age of flowering

plants is unknown.

Schenk, J. J. (2016). Consequences of secondary calibrations on divergence time
estimates. *PLoS ONE*, *11*(1). https://doi.org/10.1371/journal.pone.0148228

Smith, Stephen A., & Brown, J. W. (2018). Constructing a broadly inclusive seed
plant phylogeny. *American Journal of Botany*, *105*(3), 302–314.

Smith, Stephen A., & O'Meara, B. C. (2012). TreePL: Divergence time estimation
using penalized likelihood for large phylogenies. *Bioinformatics*, *28*(20),
2689–2690. https://doi.org/10.1093/bioinformatics/bts492

Stoltzfus, A., Lapp, H., Matasci, N., Deus, H., Sidlauskas, B., Zmasek, C. M., . . .
Jordan, G. (2013). Phylotastic! Making tree-of-life knowledge accessible, reusable
and convenient. *BMC Bioinformatics*, *14*.
https://doi.org/10.1186/1471-2105-14-158

Title, P. O., & Rabosky, D. L. (2016). Do Macrophylogenies Yield Stable
Macroevolutionary Inferences? An Example from Squamate Reptiles. *Systematic
Biology*, syw102. https://doi.org/10.1093/sysbio/syw102

Vos, R. A., Balhoff, J. P., Caravas, J. A., Holder, M. T., Lapp, H., Maddison, W. P.,
. . . others. (2012). NeXML: Rich, extensible, and verifiable representation of
comparative data and metadata. *Systematic Biology*, *61*(4), 675–689.

Webb, C. (2000). Exploring the Phylogenetic Structure of Ecological Communities :
An Example for Rain Forest Trees. *The American Naturalist*, *156*(2), 145–155.

Webb, Campbell O., Ackerly, D. D., & Kembel, S. W. (2008). Phylocom: Software for
the analysis of phylogenetic community structure and trait evolution.
*Bioinformatics*, *24*(18), 2098–2100.
https://doi.org/10.1093/bioinformatics/btn358

Webb, Campbell O., & Donoghue, M. J. (2005). Phylomatic: Tree assembly for
applied phylogenetics. *Molecular Ecology Notes*, *5*(1), 181–183.

to be formatted in the same way as the general text (double spaced and linenumbered)

568      LTT plots showing results from the cross-validation analyses of trees with branch

569  length reconstructed with data from the Barcode of Life Database (BOLD) dated using

570  PATHd8. We could construct a tree with branch lengths for all source chronograms.

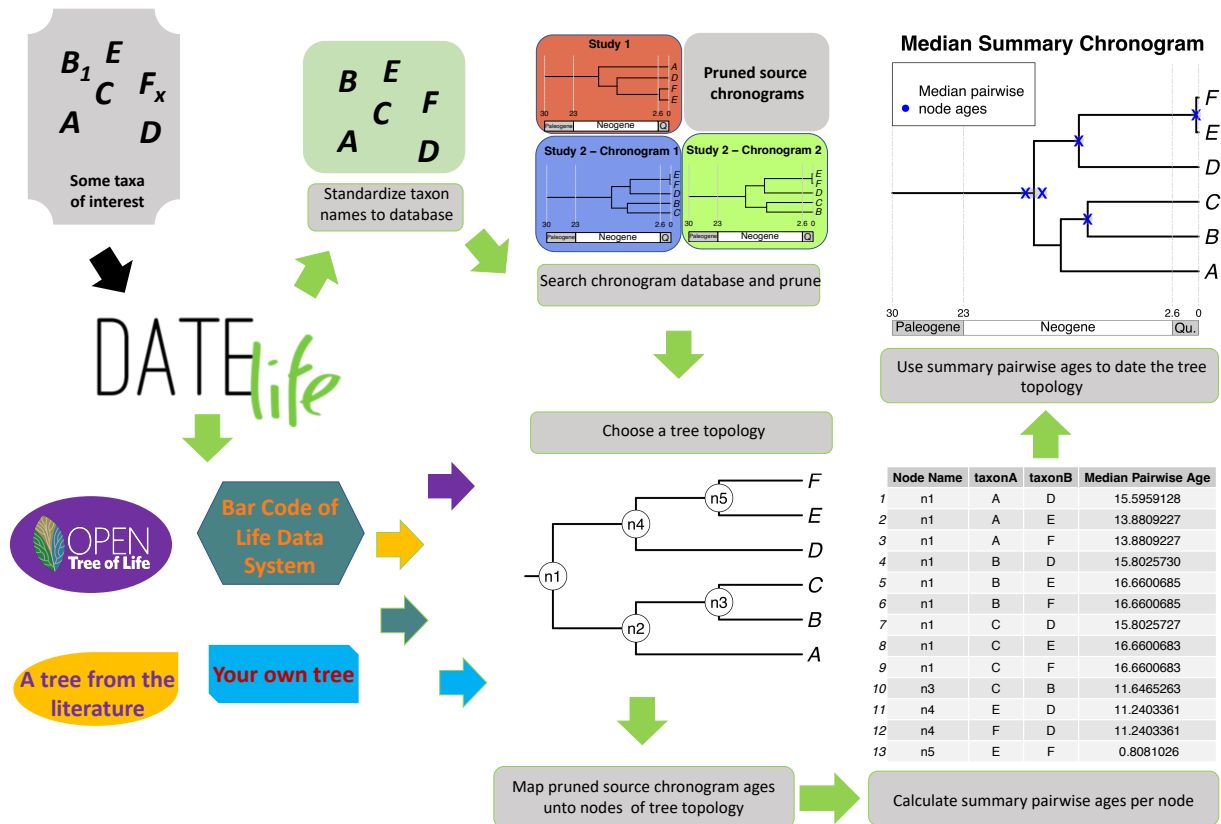571  However, dating with PATHd8 was only successful in three source chronograms shown here.

