

Phytogeographic History of the Tea Family Inferred Through High-Resolution Phylogeny and Fossils

YUJING YAN^{1,2,*}, CHARLES C. DAVIS^{2,*}, DIMITAR DIMITROV^{1,3}, ZHIHENG WANG⁴, CARSTEN RAHBK^{1,4,5,6,7}, AND MICHAEL KRABBE BORREGAARD¹

¹Center for Macroecology, Evolution and Climate, GLOBE Institute, University of Copenhagen, Universitetsparken 15, 2100 Copenhagen, Denmark;

²Department of Organismic and Evolutionary Biology, Harvard University Herbaria, 22 Divinity Ave, Cambridge, MA 02138, USA; ³Department of Natural History, University Museum of Bergen, University of Bergen, P.O. Box 7800, 5020 Bergen, Norway; ⁴Institute of Ecology, College of Urban and Environmental Sciences, Key Laboratory of Earth Surface Processes of Ministry of Education, Peking University, 5 Yiheyuan Road, Beijing 100871, China;

⁵Center for Global Mountain Biodiversity, GLOBE Institute, University of Copenhagen, Universitetsparken 15, 2100 Copenhagen, Denmark; ⁶Department of Life Sciences, Imperial College London, Silwood Park campus, Ascot SL5 7PY, UK; and ⁷Danish Institute for Advanced Study, University of Southern Denmark, 5230 Odense M, Denmark

*Correspondence to be sent to: Department of Organismic and Evolutionary Biology, Harvard University Herbaria, Cambridge, MA 02138, USA
 E-mail: africarugu@gmail.com (Y.Y.), cdavis@oeb.harvard.edu (C.C.D.).

Received 1 September 2020; reviews returned 28 May 2021; accepted 8 June 2021
 Associate Editor: Alexandre Antonelli

Abstract.—The tea family (Theaceae) has a highly unusual amphi-Pacific disjunct distribution: most extant species in the family are restricted to subtropical evergreen broadleaf forests in East Asia, while a handful of species occur exclusively in the subtropical and tropical Americas. Here, we used an approach that integrates the rich fossil evidence of this group with phylogenies in biogeographic analysis to study the processes behind this distribution pattern. We first combined genome-skimming sequencing with existing molecular data to build a robust species-level phylogeny for c.130 Theaceae species, resolving most important unclarified relationships. We then developed an empirical Bayesian method to incorporate distribution evidence from fossil specimens into historical biogeographic analyses and used this method to account for the spatiotemporal history of Theaceae fossils. We compared our method with an alternative Bayesian approach and show that it provides consistent results while significantly reduces computational demands which allows analyses of much larger data sets. Our analyses revealed a circumboreal distribution of the family from the early Cenozoic to the Miocene and inferred repeated expansions and retractions of the modeled distribution in the Northern Hemisphere, suggesting that the current Theaceae distribution could be the remnant of a larger continuous distribution associated with the boreotropical forest that has been hypothesized to occupy most of the northern latitudes in the early Cenozoic. These results contradict with studies that only considered current species distributions and showcase the necessity of integrating fossil and molecular data in phylogeny-based parametric biogeographic models to improve the reliability of inferred biogeographical events. [Biogeography; genome skimming; phylogenomics; plastid genome; Theaceae.]

The Pacific Ocean is the largest water basin on Earth and constitutes a formidable barrier for terrestrial species dispersal between tropical eastern Asia and the neotropical region in the Americas, yet several plants and animals are found in both regions. Such disjunct distributions, where a monophyletic lineage occurs on both the eastern and western edges of the Pacific basin (i.e., temperate to tropical Asia on the one side and southeastern North America and South America on the other), with no occurrences in between are referred to as “amphi-Pacific”. Amphi-Pacific distribution is exhibited by more than 100 genera and a few higher taxa within angiosperms (Steenis 1962; Thorne 1972). In spite of this relatively high number of occurrences, the origin of this distribution remains controversial. Unraveling the causes of the amphi-Pacific distribution will have key implications for our understanding of the footprint of deep-time historical processes on present plant biogeographical patterns.

Amphi-Pacific distributions have been hypothesized to be relicts of wider circumboreal distributions associated with a continuous belt of “boreotropical” evergreen forest, which is thought to have extended at middle to northern latitudes through Eurasia and America during the early Cenozoic. This forest was

most likely continuous via a North Atlantic Land Bridge and/or the Bering Land Bridge, supported by evidence from both plants and animals (Tiffney 1985a, 1985b; Lavin and Luckow 1993; Sanmartín et al. 2001; Davis et al. 2002; Condamine et al. 2013). According to paleobotanical evidence, it was replaced by mixed mesophytic forest around early Oligocene and later boreal forest in late Miocene following climate cooling, resulting in the extinction and/or southward migration of boreotropical thermophilic taxa at high latitude regions (Meseguer et al. 2015, 2018). Evidence of adaptation to more temperate climate was also found (Meseguer et al. 2018). Recent analyses have inferred such a boreotropical forest origin for several plant groups with amphi-Pacific distributions, based on dated phylogenies and ancestral range reconstructions (e.g., Antonelli et al. 2009; Li et al. 2011a; Li and Wen 2013; Fritsch et al. 2015; Xiang et al. 2016), though biogeographical analyses are inherently problematic for lineages with no current members in intervening regions. The most important alternative hypothesis proposes an origin of the group in either Eastern Asia or North America followed by one or several long-distance dispersal events, either via the Bering Land Bridge or by

sea currents across the Pacific Ocean (Wen et al. 2010; Christenhusz and Chase 2013; Wu et al. 2018).

One of the best-known plant lineages exhibiting an amphipacific distribution is the tea family (Theaceae). Placed in Ericales, the family contains three tribes, ca. nine genera (Prince 2007), and 368 accepted species according to (The plant list, 2013), comprising shrubs and trees that are mostly thermophilic species and inhabit broadleaved-evergreen forests. Within the family, all three tribes have a disjunct amphipacific distribution. *Camellia*, *Schima*, *Pyrenaria*, *Polyspora*, *Apterosperma*, and most members of *Stewartia* (including *Hartia*) and *Gordonia* are restricted to subtropical and tropical Asia, whereas ca. 20 species belonging to several morphologically distinct groups, including *Stewartia*, *Gordonia*, *Laplacea*, and *Franklinia*, are restricted to the southeastern North America and the Neotropics, with no occurrences in-between.

Several attempts have been made to understand the evolutionary history of the family using biogeographical reconstructions. A recent study based on a relatively species-sparse phylogeny of Ericales inferred the stem of the family to be of Indo-Malaysian origin, with a possible expansion into the Nearctic at ~63 Ma (Rose et al. 2018), supporting an earlier conjecture by Li et al. (2013). The crown groups for all three tribes (Stewartieae, Gordoneae, and Theaeae) were all estimated to have originated during the late Oligocene to mid-Miocene (Yu et al. 2017; Lin et al. 2019), with uncertain biogeographical origin. The most recent common ancestors of both Stewartieae (*Stewartia s.l.*) and Gordoneae may have originated in North America, with subsequent dispersal events into East Asia involving multiple species during the Miocene (Lin et al. 2019). In particular, *Gordonia s.l.* may have species on both sides of the Pacific. The relationships within this genus have been contested by several authors, and even its status as a monophyletic clade is in question (Prince and Parks 2001; Yang et al. 2004; Gunathilake et al. 2015). The Neotropical distribution of *Laplacea* in Theaeae, on the other hand, has been argued to result from a single long-distance dispersal event from Asia to South America (Li et al. 2013), though the position of key taxa involved in the disjunct distribution is uncertain.

Poorly resolved relationships within the family and the sparse sampling of relevant species may lead to underestimating rates of cladogenesis and potentially erroneous conclusions in ancestral area reconstruction (Meseguer et al. 2015). Other potentially promising approach to improve biogeographic reconstructions and to decrease their level of uncertainty is to draw on fossils that directly reveal past occurrences. Several recent studies have experimented with including fossil taxa with reliable phylogenetic positions and spatial location data into parametric biogeographic models when reconstructing range dynamics. In some cases, including fossil information has led to substantial changes in the inferred biogeographic scenarios (Mao et al. 2012; Nauheimer et al. 2012; Wood et al. 2013).

Thus, here we have made specific effort to improve not only the sampling of extant species but also of fossil species.

However, accurately placing fossils onto the phylogeny usually takes a total-evidence approach which requires: (1) fossils with very informative morphological characters, (2) a comparative and well-sampled morphological character matrix for extant species, and (3) preferably a congruent evolutionary history of sampled morphological characters and molecular evidence. Such data are not always feasible for most clades including our target group and analyses using this approach have so far been restricted to a few small clades with well-preserved fossils (Meseguer and Condamine 2020). Theaceae is known from many macrofossils dating as far back as the late Cretaceous, occurring through the Cenozoic across temperate regions in the Northern Hemisphere (Fig. 1; Grote and Dilcher 1989, 1992; Bozukov and Palamarev 1995). Interestingly, these fossils from mid to high latitude may reflect a wide and northern distribution of the group and subsequent extinctions as many of them occur far north of the distribution of extant lineages (Sanmartín and Meseguer 2016a). This hypothesis has only been proposed by paleobotanists for the genera *Gordonia* and *Schima* based on the visual appearance of fossils (Grote and Dilcher 1992; Shi et al. 2017) but has not been tested quantitatively.

In this study, we generated a robust time-calibrated species-level phylogeny of Theaceae with wide taxonomic and genomic coverage, using *de novo* sequenced and previously published chloroplast genomes and nuclear ribosomal DNA (nrDNA) sequences for 146 species in nine genera (in the broad sense) across the world. Sequencing focused particularly on problematic and undersampled groups such as *Gordonia s.l.* and *Camellia*. We developed a heuristic method to incorporate fossils into the phylogeny based on their taxonomic placement and estimated ages even when they cannot be confidently placed on the phylogeny using morphological characters. We then use an empirical-Bayesian approach to infer biogeographical history that accounts for the phylogenetic uncertainties and evaluates the effect of incorporating information from extinct taxa. To resolve the enigmatic amphipacific distribution of Theaceae, we (1) provide a new reference phylogeny for the tea family for the biogeographic analysis, (2) clarify the placement of the problematic group *Gordonia* and its relationship to *Laplacea* and *Polyspora*, and (3) evaluate the effect of incorporating fossil information for the reconstructed biogeographic history.

