

Genome-wide interrogation advances resolution of recalcitrant groups in the tree of life

Dahiana Arcila^{1,2}, Guillermo Orti¹, Richard Vari^{2,†}, Jonathan W. Armbruster³, Melanie L. J. Stiassny⁴, Kyung D. Ko¹, Mark H. Sabaj⁵, John Lundberg⁵, Liam J. Revell⁶ and Ricardo Betancur-R.^{2,7*}

Much progress has been achieved in disentangling evolutionary relationships among species in the tree of life, but some taxonomic groups remain difficult to resolve despite increasing availability of genome-scale data sets. Here we present a practical approach to studying ancient divergences in the face of high levels of conflict, based on explicit gene genealogy interrogation (GGI). We show its efficacy in resolving the controversial relationships within the largest freshwater fish radiation (Otophysi) based on newly generated DNA sequences for 1,051 loci from 225 species. Initial results using a suite of standard methodologies revealed conflicting phylogenetic signal, which supports ten alternative evolutionary histories among early otophysan lineages. By contrast, GGI revealed that the vast majority of gene genealogies supports a single tree topology grounded on morphology that was not obtained by previous molecular studies. We also reanalysed published data sets for exemplary groups with recalcitrant resolution to assess the power of this approach. GGI supports the notion that ctenophores are the earliest-branching animal lineage, and adds insight into relationships within clades of yeasts, birds and mammals. GGI opens up a promising avenue to account for incompatible signals in large data sets and to discern between estimation error and actual biological conflict explaining gene tree discordance.

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA

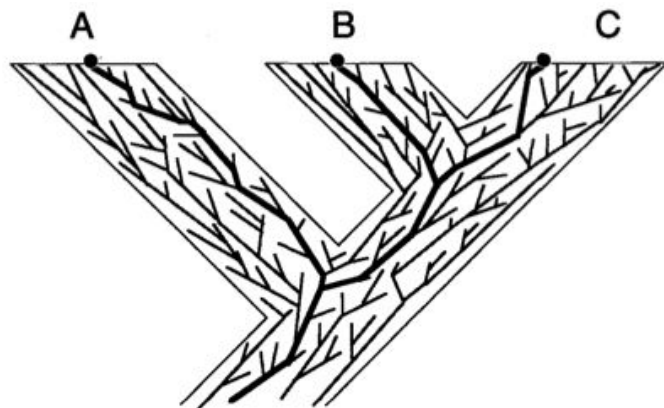


FIGURE 1. A gene tree contained within a species tree leading to three extant species: A, B, and C. Bold branches of gene tree show relationships among the sampled copies of the gene (●). Sampled copies from sister species B and C are sister copies.

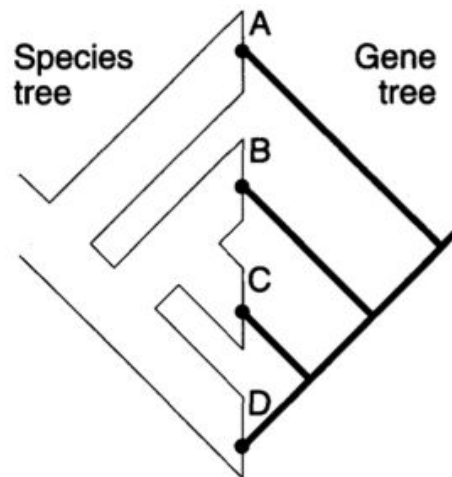


FIGURE 2. Discord between gene and species trees. At left is the species tree of four species, A, B, C, and D, and at right is the tree of a gene sampled one copy per species. Species B and C are sister species, but their gene copies are not sister copies.

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA

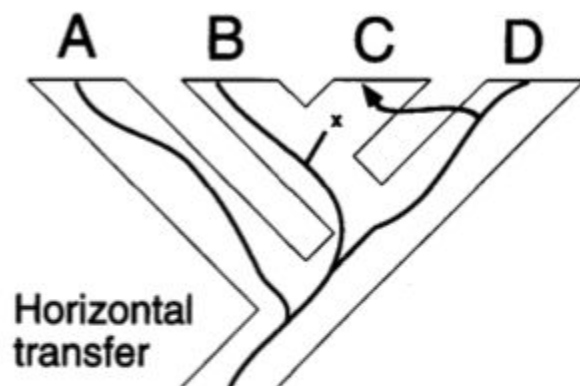


FIGURE 3. Horizontal transfer. A branch of the gene tree jumps between species lineages. If the indigenous gene copy in the receiving species lineage goes extinct or is not sampled (x), then the gene tree will disagree with the species tree, as shown in Figure 2.

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA

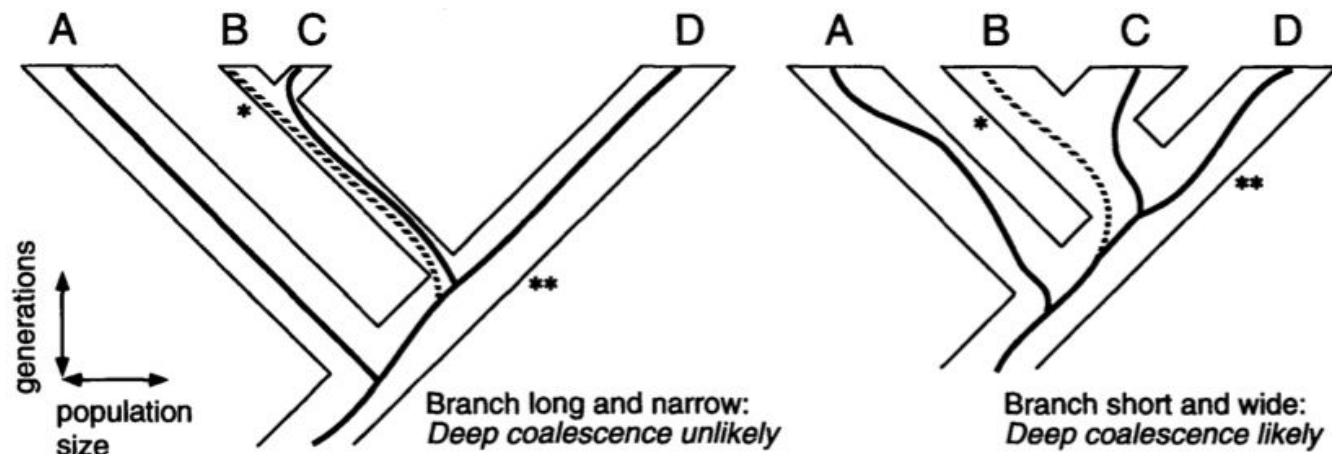


FIGURE 4. Lineage sorting (deep coalescence). Described in a time-forward sense as lineage sorting, an ancestral polymorphism at ** is retained through a lineage to the next speciation event at *, where different forms are sampled in different descendant species. Described in a time-backward sense as deep coalescence, two gene copies from species B and C meet at * but fail to coalesce until deeper than the speciation event at **, at which point the gene from C coalesces first with the gene from D. Failure to coalesce is more likely the shorter (in generations) and wider (in effective population size) the branch is between ** and *.

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA

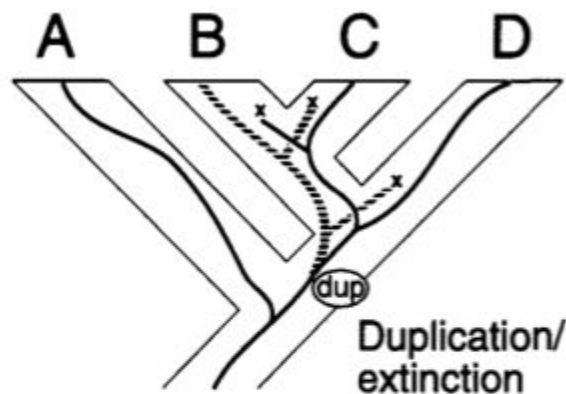


FIGURE 5. Gene duplication and extinction (or paralogous sampling). The gene is duplicated to a different locus, indicated by the dashed lines. If in descendant species one or the other locus goes extinct or is not sampled (\times), then the gene tree will disagree with the species tree, as shown in Figure 2.

GENE TREES IN SPECIES TREES

WAYNE P. MADDISON

Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, USA



FIGURE 9. Phylogeny as a cloud of gene histories. Phylogeny is more like a statistical distribution than a simple tree of discrete thin branches. It has a central tendency, but it also has a variance because of the diversity of gene trees. Gene trees that disagree with the central tendency are not wrong; rather, they are part of the diffuse pattern that is the genetic history.

