

Recent Trends in Molecular Phylogenetic Analysis: Where to Next?

CHRISTOPHER BLAIR AND ROBERT W. MURPHY

From the Department of Ecology and Evolutionary Biology, University of Toronto, Toronto, ON M5S 3B2, Canada (Blair and Murphy); and the Department of Natural History, Royal Ontario Museum, Toronto, ON M5S 2C6, Canada (Blair and Murphy).

Address correspondence to C. Blair at the address above, or e-mail: christopherblair@utoronto.ca.

The acquisition of large multilocus sequence data is providing researchers with an unprecedented amount of information to resolve difficult phylogenetic problems. With these large quantities of data comes the increasing challenge regarding the best methods of analysis. We review the current trends in molecular phylogenetic analysis, focusing specifically on the topics of multiple sequence alignment and methods of tree reconstruction. We suggest that traditional methods are inadequate for these highly heterogeneous data sets and that researchers employ newer more sophisticated search algorithms in their analyses. If we are to best extract the information present in these data sets, a sound understanding of basic phylogenetic principles combined with modern methodological techniques are necessary.

Key words: *alignment, Bayesian, mixture modeling, parsimony, species TREE*

With third-generation sequencing technology rapidly approaching, it will become more feasible to obtain large multilocus data sets to infer evolutionary relationships (Genome 10k Community of Scientists 2009). These enormous quantities of data have spawned the development of several new programs for phylogenetic inference for these highly heterogeneous data sets. From multiple sequence alignment (MSA) to species tree construction, these new methods are changing the way we gather and manipulate data and analyze and interpret results. The purpose of this paper is to provide a brief review of both traditional and more modern methods available to analyze diverse data sets in the era of multilocus phylogenetics. We start by providing a discussion of methods to infer an alignment of orthologous sequences and ways to extract information present in insertion/deletion events. This is followed by a review of the often-controversial direct optimization (DO) techniques. The last 2 sections specifically review methods of tree reconstruction, the first of which focuses mostly on concatenation; a method commonly used in multilocus studies. In contrast, species tree

methods are showing exceptional promise for these types of data but have not yet garnered sufficient attention among empiricists. As phylogenetic methods change rapidly and are often complex, we hope this review can serve as a guideline for researchers new to the field. This review will also be of use to current phylogeneticists looking for novel methods to analyze their multilocus data sets. We stress that this review is by no means exhaustive. However, we feel that the topics discussed are likely those that researchers will encounter in modern molecular systematics.

Alignment

A robust phylogenetic hypothesis is intricately linked to the quality of the MSA. The accurate determination of character homology is vital to infer evolutionary history, and thus, substantial effort has been placed on finding optimal methods to align divergent orthologous sequences. We define homology here as identity in character state due to shared evolutionary history. Due to intronic and intergenic regions, accurate assessments of base-pair homology will become even more problematic as larger multilocus data sets become available. Traditionally, an MSA is performed either manually or with the aid of different computer algorithms. Manual alignments (i.e., by eye) are often relatively straightforward when inferring relationships using highly conserved coding sequences. However, performing manual alignment of noncoding sequences can be problematic as insertion/deletion (indel) events are more common and do not necessarily occur in multiples of 3 bases.

To try to reduce the degree of subjectivity introduced into an MSA of distantly related groups, many researchers turn to computer software packages. Traditionally, the most common programs used to construct an MSA were the Clustal series (Thompson et al. 1994; Larkin et al. 2007). Clustal is a progressive method that uses a neighbor joining (NJ) guide tree to create the MSA from a series of pairwise alignments. The user can specify specific parameters before the analysis is run, including the costs of both inserting a gap

and extending a previously inserted gap. Although Clustal can be extremely helpful when working with difficult sequences, the resulting MSA is often edited manually to subjectively improve character homology. However, opponents of manual alignment methods argue that introducing researcher bias into an analysis yields the resulting alignment and subsequent phylogenetic hypothesis nonreproducible.

The acquisition of large heterogeneous data sets will result in orthologous sequences containing several distinct indel events in noncoding regions. Thus, a refinement of methods used to create an MSA is necessary if we are to accurately determine base-pair homology (Wallace et al. 2005). Over the past decade, several studies have quantified the accuracy of common MSA methods using simulated data sets (Lassmann and Sonnhammer 2002; Katoh et al. 2005; Edgar and Batzoglou 2006; Nuin et al. 2006). For example, a recent study compared the performances of 6 alignment packages to infer the correct MSA from sequences containing gaps (Golubchik et al. 2007). Gaps are inserted into an MSA to account for indel events, but the location of gap insertion is often ambiguous at best. They found that the more commonly used progressive packages like Clustal and T-Coffee (Notredame et al. 2000) performed quite poorly when indels were present. Conversely, programs such as DIALIGN-T (Morgenstern 2004) and MAFFT (Katoh et al. 2002) often recovered a more accurate MSA. However, combining results from different alignment algorithms using M-Coffee (Wallace et al. 2006) and/or removing ambiguously aligned regions with Gblocks (Talavera and Castresana 2007) may result in more robust phylogenetic hypotheses, especially when alignments differ substantially with different algorithms.

Probabilistic methods of MSA are also beginning to determine character homology with noncoding DNA. These methods rely on statistical models of sequence evolution to obtain the alignment or set of alignments with the highest probability given a set of orthologous sequences. Early models were only able to accommodate indels one base pair in length and are thus unrealistic for most molecular data sets (Thorne et al. 1991, 1992). However, more realistic indel models continue to be developed specifically for noncoding sequences (Keightley and Johnson 2004). The model used in the package MCALIGN (Keightley and Johnson 2004), for example, is parameterized by comparing the empirical distribution of indel lengths and frequencies relative to the rate of nucleotide substitution. Although the method is potentially promising, its current implementation only allows for the alignment of 2 or 3 orthologous sequences.

Gap Coding

Following the construction of an MSA for the traditional 2-step MSA + phylogeny estimation procedure, the researcher is left with the decision of how to handle the gaps inserted into the data set by the MSA algorithm to account for indel events. For most traditional maximum parsimony (MP) analyses, gaps have been either coded as

missing data (most cases) or coded as a fifth character state (Swofford et al. 1996). Both of these methods are potentially problematic in that the former completely discards relevant evolutionary information (Simmons et al. 2001; Kawakita et al. 2003), whereas the latter assumes that gaps represent independent evolutionary events; a highly unlikely scenario (Simmons and Ochoterena 2000). These issues also extended into probabilistic phylogenetic inference in that parameters were estimated without taking indel events into account.

