

Popperian Corroboration and Phylogenetics

KEVIN DE QUEIROZ*

Department of Vertebrate Zoology, National Museum of Natural History, Smithsonian Institution, Washington, DC 20560-0162, USA

*Correspondence to be sent to: Department of Vertebrate Zoology, National Museum of Natural History, Smithsonian Institution, Washington, DC 20560-0162, USA; E-mail: dequeirozk@si.edu.

Received 7 February 2014; reviews returned 17 August 2014; accepted 18 August 2014

Associate Editor: Benoit Dayrat

Starting with articles by [Bock \(1973\)](#) and [Wiley \(1975\)](#) in this journal, the field of systematic biology has a history, reviewed by [Helfenbein and DeSalle \(2005\)](#), of examining its methods in the context of the philosophy of science articulated by Karl R. Popper (e.g., [1959, 1962, 1983](#)). Two main categories of debates have emerged in this literature. In one, Popper's philosophy is assumed to be relevant, and it is used to promote some systematic methods and criticize others (e.g., [Siddall and Kluge 1997; Kluge 2001](#)), which has led to counter-arguments proposing that the criticized methods are equally compatible with Popper's philosophy (e.g., [de Queiroz and Poe 2001, 2003](#)). In a second category of debates, the relevance of Popper's philosophy to systematics has been questioned (e.g., [Rieppel 2003, 2005; Vogt 2008](#)) and defended (e.g., [Farris 2013, 2014](#)). These debates can provide insights of at least two different kinds. Systematic biologists can gain a better understanding of how their methods and practices relate to general ideas about the nature of science, while philosophers can assess how well Popper's ideas about the nature of science describe the methods and practices in a discipline other than the ones (primarily physics and astronomy) upon which those ideas were based.

As part of the continuing debates about Popperian philosophy and phylogenetics, [Farris \(2013\)](#) recently argued that [Felsenstein \(2004\)](#) was incorrect in suggesting that Popper's concept of degree of corroboration, C , relies on a Bayesian approach, that [Helfenbein and DeSalle \(2005\)](#) were misguided in adopting Felsenstein's suggestion as a justification for Bayesian approaches, and that [Rieppel \(2005\)](#) assigned an incorrect value to the term $p(e|hb)$ in the defining formula of C . Here, I wish to call attention to an article regarding Popper's philosophy and its relationship to phylogenetics ([de Queiroz 2004](#)) that is highly relevant to, but was not considered in, those disagreements. I will summarize some of the main conclusions of that article and then illustrate how the perspective developed in it and two earlier publications ([de Queiroz and Poe 2001, 2003](#)), which equates specific phylogenetic analyses with

specific components of C , clarifies issues discussed by [Felsenstein \(2004\)](#), [Helfenbein and DeSalle \(2005\)](#), [Rieppel \(2005\)](#), and [Farris \(2013\)](#).

The article in question ([de Queiroz 2004](#)) focused on the concept of test severity, which is assessed using the term $p(e|b)$ in the defining formula of C (e.g., [Popper 1983](#), p. 238), the full definition of which is $C(h|e|b) = [p(e|hb) - p(e|b)] / [p(e|hb) - p(e|h|b) + p(e|b)]$. According to Popper, the data or evidence (e) used to test a hypothesis (h) may have a certain probability given that hypothesis and the background knowledge (b), thereby offering support for the hypothesis; however, the test can only be considered severe, and therefore the support meaningful, if the data have a much lower probability given the background knowledge (b) alone—that is, in the absence of the hypothesis being tested. In other words, the hypothesis can only be considered well-corroborated if the probability of the evidence given that hypothesis in conjunction with the background knowledge, $p(e|hb)$, is substantially greater than the probability of the evidence given the background knowledge alone, $p(e|b)$. This idea is embodied in the numerator of the defining formula of C , which is the difference between these two quantities: $p(e|hb) - p(e|b)$.