Estimating phylogenetic trees from genome-scale data

Liang Liu,^{1,2} Zhenxiang Xi,³ Shaoyuan Wu,⁴ Charles C. Davis,³ and Scott V. Edwards³

Estimating phylogenetic trees from genome-scale data

Liang Liu,^{1,2} Zhenxiang Xi,³ Shaoyuan Wu,⁴ Charles C. Davis,³ and Scott V. Edwards³

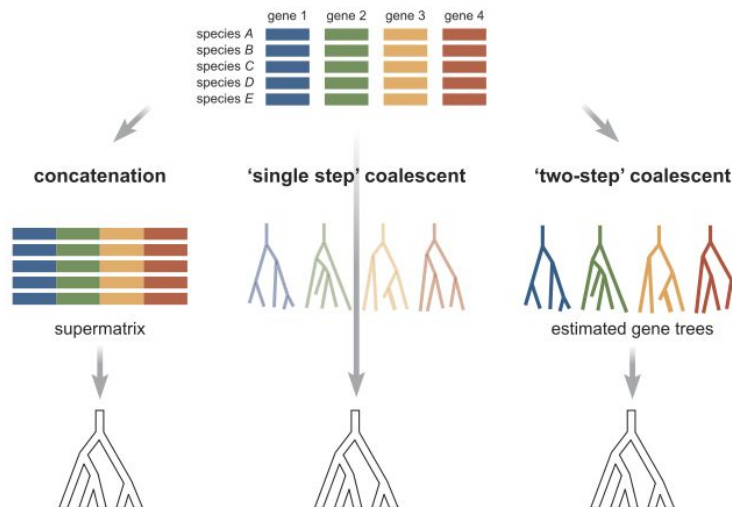


Figure 1. Schematic of the concatenation and coalescent paradigms in phylogenetics. At the top is depicted a multilocus data set consisting of five species (A–E) and four genes (1–4). On the left is indicated the classic supermatrix approach, in which all genes are concatenated to produce a single supergene, which is then subjected to phylogenetic analysis by classical or updated traditional algorithms, such as RAxML,¹⁰⁷ ExaBayes,⁸⁴ or MrBayes.¹⁰⁸ Although the resulting tree at lower left is in truth a gene tree, it is often called a species tree or phylogeny because it is the result of analysis of a complete data set. In the center is depicted a class of species tree (coalescent) methods in which both gene trees and species trees are estimated concurrently according to multilocus sequence data, priors, and a multispecies coalescent likelihood model. Examples of algorithms estimating species trees in this way include *BEAST⁴⁹ and BEST.¹⁰⁹ On the right are depicted so-called two-step species tree methods, in which gene trees are first estimated using classical approaches and then used as input data to estimate a species tree using algorithms such as MP-EST,⁴⁸ STAR,⁴⁷ STEM,⁸⁷ NJst,¹¹⁰ STELLS,⁸⁸ or ASTRAL.⁸⁹ The methods depicted typically used sets of loci each consisting of linked DNA sites.² Yet other species tree methods such as SNAPP⁹⁰ and quartet inference⁹² use unlinked SNPs rather than linked sites, and require different statistical models.

Estimating phylogenetic trees from genome-scale data

Liang Liu,^{1,2} Zhenxiang Xi,³ Shaoyuan Wu,⁴ Charles C. Davis,³ and Scott V. Edwards³

Table 1. Studies evaluating the robustness of species tree phylogenetic methods to various genetic forces and sampling schemes

Topic	Conclusions/comments ^a
General violation of multispecies coalescent model ¹⁰⁴	<ul style="list-style-type: none">• Claims the majority of multilocus sequence data sets are a poor fit to the multispecies coalescent model, although much of the violation stems from fit of substitution model or unknown sources on a minority of genes.
Gene flow ^{73,105}	<ul style="list-style-type: none">• The coalescent method is robust to low levels of gene flow• Concatenation performs poorly relative to the coalescent methods in the presence of gene flow.• Gene flow can lead to overestimation of population sizes and underestimation of species divergence times in species trees.
Sampling/mutation ^{33,74,94}	<ul style="list-style-type: none">• Increased sampling of individuals per species can significantly improve the estimation of shallow species trees.• Sampling more individuals does not significantly improve accuracy in estimating deep species trees. Adding more loci can improve the estimation of deep relationships.• Mutational variance is a major source of error in estimates of species trees.

Estimating phylogenetic trees from genome-scale data

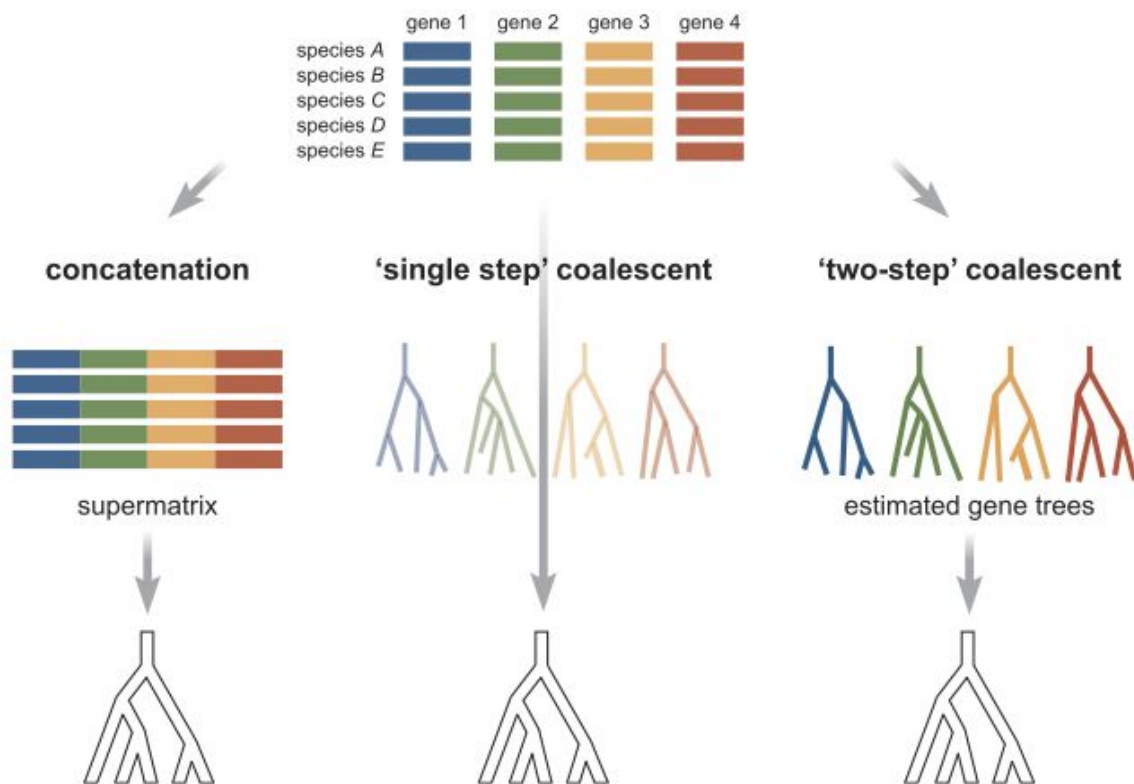
Liang Liu,^{1,2} Zhenxiang Xi,³ Shaoyuan Wu,⁴ Charles C. Davis,³ and Scott V. Edwards³

Table 1. Studies evaluating the robustness of species tree phylogenetic methods to various genetic forces and sampling schemes

Topic	Conclusions/comments ^a
Recombination ^{39,106}	<ul style="list-style-type: none">• Recombination has minor effect on species tree estimation except on extremely short species trees.• The negative effects of recombination can be easily overcome by increased sampling of alleles
Missing data ^{19,63,67}	<ul style="list-style-type: none">• Missing data can decrease the support of species tree estimates• Missing data can significantly affect the accuracy of species tree estimation• Species tree methods are “remarkably resilient” to missing data⁶⁸
Taxon sampling ¹⁵	<ul style="list-style-type: none">• Compared to concatenation, coalescent methods are more robust to poor taxon sampling
Long-branch attraction ²⁹	<ul style="list-style-type: none">• Species tree methods are more resilient to the effects of long-branch attraction than concatenation methods
Random rooting of gene trees ^{36,37}	<ul style="list-style-type: none">• Misrooting of gene trees can mimic the coalescent process
Other ⁸⁰	<ul style="list-style-type: none">• Anomalous gene trees are unlikely to pose a significant danger to empirical phylogenetic study, in part because species trees in the anomaly zone are likely to be rare.

Estimating phylogenetic trees from genome-scale data

Liang Liu,^{1,2} Zhenxiang Xi,³ Shaoyuan Wu,⁴ Charles C. Davis,³ and Scott V. Edwards³



Genome-wide interrogation advances resolution of recalcitrant groups in the tree of life

Dahiana Arcila^{1,2}, Guillermo Ortí¹, Richard Vari^{2,†}, Jonathan W. Armbruster³, Melanie L. J. Stiassny⁴, Kyung D. Ko¹, Mark H. Sabaj⁵, John Lundberg⁵, Liam J. Revell⁶ and Ricardo Betancur-R.^{2,7★}

A major assumption of these ‘summary’¹³ or ‘short-cut’¹⁴ coalescent methods is that individual gene trees accurately depict the genealogical history of fragments of the genome that independently segregate (coalescent genes, or c-genes).

Genome-wide interrogation advances resolution of recalcitrant groups in the tree of life

Dahiana Arcila^{1,2}, Guillermo Orti¹, Richard Vari^{2,†}, Jonathan W. Armbruster³, Melanie L. J. Stiassny⁴, Kyung D. Ko¹, Mark H. Sabaj⁵, John Lundberg⁵, Liam J. Revell⁶ and Ricardo Betancur-R.^{2,7*}

the analysis of short, recombination-free genes (consisting of a few hundred sites) are error-prone due to limited signal-to-noise content^{2,3,14–17}. On the other extreme, long genes or full-length transcripts with thousands of sites harbor more phylogenetic information, reducing (but not necessarily removing) stochastic error^{2,18}. Longer genes, however, are more likely to carry past recombination events, violating the assumption of a single genealogical history⁸. Both situations lead to statistical inconsistency under the multi-species coalescent,

Genome-wide interrogation advances resolution of recalcitrant groups in the tree of life

Dahiana Arcila^{1,2}, Guillermo Ortí¹, Richard Vari^{2,†}, Jonathan W. Armbruster³, Melanie L. J. Stiassny⁴, Kyung D. Ko¹, Mark H. Sabaj⁵, John Lundberg⁵, Liam J. Revell⁶ and Ricardo Betancur-R.^{2,7★}

Here, we present a phylogenomic approach that efficiently extracts the genealogical signal from short c-genes by reducing the complexity of tree space on the basis of topological constraints. This method is similar to others that place priors on gene tree topologies²⁵, but is unique in that priors are set to test specific hypotheses directly.