MATERIALS AND METHODS

Our analytical approach involved building an improved phylogeny of the Theaceae by sequencing and including in the phylogenetic analyses species that have been difficult to place. For phylogenetic inference, we

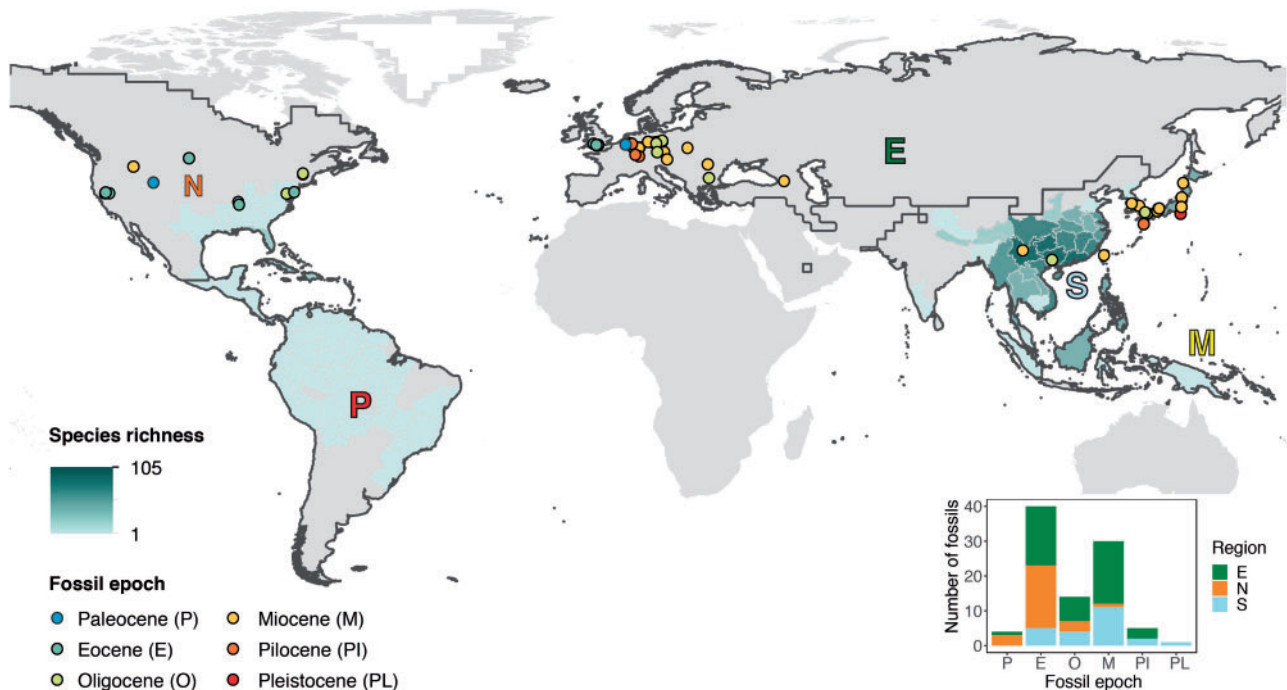


FIGURE 1. Current richness distribution pattern of Theaceae and occurrences of fossils through time and across biogeographic regions. The stack plot shows the number of reported fossils in each biogeographic region during different epoch. The capital letters on the map represent different biogeographic regions defined in the study. E = Eurasian; M = Papua-Melanesia; N = Nearctic; P = Panamanian; S = Sino-Japanese.

chose chloroplast genomes (pt, plastome) and nuclear ribosomal high-copy regions (nrDNA: the 18S rRNA-ITS1-5.8S rRNA-ITS2-26S rRNA region), because these markers have resolved relationships within the family well in previous studies and data for c.60 of species is already available (Yu et al. 2017b). We assembled target regions for 83 species *de novo* from herbarium specimens following genome skimming methods by Marinho et al. (2019). We then combined these data with existing plastid genomes and nrDNA data downloaded from GenBank and inferred a new dated phylogeny aiming to cover the entire family. We developed a novel method that allow the incorporation of information on past distributions based on the fossil records in biogeographic reconstruction.

Taxonomic Sampling, DNA Sequencing, and Data Processing

We sequenced *de novo* 83 herbarium specimens using the genome skimming method (Alamoudi et al. 2014; Dodsworth 2015; Zeng et al. 2018; Marinho et al. 2019). We focused on specimens from problematic or poorly covered taxa or regions, in particular *Gordonia s.l.* and *Camellia* (a vouchered specimen list is presented in Supplementary Table S2 available on Dryad at <http://dx.doi.org/10.5061/dryad.x0k6djhh0>). For each species, we selected the most recently collected specimen from within the native range of the species. Taxonomic assignments of specimens were based on the most recent

name identified on the specimen sheets. Combined with the GenBank data set, the resulting species coverage for all genera except *Camellia* was higher than 80% and the coverage within *Camellia* was c. 31% (according to TPL 2013) or c. 64% (according to Flora of China 2007).

Approximately 15 mg of leaf tissue was used for each DNA extraction. Both the Promega Maxwell kit (customized for herbarium specimens following the manufacturer's instructions) and a modified CTAB protocol (Doyle and Doyle 1987) were used to extract DNA from specimens. All extractions were quantified using a Qubit®3.0 Fluorometer (Life Technologies, Carlsbad, CA, USA) and quality checked using Nanodrop ND-1000 Spectrophotometer (Thermo Fisher Scientific). Samples collected at different times in the past were selected to visually assess the integrity with an Agilent Technologies 4200 TapeStation System using the Genomic DNA ScreenTape (Agilent Technologies, Santa Clara, CA, USA).

Dual-indexed libraries were prepared with 2–48 ng input genomic DNA using the Kapa DNA Hyper Plus Library Prep Kit (Kapa Biosystems) at 1/4 the recommended volume and size selected for 250–700 bp. The final libraries were pooled in equimolar ratios. The samples were sequenced either on Illumina NextSeq 500 mid output platform with 2*150-bp paired-end reads or Illumina HiSeq 2500 platform with 2*125-bp paired-end reads at the Bauer Core Facility of the Harvard University (<https://www.rc.fas.harvard.edu/odyssey-3-the-next-generation/>).

Adapters were trimmed, low-quality reads were removed, and base correction was performed using the default setting in fastp (Chen et al. 2018). Samples with >100,000 filtered reads were processed in Geneious R11 (<https://www.geneious.com>) with medium-low sensitivity and up to 10 iterations to assemble plastome and nuclear DNA sequences. For plastome assembly, we used the complete plastome of *Camellia taliensis* (NC_022264.1) as reference sequence. For nrDNA, nrDNA of *Camellia elongata* (MF171090.1) was used as reference sequence. For each region, two consensus contigs were generated using a 75% masking threshold and highest chromatogram quality. To avoid sequencing errors and potential contamination, bases with low coverage (below 6× for the plastome and 11× for nrDNA) were masked with Ns (Ripma et al. 2014). Samples with <40% high-quality bases, >60% ambiguities and an average sequencing depth <5 were excluded in the following analysis. Raw read data generated in the study are archived in the NCBI Sequence Repository Archive.

Building the Molecular Data Set from Genbank Data

We downloaded all available plastomes and nrDNA of Theaceae from GenBank to supplement our sequencing data set to achieve as complete a sampling of the family as possible at the species level. To ensure the reliability of our data, we only included plastome sequences cited in publications. As there is no widely accepted classification for Theaceae worldwide, our taxonomy followed (The plant list, 2013) and excluded taxa that were listed as unresolved. The cleaned GenBank data set included 69 complete plastomes and 41 nrDNA sequences for 70 species (Supplementary Table S1 available on Dryad).

Additionally, we downloaded chloroplast genomes and nrDNA sequences in the RefSeq data set for 19 species of various families within Ericales as outgroups. The selection of these outgroups mostly followed Yu et al. (2017b) with additional species sampled in Primulaceae and Pentaphragmaceae (Supplementary Table S1 available on Dryad).

Data Set Composition and Alignment

The final data set was assembled by combining the genome skimming data with those downloaded from GenBank (Supplementary Tables S1 and S2 available on Dryad). The data set comprised 146 species.

We constructed two subsets of these data to assess the congruence of phylogenies reconstructed from different genomic compartments, that is, plastid coding regions (CDS) and nrDNA. To assemble the CDS data set, we extracted the coding regions of all plastomes using BLAST+ (Camacho et al. 2009) against four annotated reference genomes, that is, *Camellia taliensis* (NC_022264.1), *Schima brevipedicellata* (NC_035537.1), *Euryodendron excelsum* (NC_039178.1),

and *Ardisia polysticta* (NC_021121.1) to eliminate possible annotation errors.

We used a partition strategy for the alignment to reduce the impact of sequencing and assembly errors. Sequences were roughly aligned using MAFFT v.7 (Kato et al. 2002), then partitioned based on gene position as extracted from the annotations of 49 plastomes in RefSeq database (sequence accession started with “NC_” in Supplementary Table S1 available on Dryad, including four outgroup species) and the annotations of 54 nrDNA sequences (sequence accession started with “MF” in Supplementary Table S1 available on Dryad, including 12 outgroup species). This resulted in 105 partitions for the chloroplast genome and five partitions for the nrDNA (Supplementary Table S3 available on Dryad). Each partition was aligned separately, and the resulting alignments were used to infer phylogenies with FastTree (Price et al. 2010). For each partition, we removed sequences with a threshold of 0.05 in TreeShrink (Mai and Mirarab 2018) to reduce long branch attraction artifacts. The filtered partitions were then concatenated and trimmed to exclude indel-rich positions using the “auto” setting in TrimAl (Capella-Gutiérrez et al. 2009).

Phylogenetic Analysis and Divergence Time Estimation

We concatenated the plastome and nrDNA data sets to reconstruct the full phylogeny using partitioned maximum likelihood method, which allows each data partition (plastome and nrDNA) to have different substitution models (comparable to the approach taken by Yu et al. (2017b)). The best nucleotide substitution model for each partition was determined with PartitionFinder2 (Lanfear et al. 2016), based on the corrected Akaike Information Criterion (cAIC). The partitioned maximum likelihood analysis was accomplished with the MPI version of RAXML-ng (Kozlov et al. 2019). Clade support was estimated using non-parametric bootstrap analyses with 1000 replicates and bootstrap values were mapped to the best-scoring tree. In this process, we observed several unstable taxa (rogue taxa) that affected clade support negatively (Wilkinson 1996). We qualified the influence of each species on the stability of the phylogeny using RogueNaRok (Aberer et al. 2013) (Supplementary Table S4 available on Dryad), pruned the rogue taxa and reran the above analysis using the reduced data set.

To evaluate the stability of the topology to data set selection and phylogenetic methodology, we also built phylogenies for the CDS and the nrDNA data sets separately using both maximum likelihood and a Bayesian inference framework. For the maximum likelihood analysis, we used the same parameter settings as mentioned previously, but used 200 pseudo-replicates for a bootstrap analysis of branch support. For the Bayesian analysis, we used MrBayes v3.2.6 (Huelsenbeck and Ronquist 2001). Two independent runs were conducted with four Markov chains (one

cold and three heated) for $> 3 \times 10^7$ generations, and a sampling frequency of one tree every 300 generations. Convergence was assumed when the average standard deviation of split frequencies was < 0.01 . We summarized the posterior on the best maximum likelihood tree, to facilitate comparison between the results of the two methods.

Due to the large size of the phylogeny, we used an efficient penalized likelihood method implemented in treePL (Smith and O'Meara 2012) to estimate divergence times within the trees generated by RAXML-ng. We selected three fossil calibration points for the outgroups, following Yu et al. (2017b), and four fossil calibration points for the ingroup species (Fig. 2, Supplementary Table S5 available on Dryad, Li and Wen 2013; Yu et al. 2017b). We conservatively applied all calibration constraints to the stems of the groups where respective fossils were placed. Penalized likelihood uses a smoothing parameter to accommodate rate heterogeneity. We used Random Subsample and Replicate Cross-Validation to determine the appropriate smoothing parameter for our data set. To assess uncertainty in age estimates we estimated confidence intervals on inferred ages by dating all 1000 ML bootstrap trees for the concatenated data set. Results from the dating of the bootstrapped trees were then summarized and resulting confidence intervals were visualized on the best tree using TreeAnnotator (part of the BEAST2 package, Bouckaert et al. 2014).

