# A multi-neighbor-joining approach for phylogenetic tree reconstruction and visualization

**Ana Estela A. da Silva[3], Wilfredo J.P. Villanueva[1], Helder Knidel[1], Vinícius Bonato[3], Sérgio F. dos Reis[2] and Fernando J. Von Zuben[1]**

[1]Faculdade de Engenharia Elétrica e de Computação (FEEC), UNICAMP,
Av. Albert Einstein, 400, Caixa Postal 6101, 13083-970 Campinas, SP, Brasil
[2]Instituto de Biologia, UNICAMP, Caixa Postal 6109,
13083-970 Campinas, SP, Brasil
[3]Faculdade de Ciências Matemáticas,
da Natureza e Tecnologia da Informação (FCMNTI), UNIMEP,
Piracicaba, Campus Taquaral, Rod. do Açucar, km 156,
13400-911 Piracicaba, SP, Brasil
Corresponding author: W.J.P. Villanueva
E-mail: wilfredo@dca.fee.unicamp.br

**ABSTRACT.** The computationally challenging problem of reconstructing the phylogeny of a set of contemporary data, such as DNA sequences or morphological attributes, was treated by an extended version of the neighbor-joining (NJ) algorithm. The original NJ algorithm provides a single-tree topology, after a cascade of greedy pairing decisions that tries to simultaneously optimize the minimum evolution and the least squares criteria. Given that some sub-trees are more stable than others, and that the minimum evolution tree may not be achieved by the original NJ algorithm, we propose a multi-neighbor-joining (MNJ) algorithm capable of performing multiple pairing decisions at each level of the tree reconstruction, keeping various partial solutions along the recursive execution of the NJ algorithm. The main advantages of the new recon-

struction procedure are: 1) as is the case for the original NJ algorithm, the MNJ algorithm is still a low-cost reconstruction method; 2) a further investigation of the alternative topologies may reveal stable and unstable sub-trees; 3) the chance of achieving the minimum evolution tree is greater; 4) tree topologies with very similar performances will be simultaneously presented at the output. When there are multiple unrooted tree topologies to be compared, a visualization tool is also proposed, using a radial layout to uniformly distribute the branches with the help of well-known metaheuristics used in computer science.

**Key words:** Fitch-Margoliash method, Multi-neighbor-joining algorithm, Unrooted tree visualization, Phylogenetic tree reconstruction

## INTRODUCTION

The inference of a phylogeny can be essentially characterized as an estimation task that is based solely on the information available after the occurrence of an unknown and gradual process of evolutionary differentiation and inheritance of ancestral characteristics. Most techniques are devoted to the proposition of the best estimate of the evolutionary history, in the form of an evolutionary tree. The branching order of the tree and the length of the edges provide a concise hypothetical explanation based on the data.

The process of tree reconstruction can be algorithmic or search-based. Parsimony methods (Swofford et al., 1996) and distance-based methods (Fitch and Margoliash, 1967) are examples of algorithmic procedures, and maximum likelihood methods (Felsenstein, 1988) are a well-known example of the search-based procedures.

Among the available techniques for phylogenetic inference, we have adopted and extended a distance-based method called neighbor-joining (NJ), proposed by Saitou and Nei (1987) and improved by Studier and Keppler (1988). The NJ algorithm is one of the most effective methods for reconstructing phylogenies from a matrix of pairwise evolutionary dissimilarities. The resulting tree topology is unrooted, because the branching order is time-reversible and consequently the location of the root has no influence on the quality of the solution.

The extension was denoted the multi-neighbor-joining (MNJ) algorithm, because the output of the algorithm may not be restricted to the proposition of a single final topology, given that a distinct pairing decision can also produce high-quality topologies. The occurrence of multiple pairing decisions at each step of the algorithm is governed by a threshold (a user-defined parameter).