In the context of phylogenetics, where e is a phylogenetic data set (e.g., a taxon \times character matrix) and h is a hypothesis of phylogenetic relationships (e.g., a tree), b , the background knowledge, consists of those assumptions inherent to the method (model) used to analyze the data ([de Queiroz 2004](#); see also [de Queiroz and Poe 2001, 2003](#)), including ordered versus unordered character states, rates (likelihood) or costs (parsimony) of change between states, and equal versus variable rates (likelihood) or weights (parsimony) among characters. One of the basic assumptions common to diverse phylogenetic methods is that the terminal taxa are related through common ancestry, which allows their relationships to be represented by a tree. In this context, the absence of h is the absence of the hypothesis of phylogenetic relationships under the assumption that those relationships are to be

represented by a tree, which can be represented by an unresolved (star) tree. Under this interpretation, the calculation of C , or at least its numerator, which is the most important part of the concept (the denominator is simply a “normalization factor” [Popper 1983, p. 240]), is straightforward in a likelihood framework, because both of its component terms are likelihoods: $p(e|hb)$ is likelihood of the (resolved) tree representing the hypothesis of phylogenetic relationships (normally the optimal and therefore inferred tree) and $p(e|b)$ is the likelihood of an unresolved tree for the same taxa. To the extent that parsimony methods can be expressed in a likelihood framework (e.g., Farris 1973; Felsenstein 1973; Goldman 1990; Tuffley and Steel 1997), those methods can also be used to calculate $p(e|hb)$ and $p(e|b)$, though different values would be expected for the same h and e given the differences in b (i.e., between commonly used likelihood models and those that correspond to parsimony methods).

The interpretation of Popper’s C just described thus unifies likelihood and parsimony approaches to phylogenetic inference under a common philosophical framework. It also provides a more satisfactory account of how permutation tail probability (PTP) tests (Faith and Cranston 1991) exemplify Popper’s ideas about corroboration. PTP tests have been proposed as methods for assessing test severity as described by Popper, but previous attempts to assign values to $p(e|hb)$ and $p(e|b)$ in that context (e.g., Faith 1992; Faith and Cranston 1992; Faith and Trueman 2001) have been unsatisfactory (e.g., Farris 1995; Farris et al. 2001; de Queiroz and Poe 2001; de Queiroz 2004). In the context of the interpretation of C described above, PTP tests are seen as analogous to comparisons of the likelihoods of resolved and unresolved trees (i.e., $p(e|hb)$ and $p(e|b)$): PTP tests involve the comparison of optimal tree scores with null distributions of tree scores derived from randomized data, and patterns of shared, derived character states that evolved on an unresolved tree are expected, all else being equal, to be randomly distributed among taxa (see de Queiroz [2004] for details).

Grant and Kluge (2008) disagreed with the above-described interpretation of C . According to them, “... maximum likelihood does not maximize explanatory power because it does not discern between critical evidence (severe tests) and mere data (... *contra* de Queiroz 2004).” That assertion is incorrect. Although it is true that simply identifying the maximum-likelihood (or maximum parsimony) tree for a given data set does not involve an assessment of test severity (de Queiroz and Poe 2001, 2003; de Queiroz 2004), as described above, test severity can be assessed under likelihood using a star tree (de Queiroz 2004). Such an assessment would most certainly distinguish between critical evidence (characters that have substantially higher probabilities on the optimal tree than on the star tree, such as derived character states shared by inferred sister taxa) and mere data (characters that have similar probabilities on the optimal tree and the star tree, such as autapomorphic and invariant characters). Indeed,

that is the very purpose of comparing the likelihoods of the two trees, as the test will be considered severe and thus the evidence (e) critical only if the likelihood of the optimal tree, $p(e|hb)$, is substantially higher than the likelihood of the star tree, $p(e|b)$, for the data set in question. In the remainder of this contribution, I will illustrate how the above-described interpretation of C (de Queiroz 2004) clarifies several additional issues discussed by Felsenstein (2004), Helfenbein and DeSalle (2005), Rieppel (2005), and Farris (2013).

Rieppel (2005; see also Faith [1999]) proposed that the value of $p(e|hb)$ in the defining formula of C is fixed at 1. Farris (2013) refuted that proposition by quoting a passage from Popper (1983) that describes a situation in which the value of $p(e|hb)$ is zero: when e falsifies h in the presence of b . Farris (2013) also pointed out that fixing $p(e|hb)$ at 1 would have the effect of doing away with falsifiers, maximum-likelihood estimation, and the connection between likelihood and C , which would (if it were correct) lend credence to Rieppel’s (2003) proposition that corroboration in contemporary systematics may not conform to Popper’s description of that concept. In agreement with Farris (2013; see also Farris et al. [2001]), the term $p(e|hb)$ is the likelihood of h , that is, the likelihood of the tree that corresponds to a particular phylogenetic hypothesis. As such, its value is the likelihood score of that tree (h) for the analyzed data set (e) under the assumptions of the analytical model adopted (b). For real data sets and models employed by phylogeneticists, the value of $p(e|hb)$ is never = 1, which can be confirmed by checking the tree scores of any empirical study that employs likelihood methods. Instead, those values are generally very small positive numbers (not at all close to 1), which is the reason that the tree scores are commonly presented as the negative logarithms of the probabilities.