All phylogenetic analyses were implemented on the CIPRES Science Gateway (Miller et al. 2010) or the Harvard cluster Odyssey (<https://www.rc.fas.harvard.edu/odyssey-3-the-next-generation/>).

Selection and Incorporation of Extinct Taxa

To evaluate the impact of fossil occurrences on biogeographic inference, we developed a method to directly incorporate the fossil distribution information in phylogeny-based parametric biogeographic models by using phylogenies with fossils assigned iteratively based on their taxonomic placement. We compiled a comprehensive data set of all fossils of Theaceae from the Paleobiology Database (<http://paleodb.org>), Japan Paleobiology Database (<http://jpaleodb.org/>), and published literature (see Fig. 1 for the distribution of all fossil records). Online databases do not provide confirmation of identification and might include misplaced fossils. To limit the uncertainties, we first filtered the data set to unique and most reliable records by applying the following criteria: (1) the fossil has a locality at city or provincial level, (2) the fossil is dated to a geological epoch, (3) the fossil is assigned to a modern genus, (4) the fossil is associated with a publication later than 1990, (5) the fossil is a macrofossil, and (6) there is a detailed description with comparison to modern analogs. After applying these filters, we ended with 22 unique fossil records with species names (detailed information of the fossils and arguments that support

the taxonomy are presented in Supplementary Table S7 available on Dryad).

It is often difficult to place woody plant fossils accurately on a phylogeny due to missing information on some of the morphological characters and the lack of comparative morphological character matrices for extant species. To address this issue, we implemented a heuristic approach that assigned fossils randomly to a node within the most recent clade a fossil could be assigned to with certainty (usually a genus). The fossils were added as extinct lineages terminating within the confidence interval of the deposition time of the fossil, defined as the interval from the maximum to the minimum possible age of the formation in which the fossil was found (according to the latest geological time scale, Walker et al. 2018).

To place each fossil on the phylogeny, we first identified the narrowest clade that it could be assigned to with certainty, following the taxonomic placement of the most recent publication describing the fossil (see Supplementary Table S7 available on Dryad for more details). We then selected a branch within this clade and added the fossil specimen as a side branch to it. The branch was selected among a candidate set of all branches in the clade that existed prior to the minimum age of the fossil (based on the time interval given in the literature for the youngest specimen of that species). The new side branch was added to the selected branch at a randomly chosen time, drawn from a uniform distribution along the branch (but no later than the minimum age of the fossil). The terminal age of the fossil was then drawn from a uniform distribution between the newly created node and the fossil's minimum age. Note that this procedure may cause fossils to branch off from the stem above their assigned genus, effectively increasing the crown age of the genus for subsequent fossils. For this reason, we added the fossils in order from the oldest to the youngest, making it possible to assign younger fossils to the stem nodes created by older fossils.

One potential complication with this approach is that several fossils of morphologically highly similar species, in many cases found together, may be placed in different locations by the algorithm, though the most parsimonious explanation would be that they belonged to the same radiation. Dispersing these species over the phylogeny may lead to unrealistic ancestral state reconstruction and overestimation of dispersal events. To address this issue, we subsampled the fossil data set so that only one species per geological age per area was retained. We excluded fossils found in the Sino-Japanese region, which is well represented by present taxa in space. Applying these additional criteria resulted in a final set of 10 fossils (Supplementary Table S7 available on Dryad, Fig. S4).

Biogeographic Analysis with and without Fossil Taxa

Biogeographic analyses require the a priori definition of distinct regions. We defined large regions based

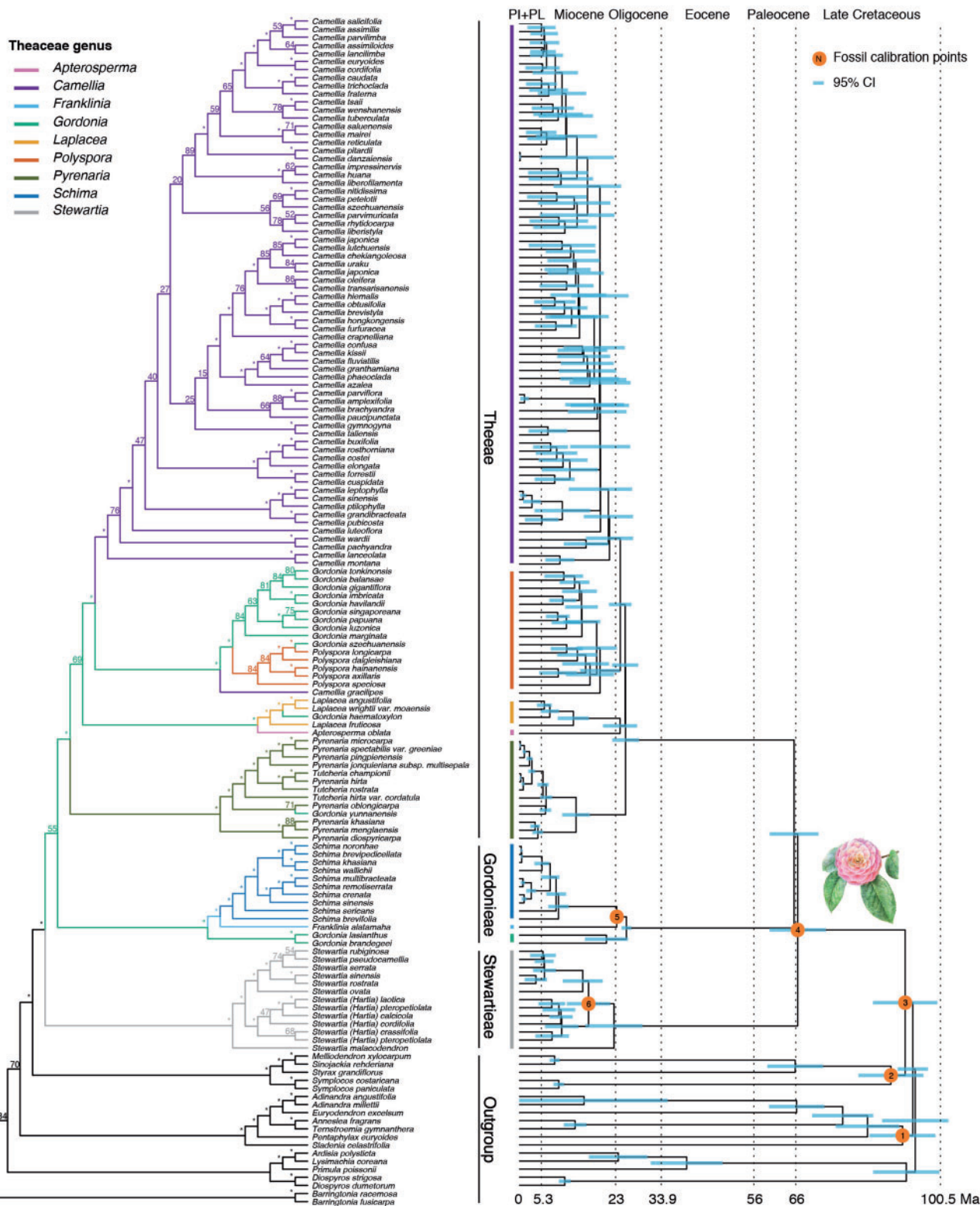


FIGURE 2. Left: Phylogeny of Theaceae inferred using the combined pt-nrDNA data set in RAxML-ng. Numbers associated with nodes indicate ML bootstrap support (BS) values. Asterisks represent nodes with BS values >90%. Right: The time-calibrated phylogeny of Theaceae. Topology is derived from the maximum likelihood tree on the left. Two species of Lecythidaceae at the root of the phylogeny are pruned as the root height is set manually in TreePL and does not influence the height of inner nodes in our study. 95% confidence intervals of the divergence time estimated in treePL are shown as blue bars at each node. Fossil calibrations are marked by orange circles with numbers corresponding to [Supplementary Table S5](#) available on Dryad. Illustration of the flower of a *Camellia japonica* cultivar by Qian Qian indicates the branch of stem Theaceae.

on the phyloregions defined by [Holt et al. \(2013\)](#), merging their regions into five areas: (1) Eurasian (E), (2) Nearctic (N), (3) Sino-Japanese (S), (4) Panamanian (P), and (5) Papua-Melanesian (M). Africa and Australia were excluded from the analysis, as there are no reliable fossils or current occurrences of Theaceae in these regions. Moreover, the age of the crown node of Theaceae was estimated to be around 60 Ma (with upper 95% highest posterior density boundary of 74.7 Ma) ([Yu et al. 2017b](#)) and is much younger than that expected for a Gondwanan origin (>80 Ma) ([Beaulieu et al. 2013](#)). Eurasian here represents the high latitude region of the Eurasian landmass, which used to be connected to the Nearctic through the North Atlantic land bridge or the Bering land bridge in historical times. The boundaries of this region were set based on the phylogenetic regionalization of [Holt et al. \(2013\)](#). Each species was assigned to one or more of these regions in ArcGIS 10.5 ([Esri Inc. 2016](#)) based on species' occurrence data collected from multiple data sources (see [Supplementary Appendix 1](#) available on Dryad), as well as from specimen tags from the species sequenced for our phylogenetic analysis. We limited the occurrences to specimen records, and cleaned the occurrences using CoordinateCleaner ([Zizka et al. 2019](#)).

The biogeographical history of the family was inferred using a likelihood-based dispersal-extinction-cladogenesis (DEC) model and a DEC model with founder events parameter (DEC+j) implemented in R package BioGeoBEARS ([Ree and Smith 2008](#); [Matzke 2013](#)). Both methods were used because DEC+J often artificially results in higher likelihood values than DEC, so the likelihoods themselves are not adequate for model selection in that case ([Ree and Sanmartín 2018](#)). The BioGeoBEARS DEC and DEC+j models allow inclusion of fossil tips in the phylogeny and inference of extinction directly from biogeographic data, which is an advantage over cladistic and event-based methods ([Matzke 2013](#); [Sanmartín and Meseguer 2016b](#)). The models were run with some key assumptions that cannot be determined *a priori* from our data set and may affect the final inference. We carried out a thorough sensitivity analyses and evaluated a larger set of feasible parameter combinations. The set of parameters that made the largest difference in terms of conclusions are presented here in the main text, the rest in the Supplementary materials available on Dryad. These parameters are whether or not to allow founder event dispersal (the DEC vs DEC+j model), the maximum number of occupied regions (set to two or three), and whether dispersal probabilities between Nearctic and Eurasian should vary over time to reflect the disappearance of the boreotropical forest corridor (yes or no). ([Mao et al. 2012](#); [Fritsch et al. 2015](#); [Meseguer et al. 2015](#); [Rose et al. 2018](#)).

