



Partition-free congruence analysis: implications for sensitivity analysis

Ward C. Wheeler^{a,*}, Martín J. Ramírez^b, Lone Aagesen^a and Susanne Schulmeister^a

^aDivision of Invertebrate Zoology, American Museum of Natural History, Central Park West at 79th St., New York, NY 10024-5192, USA; ^bDivision Aracnologia, Museo Argentino de Ciencias Naturales, Avenue. Angel Gallardo 470, Buenos Aires C1405DJR, Argentina

Accepted 19 January 2006

Abstract

A criterion is proposed to compare systematic hypotheses based on multiple sources of information under a diverse set of interpretive assumptions (i.e., sensitivity analysis of Wheeler, 1995). This metric, the Meta-Retention Index (MRI), is the retention index (RI) of Farris calculated over the set of conventional homologous qualitative characters (ordered, unordered, Sankoff, etc.) and molecular fragment characters *sensu* Wheeler (1996, 1999). The superiority of this measure to other similar measures (e.g., incongruence length difference test) comes from its independence from partition information. The only values that participate in its calculation are the minimum, maximum and observed cost (= cladogram cost) of each character. The partition (morphology, gene locus) from which the variant may have come is irrelevant. In the special cases where there is only a single data partition, this measure is equivalent to the conventional RI; and in the case where there are single fragment characters per partition (contiguous molecular loci as data sets) the measure is identical to the complement of the Rescaled Incongruence Length Difference (RILD) of Wheeler and Hayashi (1998). The MRI can serve as an optimality criterion for deciding among systematic hypotheses based on the same data, but different sets of analysis assumptions (e.g., character weights, indel costs). The MRI may lose discriminatory power in situations where a minority of highly congruent characters is given high weight. This situation can be detected and seems unlikely to occur frequently in real data sets.

© The Willi Hennig Society 2006.

Congruence is the core of phylogenetic analysis. Congruence may come in the form of agreement among characters in a combined analysis, or among the results of several independent analyses derived from varying sources of information as some vague or precise notion of topological consensus. Hierarchical pattern is thought to be reinforcing while non-hierarchical (i.e., homoplasy) is not. While it is clear for many that the most appropriate measure of congruence given a particular set of assumptions is parsimony, it is less clear how we are to compare results of different assumption sets. This point was raised by Wheeler (1995) in proposing phylogenetic sensitivity analysis. All data analyses posit weighting schemes (equal or otherwise)

and character transformation costs (ordered, unordered, or more elaborate schemes) and offer parsimonious results.

The problem, then, is how to distinguish among this abundance of parsimonious solutions. Wheeler (1995) proposed two paths. The first is to accept only those groups that are generally present over the space of assumption sets. This approach was deemed “robust” choice and allowed the investigator to cling to those groups that offer consensual support over some large fraction of parameter space. This is a highly operational approach in that the resulting scheme does not attempt to optimize anything in particular. The second approach is based on explicit calculations of character [incongruence length difference test (ILD) of Mickevich and Farris, 1981; Farris et al., 1995] or topological congruence (Present Shared Groups/Present Common Groups

*Corresponding author.

E-mail address: wheeler@amnh.org

(PSG/PCG) of Wheeler, 1995). As topological methods (in general) make no allowance for relative levels of support, Wheeler (1995) favored the character-based ILD as a meta-optimality criterion to compare the suite of parsimonious results generated by sensitivity analysis.

It was obvious almost immediately that some adjustments were needed in the ILD as it could favor highly homoplastic solutions, which was clearly undesirable. Goloboff (pers. comm.) pointed to the situation in which partitioned analyses (in a total evidence or simultaneous analysis framework) were highly homoplastic under a certain assumption set. When combined, there would be little extra homoplasy that could occur from the combination of partitions; hence the ILD measure would be low. This is not a failure of the ILD itself as it is doing exactly what it was designed for—to measure the fraction of homoplasy caused by the combination of data sets. The problem lies in its application (Swofford, 1991; Farris et al., 1995). A modification of the ILD that adjusted for the “potential” homoplasy levels (RILD, “Rescaled ILD”) was proposed by Wheeler and Hayashi (1998) to avoid this problem.