To adequately extract the phylogenetic information present in indels (gaps), researchers are beginning to look at both alternative means of character coding (Simmons and Ochoterena 2000; Simmons et al. 2007) as well as the formulation of indel models of sequence evolution (Thorne et al. 1991, 1992; Miklós et al. 2004). Simmons and Ochoterena (2000) suggested using both simple and complex coding of indels to quantify synapomorphies present between taxa. Simple indel coding assumes that each contiguous stretch of gaps with identical 5' and 3' termini comprises a single evolutionary event. A binary (presence/absence) character state matrix is then concatenated to the remaining alignment and analyzed together under any optimality criterion. For probabilistic inference, the morphological models presented in Lewis (2001) can be used for the binary matrix as a separate partition. Indel coding can be performed manually or automated using programs such as GapCoder (Young and Healy 2003) and BARCOD (Barriol 1994). Studies have shown that including the information present in gaps can: 1) dramatically increase the number of potentially phylogenetically informative characters and 2) lead to a more strongly supported phylogenetic hypothesis (Giribet and Wheeler 1999; Kawakita et al. 2003; Ogden and Rosenberg 2007b). Furthermore, indel characters tend to show reduced levels of homoplasy as compared with nucleotide characters (Simmons et al. 2001). Thus, researchers should continue to seek to maximize the evolutionary information present in indel events, especially as more noncoding sequences are used and more realistic models of evolution are developed.

Direct Optimization

An alternative to constructing an MSA prior to phylogenetic inference is to use DO procedures. DO is different from other approaches in that the alignment and phylogenetic tree are estimated simultaneously. Optimization can be performed either under a parsimony or under a probabilistic framework. The program POY (Wheeler 1996; Varón et al. 2009), for example, estimates both the phylogenetic tree and the best alignment based on the MP criterion. Previous versions of POY were also able to implement DO in a likelihood framework. Newer programs such as StatAlign (Novák et al. 2008), BALi-Phy (Suchard and Redelings 2006), and BEAST (Lunter et al. 2005) incorporate models of sequence evolution to estimate the posterior distribution of a set of trees and alignments based on Bayesian inference (BI). The BALi-Phy software shows exceptional promise in

that its models allow for nested or overlapping indel events (Redelings and Suchard 2005, 2007), whereas other methods utilize the more common TKF1 and TKF2 indel models (Thorne et al. 1991, 1992). However, joint estimation of alignment and phylogeny in a probabilistic framework is currently computationally intensive and feasible only with smaller data sets (Lunter et al. 2005). These methods also fit a single model to the data, which may not be justified with multilocus data sets.

Although DO procedures can provide an alternative to constructing a subjective MSA, there are other technical and philosophical factors that possibly confound their utility (Simmons 2004; Simmons and Ochoterena 2000). Studies have compared the performance of DO with POY versus traditional Clustal MSA + MP in PAUP* (Swofford 2002) and found the former to be inferior in most of the analyses performed (Kjer et al. 2007; Ogden and Rosenberg 2007a). However, a more recent simulation study showed that these results are due, in part, to the use of simple gap penalties in POY (Liu, Nelesen, et al. 2009). Conversely, using an affine-gap penalty in POY often leads to more accurate alignment and phylogenetic hypotheses than the traditional Clustal + MP approach. By affine penalty, we refer to the additive nonzero cost of both inserting a gap and extending a previously inserted gap. For example, under a simple gap penalty, the cost incurred would equal kC , where C equals the gap opening cost and k equals the length of the gap. Under an affine penalty, a gap of length k would cost $C_{\text{open}} + kC_{\text{extend}}$, where C_{open} equals the cost of inserting a gap and kC_{extend} is the cost of extending the inserted gap of length k . Few comprehensive simulation studies have investigated the performance of likelihood-based DO methods, but it is anticipated that the power of these methods will continue to increase as more realistic models of evolution become incorporated into the analysis (Miklós et al. 2004).

Phylogenetic Inference

Concatenation

Any method of data analysis must rely on an underlying set of assumptions. A primary objective of any scientific investigation is to minimize the number of unnecessary and unjustifiable assumptions required to explain patterns in the data. In this regard, MP analysis continues to be a popular cladistic method to reconstruct evolutionary histories. Although not statistical in the sense that the criterion does not impose a probabilistic model of sequence evolution, many cladists believe that MP outperforms model-based methods of tree reconstruction (including distance-based phenetic methods). This is due, in part, to factors such as the additional assumptions involved with phenetic and likelihood-based analyses (Siddall and Kluge 1997; Siddall 2001). Several computer programs are able to perform MP analysis using morphological and/or molecular data (Hennig86, Farris 1988; NONA/Pee-Wee, Goloboff 1994) Many of these programs are limited in their scope of

data exploration and analysis and are currently used less frequently in phylogenetic analysis than more utilitarian software packages.

Historically, the most widely used program to estimate trees under the parsimony criterion is PAUP* or phylogenetic analysis using parsimony (*and other methods) (Swofford 2002). This program has several advantages to the user, some of which include the ability to analyze molecular or morphological data or data sets containing a mixture of character types. PAUP* is also able to analyze data under different criteria including the NJ and maximum likelihood (ML) methods. Compared with several other phylogenetic software packages, the Macintosh version of PAUP* is also quite user-friendly, which has added to researchers preference for this program versus other parsimony software. Other functions in PAUP* that have made it a versatile choice for phylogenetic inference are: 1) the ability to implement a variety of different models of evolution in a distance or likelihood framework; 2) the calculation of heterogeneous patterns in data with the partition homogeneity test or incongruence length difference test (Farris et al. 1994); 3) the detection of saturation; and 4) the calculation of pairwise base differences with or without accounting for multiple hits.

Although the benefits of PAUP* are numerous, the program suffers from severe limitations when data sets get too large, especially when there is a lack of strong phylogenetic signal. Furthermore, it is well known that for implicit enumeration and branch-and-bound methods, having a data set with more than 20 taxa can take an enormous amount of time and computational power that is generally not feasible with modern computer technology (Swofford et al. 1996). Instead, performing a heuristic search, usually with a set of random addition sequences (RASs) followed by tree bisection–reconnection (TBR) branch swapping is implemented in hopes of finding a set of most parsimonious trees. However, even these heuristic methods can take an inordinate amount of time when data sets get too large (both in taxa and base pairs). This makes PAUP* an unlikely candidate to analyzing larger multilocus data sets in a parsimony framework. Furthermore, providing support for nodes via bootstrapping (Felsenstein 1985) or jackknifing (Lanyon 1985) would be a computational burden with these quantities of data.