After the conclusion of the alternative reconstructions, a visualization tool presents the resulting unrooted trees for comparison. A radial layout is adopted here, and the purpose is to maximize the distribution of the branches, independent of the configuration of the tree and the pattern of branch lengths. We then have to subsequently solve two very demanding optimization problems: reconstruct the most representative trees, given the data, and present the most informative view of the obtained trees. In fact, both problems are computationally intractable (Garey and Johnson, 1979; Foulds and Graham, 1982; Day, 1987); consequently, finding high-quality solutions, not necessarily the optimal one, is the main purpose.

## INFERRING PHYLOGENIES

Based on evolutionary attributes extracted from entities under systematic analysis, the proposal of a phylogenetic tree to best describe the sequence of evolutionary events that were responsible for the observed dissimilarity among entities is generally preceded by the selection of an appropriate method from a set of alternative computational procedures. These procedures can be classified as follows: 1) algorithmic: a specific algorithm is designed to determine the tree; 2) search-based: an optimization criterion is defined to compare the alternative tree topologies, and each tree is given by a point in a search space to be explored. The algorithmic procedures try to simultaneously perform the definition of the topology and the fulfillment of optimization criteria, in a greedy and step-by-step manner. The search-based procedures adopt a distinct paradigm, composed of two steps: one for the evaluation of topologies, taking into account evolutionary assumptions, and other for the determination of topologies with the highest evaluation among the candidates, making use of computational search strategies.

The main disadvantages of the algorithmic procedures are the possibility of getting stuck in poor local minima and the presentation of a single-tree topology. The main advantage is the reduced computational cost, with the number of hypothetical ancestors to be determined being linearly correlated with the number of entities under analysis. These methods incorporate all distance-based methods, including cluster analysis (e.g., UPGMA) and NJ. The opposite scenario can be depicted in the case of search-based procedures. The computational cost is the main disadvantage, due to the factorial increase in the number of candidate topologies with an increasing number of entities under analysis. On the other hand, the possibility of avoiding poor local minima along the search and the capability of presenting several high-quality tree topologies as the output, instead of a single one, are very desirable aspects.

Both approaches to infer a phylogeny are competitive with each other, with no clear evidence of superiority when general application fields are considered. That is why every research group devoted to systematics should consider both as relevant computational approaches to infer phylogenies. Efforts should be concentrated on alleviating the disadvantages of each one. Prado and Von Zuben (2002) have already proposed powerful metaheuristics for search-based procedures; we have sought to generalize one of the better-accepted algorithmic procedures, the NJ method.

## MULTIPLE TOPOLOGIES AND VISUALIZATION ISSUES

Taking the solution space as the set of all possible phylogenies spanning the input data, the greedy nature of the local decisions adopted by the NJ algorithm will explore only a small portion of the solution space; that is why they are so cost-effective. However, the absence of a more exploratory search strategy (with just one topology as the output) will prevent the proposal of alternative solutions with distinct topologies. The alternative topologies may fit the data nearly as well, equally well, or even better, thus alleviating the occurrence of misleading conclusions.

To improve the reliability of phylogenetic tree reconstruction, an unequivocal tendency in the study of systematics is the necessity of taking advantage of the low computational cost associated with algorithmic procedures for inferring phylogenies (Salemi and Vandamme, 2003) intimately related to the necessity of computational procedures capable of proposing multiple outputs as a result of the inference (Holmes, 2002; Holder and Lewis, 2003). So, the question is:

given that obtaining a single tree topology from the available data is fast and easy, is it possible to explore this reduced computational burden in order to produce multiple high-quality, and possibly informative, tree topologies? The most widely accepted software packages adopted around the world, like PAUP (Swofford, 2001) and Phylip (Felsenstein, 1989), do not provide such a resource, particularly when the NJ algorithm is considered.

An additional concern is the conception of an original software package devoted to the proper presentation of unrooted trees. With the new scenario involving several topologies, a computational tool with advanced graphic resources for the presentation of the resulting tree topologies for comparison is a fundamental step.

The already available software packages for presentation of unrooted trees are based on elementary rules of thumb and cannot provide a clear indication of the branching patterns in crucial circumstances (Carrizo, 2004).