Basic dispersal probabilities among regions assumed that species could move between Nearctic and Eurasia regions freely, but only disperse to southern regions from their adjacent region (M0). When evaluating the effect of allowing inter-region dispersal probabilities to

vary, we applied a time-stratified matrix with three time slices (70–35 Ma, 35–10 Ma, 10–0 Ma), where dispersal probabilities reflected the connectivity of Northern Hemisphere regions for the tea family following the vegetation reconstruction in [Meseguer et al. \(2015\)](#) (M1, [Supplementary Table S6](#) available on Dryad). We set three probability categories: 0.01 for well-separated areas, 0.5 for moderately separated areas, and 1.0 for well-connected areas. We used categories rather than an actual distance matrix in the analysis because distances among regions change continuously and are difficult to quantify.

In order to accommodate phylogenetic uncertainties, we used an empirical Bayesian approach inspired by [Nylander et al. \(2008\)](#), in which the biogeographic modeling procedure was repeated on 300 phylogenies randomly sampled from the 1000 dated bootstrap trees and the average marginal regional probabilities were summarized on the best RAxML-ng tree in BioGeoBEARS followed [Smith \(2009\)](#) and [Cai et al. \(2016\)](#). To evaluate the effect of including fossil distributions in biogeographic analysis, we assigned the two sets of fossils to each of the 300 trees using the procedure described above, ran the models with same settings, and summarized the results on the best RAxML-ng tree. We also summarized the main parameters of the different models and evaluated model performance using AICc and likelihood-ratio tests (lrt).

Averaging marginal regional probabilities from the sampled trees to the best tree only considers results of identical clades and does not necessarily reflect the statistical distribution of different evolutionary histories. To evaluate the ancestral ranges at nodes of major clades and possible historical events along branches of the sampled trees, we performed the newly developed biogeographical stochastic mappings procedure (BSM) ([Dupin et al. 2017](#)) for each tree in BioGeoBEARS simulating the biogeographic history based on the corresponding biogeographical likelihood model. The BSM realizations were summarized after 50 successful BSMs in 10,000 tries. We finally calculated the proportions of the most likely ancestral ranges for major clades and summarized the number of different events across the 300 sampled trees. See [Supplementary Figure S1](#) available on Dryad for an illustration of the workflow of our sequential empirical-Bayesian analysis.

All the biogeographical analysis was performed in R 3.5.3 ([R Core Team 2019](#)). Phylogenies were plotted using the R package ggtree ([Yu et al. 2017a](#)).

RESULTS

Phylogenomics of Theaceae

We succeeded in assembling at least 60% of the plastome for 59 samples (median sequencing depths at 29×) and the entire ribosomal sequences cluster for 72 samples (median sequencing depths at 278×) from the genome skimming data ([Supplementary Table S2](#)

available on Dryad). Two New World species (*Gordonia haematoxylon* and *Laplacea wrightii* var. *moaensis*) were represented by two samples each. As both samples grouped together with high support values for both species, only one of each was kept for the combined plastome-nrDNA data set. The proportion of parsimony-informative characters for plastid coding regions (CDs), nrDNA and concatenated plastome-nrDNA data sets were 14%, 14%, and 20%, respectively, with alignment lengths of 68,867 bp, 6142 bp, and 135,216 bp. In the phylogeny based on the combined data set, nine species were removed after being identified as rogue taxa (Supplementary Table S4 available on Dryad), for a total of 147 species in the final data set (128 ingroup species and 19 outgroup species).

All data sets recovered the genera *Pyrenaria*, *Stewartia*, and *Schima* as monophyletic, in accordance with previous molecular phylogenetic studies. Conversely, paraphyly of *Gordonia* s.l. and *Laplacea* was suggested with high support values across all the data sets. In the phylogeny based on the combined data set, only *Gordonia brandegeei* (synonym: *Laplacea grandis*) was retained in the same clade as the type species for *Gordonia*, *Gordonia lasianus*. This reduced *Gordonia* clade was strongly supported as the sister clade to *Franklinia alatamaha* + *Schima* with high bootstrap support (BS 100%). The following relationships were also supported by 100% BS: all the central American species, including *Laplacea fruticosa*, *Laplacea angustifolia*, *Laplacea portoricensis*, *Laplacea wrightii* var. *moaensis*, and *Gordonia haematoxylon* (*Laplacea haematoxylon*), formed a monophyletic clade, which was sister to *Apterosperma oblata*. *Gordonia yunnanensis* was placed in *Pyrenaria*. *Gordonia szechuanensis* was placed in a previously recognized *Polyspora* clade while all the Southeast Asian *Gordonia* species and *Gordonia balansae* formed a monophyletic clade and was sister to *Polyspora* (Fig. 2, Supplementary Fig. S2 available on Dryad).

The topologies were consistent across different phylogenetic reconstructions, and both maximum likelihood and Bayesian analysis inferred similar topologies for the CDs and nrDNA data sets (Supplementary Figs. S2 and S3 available on Dryad). The combined data set had the highest overall support and resolution. The highly supported nodes (BS > 70% and PP > 0.95) of the nrDNA phylogeny were mostly consistent with the CDs phylogeny and the combined plastome-nrDNA phylogeny, though the overall nodal support was low (average 51% BS). Despite high congruence among data sets of inferred relationships within the three tribes (Stewartieae, Gordonieae and Theeae), the relationships among them differed between data sets. This incongruence was also reflected by moderate support for Theeae and Gordonieae as sister taxa, with Stewartieae as a sister group to that clade in the combined data set. Our biogeographical analyses used an empirical Bayesian method to explicitly account for such uncertainties in the phylogenetic reconstruction.

The backbone relationships between some clades remained poorly resolved (Supplementary Figs. S2 and S3 available on Dryad). Notably, *Camellia gracilipes* grouped with the *Polyspora*+Asia *Gordonia* clade with high support in the plastid data set but was placed within *Camellia* in the analysis based on nrDNA data with very low support.

Divergence Time Estimation

For the treePL analyses, we set the maximum root depth to 125 Ma. Changing this setting did not give qualitatively different results for the internal node ages (Fig. 2 and Supplementary Table S8 available on Dryad). The stem age of Theaceae was estimated to be 92.1 Ma (95% confidence interval [CI]: 84.4–99.8) and the crown age 66.4 Ma (95% CI: 60.0–73.4). The crown age of Stewartieae (including the North American *Stewartia malacodendron*) was 22.69 Ma (95% CI: 15.1–28.7). The stem and crown ages of New World *Gordonia* were 25.6 Ma (95% CI: 24.1–26.4) and 20.8 (95% CI: 12.6–22.8). The stem and crown ages of *Laplacea*, another New World lineage, were 24.1 Ma (95% CI: 21.2–29.3) and 12.9 (95% CI: 8.7–16.1). The estimated ages were all within the 95% highest posterior density intervals of Yu et al. (2017b) and Rao et al. (2018).

Biogeographical Analysis with BioGeoBEARS

The model allowing jump dispersal (DEC+J) had higher likelihood values than the vicariance-only model (DEC) for all three phylogenetic data sets and across parameter settings (Table 1 and Supplementary Table S9 available on Dryad). However, there are arguments against model selection using likelihood-based methods in DEC-based models (Ree and Sanmartín 2018). For plants, jump dispersal is a strong assumption to take in biogeographic analysis, in that it assumes an ability to jump between remote geographical regions multiple times, and as such a vicariance-only model may be considered more parsimonious. Consequently, we present the results from both models.

The model parameters were congruent for all data sets under both M0 and M1 dispersal scenarios, and the reconstructed range dynamics were almost identical. Therefore, we focus on the results of the M0 scenario. Several observations about parameter estimates were consistent with the results of Ree and Sanmartín (2018). The dispersal rate of the DEC model was estimated to be higher than the DEC+j model for all data sets. Incorporating fossils into phylogenies increased the dispersal rate and extinction rate for DEC models, but did not influence these rates in DEC+j models. The extinction rate of DEC+j models stayed zero regardless of the number of fossils added, whereas the rate of founder events increased significantly (average 0.01 to 0.045 and 0.055 when max range was set to two) (Table 1). Similar patterns were found when the maximum number of co-occupied regions was increased

TABLE 1. Summary of ancestral range estimation using biogeographic models with maximum number of occupied ranges set to two (mean \pm standard deviation)

Data set	Model	LnL	num_params	d	e	j	AICc	AICc_wt
nfossil	M0-DEC	-64.012 \pm 2.741	2	0.004 \pm 0.001	0 \pm 0	0 \pm 0	132.121 \pm 5.482	0.185
	M0-DEC+j	-61.199 \pm 1.884	3	0.002 \pm 0.001	0 \pm 0	0.01 \pm 0.003	128.591 \pm 3.769	0.815
	M1-DEC	-63.773 \pm 2.73	2	0.004 \pm 0.001	0 \pm 0	0 \pm 0	131.641 \pm 5.459	0.195
	M1-DEC+j	-61.046 \pm 1.887	3	0.002 \pm 0.001	0 \pm 0	0.01 \pm 0.003	128.285 \pm 3.774	0.805
fossil_10	M0-DEC	-103.367 \pm 4.418	2	0.008 \pm 0.001	0.002 \pm 0.001	0 \pm 0	210.823 \pm 8.837	0.010
	M0-DEC+j	-86.777 \pm 4.288	3	0.002 \pm 0.001	0 \pm 0	0.045 \pm 0.006	179.733 \pm 8.576	0.990
	M1-DEC	-104.312 \pm 4.406	2	0.008 \pm 0.001	0.001 \pm 0.001	0 \pm 0	212.713 \pm 8.812	<0.001
	M1-DEC+j	-88.358 \pm 3.453	3	0.002 \pm 0	0 \pm 0	0.045 \pm 0.004	182.895 \pm 6.906	>0.999
fossil_22	M0-DEC	-129.707 \pm 5.347	2	0.01 \pm 0.002	0.003 \pm 0.001	0 \pm 0	263.495 \pm 10.695	<0.001
	M0-DEC+j	-109.218 \pm 6.066	3	0.002 \pm 0.001	0 \pm 0	0.055 \pm 0.004	224.601 \pm 12.132	>0.999
	M1-DEC	-131.21 \pm 5.422	2	0.01 \pm 0.002	0.003 \pm 0.001	0 \pm 0	266.501 \pm 10.844	<0.001
	M1-DEC+j	-110.292 \pm 5.915	3	0.002 \pm 0.001	0 \pm 0	0.057 \pm 0.004	226.749 \pm 11.831	>0.999

Note: d = rate of dispersal; e = rate of extinction; j = relative per-event weight of jump dispersal; LnL = log value of the likelihood.

from two to three (Supplementary Table S9 available on Dryad).

Biogeographic Reconstruction without Fossil Taxa

All models that did not include fossil taxa reconstructed similar ancestral range dynamics regardless of dispersal settings, with a discontinuous distribution of the crown node (Fig. 3a and Supplementary Figs. S5–S7 available on Dryad). These findings indicate that the dispersal matrix had a very limited impact on the ancestral state reconstruction, as also found by Chacón and Renner (2014). Very few dispersal and vicariance events were recovered and most were inferred to have occurred around the transition between boreotropical forest to mixed mesophytic forest (Fig. 4a). Most dispersal events were from Sino-Japanese regions (Fig. 3c). Analyses reconstructed the crown node of Theaceae as Nearctic+Sino-Japanese (NS) for 60–98% of the sampled phylogenies (Fig. 3a and Supplementary Table S10 available on Dryad).