Mammal madness: is the mammal tree of life not yet resolved?

Nicole M. Foley¹, Mark S. Springer² and Emma C. Teeling¹

Wodniok et al. BMC Evolutionary Biology 2011, 11:104
<http://www.biomedcentral.com/1471-2148/11/104>

RESEARCH ARTICLE

BMC
Evolutionary Biology

Open Access

Origin of land plants: Do conjugating green algae hold the key?

Sabina Wodniok^{1†}, Henner Brinkmann^{2†}, Gernot Glöckner³, Andrew J Heidel⁴, Hervé Philippe², Michael Melkonian¹ and Burkhard Becker^{1*}

THE
SOCIETY

The evolution of the Ecdysozoa

Maximilian J. Telford*, Sarah J. Bourlat, Andrew Economou,
Daniel Papillon and Omar Rota-Stabelli

ca Scripta

KUNGL.
VETENSKAPS-
AKADEMIEN
THE ROYAL SWEDISH ACADEMY OF SCIENCES

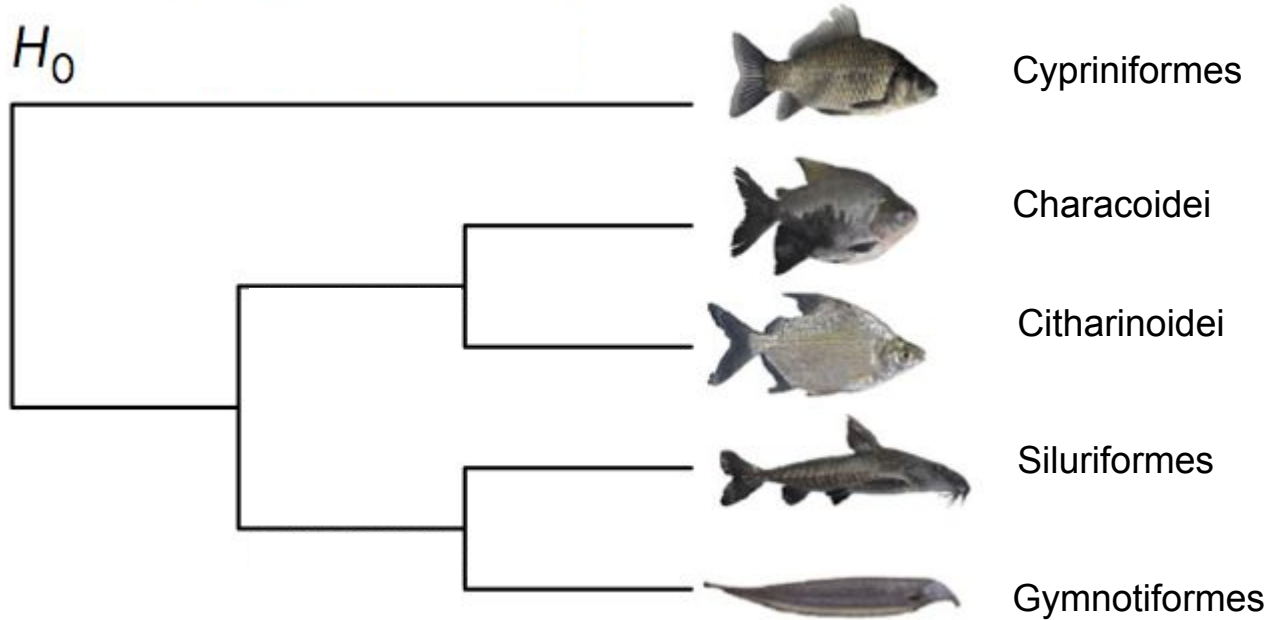


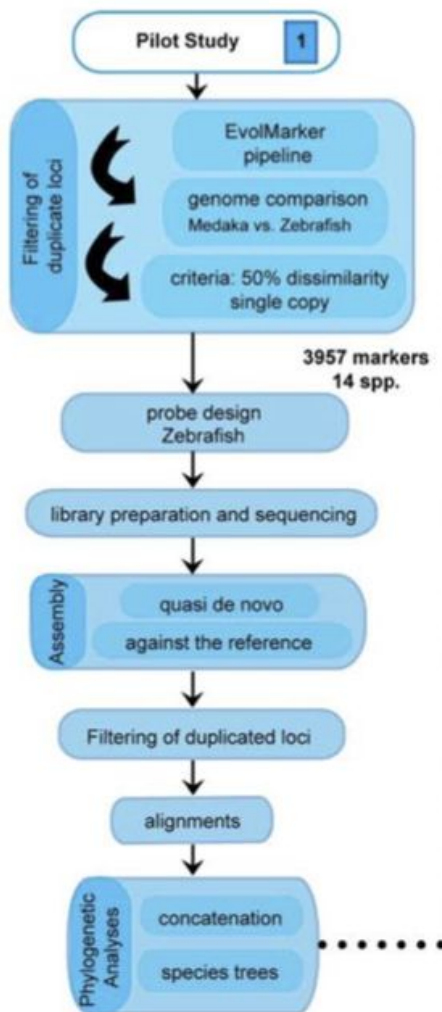
The phylogenomic forest of bird trees contains a hard polytomy at the root of Neoaves