## A BRIEF DESCRIPTION OF THE NEIGHBOR-JOINING ALGORITHM

The NJ algorithm is a method for reconstructing unrooted phylogenetic trees, using a matrix of evolutionary dissimilarities as input data. The dimension of this square matrix of pairwise distances corresponds to the number of leaves in the resulting unrooted tree topology, each one denoted a taxon or OTU (operational taxonomic unit). Initially, the $n$ taxa are neighbors, because the algorithm starts with a star tree. A sequence of agglomerative steps then follows, taking into account the minimum evolution principle (Kidd and Sgaramella-Zonta, 1971) to determine the pair of taxa to be joined, among all the $n*(n - 1)/2$ possibilities, and the Fitch-Margoliash approach (Fitch and Margoliash, 1967) to propose the branch lengths of the two new branches. At each agglomerative step, a new node is created (HTU - hypothetical taxonomic unit) to support the two additional branches, so that the star tree loses the newly joined OTUs and gains the new HTU in replacement. This iterative process is repeated until the remaining star tree has only three taxa.

In computational terms, the agglomerative process has a computational complexity of $O(n^3)$, where $n$ is the number of taxa, and it may be interpreted as a greedy strategy that tries to simultaneously satisfy the minimum evolution principle, associated with the sum of branch lengths, and the least squares criterion (Bulmer, 1991), associated with the difference between the original distance matrix and the distance matrix extracted from the obtained unrooted tree. As a consequence, taking locally best decisions towards optimizing both objectives in general cannot guarantee the construction of the global minimum evolution tree, but only a short tree whose topology may be similar (and sometimes identical) to the minimum evolution tree (Saitou and Imanishi, 1989).

Due to the greedy nature of the search, only one candidate is taken among all the candidate pairs at each step of the agglomerative process. As a consequence, the NJ algorithm is only capable of producing a single unrooted tree at the end of the execution. The first step of the agglomerative process consists of six OTUs in a star tree (Figure 1A). At this step, all the 15 pairs are considered as candidates, and the one that minimizes the total length of the resulting tree is then chosen to be joined. Let us assume that the pair 3-4 is the one selected. The new star tree for the next step is composed of the remaining 4 OTUs plus the new HTU representing the common ancestor of OTUs 3 and 4 (Figure 1B).

When the distance matrix obeys the additive property (Barthélemy and Guénoche, 1991),
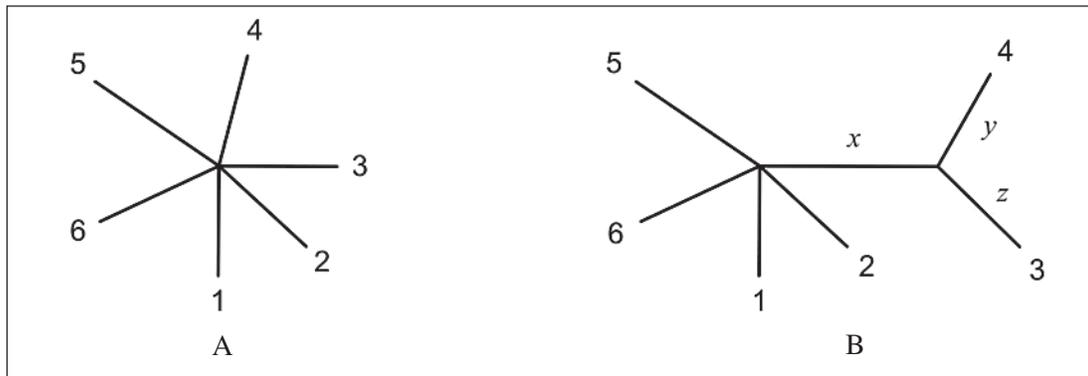
**Figure 1. A.** Unrooted tree with six operational taxonomic units (OTUs). **B.** Topology obtained with the junction of OTUs 3 and 4. The lengths of the three branches involved in the current step are determined by means of the Fitch-Margoliash method. They are represented by *x*, *y* and *z*.

the NJ algorithm satisfies a basic requirement for tree reconstruction: it necessarily finds the unique tree which perfectly represents the data (Gascuel, 1997). High-quality solutions produced by the NJ algorithm have been described in the literature (Saitou and Imanishi, 1989; Kuhner and Felsestein, 1994; Huelsenbeck, 1995). Also, theoretical results obtained by Atteson (1996) have indicated that the NJ algorithm is as efficient as possible.