A relatively new phylogenetic inference package, TNT (Tree Analysis Using New Technology; Goloboff et al. 2003, 2008) is showing exceptional promise for the analysis of large data sets (>150 taxa) under MP (Hovenkamp 2004; Giribet 2005; Meier and Ali 2005; Goloboff and Pol 2007). TNT is available for Linux, Macintosh, and Windows platforms (GUI version for Windows only) with a downloadable PowerPoint presentation that explains the basic features of the program. TNT's implementation of 4 "New Technology" search algorithms (sectorial searches, tree drifting, tree fusing, and the ratchet) provides researchers with novel and faster methods to search for optimal trees. TNT allows for the analysis of data sets that would normally require between 300 and 900 times longer with traditional

phylogenetic programs such as PAUP* (Goloboff 1999). For a comprehensive explanation of these search algorithms, the reader is directed to the original papers (Goloboff 1999; Nixon 1999). In addition to these new algorithms providing a more thorough evaluation of tree space with larger data sets, simulation studies have shown that standard RAS + TBR searches in TNT take a fraction of time as those implemented in PAUP*. Indeed it is shown that TNT is able to perform rapid parsimony analysis with a high number of taxa (~1000; Goloboff 1999), although it is not yet known how the algorithms would handle data sets composed of thousands to tens of thousands of base pairs. More simulation studies are needed to determine how TNT will handle these data sets for a greater number of taxa. However, it is highly likely that TNT will serve as a prime candidate to quickly and efficiently analyze multilocus data under the parsimony criterion.

Phylogenetic packages using probabilistic methods such as ML and BI have gained popularity due to their ability to incorporate various models of sequence evolution to best explain the data. Historically, ML methods were used as a parametric comparison with MP analysis. Over the past 2 decades, its use has diminished as newer Bayesian techniques have become available that are better able to exploit the information contained in a data set. For example, early ML methods (i.e., PAUP*) were more computationally demanding than Bayesian methods, especially when applying more than one model of evolution to the data. Indeed, most researchers who used the ML criterion to generate a phylogenetic hypothesis fit a single model of evolution to the data, even if the data were composed of many different genes and codons that are likely evolving at very different rates. Recently, newer, more sophisticated ML algorithms such as those implemented in GARLI (Zwickl 2006) and RAxML (Stamatakis 2006) have been developed for rapid likelihood analysis. RAxML in particular is showing exceptional promise for phylogenetic inference with thousands of taxa and/or base pairs. RAxML is also able to accommodate mixed model ML analysis (GARLI is currently unable to do so), fitting independent evolutionary models to alignments of both nucleotides and amino acids. With these improved algorithms, ML is again becoming a popular choice for phylogenetic analysis of exceptionally large data sets.

Bayesian methods of phylogenetic inference continue to be among the most popular choice for analyzing large complex data sets composed of different genic regions. The majority of past and current studies implement a partitioned Bayesian analysis in the program MrBayes (Ronquist and Huelsenbeck 2003) dividing the data a priori based on different functional regions such as genes, codon position in protein-coding genes, and stems and loops in ribosomal RNA (rRNA) (Nylander et al. 2004; Brandley et al. 2005). Simulation and empirical data sets have shown that partitioning the data into discrete regions and assigning different evolutionary models to each region significantly improves the marginal likelihood values for the posterior distribution (Nylander et al. 2004; Brown and Lemmon 2007). However, there has been less consensus among

researchers as to the best way to determine what the optimal partitioning strategy is for a particular data set. This problem is further exacerbated by the fact that overpartitioning a data set can be just as harmful as underpartitioning and can introduce several unnecessary additional parameters to the analysis that may or may not be justified by the data (Brown and Lemmon 2007). The use of Bayes factors (Kass and Raftery 1995) is currently the most common way to assess the optimal partition strategy for a given data set. However, there is still some degree of subjectivity in the calculation and interpretation of values. For example, some researchers advocate a significance value of 10 or greater (Brandley et al. 2005), whereas others suggest that any value above 2 is significant (Raftery 1996). This means that when comparing 2 separate models, model one versus model 2 (BF_{12}), a $2(\ln BF) > 10$ favors the former over the latter. The values used for the calculation of BF are typically the marginal likelihoods of each run, which can be calculated in software packages such as Mathematica or approximated by the harmonic mean values in MrBayes.

As sequencing technology progresses, heterogeneous data sets encompassing various genomic regions including protein-coding genes, rRNA genes, intergenic regions, and introns will become standard practice. With these diverse data sets, it will be necessary to determine the best methods to account for both the pattern as well as rate variation in nucleotide substitution present throughout the different regions. Although partitioning the data in MrBayes generally leads to improved likelihood scores, this usually comes at the cost of increased computation time required to reach convergence (Roberts et al. 2009). Consider a data set composed of 3 unlinked protein-coding genes. Assigning one model to the entire data set may be relatively computation-friendly, but it may not be philosophically justifiable based on the variance in evolutionary rate and pattern. It may be more suitable to divide the data into discrete partitions (gene and codon within gene) to fully accommodate the heterogeneous rates and patterns most likely present in the data. With 3 protein-coding genes, however, this would require 9 separate partitions, with parameters estimated separately for each partition (assuming an unlinked model). Different priors on branch lengths and evolutionary rates can also be specified for each data partition before the analysis is implemented (Brown et al. 2010). However, increasing the number of partitions and parameters increases computation time and parameter variance at an exponential rate (Roberts et al. 2009). Additional partitions also decrease the number of sites (characters) per partition, increasing the systematic error in model fitting (Brandley et al. 2005; Roberts et al. 2009). Highly partitioned data sets may also lead to errors in branch length estimation (Marshall et al. 2006) and influence convergence and mixing of the Markov chains, which may require changing the heating scheme (Beiko et al. 2006). Although certainly feasible with a small to moderate number of loci, applying these methods to large data sets may lead to erroneous conclusions due to inappropriate models and prior distributions.

