# Analysis of plastid and nuclear DNA data in plant phylogenetics—evaluation and improvement

WANG Wei[*], LI HongLei & CHEN ZhiDuan

*State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing 100093, China*

Correct combination of plastid (cp) and nuclear (nr) DNA data for plant phylogenetic reconstructions is not a new issue, but with an increasing number of nrDNA loci being used, it is of ever greater practical concern. For accurately reconstructing the phylogeny and evolutionary history of plant groups, correct treatment of phylogenetic incongruence is a vital step in the proper analysis of cpDNA and nrDNA data. We first evaluated the current status of analyzing cpDNA and nrDNA data by searching all articles published in the journal *Systematic Botany* between 2005 and 2011. Many studies combining cpDNA and nrDNA data did not rigorously assess the combinability of the data sets, or did not address in detail possible reasons for incongruence between the two data sets. By reviewing various methods, we outline a procedure to more accurately analyze and/or combine cpDNA and nrDNA data, which includes four steps: identifying significant incongruence, determining conflicting taxa, providing possible interpretations for incongruence, and reconstructing the phylogeny after treating incongruence. Particular attention is given to explanation of the cause of incongruence. We hope that our procedure will help raise awareness of the importance of rigorous analysis and help identify the cause of incongruence before combining cpDNA and nrDNA data.

**combined analysis, plastid DNA, incongruence, nuclear DNA, phylogenetics**

Tremendous progress has been made in our understanding of phylogenetic relationships at all taxonomic levels across all plant groups with developments in molecular phylogenetics. As an example, the order- and family-level phylogenetic framework of land plants has been largely resolved [1–5]. Phylogenetics has become a powerful tool and a starting point in many areas of biology, such as taxonomy, physiology, ecology, biogeography, paleobiology, genomics, and developmental genetics [5–9]. A robust and well-supported phylogenetic tree is a prerequisite for understanding and explaining many life phenomena, otherwise erroneous conclusions will be generated in an incorrect phylogenetic context.

Plastid (cp) and nuclear (nr) genomes are the most frequently used sources of genetic data for reconstructing the phylogeny of plant groups [10–14]. However, certain evolutionary events, such as gene duplication, hybridization, and lineage sorting of ancestral polymorphisms, may result in conflicting topologies based on data sets from these two genomes at all taxonomic levels [12,15,16]. It is well known that hybridization and lineage sorting of ancestral polymorphisms have occurred much more frequently in nature than previously envisioned [17,18]. The issue of combining cpDNA and nrDNA data sets is not new, but it is of greater practical concern today, with the vast amount of molecular data now available, especially as an increasing number of low-copy nuclear genes are used in plant phylogenetics [13,14,19,20]. Careful analysis of cpDNA and nrDNA data and assessment of their combinability is a vital step toward accurate reconstruction of phylogenetic relationships in a

*Corresponding author (email: wangwei1127@ibcas.ac.cn)

plant group. If significant incongruence exists, the data should not be combined without further justification, otherwise a phylogenetic tree estimated from the combined (concatenated) data may not track or may represent an over-simplification of the evolutionary history [21,22].

Various methods have been proposed to identify incongruence between phylogenies obtained from cpDNA and nrDNA data, such as tree-based comparisons [22], the parsimony-based incongruence length difference (ILD) test [23], compare-two permutation tests [24], Templeton's test [25], and the Shimodaira-Hasegawa (SH) test [26]. Recently, van der Niet & Linder [27] described a protocol as a guide on how to analyze cpDNA and nrDNA data, which included three steps: identifying incongruence and testing its significance, assessing the cause of incongruence, and reconstructing the species tree. However, it is sometimes extremely difficult to determine the cause of incongruence, particularly to distinguish hybridization from incomplete lineage sorting [28,29], which may result in similar phylogenetic patterns. Fortunately, several novel approaches have been developed for statistically distinguishing hybridization from incomplete lineage sorting [16,29].

In this paper, we first evaluate how empirical plant phylogenetic studies usually analyze cpDNA and nrDNA data at present. Surprisingly, we found that many studies directly combined cpDNA and nrDNA data sets without rigorously assessing and/or treating data combinability. We sequentially outline a procedure for correctly analyzing cpDNA and nrDNA data and list the possible methods involved in each step of the procedure. We hope that this will help raise awareness of the importance of dealing carefully with incongruence and that it may serve as a guide to help authors facing this problem.