FIGURE 1. Stylized DateLife workflow. This shows the general worflows 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 Resoultion Service (TNRS; dark gray), and using TNRS (light gray). We also performed a search using the already processed query for comparison (light blue).
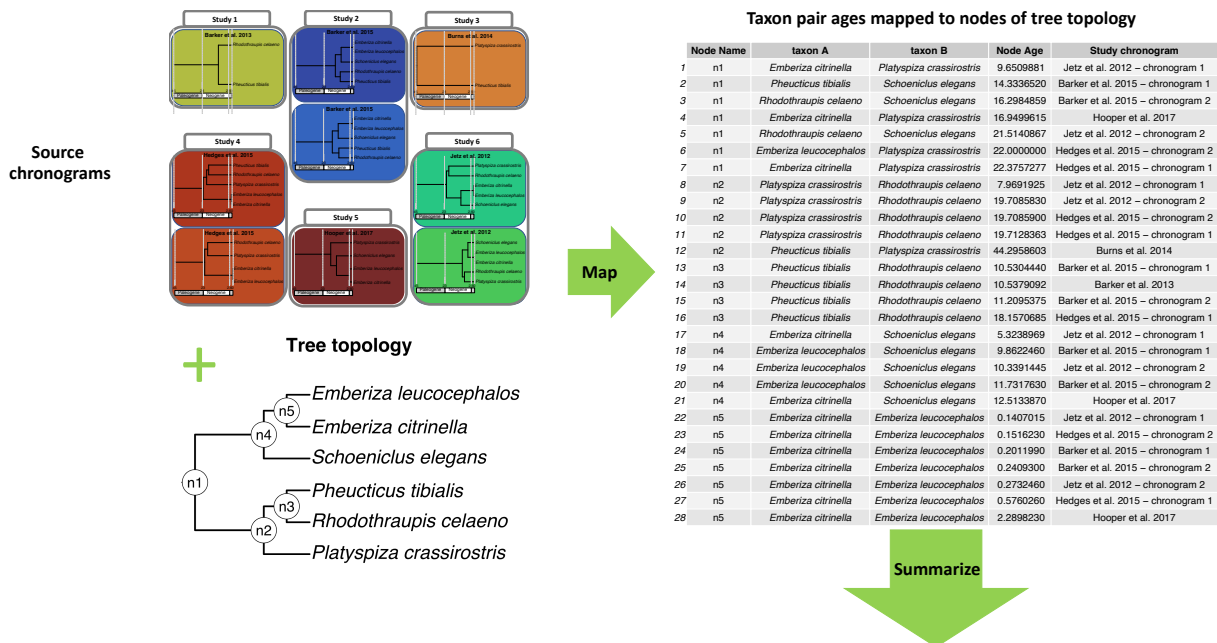
FIGURE 3. Age data results of a DateLife search of a small sample of 6 bird species within the Passeriformes. Input names were found across 9 chronograms within 6 independent studies (Barker et al. (2012), Barker et al. (2015), Burns et al. (2014), Hedges et al. (2015), Hooper and Price (2017), Jetz et al. (2012).) This revealed 28 age data points for the queried species names.