A newer and more promising approach to BI of multilocus data sets can be found in the software package BayesPhylogenies (Pagel and Meade 2004). The program implements a mixture model to account for pattern heterogeneity present in diverse data sets and assigns a model to each site in an alignment of nucleotide sequences. BayesPhylogenies uses the likelihood function to calculate how many independent substitution rate matrices, \mathbf{Q} , are needed to best explain the data without the need for a priori partitioning of the data. Each rate matrix is composed of a known model of evolution, generally corresponding to the general time reversible (GTR) model with or without gamma rate variation (Yang 1994), and parameters are estimated separately for each matrix. The GTR method is usually used for each matrix because it allows for the most free parameters. In a traditional likelihood framework, we seek to maximize the likelihood of observing the data given a specific set of parameters. This can be written as $P(\mathbf{D}|\mathbf{Q}, \mathbf{T})$, where \mathbf{D} is an alignment of nucleotide sequences, \mathbf{Q} corresponds to a model of sequence evolution, and \mathbf{T} is a phylogenetic tree. The mixture model method is able to incorporate additional models to maximize the likelihood and can be written as $P(\mathbf{D}|\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3, \dots, \mathbf{Q}_p, \mathbf{T})$, where \mathbf{Q}_1 – \mathbf{Q}_p represents how many models (rate matrices) are required to best explain the data. Simulation and empirical analyses have shown that mixture models can outperform mixed model BI (i.e., partitioning) based on increased likelihood scores, usually without the incorporation of additional unnecessary parameters (Pagel and Meade 2004, 2005).

Although employing a mixture model analysis may be more suitable to a particular data set, the researcher often does not know a priori how many independent rate matrices are required to best explain the data. One method that can be used to determine the minimum number of rate matrices is to simply repeat the analysis several times specifying a different number of \mathbf{Q} matrices with and without gamma rate variation (Γ) and separate base frequencies (π) and use model selection criteria such as Bayes factors or Akaike's information criterion (AIC) to select the best strategy (Roberts et al. 2009). The researcher may also plot the likelihood returns for each separate run versus \mathbf{Q} (with and without Γ) in order to visually assess stationarity (see e.g., Pagel and Meade 2004). Alternatively, the reversible jump (RJ) procedure (Green 1995; Huelsenbeck et al. 2004) can accommodate uncertainty by allowing the joint estimation of parameters in the posterior distribution. Along with estimating the parameters of interest such as tree topology and branch lengths, the RJ method also estimates the minimum number of rate matrices required to maximize the likelihood of observing the data. This method is advantageous as it does not to make any a priori assumptions regarding the pattern of sequence evolution.

Although computationally intensive, an additional useful benefit of employing a mixture model analysis is that the program is able to incorporate heterotachy; variation in the rate of evolution through time in a specific region of an

alignment at a specific location in the tree (Fitch and Markowitz 1970; Pagel and Meade 2008). Heterotachy is a common phenomenon in molecular evolution but is seldom incorporated into a molecular phylogenetic analysis (Lopez et al. 2002; Philippe et al. 2005). BayesPhylogenies accomplishes this by assigning more than one branch length to a specific region in a phylogenetic tree. These additional branch length sets can be specified by the user a priori or estimated directly from the data by the RJ procedure (Pagel and Meade 2008). Using the RJ procedure to estimate the likelihood of additional branch lengths has numerous benefits over assigning additional complete sets to the entire data a priori (Pagel and Meade 2008). For example, assigning additional sets forces more than one branch length for every branch of a tree, most of which are unjustified and add unnecessary parameters to the analysis. The RJ procedure, on the other hand, only assigns additional branch lengths where the data suggest that they are necessary, in turn minimizing the number of additional parameters required to best explain the data.

The benefits of employing a mixture model Bayesian analysis over a partitioned analysis are relatively straightforward. First, these models do not require a priori partitioning of the data, which is often a subjective matter. Second, mixture model methods can select the best modeling strategy a posteriori based on the given data (RJ method), eliminating the need for model selection criteria such as Bayes factors, AIC, and Bayesian Information Criteria. Third, mixture models allow for within partition pattern variability, which is quite common in an alignment of molecular sequences but seldom incorporated (Pagel and Meade 2004). Finally, mixture models have the ability to reduce node-density artifacts (i.e., underestimate of the amount of change in long branches compared with shorter branches with more nodes) often present in a data set (Venditti et al. 2008). For these reasons, as well as those presented above, researchers are beginning to incorporate new heterogeneous models in empirical data sets composed of several independent loci. Over the past few years, mixture models have been used to infer phylogeny in several groups including mammals (Roberts et al. 2009), flies (Lewis et al. 2005), and lizards (Zaldivar-Riverón et al. 2008; Schulte and Cartwright 2009). These studies find mixture model results comparable with traditional probabilistic inference methods, although likelihood scores are generally higher. However, methods to assess adequate mixing and convergence for mixture models are still in development. Additional empirical studies will allow us to draw better conclusions as to the congruence of mixture model approaches with more traditional partition-based methods.

It is predicted that mixture models will begin to be incorporated more frequently as multilocus data sets become available. This is due to a combination of computational efficiency and philosophical justifications to handle highly heterogeneous data sets. It becomes problematic to fit an inordinate number of evolutionary models to a data set containing tens to hundreds of independent loci a priori and running a partitioned analysis. These models will

not only add an enormous amount of additional parameters to the analysis but may also prove to be a poor fit to the data due to overparameterization. Estimating the number of required models directly from the data is a much more feasible and defensible goal. Consider a data set containing 100 000 bp spanning several protein coding exons, introns, and intergenic regions (soon to be a moderately sized data set). If the user wished to perform a partitioned Bayesian analysis, these data would have to be divided a priori into an unmanageable number of partitions, increasing both the number of assumptions and computational burden. Conversely, using the RJ mixture model procedure to estimate the number of required rate matrices directly from the data may show that the data can be explained by only a few different rate matrices. This in turn drastically reduces the number of required parameters in the analysis and will provide higher likelihood scores.