The ancestral areas estimated for the crown of Stewartieae were NS as well for the majority of the sampled phylogenies, while the most recent common ancestor of Theaeae and Gordonieae was Sino-Japanese region only. Within Stewartieae, *Stewartia malacodendron* was reconstructed as diverging in the Nearctic in the early-Miocene, with gene flow between New World and Old World ceasing around the mid-Miocene, after which the lineage experienced a vicariance event that left *Stewartia ovata* in the Nearctic while its sister clade diversified in the Sino-Japanese region (Fig. 3a). The most recent common ancestor of the Theaeae and Gordonieae tribes was reconstructed as expanding from the Sino-Japanese region into the Nearctic, again placing the crown of tribe Gordonieae with an NS distribution. During the late Oligocene and mid-Miocene, one descendant of Gordonieae expanded south to the Panamanian region (Fig. 3a). The distribution of the crown node of Theaeae differed between the DEC+J model and the DEC model. With jump dispersal, the crown node was placed in the Sino-Japanese region, indicating that the Neotropical distribution of the genus

Laplacea must have been the result of a later long-distance dispersal event across the Pacific Ocean (Supplementary Fig. S5 available on Dryad). A joint occupation of the Sino-Japanese+Panamaian (PS) region was inferred under the DEC model (Fig. 3a).

Biogeographic Reconstruction Including Fossil Taxa

Adding fossil taxa radically altered the reconstructed geographical distributions for the deep nodes of the family. Analysis including 10 fossil taxa placed the family in Eurasia from the early Cenozoic to the mid-Miocene (Fig. 3b). Under the M0 dispersal scenario, the DEC model inferred the ancestral range of crown Theaceae to be either Eurasian+Nearctic (EN, ~50% of the phylogenies) or Eurasian (~33% of the phylogenies). The reconstruction involved multiple vicariance and dispersal events between Eurasian and the Nearctic and dispersal events between Eurasian and the Sino-Japanese region. Eurasia was inferred as a major source region other than Sino-Japanese region (Fig. 3d). The number of events recovered was higher than reconstruction based only on extant taxa and the highest number of events happened within the period of mixed mesophytic forest (Fig. 4b). Similar results were obtained under the M1 scenario (Supplementary Fig. S6b available on Dryad).

The crown-group of Stewartieae was inferred to be originally widespread in the EN regions (supported by more than 66% of the phylogenies), followed by expansion into the Sino-Japanese and subsequent extinction in the Nearctic and Eurasia from the mid-Oligocene to mid-Miocene. The crown-group Gordonieae exhibited a similar pattern, with a slightly more complex history for the Theaeae, with an earlier Eurasian extinction around the early Oligocene with a possible temporary recolonization in the Early Miocene. The ancestral occurrence of extant crown-group Theaeae was the same as in the reconstruction without fossils. When including jump dispersal events, ~50% of the phylogenies favored Eurasia as the ancestral range, while 33% favored EN. Multiple jumping events between Eurasia and Nearctic, as well as between Eurasia and the

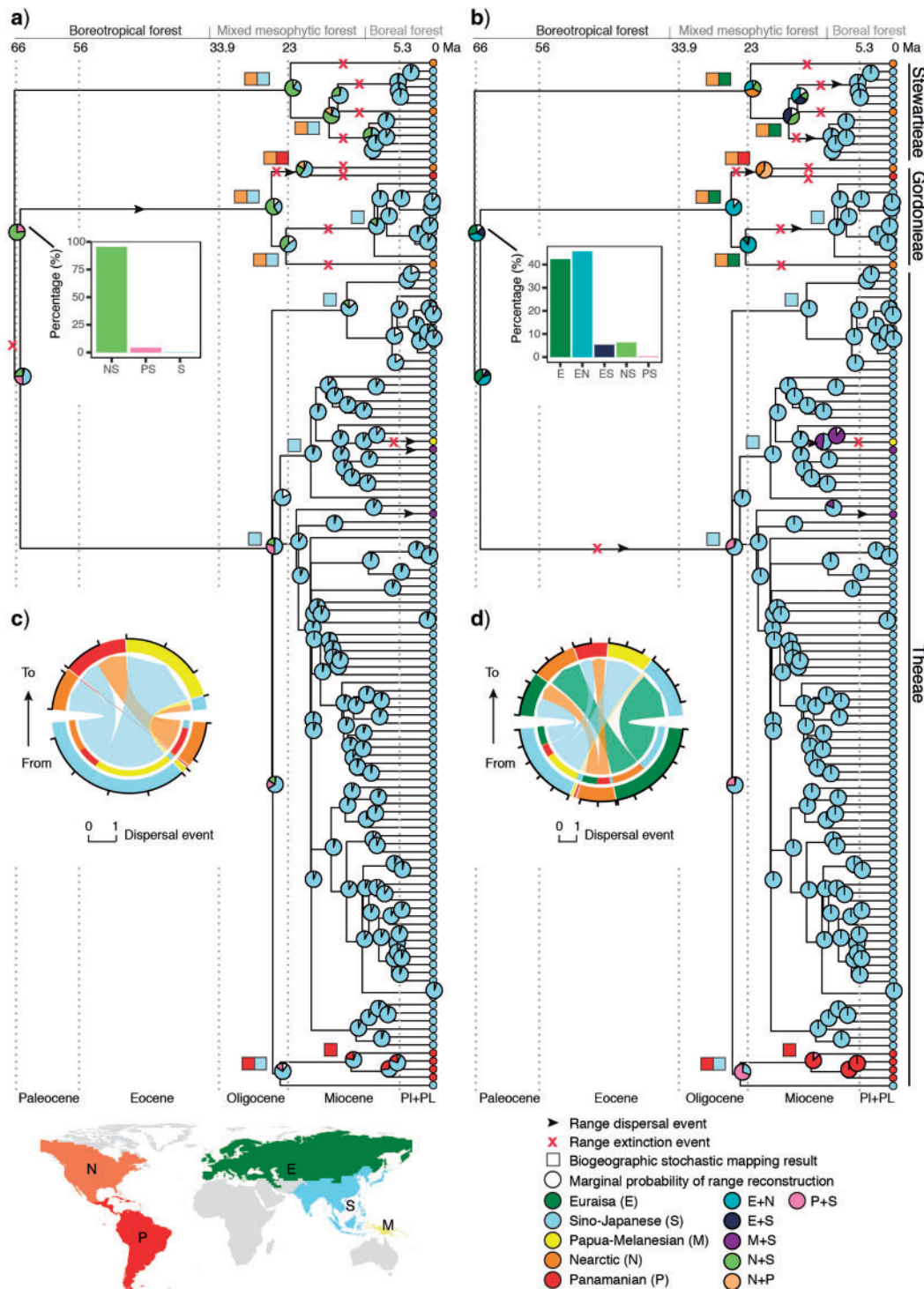


FIGURE 3. Biogeographic reconstruction of Theaceae showing the effect of incorporating fossil information into parametric biogeographic models. The reconstruction used DEC model under M0 dispersal scenario with maximum number of occupied ranges set to two over 300 phylogenies. a and b) Average marginal probability of different areas and the biogeographic stochastic mapping (BSM) result for majority of trees mapped on the time-calibrated best RAxML-ng tree based on only extant taxa (a) or based on both extant taxa and 10 fossil taxa (b). Colored circles at tips represent current ranges. Pie charts at inner nodes represent the average marginal probability of ranges across the 300 phylogenies and probabilities lower than 0.1 were combined and shown in white. The bar plots show the distribution of BSM result of the crown Theaceae across 300 phylogenies. Colored squares represent the recovered ranges of key nodes that have the highest frequencies across 300 phylogenies. Arrows and red crosses annotate the dispersal and extinction events that generated the distribution patterns. Periods of forest types are divided following Meseguer et al. (2015). c and d) Average number of reconstructed dispersals between different regions using only extant taxa (c) or using both extant and 10 extinct taxa (d). One tick mark represents one dispersal event. The lower panel show the source ranges and the upper panel show the sink ranges. The colored band within the lower panel show the sink ranges for each source range. The reconstructions using different model settings are presented in Supplementary Figs. S5–S9 available on Dryad.

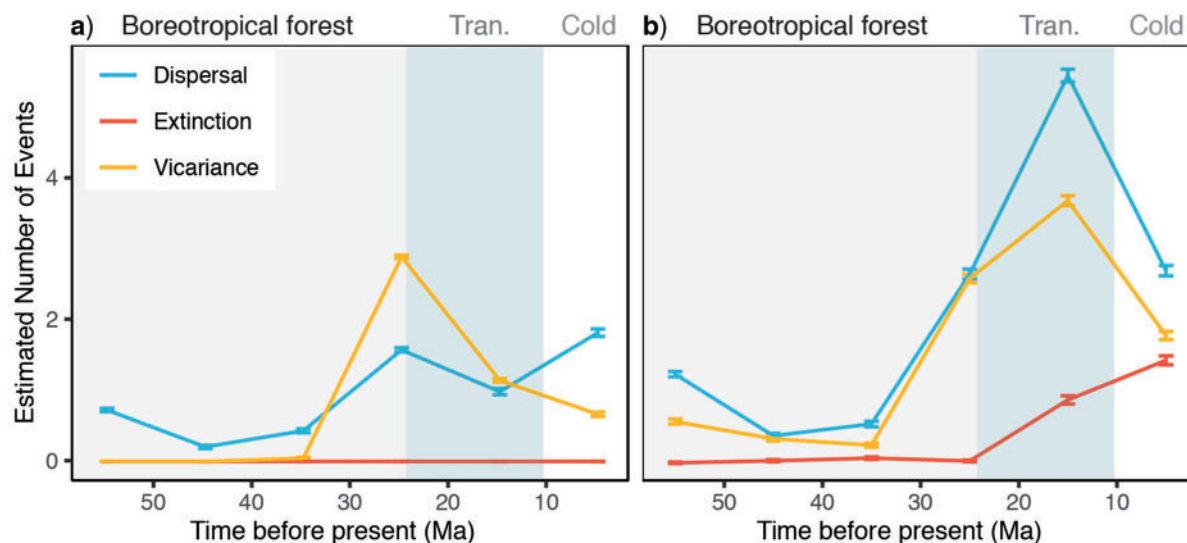


FIGURE 4. Estimated number of events through time based on biogeographic models using only extant taxa (a) and models incorporating fossil information (b). The results are the average numbers of 300 phylogenies with error bars indicate standard deviations. Shading shows the hypothesized three die out stages of the mid- to high-latitude boreotropical forest corridor followed Meseguer et al. (2015). “Trans.” represents mixed mesophytic forest and “Cold” represents boreal and temperate forest.