Two other problems remain, however. The first is epistemological and the second more operational. The calculation of the ILD (or RILD) is based explicitly on the independent and collective behavior of data partitions. This is acceptable when discussing partitioned analysis, but in the context of total or simultaneous analysis, it is somewhat antithetical to the enterprise. The fundamental tenet of such analyses is the evidential equality of characters, irrespective of origin. In order to avoid this problem, we propose a new metric based on Farris’ retention index (RI; Farris, 1989), the calculation of which is accomplished entirely within the simultaneous analysis framework, hence independent of partition information. The second problem concerns the asymptotic behavior of the proposed metric in situations where small cliques of consistent characters are assigned high weight, yielding trivial optima. The analytical and empirical discussions of the “limit” effect are presented below and should ameliorate these concerns.

The metric

The ILD is, in many ways, an outgrowth of the Consistency Index (CI) of Kluge and Farris (1969). The ILD is calculated as the sum of the minimum costs of the component partitions subtracted from the minimum cost of the combined data set, which is then divided by the combined data set cost:

$$\text{ILD} = (\text{cost Combined Data Set} - \Sigma \text{cost of Component Partitions}) / \text{cost Combined Data Set}$$

The ensemble consistency index is the sum of the minimum costs of each of the characters divided by the combined data set cost:

$$\text{CI} = \Sigma \text{Minimal cost of Characters} / \text{cost Combined Data Set}$$

Hence, the ILD is basically the complement of the CI at the data set level.

When Wheeler and Hayashi (1998) modified the ILD, it was to make it more like the ensemble RI (Farris, 1989) that adjusted for the maximum costs (worst case) of characters. The RI is defined as:

$$\text{RI} = (\Sigma \text{Maximum cost of Characters} - \text{cost Combined Data Set}) / (\Sigma \text{Maximum cost of Characters} - \Sigma \text{Minimum cost of Characters})$$

Accordingly, the RILD was defined as:

$$\text{RILD} = (\text{cost Combined Data Set} - \Sigma \text{cost of Component Partitions}) / (\Sigma \text{Maximum cost of Component Partitions} - \Sigma \text{cost of Component Partitions})$$

$$\text{I-RILD} = (\Sigma \text{Maximum cost of Component Partitions} - \text{cost Combined Data Set}) / (\Sigma \text{Maximum cost of Component Partitions} - \Sigma \text{cost of Component Partitions})$$

Both the calculations of the ILD and RILD require the comparison of the combined data cladogram cost to those of each of the component partitions. In order to have a congruence measure for a partition-free simultaneous analysis, this dependence must be removed.

Here, we propose the Meta-Retention Index (MRI) which is simply Farris’ RI (Farris, 1989) calculated over all characters, but with the sequence characters defined not as single nucleotide positions, but as contiguous, sequence fragment or locus-based entities as in Wheeler (1999; stretches of nucleotides delimited by primer positions, locus boundaries, or other criteria that do not have primary homologies). The use of fragment-based characters has been discussed earlier (Wheeler, 2001) and has the advantage of treating sequence data as the dependent strings of nucleotides they are. The term “fragment” will be used here to refer to these sequence characters of arbitrary extent.

$$\text{MRI} = (\Sigma \text{Maximum cost of Fragments} - \text{cost combined Fragments}) / (\Sigma \text{Maximum cost of Fragments} - \Sigma \text{Minimum cost of Fragments})$$

where Fragment = Character for fixed homology statements (i.e., one morphological character or a contiguous string of nucleotides).

The maximum character costs are the character costs on a bush. A point worth noting here is that the differential fragment lengths contribute to the MRI

Morphological Partition

```
test data for MRI
  abcd
TAX_A 0000
TAX_B 1000
TAX_C 1101
TAX_D 1110
TAX_E 1111
```

Molecular Partition 1 Fragment 1

```
TAX_A
  1 AA

TAX_B
  1 AG

TAX_C
  1 AAG

TAX_D
  1 A

TAX_E
  1 AGG
```

Molecular Partition 2

```
TAX_A
  1 CG

TAX_B
  1 CG

TAX_C
  1 CGC

TAX_D
  1 AC

TAX_E
  1 GCC
```

Molecular Partition 1 Fragment 2

```
TAX_A
  1 GG

TAX_B
  1 GG

TAX_C
  1 AG

TAX_D
  1 AA

TAX_E
  1 AA
```

Fig. 1. Example data for MRI calculations. The test data consist of four binary characters and two molecular partitions, the first of which contains two sequence fragments.

differently—that is, a long fragment is likely to have more impact on the MRI than a short fragment or single morphological character.