## Summary of mapped taxon pair age data

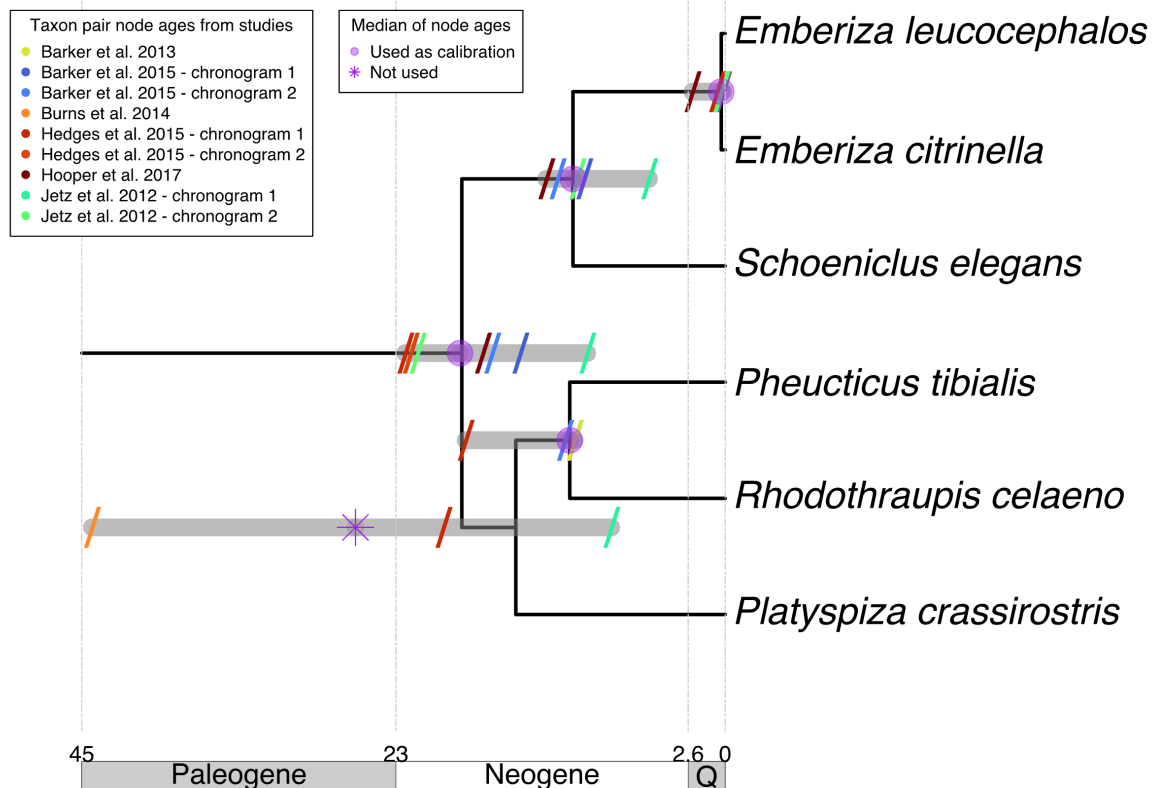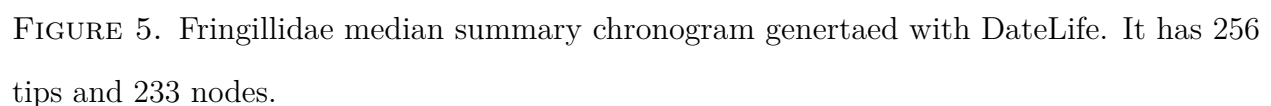| | Node Name | taxon A | taxon B | Pairwise Median Age | Node Median Age |
|---|---|---|---|---|---|
| 1 | | *Pheucticus tibialis* | *Emberiza citrinella* | 16.298486 | |
| 2 | | *Pheucticus tibialis* | *Emberiza leucocephalos* | 16.298486 | |
| 3 | | *Platyspiza crassirostris* | *Emberiza citrinella*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