Sino-Japanese region replaced the inferred vicariance in the DEC model. The crown-group of Stewartieae was estimated to have occupied Eurasia in more than 60% of the phylogenies, while that of Gordonieae was estimated to have occupied the Nearctic for most phylogenies (Fig. 3b).

Analysis on the data set including 22 fossils also favored NS for the crown-group Theaceae in 50–60% of the sampled phylogenies, with 15–26% of the samples favoring a Eurasian origin (under the assumption that no clade can occupy more than two regions at any given time; Supplementary Table S10 available on Dryad). The crown group of Stewartieae was assigned to ES with high probability. The Sino-Japanese region was inferred as a major source region, with more than four dispersal events inferred from this region to the Euraisan, on average (Supplementary Figs. S8 and S9 available on Dryad). However, including many fossils with unresolved phylogenetic position into analyses may lead to some counter-intuitive reconstructions and increase the sensitive to dispersal constraints. In the M1 dispersal scenario, only 55% of the samples succeeded in 50 realizations in 10,000 tries of BSM under the DEC model.

Relaxing the assumptions to allow occupancy of three regions simultaneously broadened the inferred occupancy of the family crown group somewhat, to either encompass the Eurasia-Nearctic-Sino-Japanese or the Eurasia-Panamanian-Sino-Japanese regions (Supplementary Table S11 available on Dryad).

DISCUSSION

We provide a combined plastome and nrDNA phylogeny for Theaceae with high taxon coverage.

Our phylogeny greatly improves the understanding of Theaceae relationships, especially of *Gordonia s.l.* and *Laplacea*. Our results demonstrate that biogeographical analyses based solely on extant species distributions substantially misrepresented the past distributional history for the clade. This was remedied by adding a small number of fossil specimens and despite the high uncertainty of their phylogenetic placement. Biogeographical analysis accounting for known fossil distributions revealed a boreotropical origin for the family, with the amphipacific disjunct distribution arising from multiple colonization events from North to South and subsequent extinction in the intervening areas of Eurasia. The biotic exchange between the Old and the New World started in the early Eocene and ended in the Miocene. The study highlights the vital importance of the targeted sequencing of problematic taxa and inclusion of fossil data for robust biogeographical analyses.

Phylogenetic Relationships of Theaceae and *Gordonia s.l.*

Both plastome and nrDNA data set support paraphyly of the problematic genus *Gordonia s.l.*, indicating that it should be divided into three clades as previously proposed by Prince and Parks (2001) and Yang et al. (2004). The *Gordonia* clade only includes the type species, *Gordonia lasianus*, and *Gordonia brandegeei* (synonym: *Laplacea grandis*). This forms a purely New World clade, distributed from the Southeastern U.S. south to Columbia. The Caribbean and South American species, which include the genus *Laplacea* and the species *Gordonia haematoxylon* (synonym: *Laplacea haematoxylon*), formed a robust clade within the Theaeae tribe. All the Asian species of *Gordonia* were placed in the genus

Polyspora of tribe Theeae, putatively placed as the sister group to *Camellia*, corroborating some previous studies (Yang et al. 2004; Li et al. 2011b; Yu et al. 2017b). These findings based on molecular evidence are also supported by recent morphological and cytogenetics studies (Gunathilake et al. 2015; Hembree et al. 2019).

The relationships among remaining genera within Theaceae were generally consistent with a recent study using the same molecular markers (Yu et al. 2017b). Within *Camellia*, the backbone remained poorly resolved, most likely reflecting numerous hybridizations and polyploidization events during the rapid radiation of this clade (Yang et al. 2013). *Camellia* species only occur within the Sino-Japanese region, so the uncertainty in the phylogenetic placement of different sections does not affect our biogeographical results.

Northern Hemisphere Origin, Dispersal, and Extinction

Including fossils in the biogeographical analysis revealed a broad mid- to high-latitude Northern Hemispheric origin of Theaceae, strongly supporting the boreotropical forest hypothesis. The different ancestral area reconstructions all inferred Eurasia to be part of the ancestral distribution for the basal node, in most model outcomes also including the Nearctic. Furthermore, the crown age was estimated to be 66.4 Ma (60.0–73.4 Ma), well within the time frame where boreotropical forest is hypothesized to have been widespread and continuous (Tiffney 1985a, 1985b; Lavin and Luckow 1993). The reconstruction revealed more range expansion and extinction events compared to analysis based only on extant taxa, and especially captured the expansion to Sino-Japanese from Eurasian after the late Eocene and the extinction in Eurasian around E-O boundary as well as mid to late Miocene.

In further support of the boreotropical hypothesis, most reconstructed dispersal events were from the Old to the New World and from north to south (i.e., Eurasian to Sino-Japanese, Sino-Japanese to Papua-Melanesian in the Old World; Nearctic to Panamanian in the New World) (Fig. 3). This supports the intriguing notion that there may be a single historical cause of this distribution pattern across many taxa. Specifically, the two tribes that included disjunctive NS distribution (i.e., Gordonieae and Stewartieae) were both estimated to have undergone vicariance between Eurasia and Nearctic, followed by dispersal to Sino-Japanese from Eurasian and subsequently the extinction in Eurasian within the late Oligocene to mid (late)-Miocene (Fig. 3).

The BioGeoBEARS model allowed species to disperse between the Old and the New World along two different pathways, through the North Atlantic Bridge and the Bering Land Bridge. In the Paleocene and Eocene, species may have taken advantage of both routes during warm intervals (Brikiatis 2014; Wen et al. 2016) though the North Atlantic Land Bridge is considered most likely in the Eocene (Tiffney and Manchester 2002; Brikiatis 2014). Although it was hypothesized that climatic cooling

after the late Eocene (approximately 35 Ma) made intercontinental dispersal for thermophilic lineages less unlikely (Tiffney and Manchester 2002), this predates the divergence time between extant Old and New World lineages for all clades in the tea family. Gene flow between the continents must have persisted during the Oligocene and possibly until the mid-Miocene, which was observed in several other Northern Hemisphere lineages and attributed to the Bering Land Bridge (Donoghue and Smith 2004; Manos and Meireles 2015).

The reasons for a sustained gene flow between the continents until the mid-Miocene could either be that the changed environment was still within the tolerance limits of the contemporary species, or, the ancestral group managed to develop traits that are better adapted to cooler environments and were able to disperse over the sea in the late Cenozoic. In *Gordonia*, the gene flow around the Oligocene-Miocene boundary supports that it once inhabited mixed mesophytic forest in the western North America in the Miocene which also agrees with paleobotanical evidence (Fig. 3, Baskin and Baskin 2016). Notably, two deciduous species (i.e., *Franklinia alatamaha* and *Stewartia malacodendron*) were involved in the disjunctive patterns, suggesting that some species may have evolved cold tolerance before the mid-Oligocene (Tiffney 1985b), as the oldest deciduous species *Stewartia malacodendron* was dated at 15.1–28.7 Ma. Thus, it is possible that some temperate disjunctions might have a tropical origin. This view is also supported by the ancestral state reconstruction analysis of Yu et al. (2017b), in which they recovered the crown Theaceae as an evergreen species. A recent study of clusioid Malpighiales using direct paleoclimate simulation methods shows similar results and indicates that some boreotropical descendants might persist through niche evolution toward temperate climate (Meseguer et al. 2018). It may have been still possible for temperate plant taxa to cross the Atlantic via the North Atlantic Land Bridge in late Miocene (Denk et al. 2010; Brikiatis 2014).

Under the current phylogenetic hypothesis, dating scheme and fossil evidence, the occurrence of *Laplacea* in Central America can only be explained by a long-distance dispersal event across the Pacific Ocean from the Sino-Japanese region around the mid-Oligocene. *Laplacea*'s sister clade is *Apterosperma* with a single species in eastern China. No fossils relate to these two groups are known. Nevertheless, we cannot rule out the boreotropical hypothesis for this disjunction, given their absence on intervening Pacific islands. Moreover, we infer several dispersals from north to south almost in the same time period of high latitude extinction (Fig. 4b), in line with observations in other clades with amphi-Pacific distribution (Thomas et al. 2017; Yang et al. 2017). One dispersal from Nearctic to Panamanian, possibly occurred along the branch of *Gordonia brangeei*. Another one is from Sino-Japanese to as far as Papua-Melanesian region, possibly within Theeae. The group crossed the Wallace line and supports the view that

Southeast Asian biodiversity includes immigrants from northwestern regions such as Indochina (De Bruyn et al. 2014). Increasing range occupancy setting from two to three found similar results (Supplementary Table S11 available on Dryad).

The Importance and Uncertainties of Using Fossils in Biogeographical Analysis

One of the most striking results of this analysis is the failure of the biogeographic reconstruction to infer a realistic clade history in the absence of incorporating fossil evidence. The importance of fossils for estimating ancestral range dynamics is increasingly becoming clear, and a growing number of studies have taken a variety of approaches to incorporate them (Crisp et al. 2011; Meseguer et al. 2015; Sanmartín and Meseguer 2016a). Yet, the great majority of currently published studies lack fossil information, a potentially concerning situation given our results.

Here, we used a sequential inference method that focuses more narrowly on evaluating the impact of adding fossils on the ancestral states of key nodes, and thus allows comparing inferences derived from different sources of data. We made the most conservative assumptions for the fossils by attaching them randomly to the phylogeny within the limits set by well-established taxonomy assignment and age constraints. We then took an empirical-Bayesian method to estimate ancestral states across the distribution of the resulting trees. The method is not computationally intensive, which allows its use on large phylogenies like that of the tea family, or even larger. We show that including even minimal fossil distribution information in biogeographical analysis may lead to substantially different results. The method uncovered a complex historical range dynamic of the Theaceae governed by past environmental change and niche evolution.

To further explore the impact of uncertain phylogenetic placement of fossils and the performance of our method, we compared our approach with an approach recently developed by Landis et al. (2020). The Landis et al. (2020) method parametrically integrates different sources of uncertainty and accounts for their effects on biogeographic reconstructions and temporal estimates using a hierarchical Bayesian framework implemented in RevBayes (Höhna et al. 2016). This method aims to find evolutionary histories that are in harmony across all available lines of evidence, effectively fitting the inferred history as closely to the data as possible. Though potentially very powerful, the computational complexity of this approach makes it too computationally demanding to be applied to data sets as large as our full data set, even on modern centralized computing clusters. Therefore, we applied the method of Landis et al. (2020) and our method to a subset of taxa from our data set that included the members of Stewartieae (with generous assistance of Michael Landis, personal communication, Supplementary

Note S1 available on Dryad). Both methods produced congruent results under the M1 time-stratified dispersal scenario, but differed slightly under the M0 dispersal scenario (Supplementary Note S1 available on Dryad). The comparison shows that the greatest advantage of including fossil data in biogeographical analysis is that they constrain ancestral states of certain clades (Landis et al. 2020), with clear effects even though only a small portion of the recorded fossils have published justifications for relationships with extant species and thus may be used in analyses. Unfortunately, we cannot confidently infer the southern boundary of the family distribution during the Paleocene and Eocene, as a large portion of the sampled trees inferred the crown group occurrences to include the Sino-Japanese region when using the data set with 22 fossil taxa (Supplementary Tables S10 and S11 available on Dryad).