The character definitions and RI calculations of standard qualitative characters (i.e., non-additive, additive, etc.) are unchanged. If a particular combination of loci (e.g., 18S rDNA and 28S rDNA) were to be treated as a series of sequence fragments, then each of these fragments would be treated homogeneously without respect to its origin. For qualitative characters, the maximum and minimum costs can be calculated directly. The situation is somewhat more complex for the sequence fragment characters where searches may be required to estimate the minimum, and other calculations for maximum character costs due to the diversity of potential character states. In the limiting situations where an analysis consisted of a single data set, the MRI would be identical to the RI over the nucleotide characters. If there were only single (fragment) characters in each partition, then the MRI would converge on $(1 - \text{RILD})$.

A synthetic example

Consider a synthetic data set with five taxa (“Tax_A” through “Tax_E”), and three sources of data (partitions): a morphological data set with four binary characters and two molecular sequences data sets, the first comprised of two fragments (perhaps exons) and the second with only a single (Fig. 1). Analyzed under two parameter regimes, one with all transformations (morphological, indel and base substitutions) equal to 1 and a second where indels cost two, yields two trees (Fig. 2). The MRI of the cladogram derived with a completely homogeneous parameter set yielded an MRI of 0.833 and that with indels twice all other transformations, 0.714 (Table 1). The MRI therefore favors the scheme based on a homogeneous assumption set. In this case, the ILD and $(1 - \text{RILD})$ do not distinguish between these alternatives; this is because the ensemble behavior of the partitions hides the character disagreement within. All analyses here and below were performed using POY ver. 2.7 (Gladstein and Wheeler, 1997) using Direct Optimization (= Optimization Alignment; Wheeler, 1996) to analyze the molecular fragments.

Real examples

To illustrate the behavior of the MRI and its performance relative to other measures, two data sets were examined. The first is an arthropod data set (Wheeler et al., 1993) and the second from a study of chelicerates (Wheeler and Hayashi, 1998). The arthropod data set contains 26 taxa (including a single extinct lineage) and 100 morphological characters and sequence data from the 18S rDNA, 28S rDNA and Polyubiquitin.

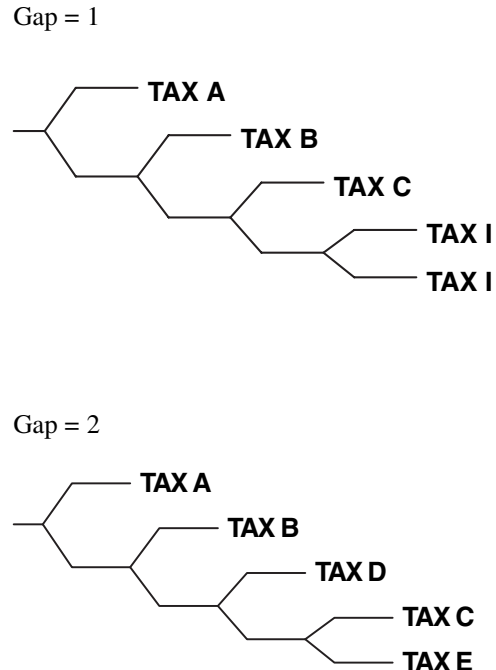


Fig. 2. Most parsimonious cladograms for the combined example data of Fig. 1. Gap = 1 is the resultant cladogram where all transformations are accorded equal weight (cost = 20), and Gap = 2 that which results when indels are up weighted to 2 (cost = 15). Cladogram drawn using CLADOS ver. 10.9 (Nixon, 1992).