# 1   Current status of cpDNA and nrDNA data analysis in plant phylogenetics

To investigate how empirical studies usually treat incongruence between cpDNA and nrDNA data sets and to explore the possible cause of incongruence between the two, we chose the journal *Systematic Botany* and examined its published articles. As an internationally famous and reputable journal that publishes research on plant molecular phylogenetics, *Systematic Botany* was considered to be representative to some extent of the present status of cpDNA and nrDNA methodology. Given that some statistical methods for distinguishing hybridization from incomplete lineage sorting have only been published relatively recently [16,27,29], we restricted our search to the period between 2011 and 2005. The articles investigated had to fulfill the following two criteria: (i) one of the goals of the paper was to reconstruct phylogenetic relationships among members of a group or to determine the systematic position of a taxon, and (ii) the paper employed both cpDNA and nrDNA data

to construct phylogenetic trees for the same taxa. A paper was excluded if the author stated that they did not combine the cpDNA and nrDNA data sets because the separate phylogenies were adequate to answer their initial questions.

With regard to data collection, if a paper contained a statement that the phylogenetic signal from cpDNA and nrDNA data was significantly incongruent, we recorded it as "incongruent" regardless of the threshold for incongruence used.

If a paper noted that an ILD test was performed in the Materials and methods, and mentioned visual inspection of the two separate bootstrap consensus trees in the Results, we recorded it as "ILD test+tree-based comparisons".

If a paper used both the ILD test and tree-based comparisons, where the ILD test indicated the cpDNA and nrDNA data were significantly incongruent but the tree-based comparisons found that incongruence was weakly supported, we recorded it as "congruent".

We identified 138 articles that combined cpDNA and nrDNA data, of which 115 (83.3%) tested for combinability and 23 (16.7%) directly combined the datasets without any test of combinability. Detailed statistics are presented in Table 1. The most commonly used method of testing for combinability was the ILD test (53/115), followed by tree-based comparisons (34/115). Among the 20 papers that used both the ILD test and tree-based comparisons, the ILD test identified significant incongruence between cpDNA and nrDNA data in 16 studies, but tree-based comparisons identified incongruence occurred in only 10 of these 16 studies.

Among the 53 studies that only used the ILD test to examine the combinability of cpDNA and nrDNA data sets, significant incongruence was identified in 16 studies, but the cpDNA and nrDNA data were still combined in 14 of these 16 studies. Tree-based comparisons identified incongruence in 15 studies, of which the cpDNA and nrDNA data were still combined in seven studies. Ten studies were indicated to be incongruent by using both the ILD test and tree-based comparisons, of which the data were still combined in seven studies.

The ILD test is the most extensively used method at present for assessing incongruence, but the threshold for incongruence (*P*-value) differs markedly among these studies, including 0.05, 0.01, 0.005, 0.002, and 0.001. Furthermore, the threshold for incongruence used in different studies for tree-based comparisons are also arbitrary, such as maximum parsimony bootstrap value (MP BS)>50%, MP BS⩾60%, MP BS⩾70%, MP BS>85%, MP BS>80% and/or posterior probabilities (PP)⩾0.95, MP BS>75% and PP>0.9, and maximum likelihood (ML) BS⩾70% or PP⩾0.95. Surprisingly, the cpDNA and nrDNA data were still combined in 29 of 46 studies in which significant incongruence was identified. The reasons usually given by the authors was that the ILD test is sensitive to differences in among-site rate variation between partitions, overall evolutionary rates, levels of noise, and the relative size of data partitions [30].

Thus, it may be problematic in that, among the clades that are in conflict between a pair of data sets, it does not differentiate between those that are weakly supported and those that potentially have different evolutionary histories. Furthermore, the ILD test may be insensitive to localized differences in the evolutionary histories of two data sets if many or several other clades are strongly supported and congruent [22]. In addition, authors also considered that total evidence could generate a phylogenetic tree with greater resolution and higher support [31–33]. Notably, statistical approaches were rarely used to explain incongruence.

# 2    A procedure for analyzing cpDNA and nrDNA data

It is surprising to note that many studies combining cpDNA and nrDNA DNA data sets did not rigorously assess and treat data combinability. We suggest researchers take more care in this important step. To correctly analyze cpDNA and nrDNA data, we here outline a procedure that includes the following four major steps: (i) identification of significant incongruence, (ii) determination of conflicting taxa, (iii) provision of possible interpretations for incongruence, and (iv) reconstruction of the phylogeny after treating incongruence.