Rieppel’s (2003) suggestion that modern systematics may not conform to Popper’s concept of corroboration rests heavily on the logical interpretation of probabilities. According to that interpretation, the contradiction or support for a hypothesis provided by empirical evidence is absolute: $p(e|hb)$ can only be either 0 or 1; it cannot be a fractional value. Under that interpretation, phylogenetic methods are indeed difficult to reconcile with C , because phylogenetic data (e.g., distributions of character states among taxa) are never either logically prohibited, $p(e|hb)=0$, or logically entailed, $p(e|hb)=1$, by a particular tree (Sober 1988; de Queiroz and Poe 2003; Vogt 2008). Popper, however, did not restrict his concept of C to the logical interpretation of probabilities. Indeed, he explicitly developed C in the context of a “formal probability calculus” that could be interpreted in various senses, including the frequency interpretation (Popper 1959, p. 320). Moreover, Popper (1959, pp. 189–205) presented a detailed analysis to support his conclusion that although probabilistic hypotheses are not falsifiable in a strict, logical sense, scientists are quite capable of deciding when such hypotheses are “empirically confirmed” or “practically falsified” (p. 191). As Popper (1959) summarized his conclusions, “Thus we can expect

to be able to refute a probabilistic hypothesis, in the sense here indicated; and we can expect this perhaps even more confidently than in the case of a non-probabilistic hypothesis" (p. 202) and "Thus our analysis shows that statistical methods are essentially hypothetical-deductive, and that they proceed by the elimination of inadequate hypotheses—as do all other methods of science" (p. 413). A major flaw in Rieppel's (2003) analysis of the relevance of Popper's philosophy to contemporary systematics is that Rieppel misinterpreted degree of corroboration as depending on the logical interpretation of probabilities (p. 268) and overlooked Popper's writings on the corroboration and refutation of probabilistic hypotheses.

Felsenstein (2004, p. 142) proposed that the calculation of $p(e|b)$ in Popper's C requires summing over all possible trees, weighting each by its prior probability. Because of this supposed reliance on prior probabilities, he concluded that C assumes a Bayesian inference framework, which is inconsistent with Popper's otherwise anti-Bayesian views (e.g., his position that assigning probabilities to hypotheses is a "mistaken solution to the problem of induction" [Popper 1983, p. 217] and his objections to the subjective interpretation of probabilities). Farris (2013) quoted several passages from Popper as evidence against Felsenstein's interpretation of $p(e|b)$ as Bayesian. He also noted the correspondence between $p(e|hb)$ and the likelihood of a tree, as well as the fact that the tree with the strongest corroboration is the maximum-likelihood tree; however, he left unanswered the question of how $p(e|b)$ is to be calculated. The interpretation of $p(e|b)$ as the likelihood of a star tree (de Queiroz 2004) answers this question directly and supports Farris' conclusion that $p(e|b)$ is not a Bayesian term. The value of $p(e|b)$ is simply the likelihood of an unresolved tree, the calculation of which requires only a data set (e) and a model of phylogenetic inference (b), including estimates of its various parameters. It involves neither the summing of probabilities over all possible trees nor the assignment of prior probabilities to trees.

Helfenbein and DeSalle (2005) suggested that Popper's foray into probability theory sent him in a direction that can be considered Bayesian/likelihoodist, inferring a Bayesian component from the presence of $p(e|hb)$ in both Popper's C and Bayes' Theorem. Farris (2013) refuted that suggestion, pointing out that Helfenbein and DeSalle confused likelihood and Bayesian approaches. Here I wish to clarify a different issue raised by Helfenbein and DeSalle (2005; see also Vogt [2008]). Although those authors did not dispute de Queiroz and Poe's (2001, 2003) proposition that Popper's concept of degree of corroboration is based on Fisher's concept of likelihood, they (see also Vogt [2008]) quoted the following passage from Popper as if it represented a contradiction: "Thus we have proved that the identification of degree of corroboration or confirmation with probability (and even with likelihood) is absurd on both formal and intuitive grounds ..." (Popper 2002, p. 407). Contrary to the implication of Helfenbein and DeSalle, this statement

in no way contradicts the idea that C is based on likelihood (see also Farris [2014]).