The 18S rDNA was split into three fragments (based on primer locations), the 28S rDNA into eight (based on structural considerations), and the Polyubiquitin left intact. The total number of characters then for this analysis was 112. Similarly, the chelicerate data set was based on the analysis of 34 taxa, morphology (93 characters); and small and large subunit sequences each with six fragments for a total of 105 characters. In both cases, the cost of each character was determined by a search for the best (minimum cost) cladogram (as there is no way at present to directly calculate this value for fragment data). The maximum cost of each fragment character was determined by calculating the cost of a bush with each fragment as the root state (sum of costs from the root state to each of the terminal states). The minimum cost of the bush over all fragments was taken as the maximum character cost (this somewhat awkward procedure will work for binary characters as well). Maximum and minimum costs were then compared with the ensemble cladogram costs to determine the MRI values. Partition searches were also performed for comparative purposes (Tables 2 and 3). Not surprisingly, the MRI and $(1 - \text{RILD})$ on one hand, and the CI and ILD on the other, behave in very similar ways. The cladograms resulting from the analyses where these metrics achieve optimal values (0.8805 at indels = 4 and $tv = ti = 1$ for arthropods, and 0.7796 at indels = 2 and

Table 1
Test data values

Indel cost	Morphology	Partition 1 fragment 1	Partition 1 fragment 2	Partition 2	Total cost	MRI
2	a b c d					
Min cost	1 1 1 1	6	2	6		
Max cost	1 2 2 2	7	4	7	20	0.714
Min cost	1 1 1 1	4	2	4		
Max cost	1 2 2 2	4	4	5	15	0.833

Indel cost	Morphology	Partition 1	Partition 2	Total cost	ILD	1 – RILD
2						
Min cost	5	9	6			
Max cost	7	11	7	20	0	1
1						
Min cost		5	6	4		
Max cost	7	7	5	15	0	1

Table 2
Congruence values for the arthropod data set

InDel & Morph	TV:TI	MRI	CI	ILD	1 – RILD
1	1	0.8619	0.9008	0.0482	0.9279
1	2	0.8536	0.8902	0.0586	0.9162
1	4	0.8414	0.8748	0.0671	0.9083
1	8	0.8257	0.8577	0.0770	0.8974
1	∞	0.8095	0.8359	0.0900	0.8857
2	1	0.8725	0.8897	0.0559	0.9311
2	2	0.8710	0.8796	0.0638	0.9272
2	4	0.8589	0.8604	0.0808	0.9132
2	8	0.8524	0.8481	0.0883	0.9085
2	∞	0.8373	0.8247	0.0984	0.9016
4	1	0.8805	0.8722	0.0651	0.9354
4	2	0.8721	0.8517	0.0760	0.9301
4	4	0.8654	0.8349	0.0874	0.9240
4	8	0.8607	0.8232	0.0948	0.9201
4	∞	0.8568	0.8113	0.1020	0.9171
8	1	0.8776	0.8402	0.0875	0.9290
8	2	0.8721	0.8219	0.1032	0.9217
8	4	0.8686	0.8087	0.1129	0.9180
8	8	0.8663	0.8006	0.1184	0.9161
8	∞	0.8643	0.7921	0.1245	0.9141
16	1	0.8746	0.8098	0.1165	0.9192
16	2	0.8761	0.8044	0.1242	0.9176
16	4	0.8696	0.7869	0.1309	0.9156
16	8	0.8685	0.7818	0.1344	0.9147
16	∞	0.8678	0.7771	0.1389	0.9133
32	1	0.8695	0.7820	0.1408	0.9116
64	1	0.8677	0.7652	0.1530	0.9096
100	1	0.8670	0.7584	0.1578	0.9089
1000	1	0.8654	0.7466	0.1663	0.9074
5000	1	0.8652	0.7455	0.1671	0.9072
10 000	1	0.8652	0.7454	0.1672	0.9072