## 2.1    Identification of significant incongruence

Tree-based comparisons can be used to visually identify incongruence between phylogenies obtained from cpDNA and nrDNA data sets [22,34]. The three most widely applied methods used in phylogenetic analyses are MP, ML and Bayesian inference (BI). The accurate alignment of each genetic region is a prerequisite for the three analytical methods. It is routine to exclude difficult-to-align regions from phylogenetic analysis. For model-based phylogenetic analysis methods (BI and ML), it is also crucial to select the best-fit model for each genetic region in the data sets. Generally, model-based phylogenetic analyses do not result in trees that substantially differ in topology from that of MP analysis [35]. In many studies, the majority of clades with MP BS≥70% and/or ML BS≥70% have PP>0.95 [36,37], and BS values of MP and ML analyses differ by less than 5% [36]. Given that different analysis methods are sometimes sensitive to different biases in the data set, Baum et al. [38] suggested that clades consistently supported in different analyses could be regarded as more robust than those supported strongly by one method but contradicted by a different method. At present, the majority of studies employ at least two phylogenetic analysis methods, most commonly MP and model-based methods (BI and/or ML) [27,35,37]. For tree-based comparisons, we propose the following thresholds as an indication of strongly supported incongru-

ence between cpDNA and nrDNA data sets: MP BS≥70% and PP≥0.95 and/or ML≥70%. The weakly supported clades are considered as potential conflicts [39] and need further examination by sampling additional molecular loci.

## 2.2    Determination of conflicting taxa in cpDNA and nrDNA trees

Given evidence of conflict, it is important to determine which taxa are involved. If only one or a few problematic taxa are involved, one can first run an MP-based ILD test with all taxa and obtain a *P*-value. The suspected problematic taxon is removed and an ILD test is re-run. If the *P*-value markedly increases, the removed taxon is considered to be in conflict between the cpDNA and nrDNA trees [36,37]. For data sets with numerous conflicting taxa, however, repeated cycles of tree comparisons, bootstrapping, pruning, and reanalysis may be impractical. In such cases, the Templeton test (also called the Wilcoxon signed-ranks (WSR) test) can be used to separately test each individual well-supported incongruent node [21]. In addition, Pelser et al. [16] designed a two-step approach to examine complex incongruence involving multiple lineages in which some also show internal incongruence. The largest mutually exclusive lineages in cpDNA and nrDNA trees are first identified by visual comparison. These lineages are then examined for the presence of strongly supported internal incongruence by evaluating branch support values and then subjecting them to ILD tests that only include the taxa of the lineage under investigation.

## 2.3    Provision of possible interpretations for incongruence

Incongruence between cpDNA and nrDNA gene trees may have a real or artificial basis, i.e., biological or artificial reasons. Artificial reasons can be easily identified, whereas biological reasons are usually complicated and need to be carefully inferred.

With failure to reconstruct the correct cpDNA or nrDNA trees, artificial reasons are responsible for incongruence, such as laboratory errors, long-branch attraction, and evolutionary saturation.

Laboratory errors can generate incorrect sequences and thereby result in incongruence. The simultaneous processing of multiple samples at one laboratory table can lead to contamination of the samples (personal observation). To identify potential contaminants, sequences from each gene can be initially analyzed using MP and/or subjected to a BLAST search against the National Center for Biotechnology Information database (http://www.ncbi.nlm.nih.gov). If different species have identical sequences, contamination may have occurred [40].

Incongruence can also be caused by long-branch attraction in one of the two data sets [27], although a long branch

observed on a phylogenetic tree may not always obscure phylogenetic signal [41]. The taxa having potential long-branch problems can first be identified by visual inspection of the phylograms with strongly incongruent lineages [16], or by using relative apparent synapomorphy analysis (RASA; http://test1.bio.psu.edu/LW/list.htm). The taxa subject to long-branch attraction are then removed and the data set is reanalyzed. If the arrangement of other taxa in the tree is changed, this indicates that long-branch attraction has occurred [42]. Increasing the number of characters may overcome some of the problems due to long-branch attraction [43]. The cpDNA and nrDNA data sets can subsequently be combined if the incongruence is caused by long-branch attraction in individual data sets. In addition, MP analysis is sensitive to long-branch attraction [44–46], whereas ML and BI analyses that implement appropriate substitution model(s) are able to largely overcome the problem [41,46,47].