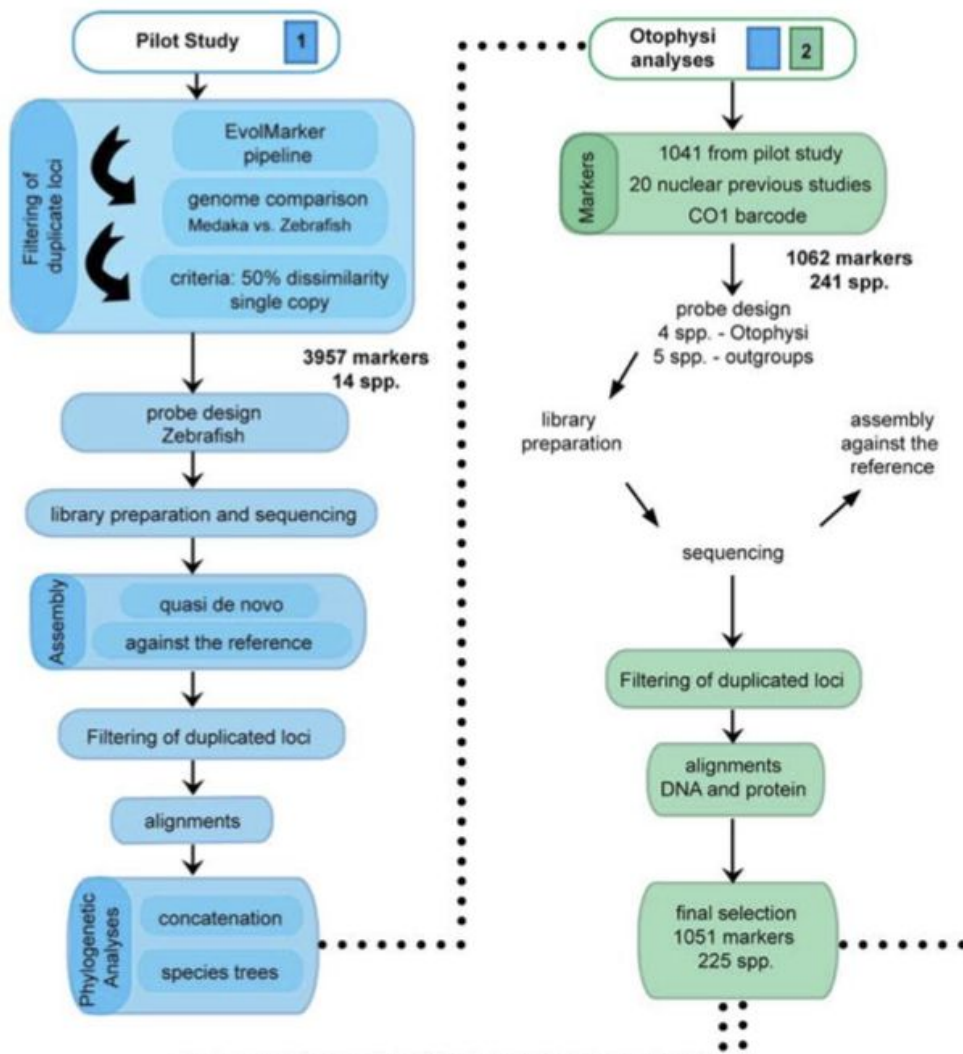
ALEXANDER SUH

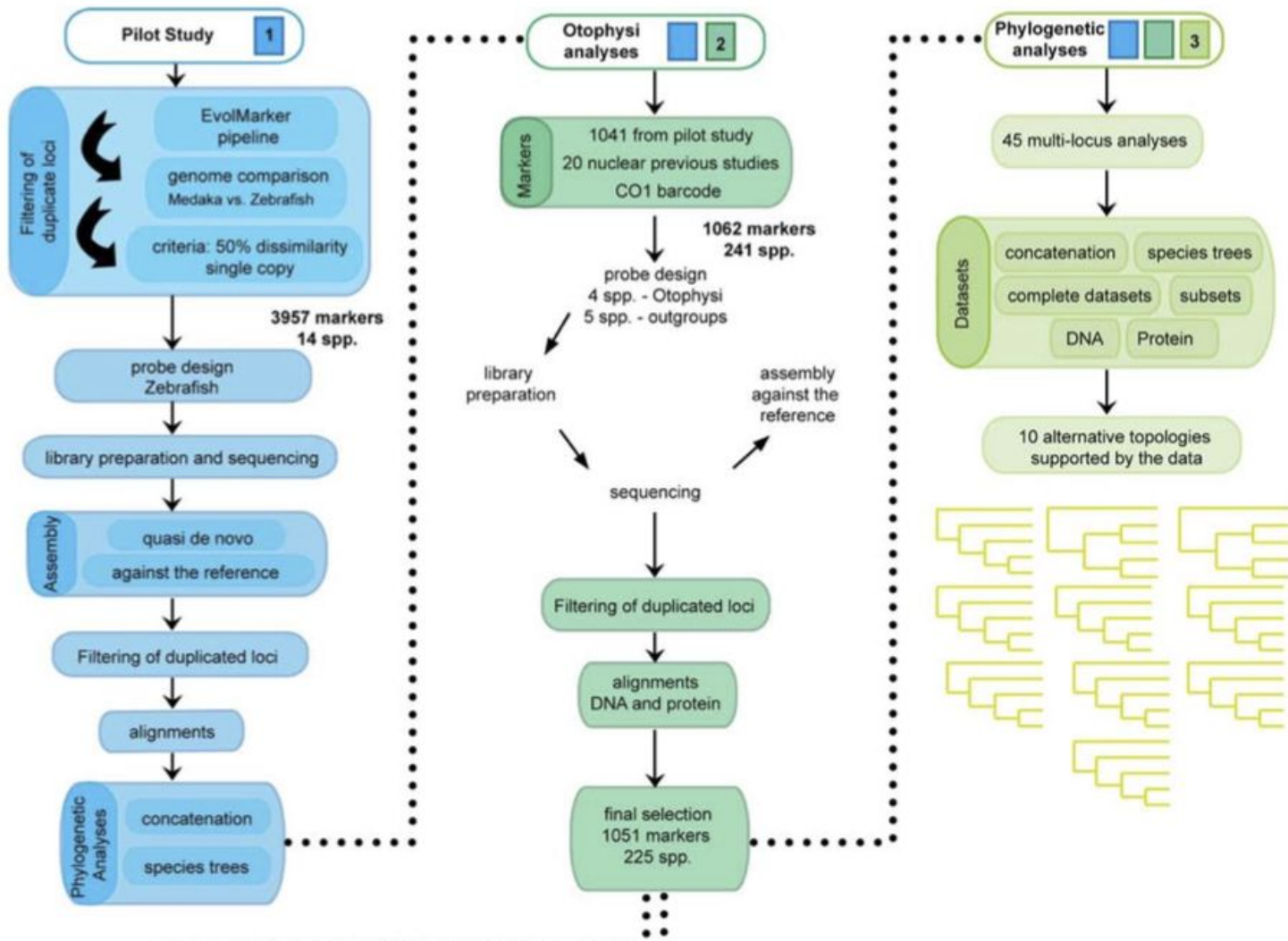
Phil. Trans. R. Soc. B (2008) 363, 1529–1537
doi:10.1098/rstb.2007.2243
Published online 11 January 2008

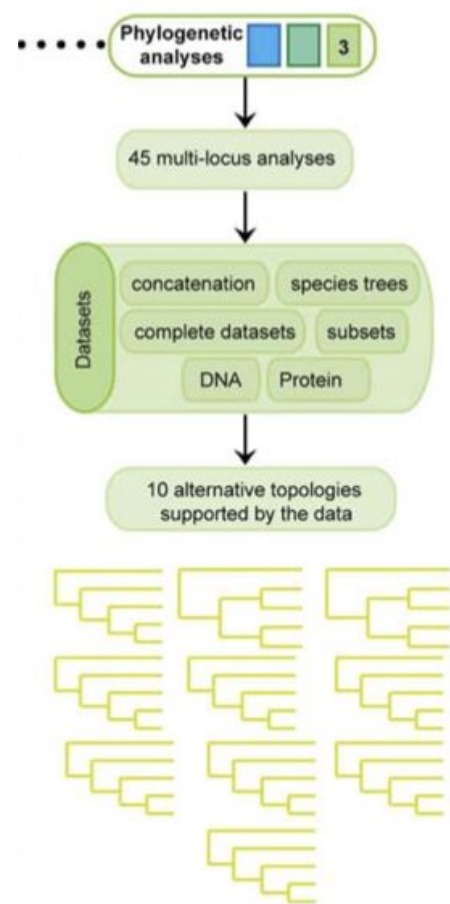
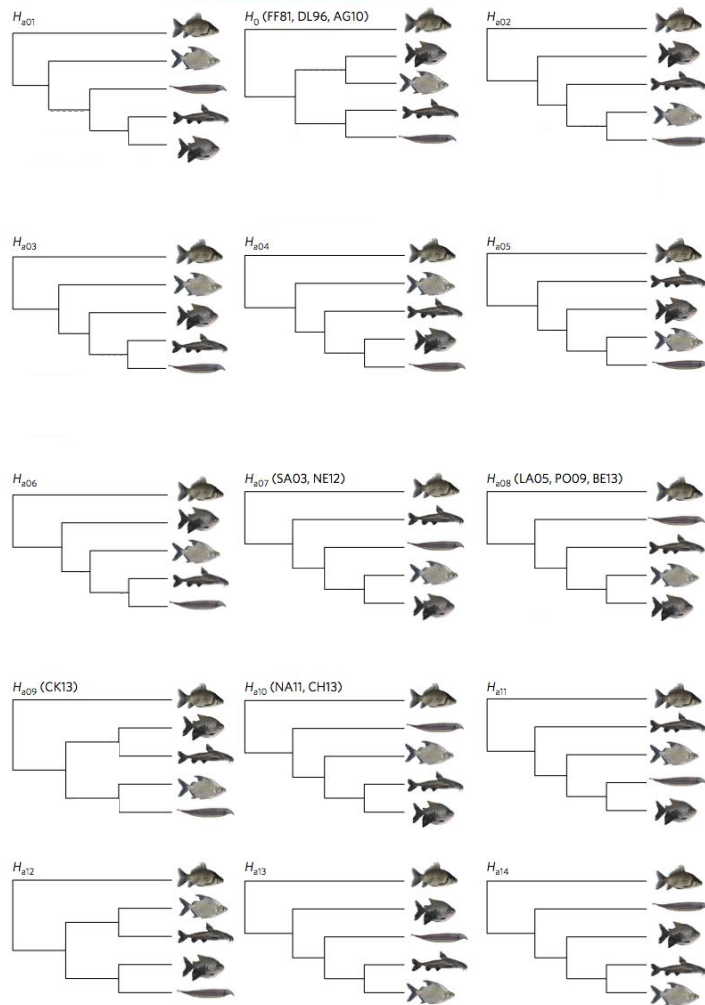
H_0











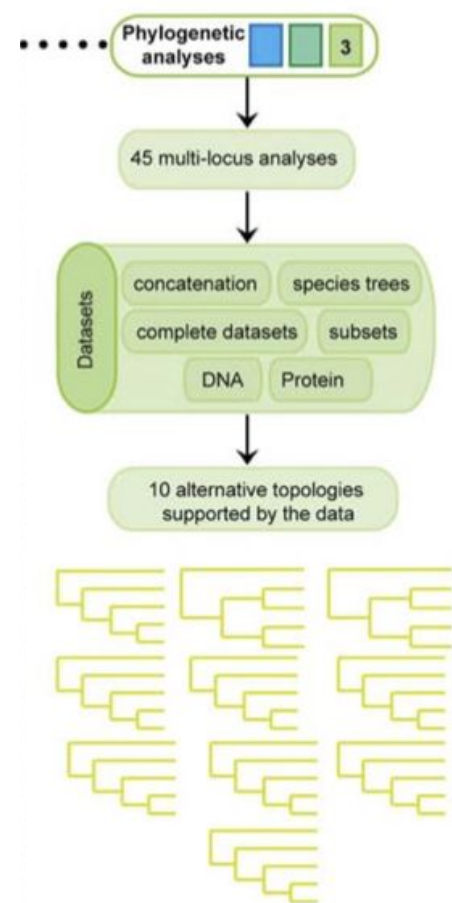
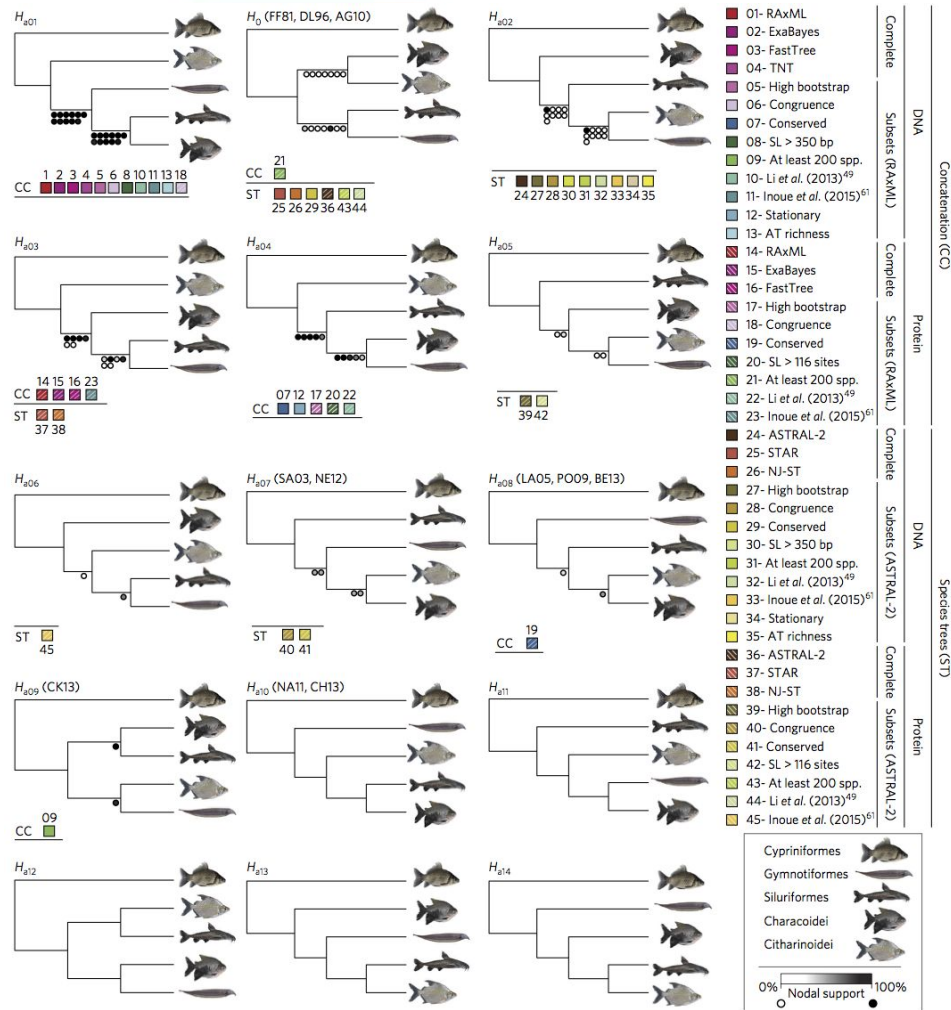