Table 3
Congruence values for the chelicerate data

InDel & Morph	TV:TI	MRI	CI	ILD	1 – RILD
1	1	0.7513	0.8261	0.0769	0.8723
1	2	0.7579	0.8194	0.0859	0.8680
1	4	0.7642	0.8089	0.0956	0.8662
1	8	0.7541	0.7919	0.1070	0.8564
1	8	0.7449	0.7716	0.1138	0.8543
2	1	0.7796	0.8167	0.1006	0.8658
2	2	0.7773	0.7988	0.1107	0.8638
2	4	0.7755	0.7826	0.1239	0.8584
2	8	0.7689	0.7679	0.1370	0.8492
2	∞	0.7571	0.7494	0.1506	0.8383
4	1	0.7701	0.7666	0.1405	0.8477
4	2	0.7688	0.7464	0.1589	0.8414
4	4	0.7615	0.7267	0.1754	0.8326
4	8	0.7563	0.7150	0.1809	0.8302
4	∞	0.7509	0.7009	0.1876	0.8277
8	1	0.7606	0.7083	0.1889	0.8306
8	2	0.7521	0.6853	0.2109	0.8191
8	4	0.7456	0.6684	0.2273	0.8104
8	8	0.7444	0.6608	0.2319	0.8098
8	∞	0.7421	0.6513	0.2383	0.8081
16	1	0.7441	0.6511	0.2407	0.8082
16	2	0.7385	0.6340	0.2557	0.8017
16	4	0.7382	0.6260	0.2633	0.8002
16	8	0.7368	0.6203	0.2679	0.7987
16	∞	0.7354	0.6143	0.2722	0.7975
32	1	0.6625	0.5988	0.2876	0.7325
64	1	0.6537	0.5734	0.3092	0.7226
100	1	0.6511	0.5638	0.3156	0.7206
1000	1	0.6467	0.5475	0.3308	0.7145
5000	1	0.6459	0.5460	0.3323	0.7137
10 000	1	0.6459	0.5458	0.3324	0.7136

tv = ti = 1 for chelicerates) are shown in Figs 3 and 4. The ILD, CI and 1 – RILD (for chelicerates) seem to show trivial best and worst values at their limits, whereas the MRI and (1 – RILD for arthropods) shows intermediate optima (Tables 2 and 3).

Limiting problems

One potential criticism of MRI as an optimality criterion is that under certain circumstances it may have a trivial optimum. Consider the synthetic test

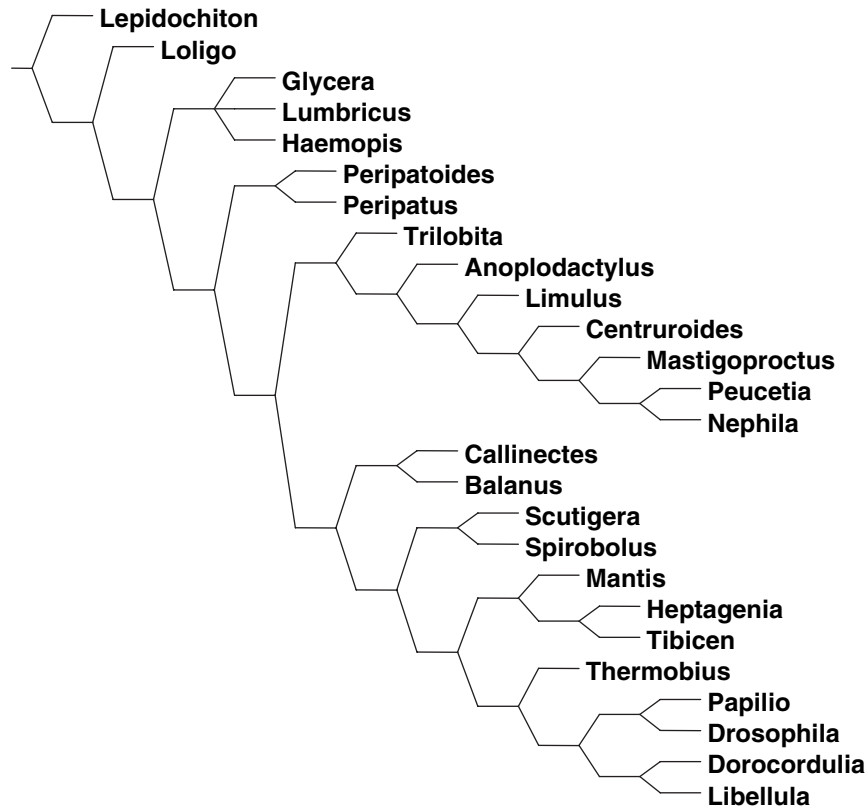


Fig. 3. Cladograms reconstructed from combined arthropod data, where indels are weighted 4, transversions 1 and transitions 1. Cladograms drawn using CLADOS ver. 10.9 (Nixon, 1992).

data used above and in Fig. 1 and Table 1. If the weight of indels (w) is increased without bound, the MRI will converge to 1 ($\lim (2w + 3)/(2w + 5)$ as $w \rightarrow \infty$ [min $3w + 12$; max $5w + 17$; combined $3w + 14$]). This will occur where there is a clique, in essence, of congruent transformations whose potential weight is unlimited (such as indels and transversions, but not transitions). Not all transformation types are subject to this, as the triangle inequality limits the range of cost differentials (Wheeler, 1993). Giribet and Wheeler (1999) touch on this as do Faith and Trueman (2001).

