

Combining and Weighting Characters and the Prior Agreement Approach Revisited



John J. Wiens; Paul T. Chippindale

Systematic Biology, Vol. 43, No. 4 (Dec., 1994), 564-566.

Stable URL:

<http://links.jstor.org/sici?sici=1063-5157%28199412%2943%3A4%3C564%3ACAWCAT%3E2.0.CO%3B2-%23>

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

Systematic Biology is published by Society of Systematic Biologists. Please contact the publisher for further permissions regarding the use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/ssbiol.html>.

Systematic Biology

©1994 Society of Systematic Biologists

JSTOR and the JSTOR logo are trademarks of JSTOR, and are Registered in the U.S. Patent and Trademark Office. For more information on JSTOR contact jstor-info@umich.edu.

©2003 JSTOR

<http://www.jstor.org/>
Thu Sep 18 10:43:10 2003

- RZHETSKY, A., AND M. NEI. 1993. Theoretical foundation of the minimum-evolution method of phylogenetic inference. *Mol. Biol. Evol.* 10:1073-1095.
- SAITOU, N. 1990. Maximum likelihood methods. *Methods Enzymol.* 183:584-598.
- STEEL, M. A., M. D. HENDY, L. A. SZÉKELY, AND P. L. ERDŐS. 1992. Spectral analysis and a closest tree method for genetic sequences. *Appl. Math. Lett.* 5:63-67.
- SWOFFORD, D. L., AND G. J. OLSEN. 1990. Phylogeny reconstruction. Pages 411-501 in *Molecular systematics* (D. M. Hillis and C. Moritz, eds.). Sinauer, Sunderland, Massachusetts.

Received 29 December 1993; accepted 25 May 1994

Syst. Biol. 43(4):564-566, 1994

Combining and Weighting Characters and the Prior Agreement Approach Revisited

JOHN J. WIENS^{1,2,3} AND PAUL T. CHIPPINDALE¹

¹*Department of Zoology, University of Texas,
Austin, Texas 78712-1064, USA*

²*Texas Memorial Museum, University of Texas,
Austin, Texas 78705, USA*

Recently, Bull et al. (1993) and de Queiroz (1993) argued for an approach to integration of diverse data sets in phylogenetic inference that we (Chippindale and Wiens, 1994) termed the prior agreement approach. This approach involves separate analyses of data sets and does not allow combination of data in certain cases in which trees from separate analyses differ substantially from each other. We argued for an approach to phylogenetic analysis of diverse data sets that involves character weighting in the context of combined analyses and discussed several potential problems of the prior agreement approach. Huelsenbeck et al. (1994) commented on our paper, and herein we attempt to clarify several points with which they took issue.

Huelsenbeck et al. (1994) stated that they did "not understand the claim that conditions under which differential weighting fails to solve the problem would cause

our approach to fail." Our point was that if a subset of the total data supports the true phylogeny, and if those characters can be identified and given the appropriate weight(s), then the true phylogeny can be recovered. If none of the data support the true phylogeny, then we fail to see how our approach or their approach could possibly find the correct tree.

Huelsenbeck et al. (1994) characterized our view of the effects of lateral gene transfer as "optimistic." We suggested that lateral gene transfer could be accommodated by downweighting characters from the transferred gene or by the addition of characters not affected by the transfer. Huelsenbeck et al. (1994) claimed that "at any given moment, weak but true phylogenetic signal from one or more data sets could be swamped by a large set of misinformative characters resulting from a lateral transfer." The main point of our paper was that sets of misleading characters such as these potentially could be downweighted so as not to "swamp" the characters that indicate the correct phylogeny. We also

³ E-mail: jwiens@mail.utexas.edu.

pointed out that lateral transfer is likely only to involve a single gene or linked genic array, and thus addition of other data should lead to recovery of the true organismal phylogeny. In this case, Huelsenbeck et al. (1994) appear to share our optimism; they suggested that the tree agreed upon by data from several genes will be the correct organismal phylogeny.

Huelsenbeck et al. (1994) stated that our "position assumes that one can always identify the more reliable (or unreliable) characters" and that (under circumstances of extreme branch length inequality) "effectively 100% of the 'phylogenetically informative' characters will support an incorrect tree." Huelsenbeck et al. (1994) did not describe how their approach could lead to the correct phylogeny if separately analyzed data sets yield alternative trees and there is no way of knowing which set of characters is more reliable. In such cases, Bull et al. (1993) preferred to entertain the alternative trees as possible phylogenies rather than risk choosing an incorrect tree from a combined analysis. We argue that the combined tree(s) should at least be recovered and (if different from the trees from the separate analyses) entertained as another alternative hypothesis. Huelsenbeck et al. (1994) did not explain how their approach would detect heterogeneity and/or find the correct phylogeny if all (or most of) the characters in all the data sets agreed on the wrong tree. They also stated that (in the case of branch length inequalities) a posteriori weighting methods are "ineffective and may even increase the likelihood of choosing the wrong tree." Although we do not know whether this claim is true or not (Huelsenbeck et al. [1994] cited no evidence to support it), successive approximations is only one weighting scheme among the many that we suggested.