Species Trees

The reconciliation of a well-supported species tree is the primary interest in systematics. Traditionally, the steps involved in a molecular phylogenetic investigation included 1) sequencing a single gene from a representative individual for each taxon under study; 2) constructing a gene tree of the sequences using available phylogenetic methods; and 3) equating the resolved gene tree with the species tree. Although these gene trees often serve as an adequate hypothesis for the species tree in older groups, discordance between gene trees can be quite common, especially for rapid speciation events and recently diverged taxa (Degnan and Rosenberg 2006, 2009; Carstens and Knowles 2007). As we are moving away from single-locus phylogenetics to more multilocus studies (Edwards 2008; Brito and Edwards 2009), gene tree heterogeneity will become a common phenomenon in molecular investigations. Extracting the phylogenetic signal from these discordant gene trees is paramount in hypothesizing the most likely species tree. Therefore, several recent methods have been proposed to estimate the most likely species tree from a set of independent gene trees (reviewed in Knowles 2009).

It is now well understood that gene tree heterogeneity is ubiquitous (Degnan and Rosenberg 2009). Factors such as horizontal gene transfer (Doolittle 1999), gene duplication (Page and Charleston 1997), recombination, hybridization and introgression (Rieseberg et al. 2000), and incomplete lineage sorting (Pamilo and Nei 1988) can result in gene trees that are discordant with other gene trees as well as the species tree (Maddison 1997). Several traditional and more recent methods have been used to deal with this discordance in multilocus data sets. For example, the most common way to incorporate information from multiple genes is to concatenate them into one supermatrix and analyze them as one under the total evidence approach (Kluge 1989). Concatenation assumes that the phylogenetic signal present in most partitions will swamp out the conflicting signal in others. Over the years, several methods have been developed to examine the contribution of individual gene partitions to the

total evidence topology. These include partitioned Bremer support (Baker and DeSalle 1997), partitioned likelihood support (Lee and Hugall 2003), and the partition addition bootstrap alteration (Struck et al. 2006). Given the total evidence topology, these methods quantify the support or conflict for each node calculated from each partition in either a parsimony or a likelihood framework. Indeed, it is the concatenation approach that has garnered the most attention from systematists over the past few decades and has proved effective in some genomic data sets (e.g., Rokas et al. 2003).

Although the concatenation method may provide more phylogenetic resolution simply by increasing the number of informative characters, information about variance in gene coalescence is lost. This has spawned an exciting new era of research into methods to 1) account for gene tree discordance in estimating a species tree and 2) move away from concatenation (Edwards et al. 2007; Kubatko and Degnan 2007; Knowles 2009).

It is currently assumed that incomplete lineage sorting due to deep coalescence is the most common mechanism leading to incongruence among gene trees (Degnan and Rosenberg 2009). Thus, most of the recent species tree methods do not account for the other potential factors outlined above (but see Kubatko 2009). Generally speaking, species tree methods are classified as either summary statistics based or probabilistic based (Knowles 2009; Liu, Yu, et al. 2009; McCormack et al. 2009). Early summary statistics approaches such as minimizing deep coalescence (Maddison and Knowles 2006) are computationally tractable but often less robust than probabilistic methods (McCormack et al. 2009). However, newer summary statistics methods that account for coalescence variance such as STAR and STEAC (Liu, Yu, et al. 2009) are showing exceptional promise to construct species trees. These methods are relatively quick and can be implemented in the statistical package R.

ML and Bayesian methods have been developed to model gene coalescence to infer the most likely species tree. These heavily parameterized models allow for the stochasticity of coalescence times between gene lineages to be integrated into species tree inference probabilistically. In a fully Bayesian framework, for example, the programs BEST (Bayesian Estimation of Species Trees; Liu 2008; Liu and Pearl 2007) and *BEAST (Heled and Drummond 2010) use multiple alleles per species to estimate the posterior distribution of gene and species trees. Prior distributions of gene trees are conditioned on the species tree in which the gene trees are embedded. As a by-product of the joint estimation of gene trees and the species tree, the programs also provide divergence times and ancestral population sizes. The latter has been shown to rely heavily on the prior distribution of $\theta = 4N_e\mu$, where N_e is the effective population size and μ is the per-site mutation rate (Liu and Pearl 2007). Although species tree inference with BEST and *BEAST have been shown to outperform concatenation, computational demand currently limits their utility when applying the methods to larger data sets (Cranston et al. 2009; Liu, Yu, et al. 2009).

Because joint estimation of gene and species trees is often intractable with large quantities of data, other ML and Bayesian methods have been employed. Bayesian concordance analysis (BCA; Ané et al. 2007), for example, uses individually calculated gene trees to infer the species tree that maximizes the bipartition concordance among each gene tree. BCA can be implemented in the program BUCKY (Ané et al. 2007) and requires a set of calculated gene trees along with a single user-defined input parameter α that quantifies the degree of gene tree incongruence. BCA has been shown to be a highly accurate alternative to the demanding algorithm used in BEST and the method is generally insensitive to values of α (Cranston et al. 2009). A promising ML method is found in the program STEM (Species Tree Estimation Using Maximum Likelihood; Kubatko et al. 2009). Like BCA, STEM uses the calculated gene trees (assuming a molecular clock) as input data to calculate the most likely species tree under coalescence. The program can also evaluate the likelihood of a user-specified input tree as well as search for a set of high likelihood topologies. Prior to analysis, the user must specify a value for θ (to convert branch lengths into coalescent units) and a per-locus evolutionary rate r . Although ML inference with STEM has been shown to improve over concatenation ML analysis (Kubatko et al. 2009), it is still unknown how robust the method is to starting parameter values.

Conclusions

We hope that this review encourages both empirical and theoretical systematists to strongly consider the status and future of phylogenetic inference methods. As multilocus data sets become the norm across laboratories, some of the most commonly employed techniques for both MSA and tree reconstruction will no longer be adequate for generating phylogenetic hypotheses. Instead, alternate and more sophisticated search algorithms are required in order to fully exploit the information contained in these large quantities of data. As highly heterogeneous data sets become available, testing the accuracy of both modern alignment algorithms and DO methods through simulation will become even more important.

For traditional phylogenetic inference, MP analysis will no doubt continue to play a role. In this regard, TNT is showing promise for dealing with difficult phylogenetic problems. Furthermore, model-based concatenation methods using mixture models in BayesPhylogenies seem promising for multilocus data sets. However, there have been few simulations to quantify the accuracy of the model compared with other methods including direct species tree inference.