When substitution rates are particularly high, phylogenetic signal may be decreased because of multiple changes that mask evolutionary history and create homoplasy [48], which may result in phylogenetic incongruence. Given that such saturation is apparent only for comparisons of highly divergent taxa, such incongruence is confined primarily to early-diverging taxa in the phylogeny and is not caused by different phylogenetic histories. In this case, data sets can subsequently be combined [49].

When cpDNA and nrDNA data may have different underlying phylogenetic histories, one must use biological reasons to explain incongruence, such as paralogy, hybridization, and incomplete lineage sorting.

If a data set includes paralogous gene copies, the actual phylogeny will partly reflect the duplication history of the gene, which is incongruent with the species divergence history. Plastid DNA markers usually lack paralogous copies because the plastid genome is a single, non-recombining locus [50–52]. Paralogy is usual for multiple-copy nuclear loci. In cases where PCR products of nuclear loci form more than one band, or double peaks or ambiguous base calls are observed in electropherograms, cloning should be attempted. Pseudogenes observed in nrDNA sequences of some taxa should first be identified and then excluded from phylogenetic analyses. Pseudogenes can be detected by searching the well-defined conserved sequence motifs and by pairwise comparison of substitution rates and need to be deleted from phylogenetic analyses [53]. Subsequently, phylogenetic analyses can be combined to obtain a preliminary tree using all sequences. If several clones of one accession occur in the same clade, one representative clone can be selected [20] or a consensus sequence is generated in which polymorphisms are coded as ambiguous characters, and this consensus sequence is used for subsequent phylogenetic analyses [16,54]. The distribution of several clones of one accession in different clades results in paralogy.

Introgression is the result of repeated backcrossing of a hybrid with one or both of its parents. Incomplete lineage sorting can occur when an ancestral species undergoes several speciation events within a short period of time [22]. If, for a given gene, the ancestral polymorphism is not fully resolved into two clades when the second speciation event takes place, the gene tree will likely differ from the species tree [55,56]. Incomplete lineage sorting is notorious for producing patterns similar to those caused by hybridization and introgression. Several methods have been proposed to distinguish between hybridization and incomplete lineage sorting:

(i) Morphological intermediacy.   Hybrids contain a combination of different genotypes and accordingly may display some phenotypic traits that are intermediate between their parental taxa [57]. Thus, morphological intermediacy is widely used as evidence for hybridization [36,58,59]. However, intermediate characters can also arise from convergent morphological evolution or from the existence of ancestral populations from which the two species diverged [60]. Moreover, gene expression in hybrid genotypes may be complicated and the resulting phenotypes may show intermediacy, resemble either parental species, or even exhibit novel character states. Morphology alone can therefore be misleading as evidence of hybridization, as has been demonstrated in recent molecular studies [61,62]. Morphological evidence for hybridization represents only a probable hypothesis. Nevertheless, if the hybrid nature of a species is supported by artificial and/or field experiments, this can be considered as robust evidence for hybridization [39,63,64].

(ii) Comparison of distribution, phenology, and habitats. From inspection of cpDNA and nrDNA trees, one can identify two hypothetical parents if hybridization is responsible for incongruence. Plastid capture is much more likely to take place than nuclear capture, owing to maternal inheritance [65,66], lack of linkage relative to nuclear gene selection [67], and smaller effective population size because of clonal inheritance [68]. If the putative hybrid grows sympatrically with its hypothetical mother, hybridization can be considered to be a factor [53,63]. This is, however, not always the case, as has been shown in *Cornus* [69]. Instead, phenology and habitat can be used to identify possible hybridization events. For example, *Cornus eydeana* QY Xiang & YM Shui was sister to the *Cornus mas* L.-*Cornus officinalis* Seib. & Zucc. clade in the cpDNA tree, whereas *Cornus eydeana* and *Cornus chinensis* Wangerin formed a clade in the nrDNA tree; Xiang et al. [69] hypothesized that the conflict in the position of *Cornus eydeana* in the cpDNA and nrDNA trees was due to cpDNA lineage sorting because the flowering time and habitats of *Cornus eydeana* and the *Cornus mas-Cornus officinalis* clade are non-overlapping, although both are distributed in eastern Asia. However, dispersal in combination with extinction in the parental distribution area, and/or novel phenology or a novel

habitat may invalidate this approach.