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

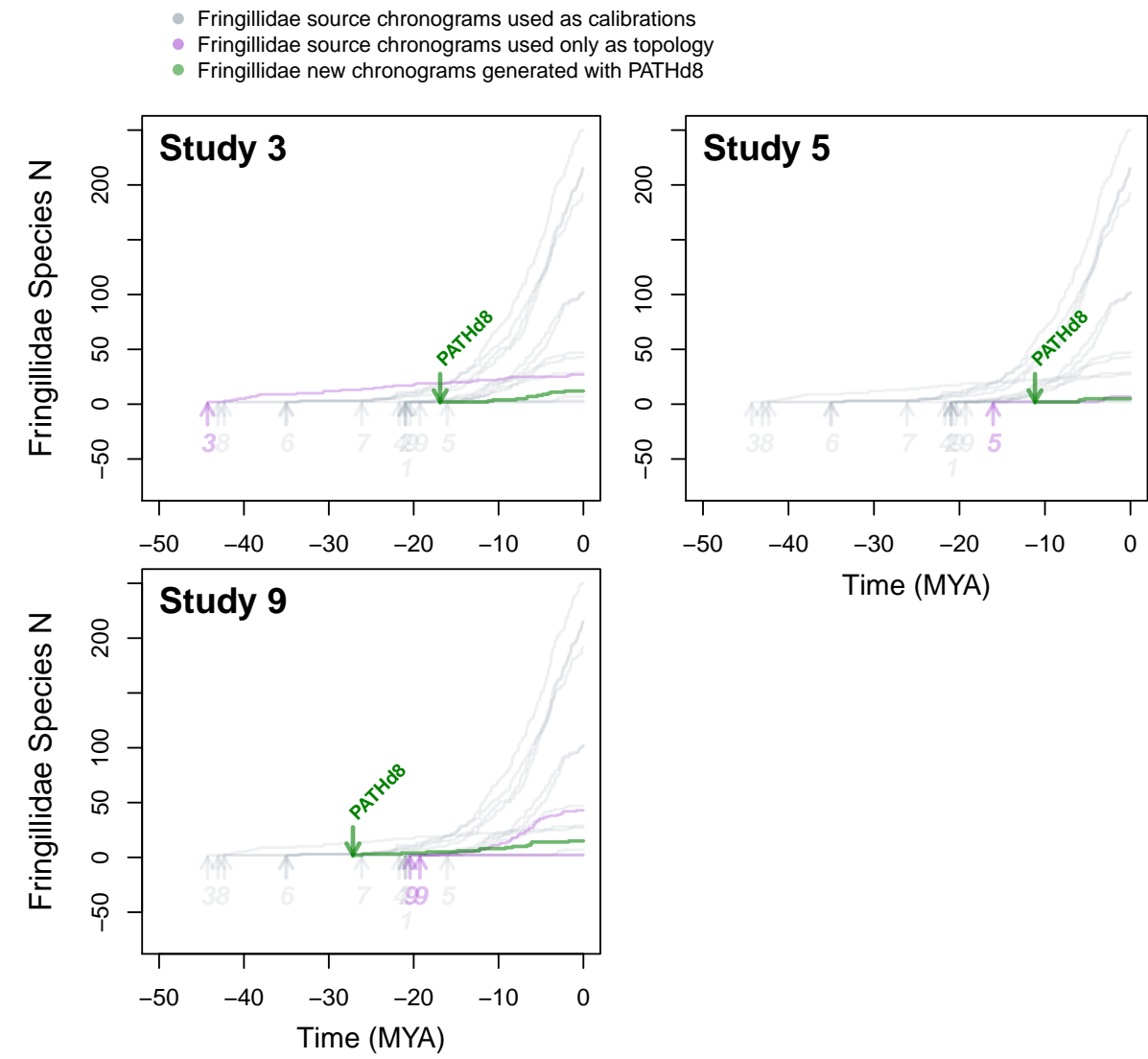FIGURE 5. Fringillidae median summary chronogram genertaed with DateLife. It has 256 tips and 233 nodes.

FIGURE 6