Several methods outlined in this review are also showing promise for analyzing phylogenomic data sets, including the alignment algorithms in the programs Multi-LAGAN (Brudno et al. 2003) and MAUVE (Darling et al. 2004) and the species tree methods found in programs such as BCA, STEM, and STEAC. The field of molecular systematics has always been a rapidly evolving discipline

drawing from expertise across diverse backgrounds. With new sequencing technologies upon us, and new methods of analysis, simulation studies are vital if we are to ascertain the best approach to phylogenetic inference. This simulation-based research has thus been a necessary and rapidly progressing component of multilocus phylogenetics and much more work is needed. Furthermore, little work has been done to quantify the application of new methods to phylogeographic studies. We are approaching an exciting new era of molecular systematics and as more heterogeneous data sets begin to accumulate, the efficacy of current methodology is sure to be put to the test.

Funding

Natural Sciences and Engineering Research Council of Canada Discovery Grant A3148 to R.W.M.

Acknowledgments

We thank the graduate students and faculty of the Royal Ontario Museum for their discussions and insightful criticisms. We thank Ida M. Conflitti and 3 anonymous reviewers for their comments to improve the manuscript.

References

- Ané C, Larget B, Baum DA, Smith SD, Rokas A. 2007. Bayesian estimation of concordance among gene trees. *Mol Biol Evol.* 24:412–426.
- Baker RH, DeSalle R. 1997. Multiple sources of character information and the phylogeny of Hawaiian drosophilids. *Syst Biol.* 46:654–673.
- Barriel V. 1994. Molecular phylogenies and nucleotide insertion-deletion. *CR Acad Sci III.* 317:693–701.
- Beiko RG, Keith JM, Harlow TJ, Ragan MA. 2006. Searching for convergence in phylogenetic Markov chain Monte Carlo. *Syst Biol.* 55:553–565.
- Brandley MC, Schmitz A, Reeder TW. 2005. Partitioned Bayesian analyses, partition choice, and the phylogenetic relationships of scincid lizards. *Syst Biol.* 54:373–390.
- Brito P, Edwards SV. 2009. Multilocus phylogeography and phylogenetics using sequence-based markers. *Genetica.* 135:439–455.
- Brown JM, Hedtke SM, Lemmon AR, Lemmon EM. 2010. When trees grow too long: investigating the causes of highly inaccurate Bayesian branch-length estimates. *Syst Biol.* 59:145–161.
- Brown JM, Lemmon AR. 2007. The importance of data partitioning and the utility of Bayes factors in Bayesian phylogenetics. *Syst Biol.* 56:643–655.
- Brudno M, Do CB, Cooper GM, Kim MF, Davydov E, NISC Comparative Sequencing Program, Green ED, Sidow A, Batzoglou S. 2003. LAGAN and Multi-LAGAN: efficient tools for large-scale multiple alignment of genomic DNA. *Genome Res.* 13:721–731.
- Carstens BC, Knowles LL. 2007. Estimating species phylogeny from gene-tree probabilities despite incomplete lineage sorting: an example from *Melanoplus* grasshoppers. *Syst Biol.* 56:400–411.
- Cranston KA, Hurwitz B, Ware D, Stein L, Wing RA. 2009. Species trees from highly incongruent gene trees in rice. *Syst Biol.* 58:489–500.
- Darling ACE, Mau B, Blatter FR, Perna NT. 2004. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res.* 14:1394–1403.
- Degnan JH, Rosenberg NA. 2006. Discordance of species trees with their most likely gene trees. *PLoS Genet.* 2:762–768.