The quoted statement was made in the context of a proposed "... mathematical refutation of all those theories of induction which identify the degree to which a statement is supported or confirmed or corroborated by empirical tests with its degree of probability in the sense of the calculus [axioms] of probability" (Popper 2002, p. 405). Popper's refutation of those theories is highly consistent with Fisher's motivation for proposing the concept of likelihood—that is, to avoid the idea that the support for a hypothesis is to be construed as the probability of that hypothesis (see de Queiroz and Poe 2001, p. 309). The likelihood of a hypothesis is not the probability of the hypothesis given the evidence, $p(h|e)$, but the probability of the evidence given the hypothesis (and the background knowledge), $p(e|hb)$. In Fisher's (1925, pp. 10–11) own words, "... the mathematical concept of probability is, in most cases, inadequate to express our mental confidence or diffidence in making such inferences ... the mathematical quantity which appears to be appropriate for measuring our order of preference among different possible populations does not in fact obey the laws of probability. To distinguish it from probability, I have used the term 'Likelihood' to designate this quantity."

Despite the basis of C in likelihood, Popper emphasized that $p(e|hb)$ alone is often inadequate for assessing the degree of corroboration of a hypothesis, particularly in cases involving small samples (e.g., 1959, pp. 413–414). This is the reason why C is not simply the likelihood of the hypothesis, $p(e|hb)$, but instead is defined as the (normalized) difference between the likelihood of the hypothesis in conjunction with the background knowledge and the likelihood of the background knowledge alone: $p(e|hb) - p(e|b)$. Advocates of likelihood-based inference commonly use the (ln transformed) likelihood ratio, termed *support*, S , by Edwards (1972), to compare the relative support for alternative hypotheses by a given set of data. When the hypotheses being compared are nested, the likelihoods are $p(e|hb)$ and $p(e|b)$, the same two quantities in the numerator of Popper's C (de Queiroz 2004). Thus, although degree of corroboration is not to be identified with the probability of a hypothesis, $p(h|e)$, or even with its likelihood, $p(e|hb)$, as noted by Popper in the passage quoted by Helfenbein and DeSalle, C is clearly based on probabilities in general, and on likelihoods in particular. Moreover, degree of corroboration, $p(e|hb) - p(e|b)$, is closely analogous to the likelihood ratio of nested hypotheses, $p(e|hb)/p(e|b)$.

Helfenbein and DeSalle (2005) discussed Popper's writings about probability in general, and those concerning degree of corroboration in particular, as if they were strangely incongruent elements that Popper was somehow trapped into introducing: "Given Popper's outlook as detailed in the main text of *The Logic of Scientific Discovery* written in the 1930s, before the appendixes from which the two quotes above were taken, it is odd that he would rely so heavily on any probabilistic

approach to the logic of science. Perhaps Popper found himself on a slope: needing a complement to falsification, he wrought corroboration; needing an explanation of corroboration, one distinct from ‘truth confirmation,’ he slid down into probabilistic reasoning, ending up with the corroboration formula discussed above” (p. 277).

Contrary to Helfenbein and DeSalle’s unsubstantiated speculations, Popper’s writings about probability and degree of corroboration represent well-reasoned extensions of his fundamental ideas about the logic of science. Even in the original *Logik der Forschung* (1934), Popper devoted an entire chapter (8) and substantial parts of two others (9 and 10) as well as seven appendices (i–vii) to probability (*The Logic of Scientific Discovery* [Popper 1959] is a translation of *Logik der Forschung* [Popper 1934]; additions to the original text [various footnotes and xii new appendices] are flagged with asterisks). Moreover, in that same edition Popper also noted the successes achieved by physics using probabilistic hypotheses (Chapter 8, section 65), and he discussed the concepts of corroboration and degree of corroboration in qualitative (rather than mathematical) terms (Chapter 10, especially section 82). Thus, it is not at all surprising that Popper would develop those ideas further in subsequent publications.