Figure 1 | Null morphological hypothesis (H_0) and all 14 possible alternative trees for the five major lineages in Otophysi. Previous studies supporting each

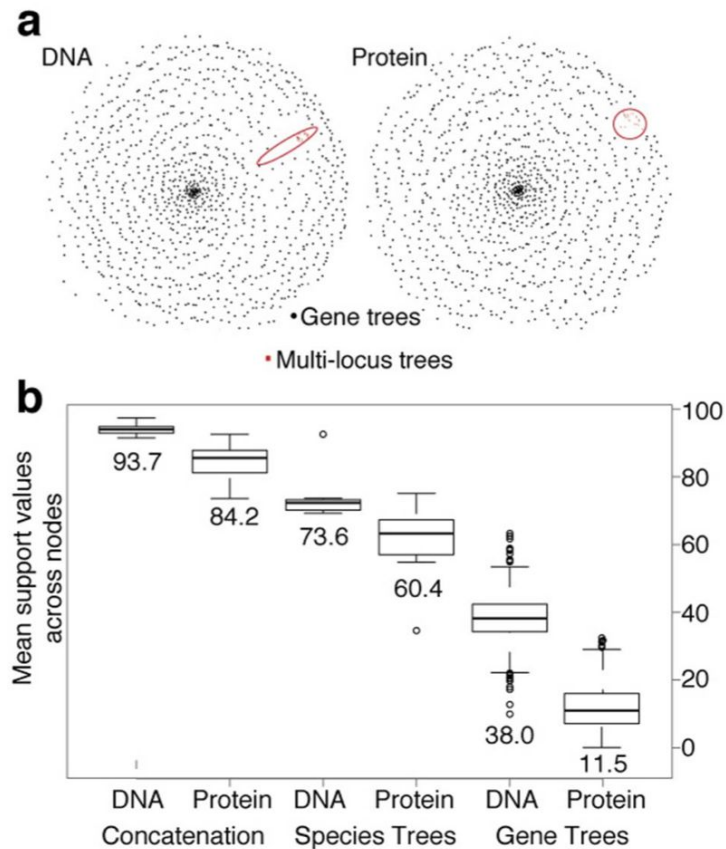
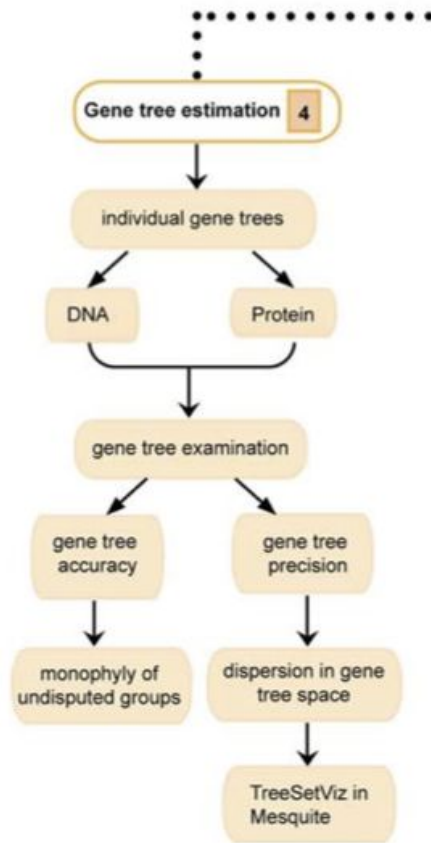
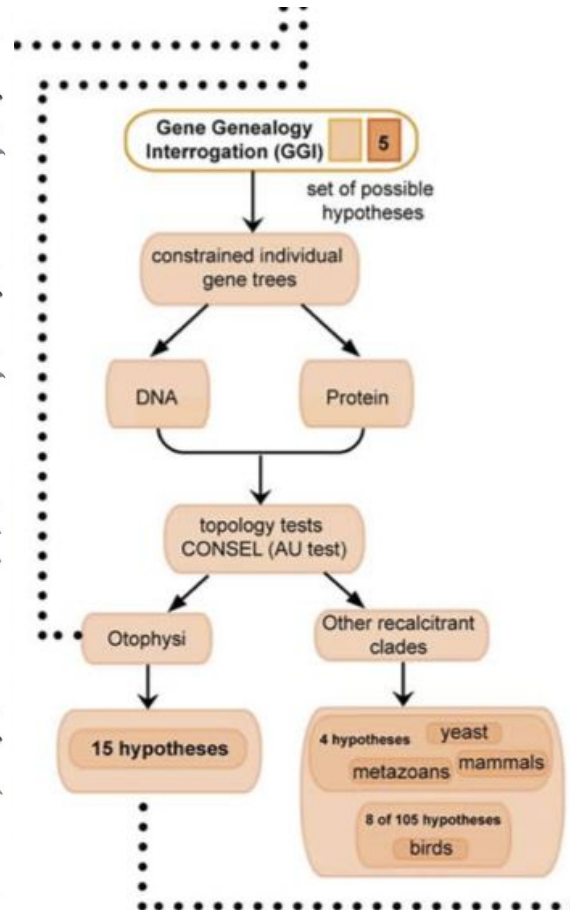
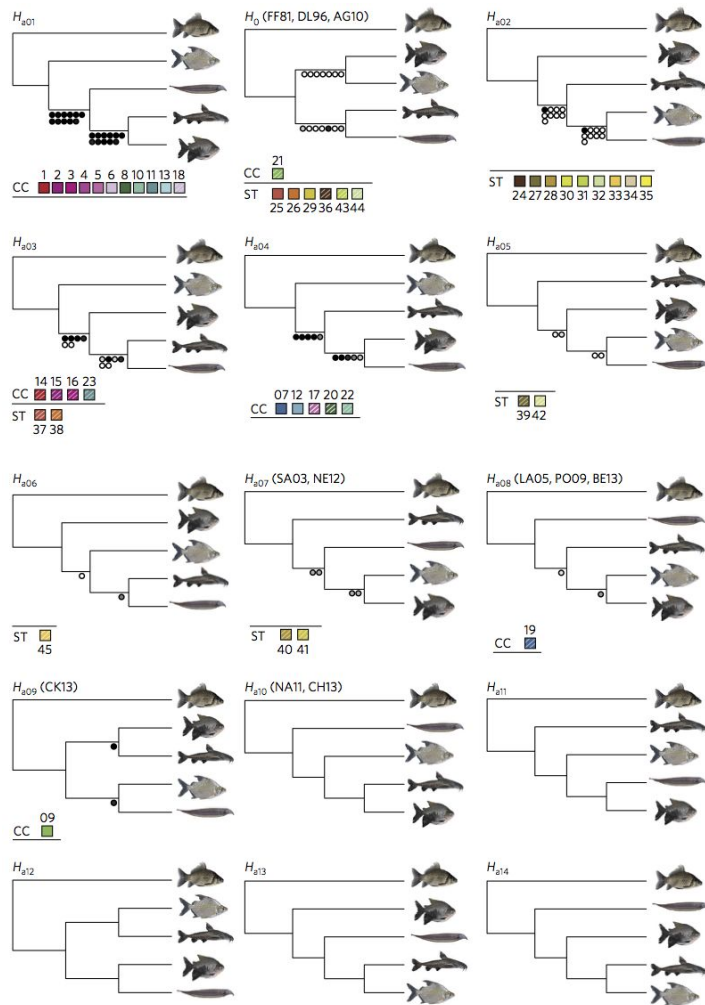
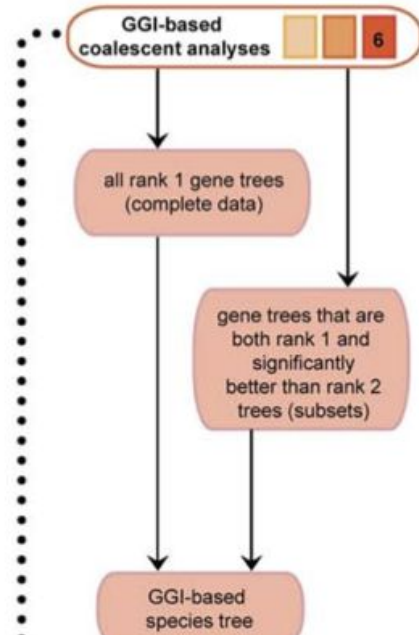
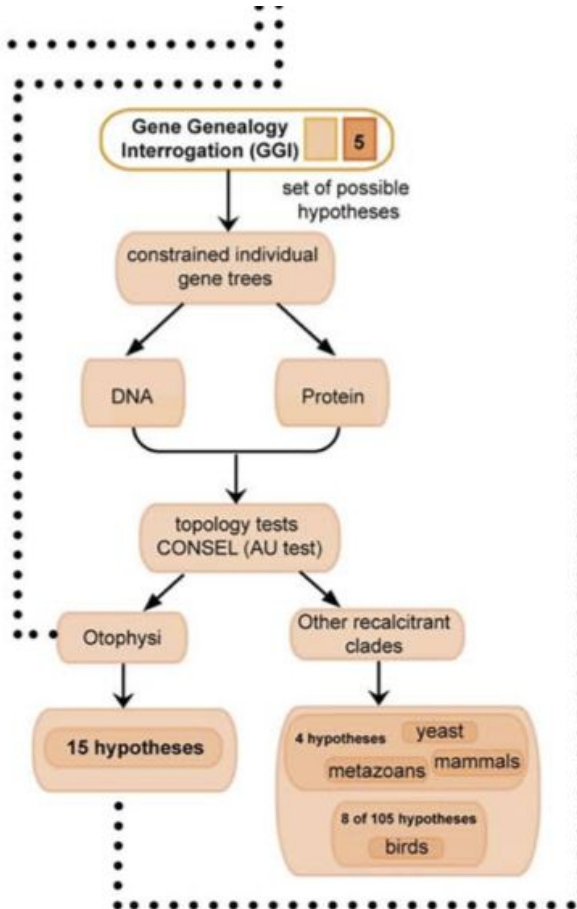
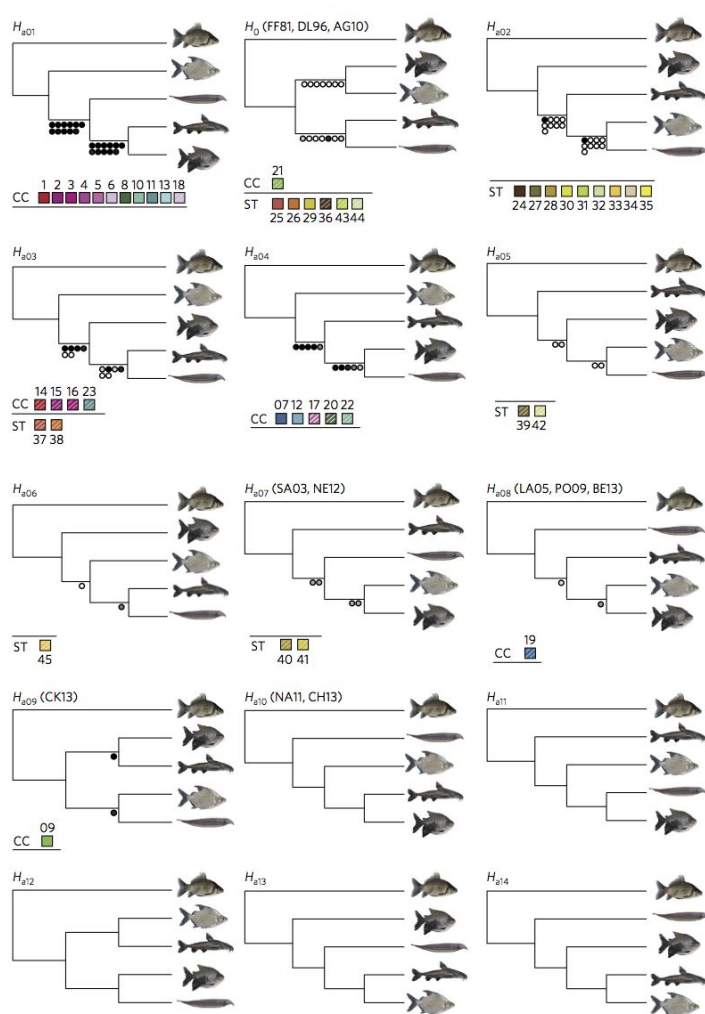
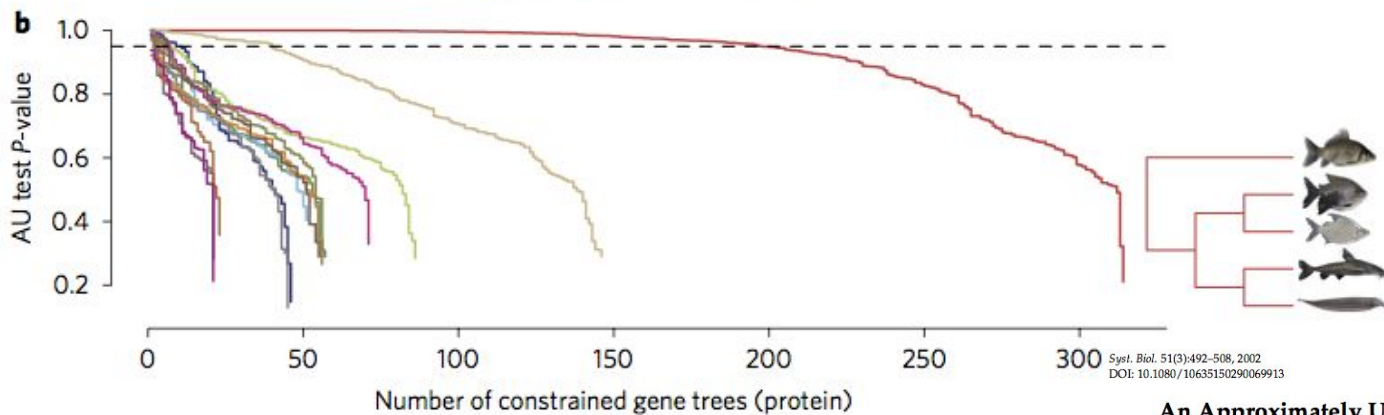