## MULTIPLE TOPOLOGIES AND THE MULTI-NEIGHBOR-JOINING ALGORITHM

Given that the NJ algorithm will always produce a single-tree topology and that the obtained tree may not be the minimum evolution tree, though it is a short tree, three approaches can be considered to improve the capability of achieving the minimum evolution tree, or a proposal close to it.

The first one considers some strategies for local rearrangements aiming at searching for the minimum evolution tree in the neighborhood of the obtained NJ tree (Rzhetsky and Nei, 1993). A bootstrap procedure is the essence of the second approach, which operates at the level of the distance matrix (Rzhetsky and Nei, 1994). The third approach is the one that we have adopted, which tries to diversify the local decision at each agglomerative step.

Essentially, instead of always selecting a single pair to be joined, the one that supposes the minimum evolution at that specific state of the topology reconstruction, alternative pairs that do not promote an increase in the total length of the tree above a specified threshold will also be considered, giving rise to alternative topologies from them on. Consequently, as the NJ algorithm has to perform *n*-3 junctions, we have *n*-3 possibilities of generating alternative topologies as close as desired to the minimum evolution principle. The newly generated alternative topologies will have the power to generate their own alternative topologies at the subsequent pairing decisions. Distinct sequences of pairing may produce the same final topology, indicating that only a subset of the alternative topologies will be effectively distinct from each other.

The work by Pearson et al. (1999) can be considered the first attempt towards the proposition of multiple tree topologies under the recursive application of the NJ algorithm. How-

———————————————

ever, their approach is conceptually incompatible with ours, because the main purpose was to obtain the most distinct tree topologies, and not necessarily the ones most capable of attending the optimization criteria, no matter the similarity in topology. They did not neglect the optimization criteria, but they incorporated a significant role towards the maintenance of diversity among the obtained topologies.

Instead of having the main purpose of achieving diversity in topology, our aim has been solely to induce the detection of alternative high-quality topologies. Whether the multiple solutions, including the one associated with the original NJ algorithm, will or will not have similar topologies is a further inquiry directly related to the nature of the ancestral dependencies indicated by the input data.

## VISUALIZATION DEVICE FOR UNROOTED TREES

The visualization tool is able to recognize trees codified in the Newick format - the most popular format adopted to represent phylogenetic trees. The MNJ algorithm adopts this format to express each alternative topology.

A large number of aspects may be involved in the characterization of a good visualization, including subjective and objective ones. Here we consider the following objective aspects: number of crossing branches, angle between branches, and alignment involving descendent branches and ancestral branches. A fitness function is then derived, taking these aspects into consideration. So, the resulting tree will tend to have no crossing branches and a homogeneous distribution of branches along the 2-D plane.

After modeling the visualization task as an optimization problem with a fitness function, the attributes to be optimized are associated with the angles between consecutive branches. The search space is a Euclidean space with a multimodal fitness function, and each angle is a component of a vector (a point in the Euclidean space). In mathematical terms, the purpose is to find a point in the search space with a high value of the fitness function, not necessarily the maximum one.

A population-based search strategy called evolution strategies (Beyer, 2001) is then applied to solve the resulting optimization problem. There will be one optimization problem for each tree to be visualized, and the computational cost allows for application to topologies with hundreds or thousands of OTUs.

## MORPHOLOGICAL DATA AS A CASE STUDY

Description of patterns of variation in morphological or genetic characters within and among populations is fundamental to define the boundaries of independent evolutionary units in nature. An important initial step in recognizing such evolutionary units is the identification of groups of populations that share morphological traits and geographic continuity over geographic space (Carleton, 1988; Myers, 1989; Patton and Smith, 1990; Patton and da Silva, 1997).