(iii) Counting the minimum number of evolution events. Lineage sorting can be distinguished from hybridization by comparing the minimum numbers of evolution events presumed necessary to attain the observed pattern of incongruence [27]. For lineage sorting, the minimum number of multiple lineages that survive through any particular branch segment of a tree has been assessed using GeneTree [70,71]. For hybridization, the minimum number of dispersal and/or extinction events needs to be postulated for probable hybridization between putative parental ancestors based on their present distribution ranges.

(iv) Minimum genetic distance. If incomplete lineage sorting is responsible for incongruence, the similar sequences will have coalesced before the speciation event. If hybridization is responsible for incongruence, the similar sequences from different species could have coalesced either before or after the speciation event. Joly et al. [29] describe a parametric approach for statistically distinguishing some hybridization events from incomplete lineage sorting scenarios based on minimum genetic distances.

(v) Coalescence-based methods. Based on coalescent theory, ancestral polymorphisms are likely to coalesce within approximately 5 $N_e$ generations ($N_e$ being the effective population size) [72,73]. Thus, congruence between gene trees and species trees is highly probable. If incongruence is to be explained by incomplete lineage sorting, one can calculate the assumed minimum $N_e$. If the assumed $N_e$ is much higher than that observed in nature, then incomplete lineage sorting can be excluded, and hybridization is supported as the most likely explanation for the observed incongruence [16].

## 2.4 Reconstruction of the phylogeny after treating incongruence

After conflicting taxa are identified and possible interpretations for incongruence are given, taxa responsible for the conflict are usually removed before the combined analysis of cpDNA and nrDNA data sets is carried out [36,37,74]. This method may be problematic, however, in that the placement of the conflicting taxa cannot be indicated in the larger tree [22]. If hybridization is responsible for the incongruence, van der Niet & Linder grafted subsequently the incongruent taxa onto the tree obtained from the combined analysis [27], whereas Pelser et al. performed the combined analysis of cpDNA and nrDNA data sets by recoding the incongruent taxa twice: once as a cpDNA-only accession (nrDNA characters were scored as missing) and once as a nrDNA-only accession (cpDNA characters were scored as missing) [16]. If nrDNA lineage sorting is responsible for the incongruence, the combined analysis of cpDNA and nrDNA data sets can be carried out by recoding nrDNA characters as missing.

## 3 Conclusion

If cpDNA and nrDNA data do indeed reflect different evolutionary histories, their data sets may result in different topologies, and a phylogenetic tree estimated from the simple combined data set would produce an incorrect estimate of the phylogeny or may sometimes represent an oversimplified version of the genetic history. To more accurately reconstruct the phylogeny and evolutionary history of plant groups, the combined analysis of cpDNA and nrDNA data sets must be done with caution, and if incongruence between the data sets exists, its possible causes should be addressed in detail.

1   APG III. An update of the angiosperm phylogeny group classification for the orders and families of flowering plants: APG III. Bot J Linn Soc, 2009, 161: 105–121

2   Cox C J, Goffinet B, Wickett N J, Boles S B, Shaw A J. Moss diversity: a molecular phylogenetic analysis of genera. Phytotaxa, 2010, 9: 175–195

3   Finet C, Timme RE, Delwiche CF, Marlétaz F. Multigene phylogeny of the green lineage reveals the origin and diversification of land plants. Curr Biol, 2010, 20: 2217–2222

4   Lehtonen S. Towards resolving the complete fern Tree of Life. PLoS ONE, 2011, 6: e24851

5   Soltis DE, Smith SA, Cellinese N, Wurdack KJ, Tank DC, Brockington SF, Refulio-Rodriguez NF, Walker JB, Moore MJ, Carlsward BS, Bell CD, Latvis M, Crawley S, Black C, Diouf D, Xi Z, Rushworth CA, Gitzendanner MA, Sytsma KJ, Qiu YL, Hilu KW, Davis CC, Sanderson MJ, Beaman RS, Olmstead RG, Judd WS, Donoghue MJ, Soltis PS. Angiosperm phylogeny: 17 genes, 640 taxa. Am J Bot, 2011, 98: 704–730

6   Soltis PS, Soltis DE, Chase MW. Angiosperm phylogeny inferred from multiple genes as a tool for comparative biology. Nature, 1999, 402: 402–404