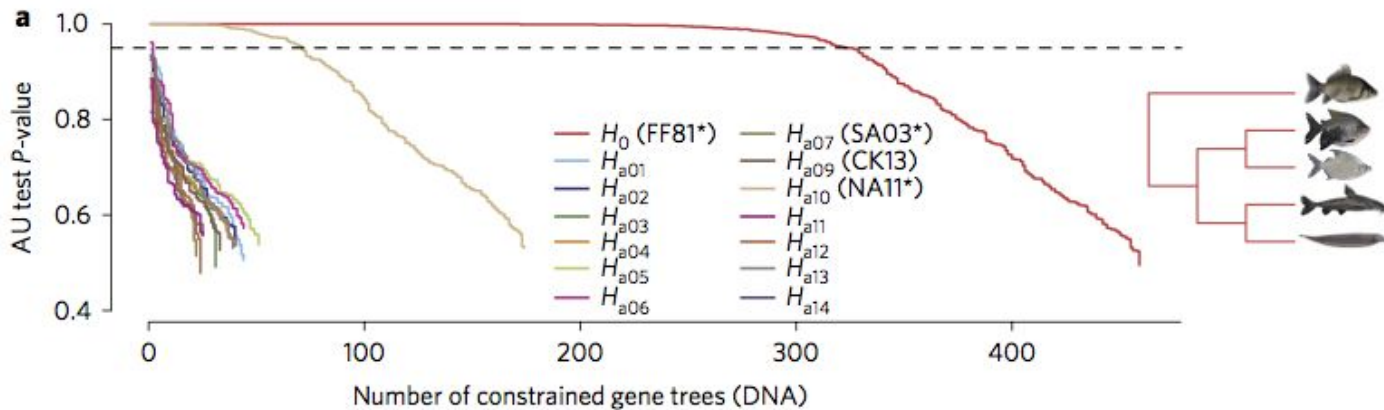


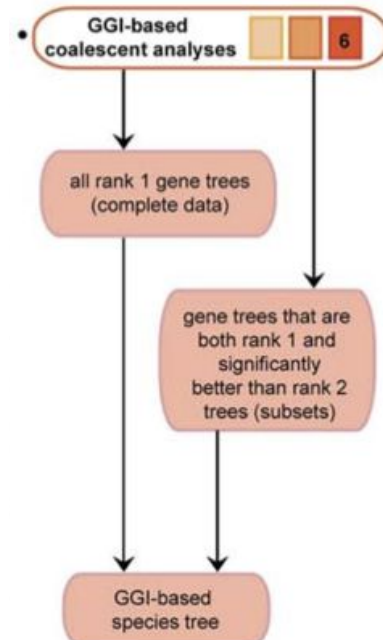
Figure S2. Assessments of phylogenetic precision for multi-locus trees and individual gene trees | **a**, dispersion in multidimensional tree space based on un-weighted Robinson-Foulds distances; **b**, plots of mean support values across all nodes for each tree obtained with different methods (ExaBayes: posterior probabilities; all other analyses: bootstrap values).