Nor should it be thought that Popper’s formal definition of degree of corroboration was incidental based on its appearance in an appendix. The material in that appendix (*ix of the *Logic of Scientific Discovery*) was originally published in three articles in *The British Journal for the Philosophy of Science* (Popper 1954, 1957, 1958). Moreover, Popper discussed his formal definition of degree of corroboration in the main texts of *Conjectures and Refutations* (Popper 1962, Chapter 11) and *Realism and the Aim of Science* (Popper 1983, Chapter 4). Given the widespread use of probabilistic approaches in science, it is not at all surprising that Popper would develop a philosophy of science that encompasses such approaches. In Popper’s own words, “... I now think that it is possible to define ‘degree of corroboration’ in such a way that we can compare degrees of corroboration [of alternative hypotheses]. Moreover, this definition makes it even possible to attribute numerical degrees of corroboration to statistical hypotheses” (Popper 1959, p. 268).

The arguments presented here and previously (de Queiroz and Poe 2001, 2003; de Queiroz 2004) bear on both of the main debates concerning Karl Popper’s philosophy of science in the discipline of systematic biology. With regard to disagreements about the relevance of Popper’s philosophy to contemporary phylogenetics, they demonstrate that Popper’s concept of degree of corroboration most certainly is relevant in that it is exemplified by methods developed and used by phylogeneticists. As argued here and previously (de Queiroz 2004), PTP tests (Faith and Cranston 1991) and tree length distribution skewness tests (Hillis 1991; Hillis and Huelsenbeck 1992) embody the spirit of Popper’s concept of degree of corroboration by assessing test

severity, and analogous likelihood-based resolution tests (e.g., Ota et al. 1999, 2000) involve direct comparisons of $p(e|hb)$ and $p(e|b)$ —the two fundamental components of C. This conclusion bears in turn on the debate concerning the consistency of parsimony versus likelihood methods with Popper’s philosophy. It reveals that any optimality-criterion-based phylogenetic method can be reconciled with Popper’s concept of degree of corroboration by employing methods that evaluate test severity, such as PTP, skewness, and resolution tests. Consequently, there is little point in attempting to use Popper’s philosophy to argue for the superiority of one optimality criterion over another. In any case, the evaluation of test severity based on $p(e|b)$ is a critical component of the logic of science described by Popper, and the realization that in phylogenetics $p(e|b)$ corresponds to the likelihood of an unresolved tree (de Queiroz 2004) is important for understanding how phylogenetics relates to Popperian philosophy.

ACKNOWLEDGMENTS

F. Anderson, B. Dayrat, and five anonymous reviewers provided comments on earlier versions of this article.

REFERENCES

- Bock W.J. 1973. Philosophical foundations of classical evolutionary classification. *Syst. Zool.* 22:375–392.
- de Queiroz K. 2004. The measurement of test severity, significance tests for resolution, and a unified philosophy of phylogenetic inference. *Zool. Scr.* 33:463–473.
- de Queiroz K., Poe S. 2001. Philosophy and phylogenetic inference: a comparison of likelihood and parsimony methods in the context of Karl Popper’s writings on corroboration. *Syst. Biol.* 50:305–321.
- de Queiroz K., Poe S. 2003. Failed refutations: further comments on parsimony and likelihood methods and their relationship to Popper’s degree of corroboration. *Syst. Biol.* 52:352–367.
- Edwards A.W.F. 1972. *Likelihood*. Cambridge (UK): Cambridge University Press.
- Faith D.P. 1992. On corroboration: a reply to Carpenter. *Cladistics* 8:265–273.
- Faith D.P. 1999. Review of *Error and the growth of experimental knowledge*. *Syst. Biol.* 48:675–679.
- Faith D.P., Cranston P.S. 1991. Could a cladogram this short have arisen by chance alone?: On permutation tests for cladistic structure. *Cladistics* 7:1–28.
- Faith D.P., Cranston P.S. 1992. Probability, parsimony, and Popper. *Syst. Biol.* 41:252–257.
- Faith D.P., Trueman J.W.H. 2001. Towards an inclusive philosophy of phylogenetic inference. *Syst. Biol.* 50:331–350.
- Farris J.S. 1973. A probability model for inferring evolutionary trees. *Syst. Zool.* 22:250–256.
- Farris J.S. 1995. Conjectures and refutations. *Cladistics* 11:105–118.
- Farris J.S. 2013. Popper: not Bayes or Rieppel. *Cladistics* 29:230–232.
- Farris J.S. 2014. Popper with probability. *Cladistics* 30:5–7.
- Farris J.S., Kluge A.G., Carpenter J.M. 2001. Popper and likelihood versus “Popper*^o”. *Syst. Biol.* 50:438–444.
- Felsenstein J. 1973. Maximum likelihood and minimum-steps methods for estimating evolutionary trees from data on discrete characters. *Syst. Zool.* 22:240–249.
- Felsenstein J. 2004. *Inferring phylogenies*. Sunderland (MA): Sinauer.
- Fisher R.A. 1925. *Statistical methods for research workers*. Edinburgh (Scotland): Oliver and Boyd.