Huelsenbeck et al. (1994) strongly criticized our two hypothetical examples. Our first example (Chippindale and Wiens, 1994: fig. 2) was meant to illustrate the problem of incongruence between trees from separate and combined analyses (i.e., novel alternative hypotheses may be

missed completely if the data are not combined), using an example similar to the one presented by Barrett et al. (1991). We simply increased the number of characters from Barrett et al.'s (1991) example to show that nonparametric bootstrapping (Felsenstein, 1985) was not necessarily a solution for avoiding such problems. Huelsenbeck et al. (1994) criticized us for declaring the combined tree to be correct in this hypothetical example (perhaps reasonably), but whether the combined tree actually is the true phylogeny in this case (it certainly may be) is not critical to our argument. Furthermore, we argued that such cases of incongruence between the trees from separately analyzed and combined data are common in real data sets (Chippindale and Wiens, 1994: table 1); we believe that this is more relevant than whether or not we can construct a simple evolutionary model to explain the phenomenon. We also pointed out that the simulations of Bull et al. (1993) showed that accuracy generally increases with increasing numbers of characters, leading us to believe that the combined tree (which uses the most characters) is most likely to be the correct one (barring inconsistency). The point of our second example was that the methods suggested by Bull et al. (1993) for detection of heterogeneity (e.g., nonparametric bootstrapping [Felsenstein, 1985]; T-PTP test [Faith, 1991]) could ignore extremely disparate amounts of phylogenetic evidence present in separately analyzed data sets. Huelsenbeck et al. (1994) made much of the fact that we cited no specific evolutionary model for generating our hypothetical data. Yet, an example similar to the one in our (1994) figure 3 presumably could be generated (for a large number of characters) using the same model as in the simulations in figure 4 of Bull et al. (1993), with a ratio of roughly 95% consistent characters to 5% inconsistent characters (vs. the 50:50 or 20:80 used in their study). We do not see what makes our choice of combinations of character types "whimsical" and theirs "realistic." Our hypothetical examples illustrated a few of the general problems of dealing only with subsets of the total data,

problems that de Queiroz (1993), Bull et al. (1993), and Huelsenbeck et al. (1994) have not addressed to our satisfaction.

We laud the attempt by Bull et al. (1993) to examine the consequences of combining data versus separate analyses on the accuracy of phylogeny estimation, and we look forward to seeing the results of their future research. However, these authors acknowledged that they have not yet (1) recommended a specific method to detect heterogeneity, (2) shown whether heterogeneity exists in nature, and (3) shown that its recognition (and the resulting prevention of data combination) will improve phylogeny reconstruction. We believe that data combination is a crucial part of the analysis of diverse data sets because (1) accuracy in phylogenetic inference generally increases with increasing numbers of characters, and the combined analysis will simultaneously incorporate the maximum number of characters; (2) most characters (considered across a given genome) should share a common phylogenetic history, so that misleading characters should (overall) be outnumbered by characters that reflect the true organismal relationships; and (3) combination of data can reveal interactions among sets of characters, whereas these interactions may be hidden by considering only subsets of the total data. We believe that systematists should not be discouraged from combining their data, and we suggest that a more useful direction for future studies of data integration might be

to emphasize consideration of the processes of character evolution (through character weighting) rather than agreement among partitions of the data.

ACKNOWLEDGMENTS

We thank David Hillis and David Swofford for helpful discussion of this manuscript.

REFERENCES

- BARRETT, M., M. J. DONOGHUE, AND E. SOBER. 1991. Against consensus. *Syst. Zool.* 40:486-493.
- BULL, J. J., J. P. HUELSENBECK, C. W. CUNNINGHAM, D. L. SWOFFORD, AND P. J. WADDELL. 1993. Partitioning and combining data in phylogenetic analysis. *Syst. Biol.* 42:384-397.
- CHIPPINDALE, P. T., AND J. J. WIENS. 1994. Weighting, partitioning, and combining characters in phylogenetic analysis. *Syst. Biol.* 43:278-287.
- DE QUEIROZ, A. 1993. For consensus (sometimes). *Syst. Biol.* 42:368-372.
- FAITH, D. P. 1991. Cladistic permutation tests for monophyly and nonmonophyly. *Syst. Zool.* 40:366-375.
- FELSENSTEIN, J. 1985. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution* 39:783-791.
- HUELSENBECK, J. P., D. L. SWOFFORD, C. W. CUNNINGHAM, J. J. BULL, AND P. J. WADDELL. 1994. Is character weighting a panacea for the problem of data heterogeneity in phylogenetic analysis? *Syst. Biol.* 43:288-291.

Received 15 February 1994; accepted 10 June 1994

Note added in proof.—In recent talks at the annual meetings of the Society of Systematic Biologists (Athens, Georgia, 16 June 1994), Swofford and Cunningham, Bull, and Huelsenbeck suggested that the Mickevich-Farris Index (1981, *Syst. Zool.*, 30:351-370) may be useful in tests of heterogeneity.