In systematic biology, information that allows recognition of such units has classically been derived from analysis of variation in the shape of morphological structures. Here we analyze geographic variation in cranial shape in 22 populations of the rodent species *Thrichomys apereoides*, sampled from northeastern, southeastern and central Brazil (see Figure 2). Shape variation was described using partial warps, which are variables derived from the formalism of

geometric morphometrics, giving rise to a distance matrix to be considered as input data for the MNJ algorithm. The data used to construct the shape variables are three-dimensional coordinates of points, which are landmarks that are defined for the skull of *Thrichomys apereoides*.
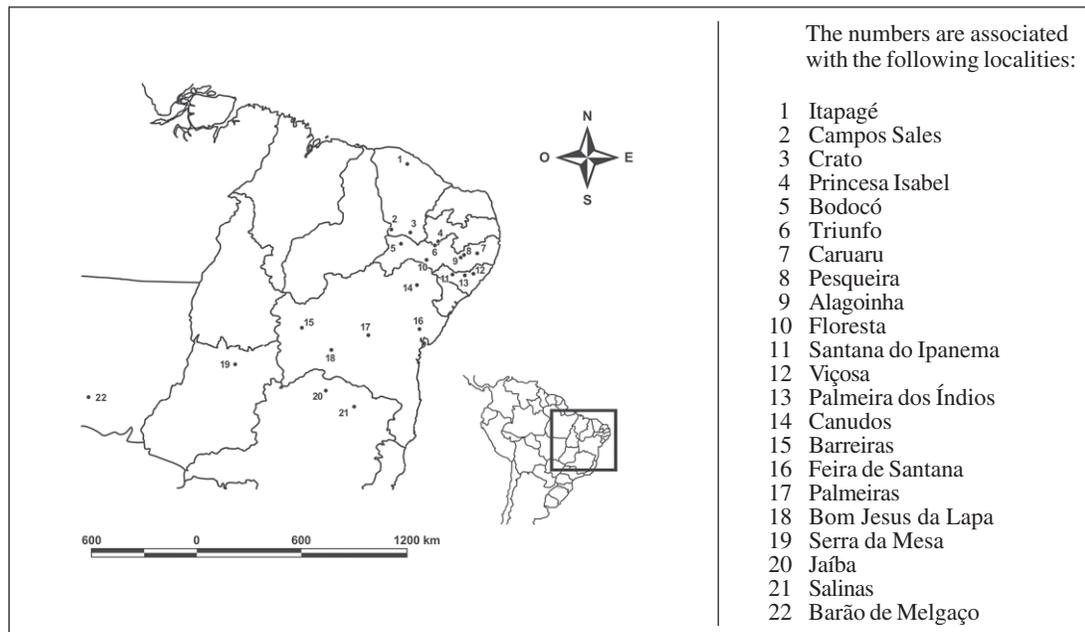


The numbers are associated with the following localities:

1   Itapagé
2   Campos Sales
3   Crato
4   Princesa Isabel
5   Bodocó
6   Triunfo
7   Caruaru
8   Pesqueira
9   Alagoinha
10  Floresta
11  Santana do Ipanema
12  Viçosa
13  Palmeira dos Índios
14  Canudos
15  Barreiras
16  Feira de Santana
17  Palmeiras
18  Bom Jesus da Lapa
19  Serra da Mesa
20  Jaíba
21  Salinas
22  Barão de Melgaço

**Figure 2.** Locations from which the samples were collected.

## RESULTS

Comparative results involving our visualization tool and alternative proposals available in the literature will not be presented here, though exhaustive scenarios have been considered to validate our proposal. Based on the fitness criterion adopted, our population-based and exploratory search algorithm is capable of maintaining a high level of performance, even for unbalanced tree topologies. The global search for a maximum dispersion of branches with arbitrary length is capable of producing a proper spatial configuration, even in the most challenging instances. Figure 3 presents an example of the output provided by the visualization tool for the unrooted tree associated with the single output of the original NJ algorithm, together with two sub-trees extracted from alternative topologies produced by the MNJ algorithm. The case study involves the morphological data set described in the "Morphological data as a case study" section, composed of 22 OTUs.