- Degnan JH, Rosenberg NA. 2009. Gene tree discordance, phylogenetic inference, and the multispecies coalescent. *Trends Ecol Evol.* 24:332–340.
- Doolittle WF. 1999. Lateral genomics. *Trends Cell Biol.* 9:M5–M8.
- Edgar RC, Batzoglou S. 2006. Multiple sequence alignment. *Curr Opin Struct Biol.* 16:368–373.
- Edwards SV. 2008. Is a new and general theory of molecular systematics emerging? *Evolution.* 63:1–19.
- Edwards SV, Liu L, Pearl DK. 2007. High-resolution species trees without concatenation. *Proc Natl Acad Sci U S A.* 104:5936–5941.
- Farris JS. 1988. Hennig86, version 1.5. Program and documentation. Port Jefferson Station (NY): Distributed by the Author.
- Farris JS, Källersjö M, Kluge AG, Bult C. 1994. Testing the significance of incongruence. *Cladistics.* 10:315–319.
- Felsenstein J. 1985. Confidence limits on phylogenies: an approach using the bootstrap. *Evolution.* 39:783–791.
- Fitch WM, Markowitz E. 1970. An improved method for determining codon variability in a gene and its application to the rate of fixation of mutations in evolution. *Biochem Genet.* 4:579–593.
- Genome 10k Community of Scientists. 2009. Genome 10k: a proposal to obtain whole-genome sequence for 10 000 vertebrate species. *J Hered.* 100:659–674.
- Giribet G. 2005. A review of: “TNT: tree analysis using New Technology.”. *Syst Biol.* 54:176–178.
- Giribet G, Wheeler WC. 1999. On gaps. *Mol Phylogenet Evol.* 13:132–143.
- Goloboff PA. 1994. Nona: A Tree-searching program. Program and documentation [Internet]. Available from: <http://www.zmuc.dk/public/phylogeny/Nona-PeeWee>. Published by the author, Tucumán, Argentina.
- Goloboff PA. 1999. Analyzing large datasets in reasonable times: solutions for composite optima. *Cladistics.* 15:415–428.
- Goloboff PA, Farris J, Nixon K. 2003. T.N.T.: Tree Analysis Using New Technology. Program and documentation [Internet]. Available from: <http://www.zmuc.dk/public/phylogeny/tnt>.
- Goloboff PA, Pol D. 2007. On divide-and-conquer strategies for parsimony analysis of large data sets: rec-I-dcm3 vs TNT. *Syst Biol.* 56:485–495.
- Goloboff PA, Farris JS, Nixon KC. 2008. TNT, a free program for phylogenetic analysis. *Cladistics.* 24:774–786.
- Golubchik T, Wise MJ, Eastale S, Jermini LS. 2007. Mind the gaps: evidence of bias in estimates of multiple sequence alignments. *Mol Biol Evol.* 24:2433–2442.
- Green PJ. 1995. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika.* 82:711–732.
- Heled J, Drummond AJ. 2010. Bayesian inference of species trees from multilocus data. *Mol Biol Evol.* 27:570–580.
- Hovenkamp P. 2004. Review of TNT—tree analysis using New Technology, ver. 1.0. *Cladistics.* 20:378–383.
- Huelsenbeck JP, Larget B, Alfaro ME. 2004. Bayesian phylogenetic model selection using reversible jump Markov chain Monte Carlo. *Mol Biol Evol.* 21:1123–1133.
- Kass RE, Raftery AE. 1995. Bayes factors. *J Am Stat Assoc.* 90:773–795.
- Katoh K, Kuma K, Toh H, Miyata T. 2005. MAFFT version 5: improvement in accuracy of multiple sequence alignment. *Nucleic Acids Res.* 33:511–518.
- Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.
- Kawakita A, Sota T, Ascher JS, Ito M, Tanaka H, Kato M. 2003. Evolution and phylogenetic utility of alignment gaps within intron sequences of three nuclear genes in bumble bees (*Bombus*). *Mol Biol Evol.* 20:87–92.
- Keightley PD, Johnson T. 2004. MCALIGN: stochastic alignment of noncoding DNA sequences based on an evolutionary model of sequence evolution. *Genome Res.* 14:442–450.
- Kjer KM, Gillespie JJ, Ober KA. 2007. Opinions on multiple sequence alignment, and an empirical comparison of repeatability and accuracy between POY and structural alignment. *Syst Biol.* 56:133–146.
- Kluge AG. 1989. A concern for evidence and a phylogenetic hypothesis for relationships among *Epicrates* (Boidea, Serpentes). *Syst Zool.* 38:7–25.
- Knowles LL. 2009. Estimating species trees: methods of phylogenetic analysis when there is incongruence across genes. *Syst Biol.* 58:463–467.
- Kubatko LS. 2009. Identifying hybridization events in the presence of coalescence via model selection. *Syst Biol.* 58:478–488.
- Kubatko LS, Carstens BC, Knowles LL. 2009. STEM: species tree estimation using maximum likelihood for gene trees under coalescence. *Bioinformatics.* 25:971–973.
- Lanyon SM. 1985. Detecting internal inconsistencies in distance data. *Syst Zool.* 38:7–25.
- Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, et al. 2007. Clustal W and Clustal X version 2.0. *Bioinformatics.* 23:2947–2948.
- Lassmann T, Sonnhammer ELL. 2002. Quality assessment of multiple alignment programs. *FEBS Lett.* 529:126–130.
- Lee MS, Huggall AF. 2003. Partitioned likelihood support and the evaluation of data set conflict. *Syst Biol.* 52:15–22.
- Lewis PO. 2001. A likelihood approach to estimating phylogeny from discrete morphological character data. *Syst Biol.* 50:913–925.
- Lewis RL, Beckenbach AT, Mooers AO. 2005. The phylogeny of the subgroups within the *melanogaster* species group: likelihood tests on *COI* and *COII* sequences and a Bayesian estimate of phylogeny. *Mol Phylogenet Evol.* 37:15–24.
- Liu K, Nelesen S, Raghavan S, Linder CR, Warnow T. 2009. Barking up the wrong treelength: the impact of gap penalty on alignment and tree accuracy. *IEEE/ACM Trans Comput Biol Bioinform.* 6:7–21.
- Liu L. 2008. BEST: Bayesian estimation of species trees under the coalescent model. *Bioinformatics.* 24:2542–2543.
- Liu L, Pearl DK. 2007. Species trees from gene trees: reconstructing Bayesian posterior distributions of a species phylogeny using estimated gene tree distributions. *Syst Biol.* 56:504–514.
- Liu L, Yu L, Pearl DK, Edwards SV. 2009. Estimating species phylogenies using coalescence times among sequences. *Syst Biol.* 58:468–477.
- Lopez P, Casane D, Philippe H. 2002. Heterotachy, an important process of protein evolution. *Mol Biol Evol.* 19:1–7.
- Lunter G, Miklós I, Drummond A, Jensen JL, Hein J. 2005. Bayesian coestimation of phylogeny and sequence alignment. *BMC Bioinformatics.* 6:83.
- Maddison WP. 1997. Gene trees in species trees. *Syst Biol.* 46:523–536.
- Maddison WP, Knowles LL. 2006. Inferring phylogeny despite incomplete lineage sorting. *Syst Biol.* 55:21–30.
- Marshall DC, Simon C, Buckley TR. 2006. Accurate branch length estimation in partitioned Bayesian analyses requires accommodation of among-partition rate variation and attention to branch length priors. *Syst Biol.* 55:993–1003.
- McCormack JE, Huang H, Knowles LL. 2009. Maximum likelihood estimates of species trees: how accuracy of phylogenetic inference depends upon the divergence history and sampling design. *Syst Biol.* 58:501–508.
- Meier R, Ali F. 2005. The newest kid on the parsimony block: TNT (Tree analysis using new technology). *Syst Entomol Am.* 30:179–182.
- Miklós I, Lunter GA, Holmes I. 2004. A “long indel” model for evolutionary sequence alignment. *Mol Biol Evol.* 21:529–540.
- Morgenstern B. 2004. DIALIGN: multiple DNA and protein sequence alignment at BiBiServ. *Nucleic Acids Res.* 32:W33–W36.
- Nixon KC. 1999. The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics.* 15:407–414.