In addition, we tested manually minimizing state space in the biogeographic reconstructions on the extant species phylogeny allowing only single or adjacent regions. We recovered similar marginal probability for different states of key nodes as that of analysis including 10 fossils (Supplementary Fig. S10 available on Dryad), supporting the boreotropical hypothesis. However, when we conducted BSM to simulate the range evolution process, the realizations failed with a warning of complex history on a branch and disallowed necessary intermediate states. This test again stresses the importance of increasing sampling over setting stricter priors in modeling biogeographic processes.

Ideally, fossils provide information on traits, time, and distribution. Apart from directly incorporating fossil distributions into parametric biogeographic models, other methods could be explored to improve the understanding of movements between regions. For example, using fossil occurrences to model the historical niche of clades (e.g., Meseguer et al. 2015).

CONCLUSIONS

Our study provides insights into the origin of the amphipacific disjunctive distribution of Theaceae and its historical latitudinal range dynamics. We argue for the importance of incorporating fossil information in phylogeny-based biogeographical analysis and we provide a novel method to incorporate such information in biogeographic analyses. We show that even randomly associating fossils to extant phylogenies using limited constraints such as age and genus-level taxonomic information appears to lead to more accurate results for clades where extinction rate was high or spatially biased. Although inference of extinction and dispersal rates was not strongly influenced by including fossils in the biogeographic parametric models, this allowed us to demonstrate repetitive expansion and extraction of the putative distribution of the tea family from mid to high latitude regions in the Northern Hemisphere.

SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository:
<http://dx.doi.org/10.5061/dryad.x0k6djh0>.

ACKNOWLEDGMENTS

We thank the Harvard University Herbaria and the New York Botanical Garden who generously provided the material for the study. We thank the Bauer Core Facility of the Harvard University for providing technical support during the laboratory process. We thank Sen Li, Petter Marki, Liming Cai, Elizabeth Spriggs, Xiaoshan Duan, and Camille Desisto for helping with bioinformatics and wet lab work. We thank three anonymous reviewers for extremely insightful comments on the manuscript. The computations in this paper were partly run on the FASRC Odyssey cluster supported by the FAS Division of Science Research Computing Group at Harvard University and partly run on the CIPRES Science Gateway.

FUNDING

This work was supported by the Danish National Research Foundation (DNRF96 to Y.Y., C.R., and M.K.B.); the Chinese Scholarship Council (No. 201606010394 to Y.Y.); the Norwegian Metacenter for Computational Science (NOTUR; project NN9601K to D.D.); Harvard University (setup funding to C.C.D.), and a Carlsberg Young Researcher Award (CF19-0695 to M.K.B.).

REFERENCES

- Aberer A.J., Krompass D., Stamatakis A. 2013. Pruning rogue taxa improves phylogenetic accuracy: an efficient algorithm and webservice. *Syst. Biol.* 62:162–6.
- Alamoudi E.F., Khalil W.K.B., Ghaly I.S., Hassan N.H.A., Ahmed E.S. 2014. Nanoparticles from *Costus speciosus* extract improves the antidiabetic and antilipidemic effects against STZ-induced diabetes mellitus in albino rats. *Int. J. Pharm. Sci. Rev. Res.* 29:279–288.
- Antonelli A., Nylander J.A.A., Persson C., Sanmartin I. 2009. Tracing the impact of the Andean uplift on Neotropical plant evolution. *Proc. Natl. Acad. Sci. USA*. 106:9749–9754.
- Baskin J.M., Baskin C.C. 2016. Origins and relationships of the mixed mesophytic forest of Oregon–Idaho, China, and Kentucky: review and synthesis. *Ann. Missouri Bot. Gard.* 101:525–552.
- Beaulieu J.M., Tank D.C., Donoghue M.J. 2013. A Southern Hemisphere origin for campanulid angiosperms, with traces of the break-up of Gondwana. *BMC Evol. Biol.* 13:80.
- Bouckaert R., Heled J., Kühnert D., Vaughan T., Wu C.-H., Xie D., Suchard M.A., Rambaut A., Drummond A.J. 2014. BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* 10:e1003537.
- Bozukov V., Palamarev E. 1995. On the tertiary history of the Theaceae in Bulgaria. *Flora Mediterr.* 5:177–190.
- Brikiatis L. 2014. The deer, thulean and beringia routes: Key concepts for understanding early Cenozoic biogeography. *J. Biogeogr.* 41:1036–1054.
- De Bruyn M., Stelbrink B., Morley R.J., Hall R., Carvalho G.R., Cannon C.H., Van Den Bergh G., Meijaard E., Metcalfe I., Boitani L., Maiorano L., Shoup R., Von Rintelen T. 2014. Borneo and Indochina are major evolutionary hotspots for Southeast Asian biodiversity. *Syst. Biol.* 63:879–901.
- Cai L., Xi Z., Peterson K., Rushworth C., Beaulieu J., Davis C.C. 2016. Phylogeny of Elatinaceae and the tropical Gondwanan origin of the Centroplacaceae (Malpighiaceae, Elatinaceae) clade. *PLoS One*. 11:1–21.
- Camacho C., Coulouris G., Avagyan V., Ma N., Papadopoulos J., Bealer K., Madden T.L. 2009. BLAST+: architecture and applications. *BMC Bioinformatics*. 10:421.
- Capella-Gutiérrez S., Silla-Martínez J.M., Gabaldón T. 2009. trimAl: A tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*. 25:1972–1973.
- Chacón J., Renner S.S. 2014. Assessing model sensitivity in ancestral area reconstruction using Lagrange: a case study using the Colchicaceae family. *J. Biogeogr.* 41:1414–1427.
- Chen S., Zhou Y., Chen Y., Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*. 34:i884–i890.
- Christenhusz M.J.M., Chase M.W. 2013. Biogeographical patterns of plants in the Neotropics—dispersal rather than plate tectonics is most explanatory. *Bot. J. Linn. Soc.* 171:277–286.
- Condamine F.L., Sperling F.A.H., Kergoat G.J. 2013. Global biogeographical pattern of swallowtail diversification demonstrates alternative colonization routes in the Northern and Southern hemispheres. *J. Biogeogr.* 40:9–23.
- Crisp M.D., Trewick S.A., Cook L.G. 2011. Hypothesis testing in biogeography. *Trends Ecol. Evol.* 26:66–72.
- Davis C.C., Bell C.D., Mathews S., Donoghue M.J. 2002. Laurasian migration explains Gondwanan disjunctions: evidence from Malpighiaceae. *Proc. Natl. Acad. Sci. USA*. 99:6833–6837.
- Denk T., Grímsson F., Zetter R. 2010. Episodic migration of oaks to Iceland: Evidence for a north Atlantic “land bridge” in the latest Miocene. *Am. J. Bot.* 97:276–287.
- Dodsworth S. 2015. Genome skimming for next-generation biodiversity analysis. *Trends Plant Sci.* 20:525–527.
- Donoghue M.J., Smith S.A. 2004. Patterns in the assembly of temperate forests around the Northern Hemisphere. *Philos. Trans. R. Soc. B Biol. Sci.* 359:1633–1644.
- Dupin J., Matzke N.J., Särkinen T., Knapp S., Olmstead R.G., Bohs L., Smith S.D. 2017. Bayesian estimation of the global biogeographical history of the Solanaceae. *J. Biogeogr.* 44:887–899.
- ESRI. 2016. ArcGIS 10.4. Redlands (CA): ESRI. Available from: <http://www.esri.com>.
- Fritsch P.W., Manchester S.R., Stone R.D., Cruz B.C., Almeda F. 2015. Northern Hemisphere origins of the amphi-Pacific tropical plant family Symplocaceae. *J. Biogeogr.* 42:891–901.
- Grote P., Dilcher D. 1989. Investigations of angiosperms from the Eocene of North America: a new genus of Theaceae based on fruit and seed remains. *Bot. Gaz.* 150:190–206.
- Grote P.J., Dilcher D.L. 1992. Fruits and Seeds of Tribe Gordoniaceae (Theaceae) from the Eocene of North America. *Am. J. Bot.* 79:744–753.
- Gunathilake L.A.A.H., Prince J.S., Whitlock B.A. 2015. Seed coat micromorphology of *Gordonia* sensu lato (including Polyspora and Laplacea; Theaceae). *Brittonia*. 67:68–78.
- Hembree W.G., Ranney T.G., Jackson B.E., Weathington M. 2019. Cytogenetics, ploidy, and genome sizes of *Camellia* and related genera. *HortScience*. 54:1124–1142.
- Holt B.G., Lessard J.-P., Borregaard M.K., Fritz S.A., Araujo M.B., Dimitrov D., Fabre P.-H., Graham C.H., Graves G.R., Jonsson K.A., Nogues-Bravo D., Wang Z., Whittaker R.J., Fjeldsa J., Rahbek C. 2013. An update of Wallace’s Zoogeographic Regions of the World. *Science* (80-). 339:74–78.
- Höhna, Landis, Heath, Boussau, Lartillot, Moore, Huelsenbeck, Ronquist. 2016. RevBayes: Bayesian phylogenetic inference using graphical models and an interactive model-specification language. *Syst. Biol.* 65:726–736.
- Huelsenbeck J.P., Ronquist F. 2001. MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics*. 17:754–755.
- Katoh K., Misawa K., Kuma K., Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–66.
- Kozlov A.M., Darriba D., Flouri T., Morel B., Stamatakis A. 2019. RAXML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* 35:4453–4455.