Starting with the pairwise distance matrix, the purpose here is to determine stable and unstable sub-trees when multiple alternatives are considered to explain the input data (Felsenstein, 2004). Within an empirically defined threshold of 0.02% of admissible increase in the total size of the tree (when compared with the best pairing decision at that moment), multiple pairing decisions were taken, leading to five distinct topologies as the final result. The single output of
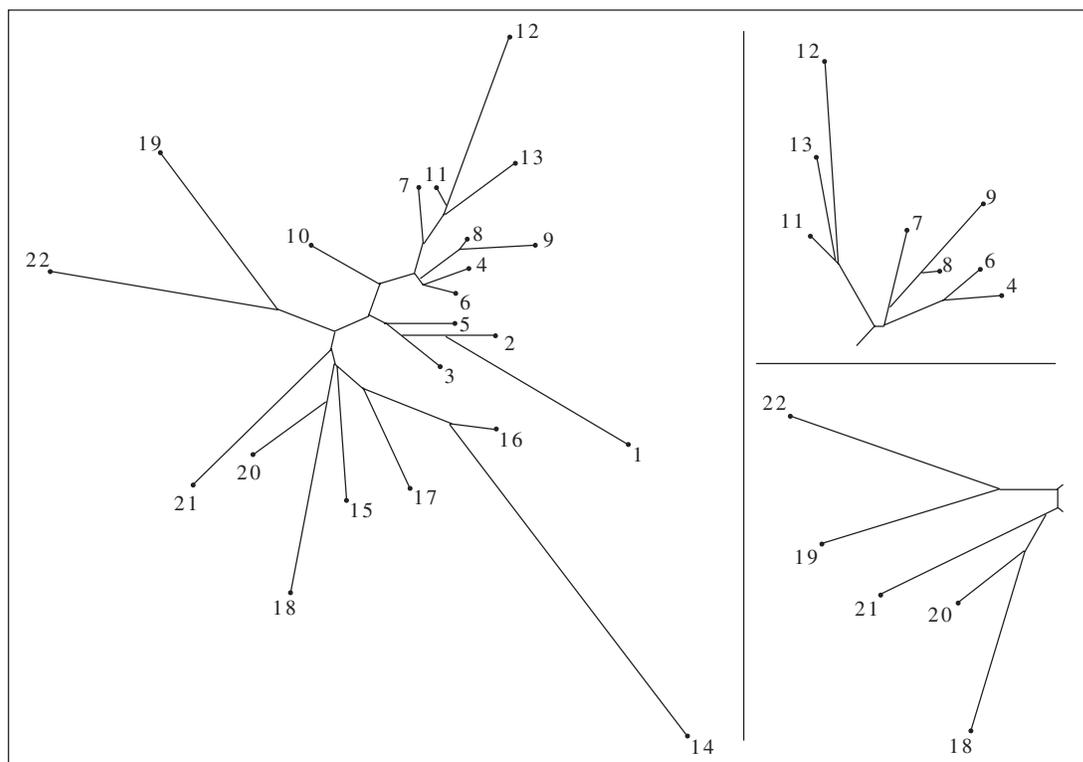
**Figure 3.** The tree produced by the original neighbor joining (NJ) algorithm (at the left), and two alternative proposals for unstable sub-trees (at the right), both extracted from topologies produced by the multi-neighbor joining (MNJ) algorithm.

the original NJ algorithm (see Figure 3, on the left) is guaranteed to be one of them, which means that the MNJ algorithm is formally a generalization of the NJ algorithm.