7   Yang Y, Zhou ZK. New insights into the species problem. Sci China Life Sci, 2010, 53: 964–972

8   Shen X, Li X, Sha Z, Yan B, Xu Q. Complete mitochondrial genome of the Japanese snapping shrimp *Alpheus japonicus* (Crustacea: Decapoda: Caridea): gene rearrangement and phylogeny within Caridea. Sci China Life Sci, 2012, 55: 591–598

9   Zheng H, Zhuang W. Four new species of the genus *Hymenoscyphus* (fungi) based on morphology and molecular data. Sci China Life Sci, 2013, 56: 90–100

10   Li X, Gao H, Wang Y, Song J, Henry R, Wu H, Hu Z, Yao H, Luo H, Luo K, Pan H, Chen S. Complete chloroplast genome sequence of *Magnolia grandiflora* and comparative analysis with related species. Sci China Life Sci, 2013, 56: 189–198

11   Liu H, Zhang X, Chen Z, Dong S, Qiu Y. Polyphyly of the fern family Tectariaceae sensu Ching: insights from cpDNA sequence data. Sci China Ser C-Life Sci, 2007, 50: 789–798

12   Rieseberg LH, Soltis DE. Phylogenetic consequences of cytoplasmic gene flow in plants. Evol Trends Plant, 1991, 5: 65–83

13   Sang T. Utility of low-copy nuclear gene sequences in plant phylogenetics. Crit Rev Biochem Molec Biol, 2002, 37: 121–147

14   Small RL, Cronn RC, Wendel JF. Use of nuclear genes for phylogeny

reconstruction in plants. Aust Syst Biol, 2004, 17: 145–170

15 Doyle JJ. Gene trees and species trees: molecular systematics as one character taxonomy. Syst Biol, 1992, 17: 144–163

16 Pelser PB, Kennedy AH, Tepe EJ, Shidler JB, Nordenstam B, Kadereit JW, Watson LE. Patterns and causes of incongruence between plastid and nuclear Senecioneae (Asteraceae) phylogenies. Am J Bot, 2010, 97: 856–873

17 Soltis DE, Kuzoff RK. Discordance between nuclear and chloroplast phylogenies in the *Heuchera* group (Saxifragaceae). Evolution, 1995, 49: 727–742

18 Zhaxybayeva O, Lapierre P, Gogarten JP. Genome mosaicism and organismal lineages. Trends Genet, 2004, 20: 254–260

19 Peng D, Wang XQ. Reticulate evolution in *Thuja* inferred from multiple gene sequences: implications for the study of biogeographical disjunction between eastern Asia and North America. Molec Phylog Evol, 2008, 47: 1190–1202

20 Joly S, Heenan PB, Lockhart PJ. A Pleistocene inter-tribal allopolyploidization event precedes the species radiation of *Pachycladon* (Brassicaceae) in New Zealand. Molec Phylog Evol, 2009, 51: 365–372

21 Mason-Gamer RJ, Kellogg EA. Testing for phylogenetic conflict among molecular data sets in the tribe Triticeae (Gramineae). Syst Biol, 1996, 45: 524–545

22 Wiens JJ. Combining data sets with different phylogenetic histories. Syst Biol, 1998, 47: 568–581

23 Farris JS, Källersjö M, Kluge A G, Bult C. Testing significance of incongruence. Cladistics, 1994, 10: 315–319

24 Swofford DL. PAUP*: Phylogenetic Analysis Using Parsimony, Beta Test Version 40. Sunderland: Sinauer Associates, 1995

25 Templeton AR. Phylogenetic inference from restriction endonuclease cleavage site maps with particular reference to the evolution of humans and the apes. Evolution, 1983, 37: 221–244

26 Shimodaira H, Hasegawa M. Multiple comparisons of log-likelihoods with applications to phylogenetic inference. Molec Biol Evol, 1999, 16: 1114–1116

27 van der Niet T, Linder HP. Dealing with incongruence in the quest for the species tree: a case study from the orchid genus *Satyrium*. Molec Phylog Evol, 2008, 47: 154–174

28 Holder MT, Anderson JA, Holloway AK. Difficulties in detecting hybridization. Syst Biol, 2001, 50: 978–982

29 Joly S, McLenachan PA, Lockhart PJ. A statistical approach for distinguishing hybridization and incomplete lineage sorting. Am Nat, 2009, 174: E54–E70

30 Hipp AL, Hall JC, Sytsma KJ. Congruence versus phylogenetic accuracy: revisiting the incongruence length difference test. Syst Biol, 2004, 53: 81–89