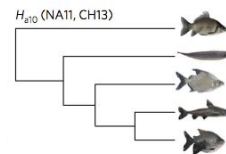
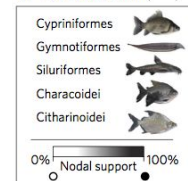
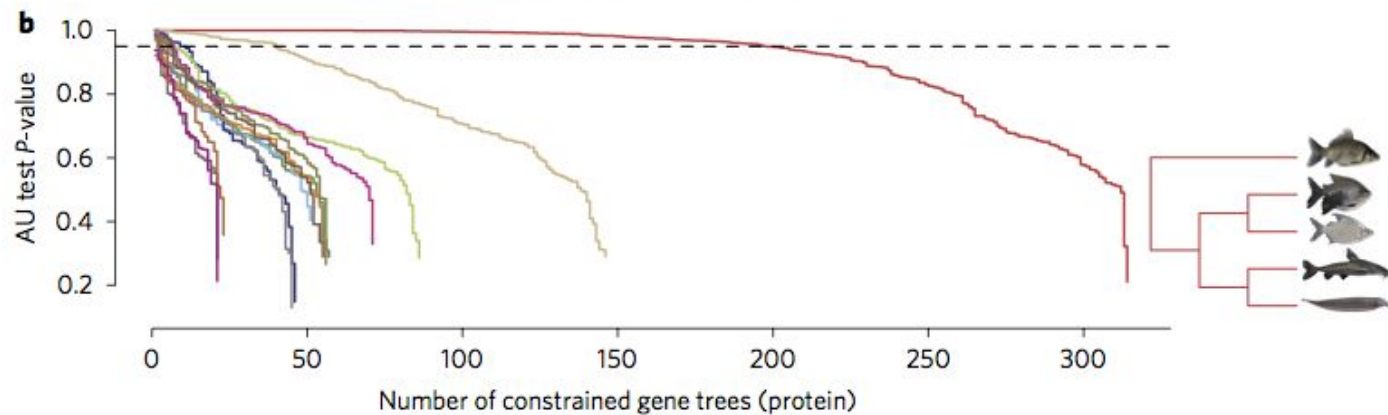
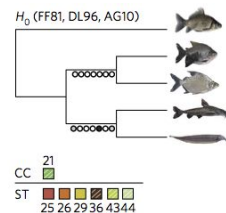
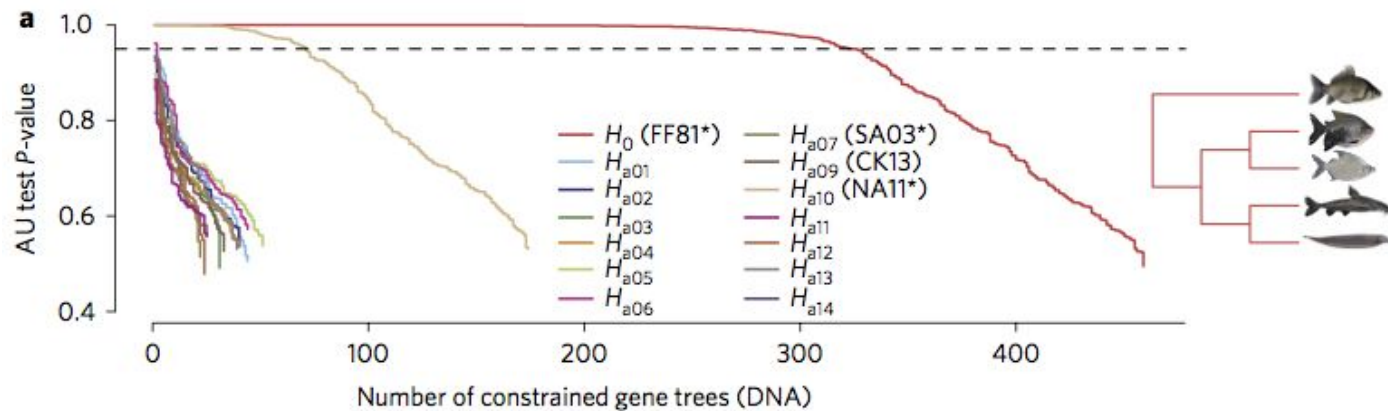
Syst. Biol. 51(3):492–508, 2002
 DOI: 10.1080/10635150290069913

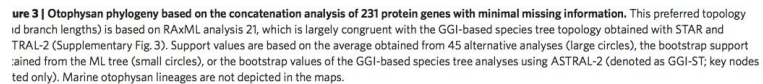
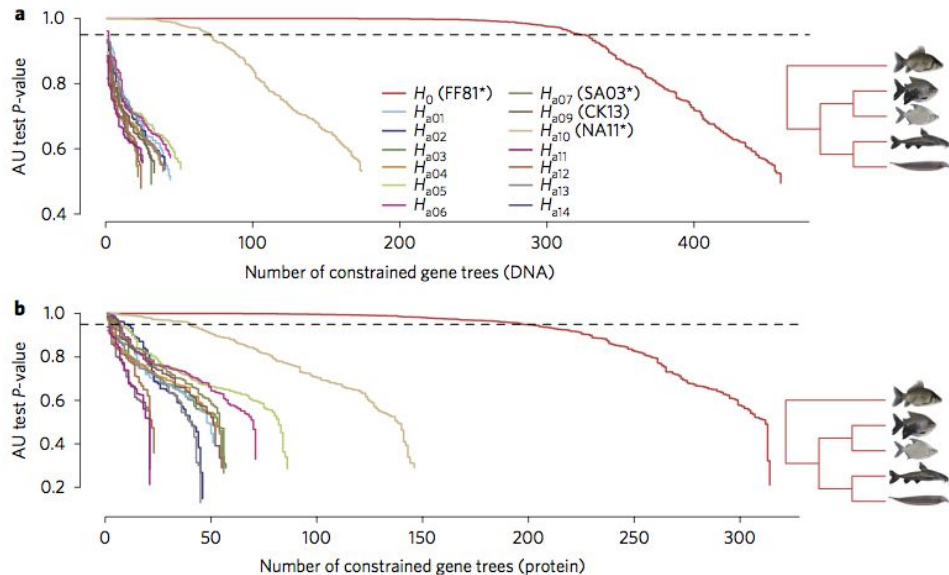


An Approximately Unbiased Test of Phylogenetic Tree Selection

HIDETOSHI SHIMODAIRA

Institute of Statistical Mathematics, 4-6-7 Minami-Azabu, Minatoku, Tokyo 106-8569, Japan;
 E-mail: shimo@ism.ac.jp





ILS is not the main problem

Coalescent theory predicts that phylogenetic histories of lineages evolving under a combination of short internal branches and large effective population sizes are prone to high incidence of ILS⁸. It has been demonstrated that for five or more lineages such conditions can generate gene trees with topologies that differ from the underlying species phylogeny with highest probability^{36,37}. When the evolutionary history of a clade falls within this so-called anomaly zone⁸, simply adopting the most frequent gene tree as a surrogate for the species phylogeny (the democratic vote procedure) is positively misleading.

To account for this possibility, as the genuine backbone tree of inter-ordinal relationships must be 1 of 15 possibilities (enumeration of all possibilities for an unrooted tree of five taxa), we used the GGI trees selected by the topology tests (the preferred constrained gene trees optimized by ML) as input for summary coalescent analyses. For this test we employed both DNA- and protein-based trees in combination with two different species-tree methods. We also applied two alternative approaches for sampling GGI trees, one using all rank 1 trees (complete data with 1,051 genes) and another using only the set of rank 1 trees that are significantly better than the alternatives ($P < 0.05$; a subset of 397 DNA trees and 275 protein trees; Supplementary Table 3). Of the eight species-tree analyses conducted, all converged on the H_0 tree, with each backbone node receiving 100% bootstrap support. Finally, an adapted version of the GGI-based coalescent method that uses constrained topologies in combination with unconstrained gene trees also supports the H_0 tree (Supplementary Information).

Our results suggest that the evolutionary history of major otophysan lineages is not trapped in the anomaly zone. In fact, these analyses identify only a minor proportion of gene trees that are significantly discordant with the inferred species phylogeny (17.7–28.4%, most supporting H_{a10}), suggesting that other sources of error rather than ILS are likely the main cause of incongruence. Gene tree estimation error may be biasing summary coalescent approaches, but the causes for discrepancy between coalescent and concatenation results are unclear. For two hypotheses (H_0 and H_{a03}), some concatenation and species tree methods converge, but more often they seem to produce non-overlapping sets of results (Fig. 1). We were unable to isolate any single factor as the principal explanation for discordance in multi-locus analysis. Possibilities include the combination of slight model misspecifications interacting in analyses of large data sets and amplifying systematic biases, or processes such as horizontal gene transfer or duplication/extinction affecting some of the sampled genes³⁸. What is perhaps most surprising is the observation that the most common topology from concatenation is incongruent with our GGI tree, even in the absence of evidence for substantial ILS. An investigation of factors that could account for this pattern would be a fruitful subject of future theoretical and analytical studies. In summary, the coalescent analyses using GGI trees resolve with high confidence the branching order of major otophysan groups (Supplementary Fig. 3), a result that is fully congruent with the morphological hypothesis (H_0)²⁷, thereby reconciling a long history of molecular and morphological conflict.