A careful analysis of the four alternative topologies (not presented here as a whole), also plotted by means of the proposed visualization tool, leads to the conclusion that the sub-trees 4-6-7-8-9-11-12-13 and 18-20-21 are unstable, in the sense that the alternative topologies propose distinct configurations for these sub-trees (see the two sub-trees presented in Figure 3, on the right). Basically, two most significant distinctions can be extracted from the visual comparison involving the original topology and the alternative ones:

- OTU 7 is clustered together with OTUs 8 and 9, with the other OTUs in the sub-tree reallocated in various configurations inside the same sub-tree;
- OTU 21 is first clustered with OTUs 18 and 20 before being clustered with the remaining OTUs.

Given that our intent here is to identify groups of populations that share morphological traits and geographic continuity over geographic space, the alternative topologies are very suggestive. Based on Figure 2, putting OTU 7 closer to OTUs 8 and 9, and OTU 21 closer to OTUs 18 and 20, is akin to the hypotheses of geographic continuity.

We are not arbitrarily looking for distinct explanations to better support a set of hypothesis, simply because the alternative topologies are still high-quality solutions in terms of the minimum evolution and the least square criteria. So, they represent a reliable and optional perspective extracted from the same input data.

## CONCLUSIONS

In phylogenetic reconstruction, making sense of a portion of the huge amount of data being generated is very desirable. Multiple high-quality views of the evolutionary relationships may give rise to a broader perspective for data analysis, opening the possibility of supporting new insights into biological systems.

Though preliminary, our results are promising in the sense that additional hypotheses may be raised, based on very similar or even strongly diverse tree topologies. With a powerful visualization tool attached to an MNJ algorithm for phylogenetic reconstruction, explicitly designed to propose a broader set of high-quality candidate solutions, experts in biological systematics will gain an additional tool to infer phylogenies.

As further research topics, the inclusion of similarity measures to enrich the MNJ decision process and the extension of the proposed methodology to other algorithmic procedures, like parsimony methods, are being investigated.

## ACKNOWLEDGMENTS

## REFERENCES

**Atteson, K.** (1996). An Analysis of the Performance of the Neighbor-Joining Method of Phylogeny Reconstruction. *DIMACS Workshop on Mathematical Hierarchies and Biology*. Rutgers University, New Brunswick, NJ, USA.

**Barthélemy, J.P.** and **Guénoche, A.** (1991). *Trees and Proximity Representations*. Wiley, Chichester, England.

**Beyer, H.-G.** (2001). *Theory of Evolution Strategies*. Springer-Verlag, Heidelberg, Germany.

**Bulmer, M.** (1991). Use of the Method of Generalized Least Squares in Reconstructing Phylogenies from Sequence *Data. Mol. Biol. Evol. 8*: 868-883.

**Carleton, M.D.** (1988). Systematics and Evolution. In: *Advances in the Study of Peromyscus* (Kirkland Jr., G.L. and Layne, J.N., eds.). Texas Tech University Press, Lubbock, TX, USA, pp. 7-140.

**Carrizo, S.** (2004). Phylogenetic Trees: an Information Visualization Perspective. *Proceedings of the Second Conference on Asia-Pacific Bioinformatics*, Dunedin, New Zeland, *29*: 315-320.

**Day, W.H.E.** (1987). Computational complexity of inferring phylogenies by dissimilarity matrices. *Bull. Math. Biol. 49*: 461-467.

**Felsenstein, J.** (1988). Phylogenies from molecular sequences: inference and reliability. *Annu. Rev. Genet. 22*: 521-565.

**Felsenstein, J.** (1989). Phylip: phylogeny inference package (version 3.2). *Cladistics 5*: 164-166.

**Felsenstein, J.** (2004). *Inferring Phylogenies*. Sinauer Associates, Sunderland, MA, USA.

**Fitch, W.M.** and **Margoliash, E.** (1967). Construction of phylogenetic trees. *Science 155*: 279-284.

**Foulds, L.R.** and **Graham, R.L.** (1982). The Steiner problem in phylogeny is NP-complete. *Adv. Appl. Math. 3*: 43-49.

**Garey, M.R.** and **Johnson, D.** (1979). *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman, New York, NY, USA.