31 Kluge AG. A concern for evidence and a phylogenetic hypothesis of relationships among Epicrates (Boidae, Serpentes). Syst Zool, 1989, 38: 7–25

32 Nixon KC, Carpenter JM. On simultaneous analysis. Cladistics, 1996, 12: 221–241

33 Walker JB, Sytsma KJ. Staminal evolution in the genus *Salvia* (Lamiaceae): molecular phylogenetic evidence for multiple origins of the staminal lever. Ann Bot, 2007, 100: 375–391

34 Zhang J, Wang J, Xia T, Zhou S. DNA barcoding: species delimitation in tree peonies. Sci China Ser C-Life Sci, 2009, 52: 568–578

35 Rindal E, Brower AVZ. Do model-based phylogenetic analyses perform better than parsimony? A test with empirical data. Cladistics, 2011, 27: 331–334

36 Yi T, Miller AJ, Wen J. Phylogenetic and biogeographic diversification of *Rhus* (Anacardiaceae) in the Northern Hemisphere. Molec Phylog Evol, 2004, 33: 861–879

37 Wang W, Wang HC, Chen ZD. Phylogeny and morphological evolution of tribe Menispermeae (Menispermaceae) inferred from chloroplast and nuclear sequences. Perspect Plant Ecol Evol Syst, 2007, 8: 141–154

38 Baum DA, Sytsma KJ, Hoch PC. A phylogenetic analysis of *Epilobium* (Onagraceae) based on nuclear ribosomal DNA sequence. Syst Biol, 1994, 19: 363–388

39 Willyard A, Wallace LE, Wagner WL, Weller SG, Sakai AK, Nepokroeff M. Estimating the species tree for Hawaiian *Schiedea* (Caryophyllaceae) from multiple loci in the presence of reticulate evolution. Molec Phylog Evol, 2011, 60: 29–48

40 Wiens JJ, Kuczynski CA, Stephens PR. Discordant mitochondrial and nuclear gene phylogenies in emydid turtles: implications for speciation and conservation. Biol J Linn Soc, 2010, 99: 445–461

41 Xiang JQ, Moody ML, Soltis DE, Fan C, Soltis PS. Relationships within Cornales and circumscription of Cornaceae—*matK* and *rbcL* sequence data and effects of outgroups and long branches. Molec Phylog Evol, 2002, 24: 35–57

42 Whiting MF. Long-branch distraction and the strepsiptera. Syst Biol, 1998, 47: 134–138

43 Lyons-Weiler J, Hoelzer GA, Tausch RJ. Relative apparent synapomorphy analysis (RASA) I: the statistical measurement of phylogenetic signal. Molec Biol Evol, 1996, 13: 749–757

44 Felsenstein J. Cases in which parsimony or compatibility methods will be positively misleading. Syst Zool, 1978, 27: 401–410

45 Swofford DL, Waddell PJ, Huelsenbeck JP, Foster PG, Lewis PO, Rogers JS. Bias in phylogenetic estimation and its relevance to the choice between parsimony and likelihood methods. Syst Biol, 2001, 50: 525–539

46 Swofford DL, Olsen GJ, Waddell PJ, Hillis DM. Phylogenetic inference. In: Hillis DH, Moritz C, Mable BK, eds. Molecular Systematics. Sunderland: Sinauer Associates, 1996. 407–514

47 Felsenstein J. Evolutionary trees from DNA sequence: a maximum likelihood approach. J Molec Evol, 1981, 17: 368–376

48 DeSalle R, Freedman T, Prager EM, Wilson AC. Tempo and mode sequence evolution in mitochondrial DNA of Hawaiian *Drosophila*. J Molec Evol, 1987, 26: 157–164

49 Baker RH, DeSalle R. Multiple sources of character information and the phylogeny of Hawaiian drosophilids. Syst Biol, 1997, 46: 654–673

50 Palmer JD, Jansen RK, Michaels HJ, Chase MW, Manhart JR. Chloroplast DNA variation and plant phylogeny. Ann Missouri Bot Gard, 1988, 75: 1180–1206

51 Soltis DE, Soltis PS, Doyle JJ. Molecular Systematics of Plants II. Boston: Kluwer Academic Publishers, 1998

52 Soltis DE, Moore MJ, Burleigh JG, Bell CD, Soltis PS. Assembling the angiosperm tree of life: progress and further prospects. Ann Missouri Bot Gard, 2010, 97: 514–526