This situation can be demonstrated in a simple analytical cartoon. Consider a scenario with four binary characters, distributed as in Fig. 5. Where n_i = the number of characters of type i , there are $n_1 + n_2 + n_3 + n_4$ characters, and $n_1 > n_2$, $n_3 > n_4$, and $n_1 + n_4 > n_2 + n_3$. Furthermore, let there be two weights, a and b , attached to the characters where n_1 and n_2 receive weight a and $n_3 + n_4$ receive weight b . This situation describes two groups of characters each of which conflicts with each other and internally. The cost of the shortest tree would then be:

$$C_{\text{combined}} = a n_1 + b n_4 + 2a n_2 + 2b n_3$$

The sum of minimum and maximum possible character costs would be:

$$\Sigma C_{\text{minimum}} = a n_1 + b n_4 + a n_2 + b n_3$$

$$\Sigma C_{\text{maximum}} = 2a n_1 + 2b n_4 + 2a n_2 + 2b n_3$$

The MRI (or RI in this case) would be:

$$\text{MRI} = (\Sigma C_{\text{maximum}} - C_{\text{combined}}) / (\Sigma C_{\text{maximum}} - \Sigma C_{\text{minimum}})$$

Or

$$\text{MRI} = (a n_1 + b n_4) / (a n_1 + b n_4 + a n_2 + b n_3)$$

Hence

$$\lim \text{MRI} = n_1 / (n_1 + n_2) \text{ as } a \rightarrow \infty.$$

As $n_1 > n_2$, MRI must vary between 0.5 and 1.0. The only case in which the MRI will limit to 1 would be when n_2 is effectively zero (= no homoplasy). Any conflict among characters with weight a will decrease the limiting MRI.

Returning to the real data sets discussed above, it is clear that this effect does not cloud the sensitivity