**Gascuel, O.** (1997). BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Mol. Biol. Evol. 14*: 685-695.

**Holder, M.T.** and **Lewis, P.O.** (2003). Phylogeny estimation: traditional and Bayesian approaches. *Nat. Rev. Genet. 4*: 275-284.

**Holmes, S.P.** (2002). Statistics for phylogenetic trees. *Theor. Popul. Biol. 63*: 17-32.

**Huelsenbeck, P.J.** (1995). Performance of phylogenetic methods in simulation. *Syst. Biol. 44*: 17-48.

**Kidd, K.K.** and **Sgaramella-Zonta, L.A.** (1971). Phylogenetic analysis: concepts and methods. *Am. J. Hum. Genet. 23*: 235-252.

**Kuhner, M.K.** and **Felsenstein, J.** (1994). A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol. Biol. Evol. 11*: 459-468.

**Myers, P., Patton, J.L.** and **Smith, M.F.** (1989). Revision of the Boliviensis Group of Akodon (Muridae: Sigmodontinae), with Emphasis on Perú and Bolivia. *Misc. Publ. Mus. Zool. Univ. Mich 177*: 1-105.

**Patton, J.L.** and **da Silva, N.F.** (1997). Definition of species of pouched four-eyed opossums (Didelphidae, Philander). *J. Mammal. 78*: 90-102.

**Patton, J.L.** and **Smith, M.F.** (1990). Evolutionary Dynamics of Thomomys Bottae Pocket Gophers, with Emphasis on California Populations. *UC Publ. Zool. 123*: 1-161.

**Pearson, W.R., Robins, G.** and **Zhang, T.** (1999). Generalized neighbor-joining: more reliable phylogenetic tree reconstruction. *J. Mol. Evol. 16*: 806-816.

**Prado, O.G.** and **Von Zuben, F.J.** (2002). A step-by-step description of a multi-purpose evolutionary algorithm for phylogenetic tree reconstruction. In: *Proceedings of the Late-Breaking Papers at 2002 Genetic and Evolutionary Computation Conference* (*GECCO-2002*), (Cantú-Paz, E., ed.). New York City, NY, USA, pp. 377-383.

**Rzhetsky, A.** and **Nei, M.** (1993). Theoretical foundation of the minimum-evolution method of phylogenetic inference. *Mol. Biol. Evol. 10*: 1073-1095.

**Rzhetsky, A.** and **Nei, M.** (1994). METREE: a program package for inferring and testing minimum-evolution trees. *Comput. Appl. Biosci. 10*: 409-412.

**Saitou, N.** and **Imanishi, M.** (1989). Relative efficiencies of the fitch-margoliash, maximum-parsimony, maximum-likelihood, minimum-evolution, and neighbor-joining methods of phylogenetic reconstructions in obtaining the correct tree. *Mol. Biol. Evol. 6*: 514-525.

**Saitou, N.** and **Nei, M.** (1987). The neighbor-joining method: a new method for reconstruction of phylogenetic trees. *Mol. Biol. Evol. 4*: 406-425.

**Salemi, M.** and **Vandamme, A.M.** (2003). *The Phylogenetic Handbook - A Practical Approach to DNA and Protein Phylogeny*. Cambridge University Press, Cambridge, UK.

**Studier, J.A.** and **Keppler, K.J.** (1988). A note on the neighbor-joining method of Saitou and Nei. *Mol. Biol. Evol. 5*: 729-731.

**Swofford, D.L.** (2001). *PAUP: Phylogenetic Analysis Using Parsimony and Related Methods*. Version 4.0. Sinauer Associates, Sunderland, MA, USA.

**Swofford, D.L., Olsen, G.J., Waddell, P.J.** and **Hillis, D.M.** (1996). Phylogenetic inference. In: *Molecular Systematics* (Hillis, D.M., Moritz, C. and Mable, B.K., eds.). 2nd edn. Sinauer Associates, Sunderland, MA, USA, pp. 407-543.