- Landis M.J., Eaton D.A.R., Clement W.L., Park B., Spriggs E.L., Sweeney P.W., Edwards E.J., Donoghue M.J. 2020. Joint phylogenetic estimation of geographic movements and biome shifts during the global diversification of viburnum. *Syst. Biol.* 70:67–85.
- Lanfear R., Frandsen P.B., Wright A.M., Senfeld T., Calcott B. 2016. PartitionFinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. *Mol. Biol. Evol.* 34(3):772–773.
- Lavin M., Luckow M. 1993. Origins and relationships of tropical North America in the context of the boreotropics hypothesis. *Am. J. Bot.* 80:1–14.
- Li J., Del Tredici P., Yang S., Donoghue M.J. 2002. Phylogenetic relationships and biogeography of *Stewartia* (Camellioideae, Theaceae) inferred from nuclear ribosomal DNA ITS sequences. *Rhodora*. 104:117–133.
- Li L., Li J., Rohwer J.G., van der Werff H., Wang Z.H., Li H.W. 2011a. Molecular phylogenetic analysis of the Persea group (Lauraceae) and its biogeographic implications on the evolution of tropical and subtropical Amphi-Pacific disjunctions. *Am. J. Bot.* 98: 1520–1536.
- Li R., Wen J. 2013. Phylogeny and biogeography of *Dendropanax* (Araliaceae), an Amphi-Pacific disjunct genus between tropical/subtropical Asia and the Neotropics. *Syst. Bot.* 38:536–551.
- Li R., Yang J.B., Yang S.X., Li D.Z. 2011b. Phylogeny and taxonomy of the *Pyrenaria* complex (Theaceae) based on nuclear ribosomal ITS sequences. *Nord. J. Bot.* 29:780–787.
- Li Y., Awasthi N., Yang J., Li C. Sen. 2013. Fruits of *Schima* (Theaceae) and seeds of *Toddalia* (Rutaceae) from the Miocene of Yunnan Province, China. *Rev. Palaeobot. Palynol.* 193:119–127.
- Lin H.-Y., Hao Y.-J., Li J.-H., Fu C.-X., Soltis P.S., Soltis D.E., Zhao Y.-P. 2019. Phylogenomic conflict resulting from ancient introgression following species diversification in *Stewartia* s.l. (Theaceae). *Mol. Phylogenet. Evol.* 135:1–11.
- Mai U., Mirarab S. 2018. TreeShrink: Fast and accurate detection of outlier long branches in collections of phylogenetic trees. *BMC Genomics*. 19(Suppl 5):272.
- Manos P.S., Meireles J.E. 2015. Biogeographic analysis of the woody plants of the Southern Appalachians: implications for the origins of a regional flora. *Am. J. Bot.* 102:780–804.
- Mao K., Milne R.L., Zhang L., Peng Y., Liu J., Thomas P., Mill R.R., S. Renner S. 2012. Distribution of living Cupressaceae reflects the breakup of Pangea. *Proc. Natl. Acad. Sci. USA*. 109:7793–7798.
- Marinho L.C., Cai L., Duan X., Ruhfel B.R., Fiaschi P., Amorim A.M., van den Berg C., Davis C.C. 2019. Plastomes resolve generic limits within tribe Clusiaceae (Clusiaceae) and reveal the new genus *Arawakia*. *Mol. Phylogenet. Evol.* 134:142–151.
- Matzke N.J. 2013. BioGeoBEARS: BioGeography with Bayesian (and Likelihood) Evolutionary Analysis in R Scripts.
- Meseguer A.S., Condamine F.L. 2020. Ancient tropical extinctions contributed to the latitudinal diversity gradient. *Evolution* 74:1966–1987.
- Meseguer A.S., Lobo J.M., Cornuault J., Beerling D., Ruhfel B.R., Davis C.C., Jousset E., Sanmartín I. 2018. Reconstructing deep-time palaeoclimate legacies in the clusioid Malpighiales unveils their role in the evolution and extinction of the boreotropical flora. *Glob. Ecol. Biogeogr.* 27:616–628.
- Meseguer A.S., Lobo J.M., Ree R., Beerling D.J., Sanmartín I. 2015. Integrating fossils, phylogenies, and niche models into biogeography to reveal ancient evolutionary history: the case of *Hypericum* (Hypericaceae). *Syst. Biol.* 64:215–232.
- Miller M.A., Pfeiffer W., Schwartz T. 2010. Creating the CIPRES Science Gateway for inference of large phylogenetic trees. 2010 Gateway. *Comput. Environ. Work.* 11:1–8.
- Min T., Bartholomew B. 2007. Theaceae. *Flora of China*. p. 366–478.
- Nauheimer L., Metzler D., Renner S.S. 2012. Global history of the ancient monocot family Araceae inferred with models. *New Phytol.* 195:938–950.
- Nylander J.A.A., Olsson U., ALSTRÖM P., Sanmartín I. 2008. Accounting for phylogenetic uncertainty in biogeography: a bayesian approach to dispersal-vicariance analysis of the thrushes (Aves: Turdus). *Syst. Biol.* 57:257–268.
- Price M.N., Dehal P.S., Arkin A.P. 2010. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS One*. 5:e9490.
- Prince L. 2007. A brief nomenclatural review of genera and tribes in Theaceae. *Aliso*. 24:105–121.
- Prince L.M. 2002. Circumscription and biogeographic patterns in the Eastern North American-East Asian Genus *Stewartia* (Theaceae: Stewartiaceae): insight from Chloroplast and Nuclear DNA Sequence Data. *Castanea*. 67:290–301.
- Prince L.M., Parks C.R. 2001. Phylogenetic relationships of Theaceae inferred from chloroplast DNA sequence data. *Am. J. Bot.* 88:2309–2320.
- R Core Team. 2019. R: a language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing. Available from: <https://www.R-project.org>. Accessed February 8, 2019.
- Rao M., Steinbauer M.J., Xiang X., Zhang M., Mi X., Zhang J., Ma K., Svenning J.C. 2018. Environmental and evolutionary drivers of diversity patterns in the tea family (Theaceae s.s.) across China. *Ecol. Evol.* 11:663–11676.
- Ree R.H., Sanmartín I. 2018. Conceptual and statistical problems with the DEC+J model of founder-event speciation and its comparison with DEC via model selection. *J. Biogeogr.* 45:741–749.
- Ree R.H., Smith S.A. 2008. Maximum likelihood inference of geographic range evolution by dispersal, local extinction, and cladogenesis. *Syst. Biol.* 57:4–14.
- Ripma L.A., Simpson M.G., Hasenstab-Lehman K. 2014. Geneious! Simplified genome skimming methods for phylogenetic systematic studies: a case study in *Oreocarya* (Boraginaceae). *Appl. Plant Sci.* 2:1400062.
- Rose J.P., Kleist T.J., Löfstrand S.D., Drew B.T., Schönenberger J., Sytsma K.J. 2018. Phylogeny, historical biogeography, and diversification of angiosperm order Ericales suggest ancient Neotropical and East Asian connections. *Mol. Phylogenet. Evol.* 122:59–79.
- Sanmartín I., Enghoff H., Ronquist F. 2001. Patterns of animal dispersal, vicariance and diversification in the Holarctic. *Biol. J. Linn. Soc.* 73:345–390.
- Sanmartín I., Meseguer A.S. 2016a. Extinction in phylogenetics and biogeography: From timetrees to patterns of biotic assemblage. *Front. Genet.* 7:1–17.
- Sanmartín I., Meseguer A.S. 2016b. Extinction in phylogenetics and biogeography: From timetrees to patterns of biotic assemblage. *Front. Genet.* 7:1–17.
- Smith S.A. 2009. Taking into account phylogenetic and divergence-time uncertainty in a parametric biogeographical analysis of the Northern Hemisphere plant clade Caprifoliaceae. *J. Biogeogr.* 36:2324–2337.
- Smith S.A., O'Meara B.C. 2012. treePL: divergence time estimation using penalized likelihood for large phylogenies. *Bioinformatics*. 28:2689–2690.
- Steenis C.G.G.J. van. 1962. The land-bridge theory in botany with particular reference to tropical plants. *Blumea Tijdschr. voor Syst. en Geogr. der Planten*. 11:235–372.
- The Plant List. 2013. Version 1.1. Published on the Internet; Available from: <http://www.theplantlist.org/>. Accessed December 11, 2017.
- Thomas D.C., Tang C.C., Saunders R.M.K. 2017. Historical biogeography of Goniothalamus and Annonaceae tribe Annoneae: dispersal–vicariance patterns in tropical Asia and intercontinental tropical disjunctions revisited. *J. Biogeogr.* 44:2862–2876.
- Thorne R. 1972. Major Disjunctions in the geographic ranges of seed plants. *Q. Rev. Biol.* 90:365–411.
- Tiffney B. 1985a. The Eocene North Atlantic Land Bridge: its importance in Tertiary and modern phytogeography of the Northern Hemisphere. *J. Arnold Arbor.* 66:243–273.
- Tiffney B.H. 1985b. Perspectives on the origin of the floristic similarity between Eastern Asia and Eastern North America. *J. Arnold Arboretum*. 66:73–94.
- Tiffney B.H., Manchester S.R. 2002. The use of geological and paleontological evidence in evaluating plant phylogeographic hypotheses in the Northern Hemisphere Tertiary. *Int. J. Plant Sci.* 162:S3–S17.
- Walker D., Geissman J., Compilers. 2018. GSA Geologic time scale v. 5.0. *Geol. Soc. Am.* 204:59425.
- Wen J., Ickert-Bond S., Nie Z.-L., Li R. 2010. Timing and modes of evolution of Eastern Asian–North American biogeographic disjunctions in seed plants. In: Long M., Gu H., Zhou Z., editors.

- Darwin's heritage today: Proceedings of the Darwin 200 Beijing International Conference. Beijing: Higher Education Press. p. 252–269.
- Wen J., Nie Z.-L., Ickert-Bond S.M. 2016. Intercontinental disjunctions between eastern Asia and western North America in vascular plants highlight the biogeographic importance of the Bering land bridge from late Cretaceous to Neogene. *J. Syst. Evol.* 54:469–490.
- Wood H.M., Matzke N.J., Gillespie R.G., Griswold C.E. 2013. Treating fossils as terminal taxa in divergence time estimation reveals ancient vicariance patterns in the palpimanoid spiders. *Syst. Biol.* 62:264–284.
- Wu Z.Y., Liu J., Provan J., Wang H., Chen C.J., Cadotte M.W., Luo Y.H., Amorim B.S., Li D.Z., Milne R.I. 2018. Testing Darwin's transoceanic dispersal hypothesis for the inland nettle family (Urticaceae). *Ecol. Lett.* 21:1515–1529.
- Xiang X.G., Mi X.C., Zhou H.L., Li J.W., Chung S.W., Li D.Z., Huang W.C., Jin W.T., Li Z.Y., Huang L.Q., Jin X.H. 2016. Biogeographical diversification of mainland Asian *Dendrobium* (Orchidaceae) and its implications for the historical dynamics of evergreen broad-leaved forests. *J. Biogeogr.* 43:1310–1323.
- Yang J.B., Yang S.X., Li H.T., Yang J., Li D.Z. 2013. Comparative chloroplast genomes of camellia species. *PLoS One.* 8(8):e73053.
- Yang M.Q., Li D.Z., Wen J., Yi T.S. 2017. Phylogeny and biogeography of the amphi-Pacific genus *Aphananthe*. *PLoS One.* 12:1–18.
- Yang S.X., Yang J.B., Lei L.G., Li D.Z., Yoshino H., Ikeda T. 2004. Reassessing the relationships between *Gordonia* and *Polyspora* (Theaceae) based on the combined analyses of molecular data from the nuclear, plastid and mitochondrial genomes. *Plant Syst. Evol.* 248:45–55.
- Yu G., Smith D.K., Zhu H., Guan Y., Lam T.T.Y. 2017a. ggtree: an R package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol. Evol.* 8:28–36.
- Yu X.Q., Gao L.M., Soltis D.E., Soltis P.S., Yang J.B., Fang L., Yang S.X., Li D.Z. 2017b. Insights into the historical assembly of East Asian subtropical evergreen broadleaved forests revealed by the temporal history of the tea family. *New Phytol.* 215:1235–1248.
- Zeng C.X., Hollingsworth P.M., Yang J., He Z.S., Zhang Z.R., Li D.Z. 2018. Genome skimming herbarium specimens for DNA barcoding and phylogenomics. *Plant Methods.* 14:43.
- Zizka A., Silvestro D., Andermann T., Azevedo J., Duarte Ritter C., Edler D., Farooq H., Herdean A., Ariza M., Scharn R., Svantesson S., Wengström N., Zizka V., Antonelli A. 2019. CoordinateCleaner: Standardized cleaning of occurrence records from biological collection databases. *Methods Ecol. Evol.* 10:744–751.