Indel = 2, TV = 1, TI = 1

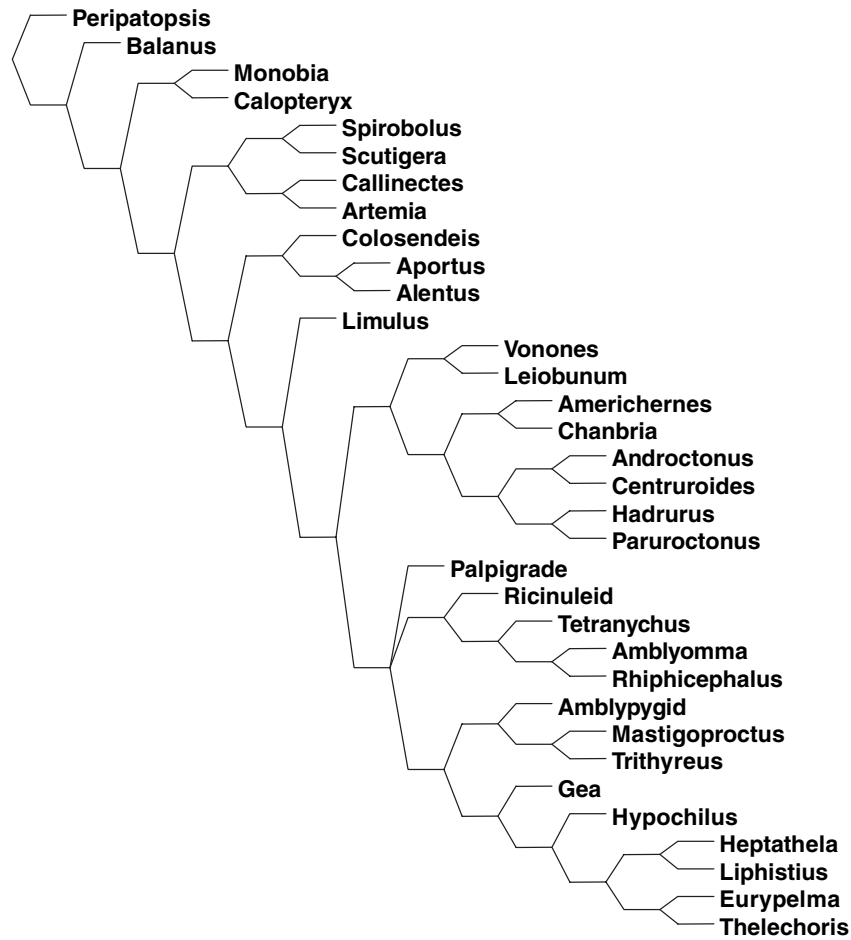


Fig. 4. Cladograms reconstructed from combined chelicerate data; indels are weighted 2, transversions 1 and transitions 1. Cladograms drawn using CLADOS ver. 10.9 (Nixon, 1992).

	1234
weight	aabb
One	0000
Two	0000
Three	0110
Four	1111
Five	1001

Fig. 5. Example data.

analysis of the arthropod or chelicerate data. As indel costs are increased from 1 to 10 000 times the cost of base transformations, the MRI does settle down to a limiting value, but this is not the global optimum for these data. In both cases, intermediate (neither minimum nor maximum) indel costs achieve optimal MRI values. Clearly, this is due to the presence of homoplasy in the optimization of cost variation. If there had been

complete agreement, the MRI would have tended towards unity (Fig. 6).

Discussion

There are several points about the MRI that should be made clear. First is that the delineation of fragment characters will affect the MRI measure. There are many unambiguous sorts of delineation that can be made, such as coding regions, exon/intron boundaries, primer locations or chromosomal location. Others may be less obvious (e.g., secondary structure). In cases where these delineations are not obvious, it would seem wise to examine this effect on the MRI values calculated. Secondly, the distinctions among parameter sets made by MRI values may be small (< 1%). Given that the MRI is an optimality criterion, this is not an issue *per se*, optimal is optimal after all, as with cladogram search

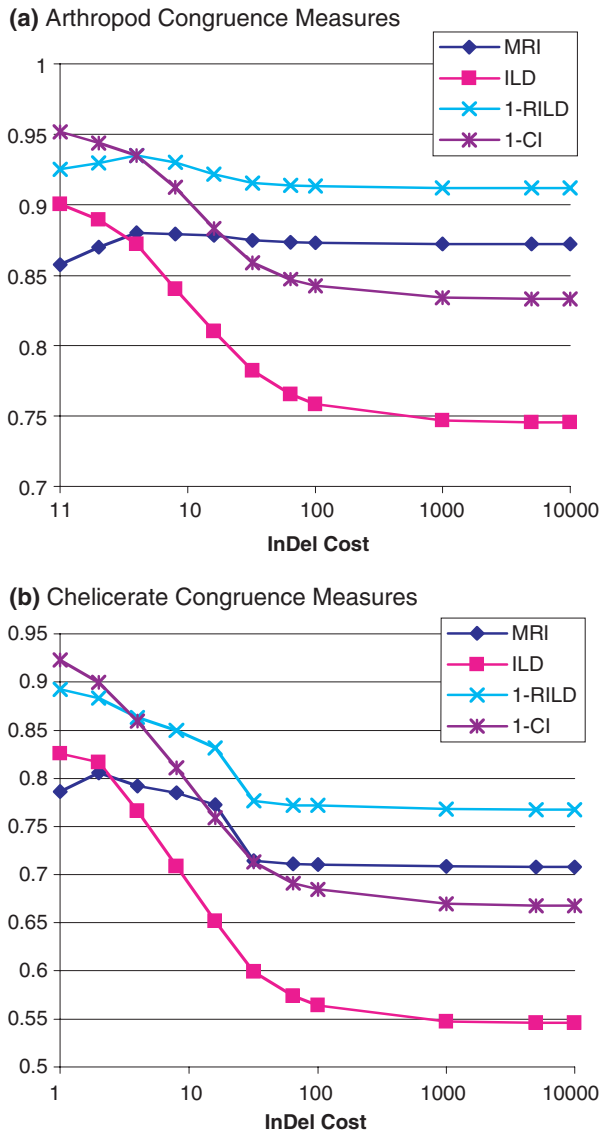


Fig. 6. MRI, ILD, 1-CI and 1 – RILD values as a function of increasing Indel cost for (a) arthropod data and (b) chelicerate data.