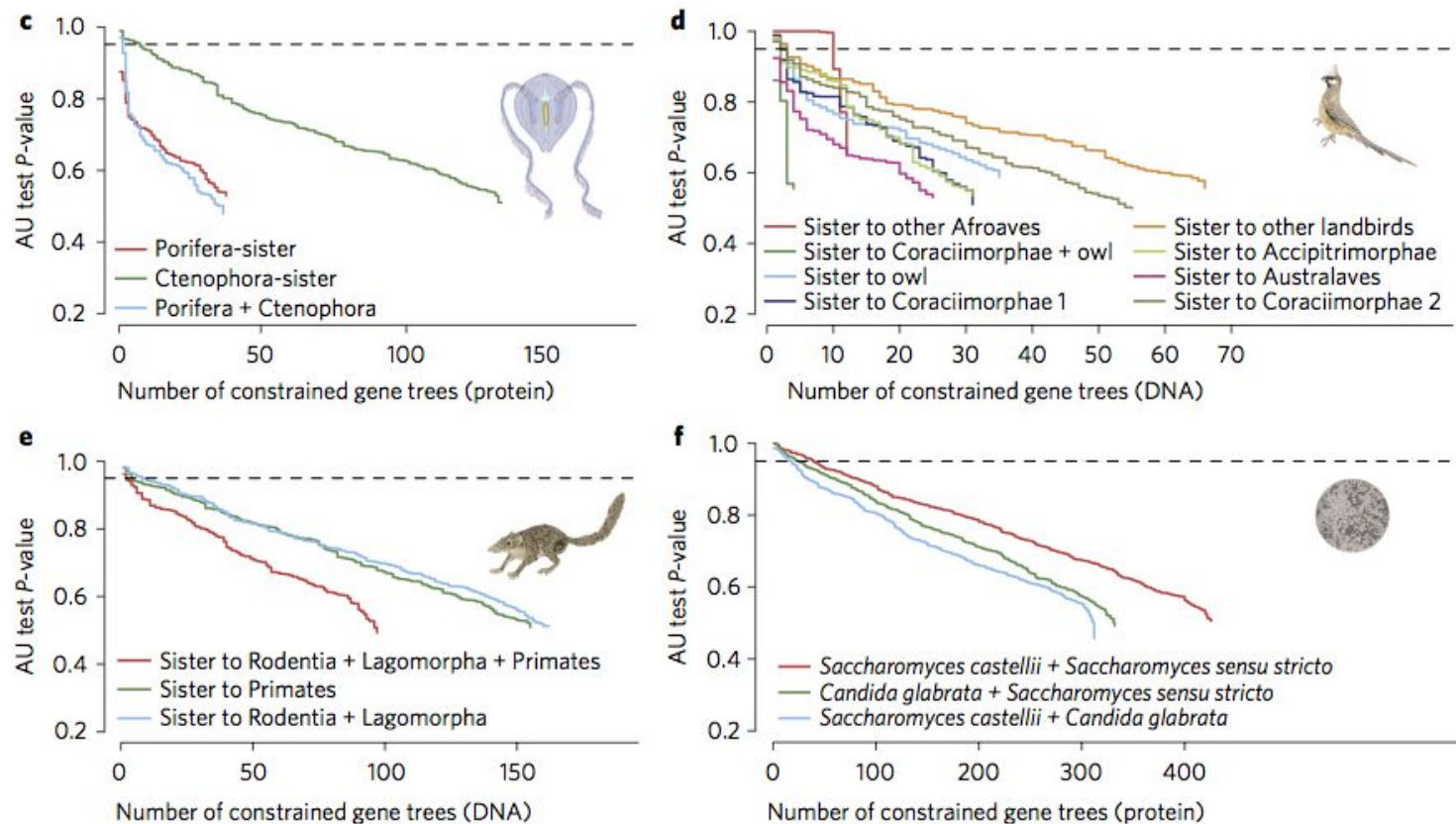
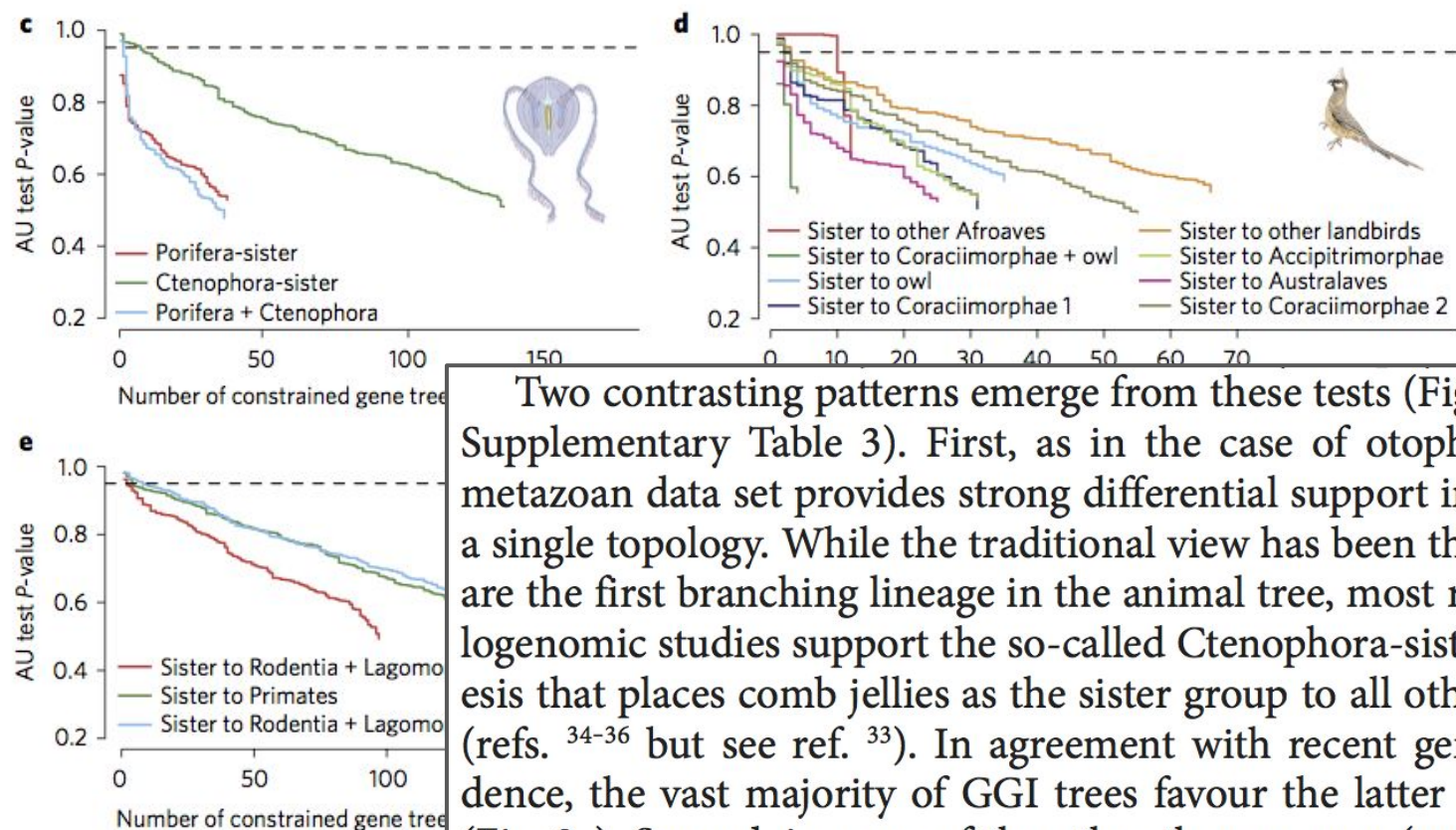


Figure 2 | Gene genealogy interrogation (GGI) applied to phylogenomic data sets to test alternative hypotheses. Lines represent the cumulative number of genes (x axes) supporting each hypothesis with highest probability (rank 1) and their associated P-values (y axes) according to the approximately unbiased (AU) topology test. Values above the dashed line indicate all rank 1 hypotheses that are significantly better than the alternatives ($P < 0.05$), whereas those below the dashed line are also rank 1 but without statistical significance. **a,b**, Otophysi (*see Supplementary Table 1). **c**, Metazoans. **d**, Neoaves (mousebird). **e**, Eutherian mammals (tree shrew). **f**, Yeast phylogeny.



Two contrasting patterns emerge from these tests (Fig. 2c–f and Supplementary Table 3). First, as in the case of otophysans, the metazoan data set provides strong differential support in favour of a single topology. While the traditional view has been that sponges are the first branching lineage in the animal tree, most recent phylogenomic studies support the so-called Ctenophora-sister hypothesis that places comb jellies as the sister group to all other animals (refs. ^{34–36} but see ref. ³³). In agreement with recent genomic evidence, the vast majority of GGI trees favour the latter hypothesis (Fig. 2c). Second, in none of the other three groups (yeasts, mammals and Neoaves) does GGI select a particular tree topology over another with overwhelming support (Fig. 2d–f), indicating that genealogical conflict in these groups is substantial.

Figure 2 | Gene genealogy interrogation (GGI) app of genes (x axes) supporting each hypothesis with unbiased (AU) topology test. Values above the dashed line are rank 1, whereas those below the dashed line are also rank 1. **d**, Neoaves (mousebird). **e**, Eutherian mammals (tree shrew). **f**, Yeast phylogeny.