53 Schneider JV, Schulte K, Aguilar JF, Huertas ML. Molecular evidence for hybridization and introgression in the neotropical coastal desert-endemic *Palaua* (Malveae, Malvaceae). Molec Phylog Evol, 2011, 60: 373–384

54 Pelser PB, Nordenstam B, Kadereit JW, Watson LE. An ITS phylogeny of Tribe Senecioneae (Asteraceae) and a new delimitation of *Senecio* L. Taxon, 2007, 56: 1077–1104

55 Pamilo P, Nei M. Relationship between gene trees and species trees. Molec Biol Evol, 1988, 5: 568–583

56 Tajima F. Evolutionary relationship of DNA sequences in finite populations. Genetics, 1983, 105: 437–460

57 Lexer C, Joseph J, van Loo M, Prenner G, Heinze B, Chase MW, Kirkup D. The use of digital image-based morphometrics to study the phenotypic mosaic in taxa with porous genomes. Taxon, 2009, 58: 349–364

58 Rieseberg LH. Hybrid origins of plant species. Annu Rev Ecol Syst, 1997, 28: 289–359

59 Hardig TM, Brunsfeld SJ, Fritz RS, Morgan M, Orians CM. Morphological and molecular evidence for hybridization and introgression in a willow (*Salix*) hybrid zone. Molec Ecol, 2000, 9: 9–24

60 Dobzhansky TH. Genetics and the Origin of Species. New York: Columbia University Press, 1941

61 Lihová J, Kučera J, Perný M, Marhold K. Hybridization between two polyploid *Cardamine* (*Brassicaceae*) species in northwestern Spain: discordance between morphological and genetic variation patterns. Ann Bot, 2007, 99: 1083–1096

62 Ortego J, Bonal R. Natural hybridisation between kermes (*Quercus*

*coccifera* L) and holm oaks (*Q ilex* L) revealed by microsatellite markers. Plant Biol, 2010, 12: 234–238

63  Yi T, Miller AJ, Wen J. Phylogeny of *Rhus* (Anacardiaceae) based on sequences of nuclear *Nia*-i3 intron and chloroplast *trnC-trnD*. Syst Bot, 2007, 32: 379–391

64  Ramdhani S, Barker N, Cowling RM. Revisiting monophyly in *Haworthia* Duval (Asphodelaceae): incongruence, hybridization and contemporary speciation. Taxon, 2011, 60: 1001–1014

65  Rieseberg LH, Whitton J, Linder CR. Molecular marker incongruence in plant hybrid zones and phylogenetic trees. Acta Bot Neerl, 1996, 45: 243–262

66  Chan KM, Levin SA. Leaky prezygotic isolation and porous genomes: rapid introgression of maternally inherited DNA. Evolution, 2005, 59: 720–729

67  Funk DJ, Omland KE. Species-level paraphyly and polyphyly: frequency, causes, and consequences, with insights from animal mitochondrial DNA. Annu Rev Ecol Evol Syst, 2003, 34: 397–423

68  Petit RJ, Kremer A, Wagner DB. Geographic structure of chloroplast DNA polymorphisms in European oaks. Theor Appl Genet, 1993, 87: 122–128

69  Xiang QY, Manchester SR, Thomas DT, Zhang W, Fan C. Phylogeny, biogeography, and molecular dating of cornelian cherries (*Cornus*, Cornaceae): tracking Tertiary plant migration. Evolution, 2005, 59: 1685–1700

70  Page RDM. GeneTree: comparing gene and species phylogenies using reconciled trees. Bioinformatics, 1998, 14: 819–820

71  Page RDM, Charleston MA. From gene to organismal phylogeny: reconciled trees and the gene tree/species tree problem. Molec Phylog Evol, 1997, 7: 231–240

72  Rosenberg NA. The shapes of neutral gene genealogies in two species: probabilities of monophyly, paraphyly, and polyphyly in a coalescent model. Evolution, 2003, 57: 1465–1477

73  Degnan JH, Rosenberg NA. Gene tree discordance, phylogenetic inference and the multispecies coalescent. Trends Ecol Evol, 2009, 24: 332–340

74  Rodrigo AG, Kelly-Borges M, Bergquist PR, Bergquist PL. A randomization test of the null hypothesis that two cladograms are sample estimates of a parametric phylogenetic tree. New Zeal J Bot, 1993, 31: 257–268