- Goldman N. 1990. Maximum likelihood inference of phylogenetic trees, with special reference to a Poisson process model of DNA substitution and to parsimony analysis. *Syst. Zool.* 39: 345–361.
- Grant T., Kluge A.G. 2008. Clade support measures and their adequacy. *Cladistics* 24:1051–1064.
- Helfenbein K.G., DeSalle R. 2005. Falsifications and corroborations: Karl Popper's influence on systematics. *Mol. Phylogenet. Evol.* 35:271–280.
- Hillis D.M. 1991. Discriminating between phylogenetic signal and random noise in DNA sequences. In: Miyamoto M., Cracraft J., editors. *Phylogenetic analysis of DNA sequences*. New York: Oxford University Press. p. 278–294.
- Hillis D.M., Huelsenbeck J.P. 1992. Signal, noise, and reliability in molecular phylogenetic analyses. *J. Hered.* 83:189–195.
- Kluge A.G. 2001. Philosophical conjectures and their refutation. *Syst. Biol.* 50:322–330.
- Ota R., Waddell P.J., Hasegawa M., Shimodaira H., Kishino H. 2000. Appropriate likelihood ratio tests and marginal distributions for evolutionary tree models with constraints on parameters. *Mol. Biol. Evol.* 17:798–803.
- Ota R., Waddell P.J., Kishino H. 1999. Statistical distribution for testing the resolved tree against [the] star tree. *Proceedings of the Annual Joint Conference of the Japanese Biometrics and Applied Statistics Societies*. Tokyo: Sinfonica. p. 15–20.
- Popper K.R. 1934. *Logik der Forschung*. Vienna: Julius Springer.
- Popper K.R. 1954. Degree of confirmation. *Br. J. Philos. Sci.* 5:143–149.
- Popper K.R. 1957. A second note on degree of confirmation. *Br. J. Philos. Sci.* 7:350–353.
- Popper K.R. 1958. A third note on degree of corroboration or confirmation. *Br. J. Philos. Sci.* 8:294–302.
- Popper K.R. 1959. *The logic of scientific discovery*. New York: Basic Books.
- Popper K.R. 1962. *Conjectures and refutations*. New York: Basic Books.
- Popper K.R. 1983. *Realism and the aim of science*. London: Routledge.
- Popper K.R. 2002. *The logic of scientific discovery*. New York: Routledge Classics.
- Rieppel O. 2003. Popper and systematics. *Syst. Biol.* 52:259–271.
- Rieppel O. 2005. The philosophy of total evidence and its relevance for phylogenetic inference. *Pap. Avulsos Zool.* 45:77–89.
- Siddall M.E., Kluge A.G. 1997. Probabilism and phylogenetic inference. *Cladistics* 13:313–336.
- Sober E. 1988. *Reconstructing the past. Parsimony, evolution, and inference*. Cambridge (MA): MIT Press.
- Tuffley C., Steel M. 1997. Links between maximum likelihood and maximum parsimony under a simple model of site substitution. *Bull. Math. Biol.* 59:581–607.
- Vogt L. 2008. The unfalsifiability of cladograms and its consequences. *Cladistics* 24:62–73.
- Wiley E.O. 1975. Karl R. Popper, systematics, and classification: a reply to Walter Bock and other evolutionary taxonomists. *Syst. Zool.* 24:233–243.