- Notredame C, Higgins DG, Heringa J. 2000. T-Coffee: a novel method for fast and accurate multiple sequence alignment. *J Mol Biol.* 302:205–217.
- Novák A, Miklós I, Lyngso R, Hein J. 2008. StatAlign: an extendable software package for joint Bayesian estimation of alignments and evolutionary trees. *Bioinformatics.* 24:2403–2404.
- Nuin PAS, Wang Z, Tillier ERM. 2006. The accuracy of several multiple sequence alignment programs for proteins. *BMC Bioinformatics.* 7:471.
- Nylander JAA, Ronquist F, Huelsenbeck JP, Nieves-Aldrey J. 2004. Bayesian phylogenetic analysis of combined data. *Syst Biol.* 53:47–67.
- Ogden TH, Rosenberg MS. 2007a. Alignment and topological accuracy of the direct optimization approach via POY and traditional phylogenetics via ClustalW + PAUP*. *Syst Biol.* 56:182–193.
- Ogden TH, Rosenberg MS. 2007b. How should gaps be treated in parsimony? A comparison of approaches using simulation. *Mol Phylogenet Evol.* 42:817–826.
- Page RDM, Charleston MA. 1997. From gene to organismal phylogeny: reconciled trees and the gene tree species tree problem. *Mol Phylogenet Evol.* 7:231–240.
- Pagel M, Meade A. 2004. A phylogenetic mixture model for detecting pattern-heterogeneity in gene sequence or character-state data. *Syst Biol.* 53:571–581.
- Pagel M, Meade A. 2005. Mixture models in phylogenetic inference. In: Gascuel O, editor. *Mathematics of evolution and phylogeny.* Oxford: Clarendon Press. p. 121–142.
- Pagel M, Meade A. 2008. Modelling heterotachy in phylogenetic inference by reversible-jump Markov chain Monte Carlo. *Philos Trans R Soc B Biol Sci.* 363:3955–3964.
- Pamilo P, Nei M. 1988. Relationships between gene trees and species trees. *Mol Biol Evol.* 5:568–583.
- Philippe H, Zhou Y, Brinkmann H, Rodrigue N, Delsuc F. 2005. Heterotachy and long-branch attraction in phylogenetics. *BMC Evol Biol.* 5:50.
- Raftery AE. 1996. Hypothesis testing and model selection. In: Gilks WR, Richardson S, Spiegelhalter DJ, editors. *Markov chain Monte Carlo in practice.* London: Chapman Hall. p. 163–188.
- Redelings BD, Suchard MA. 2005. Joint Bayesian estimation of alignment and phylogeny. *Syst Biol.* 54:401–418.
- Redelings BD, Suchard MA. 2007. Incorporating indel information into phylogeny estimation for rapidly emerging pathogens. *BMC Evol Biol.* 7:40.
- Rieseberg LH, Baird SJ, Gardner KA. 2000. Hybridization, introgression, and linkage evolution. *Plant Mol Biol.* 42:205–224.
- Roberts TE, Sargis EJ, Olson LE. 2009. Networks, trees, and treeshrews: assessing support and identifying conflict with multiple loci and a problematic root. *Syst Biol.* 58:257–270.
- Rokas A, Williams BL, King N, Carroll SB. 2003. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature.* 425:798–804.
- Ronquist F, Huelsenbeck JP. 2003. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics.* 19:1572–1574.
- Schulte JAIL, Cartwright EM. 2009. Phylogenetic relationships among iguanian lizards using alternative partitioning methods and *TSHZ1*: a new phylogenetic marker for reptiles. *Mol Phylogenet Evol.* 50:391–396.
- Siddall ME. 2001. Philosophy and phylogenetic inference: a comparison of likelihood and parsimony methods in the context of Karl Popper's writings on corroboration. *Cladistics.* 17:395–399.
- Siddall ME, Kluge AG. 1997. Probabilism and phylogenetic inference. *Cladistics.* 13:313–336.
- Simmons MP. 2004. Independence of alignment and tree search. *Mol Phylogenet Evol.* 31:874–879.
- Simmons MP, Müller K, Norton AP. 2007. The relative performance of indel-coding methods in simulation. *Mol Phylogenet Evol.* 44:724–740.
- Simmons MP, Ochoterena H. 2000. Gaps as characters in sequence-based phylogenetic analysis. *Syst Biol.* 49:369–381.
- Simmons MP, Ochoterena H, Carr TG. 2001. Incorporation, relative homoplasy, and effect of gap characters in sequence-based phylogenetic analyses. *Syst Biol.* 50:454–462.
- Stamatakis A. 2006. RAxML-VI-HPC: maximum-likelihood based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics.* 22:2688–2690.
- Struck TH, Purschke G, Halaných KM. 2006. Phylogeny of Eunicida (Annelida) and exploring data congruence using a partition addition bootstrap alteration (PABA) approach. *Syst Biol.* 55:1–20.
- Suchard MA, Redelings BD. 2006. Bali-Phy: simultaneous Bayesian inference of alignment and phylogeny. *Bioinformatics.* 22:2047–2048.
- Swofford DL. 2002. PAUP*: phylogenetic analysis using parsimony (* and other methods), version 4. Sunderland (MA): Sinauer Associates.
- Swofford DL, Olsen GJ, Waddell PJ, Hillis DM. 1996. Phylogenetic inference. In: Hillis DM, Moritz C, Mable BK, editors. *Molecular systematics.* 2nd ed. Sunderland (MA): Sinauer Associates. p. 407–514.
- Talavera G, Castresana J. 2007. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Syst Biol.* 56:564–577.
- Thompson JD, Higgins DG, Gibson TJ. 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* 22:4673–4680.
- Thorne JL, Kishino H, Felsenstein J. 1991. An evolutionary model for maximum likelihood alignment of DNA sequences. *J Mol Evol.* 33:114–124.
- Thorne JL, Kishino H, Felsenstein J. 1992. Inching towards reality—an improved likelihood model of sequence evolution. *J Mol Evol.* 34:3–16.
- Varón A, Vinh LS, Wheeler WC. 2009. POY version 4: phylogenetic analysis using dynamic homologies. *Cladistics.* 26:72–85.
- Venditti C, Meade A, Pagel M. 2008. Phylogenetic mixture models can reduce node-density artifacts. *Syst Biol.* 57:286–293.
- Wallace IM, Blackshields G, Higgins DG. 2005. Multiple sequence alignments. *Curr Opin Struct Biol.* 15:261–266.
- Wallace IM, O'Sullivan O, Higgins DG, Notredame C. 2006. M-Coffee: combining multiple sequence alignment methods with T-Coffee. *Nucleic Acids Res.* 34:1692–1699.
- Wheeler W. 1996. Optimization alignment: the end of multiple sequence alignment in phylogenetics. *Cladistics.* 12:1–9.
- Yang Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *J Mol Evol.* 39:306–314.
- Young ND, Healy J. 2003. GapCoder automates the use of indel characters in phylogenetic analysis. *BMC Bioinformatics.* 4:6.
- Zaldivar-Riverón A, Nieto-Montes de Oca A, Manríquez-Morán N, Reeder TW. 2008. Phylogenetic affinities of the rare and enigmatic limb-reduced *Anehyropsis* (Reptilia: Squamata) as inferred with mitochondrial 16S rRNA sequence data. *J Herpetol.* 42:303–311.
- Zwickl DJ. 2006. Genetic algorithm approaches for the phylogenetic analysis of large biological sequence datasets under the maximum likelihood criterion [PhD dissertation]. [Austin (TX)]: The University of Texas at Austin.

Received April 21, 2010; Revised June 11, 2010;
Accepted July 12, 2010

Corresponding Editor: William Murphy