procedures. Some notion of robustness, however, may be in order. In situations where many widely divergent hypotheses are supported in different areas of parameter space, and MRI distinctions are small, the stability of the result to the addition of new data may be slight. On the other hand, there may be small, but consistent favoring of a hypothesis, which would yield stable solutions. As discussed in Wheeler (1995), these situations are worthy of examination. Robustness arguments can also come into play with the heuristic cladogram cost calculations of Direct Optimization. Given that these costs are upper bounds, small differences in MRI values may suffer additional variability based on the aggressiveness of searches and cost heuristics.

Conclusions

A partition-free metric can be constructed to distinguish among most parsimonious solutions of a given data set based on different parameter assumptions. The metric proposed here (MRI), may in some circumstances offer trivial optima, but these cases are likely to be infrequent and can be detected. Hence, the MRI has the potential to be a useful optimality criterion for phylogenetic sensitivity analysis.

Acknowledgments

Many of these ideas were spawned in discussion with John Wenzel. Thanks to Rob Asher, Taran Grant, Cyrille D’Haese, and Jan DeLaet for discussion of these ideas, a Postdoctoral Fellowship from CONICET, a Fessenden Research Fellowship from AMNH (MJR), and grant support from the National Aeronautics and Space Administration and National Science Foundation. Comments from an anonymous reviewer also improved the manuscript.

References

- Faith, D.P., Trueman, J., 2001. Towards an inclusive philosophy for phylogenetic systematics. *Syst. Biol.* 50, 331–350.
- Farris, J.S., 1989. The retention index and homoplasy excess. *Syst. Zool.* 38, 406–407.
- Farris, J.S., Källersjö, M., Kluge, A.G., Bult, C., 1995. Constructing a significance test for incongruence. *Syst. Biol.* 44, 570–572.
- Giribet, G., Wheeler, W.C., 1999. On gaps. *Mol. Phyl. Evol.* 13, 132–143.
- Gladstein, D.S. & Wheeler, W.C., 1997. POY: The Optimization of Alignment Characters. Program and Documentation. New York. Available at ftp.amnh.org/pub/molecular. Documentation by D. Janies and W.C. Wheeler.
- Kluge, A., Farris, J.S., 1969. Quantitative phyletics and the evolution of anurans. *Syst. Zool.* 18, 1–32.
- Mickevich, M.F., Farris, J.S., 1981. The implications of congruence in *Menidia*. *Syst. Zool.* 30, 351–370.
- Nixon, K.N., 1992. Clados, version 10.9. Program and Documentation. Trumansburg, NY.
- Swofford, D.L., 1991. When are phylogeny estimates from molecular and morphological data incongruent? In: Miyamoto, M.M., Cracraft, J. (Eds.), *Phylogenetic Analysis of DNA Sequences*, pp. 295–333. Oxford University Press, New York.
- Wheeler, W.C., 1993. The triangle inequality and character analysis. *Mol. Biol. Evol.* 10, 707–712.
- Wheeler, W.C., 1995. Sequence alignment, parameter sensitivity, and the phylogenetic analysis of molecular data. *Syst. Biol.* 44, 321–331.
- Wheeler, W.C., 1996. Optimization alignment: the end of multiple sequence alignment in phylogenetics? *Cladistics* 12, 1–9.
- Wheeler, W.C., 1999. Fixed character states and the optimization of molecular sequence data. *Cladistics* 15, 379–385.
- Wheeler, W.C., 2001. Homology and the optimization of DNA sequence data. *Cladistics* 17, S3–S11.
- Wheeler, W.C., Hayashi, C.Y., 1998. The phylogeny of the extant chelicerate orders. *Cladistics* 14, 173–192.
- Wheeler, W.C., Cartwright, P., Hayashi, C.Y., 1993. Arthropod phylogeny: a combined approach. *Cladistics* 9, 1–39.